

# Intermediate Logic

An Open Introduction



F21 $\alpha$

# **Intermediate Logic**

## *The Open Logic Project*

### **Instigator**

Richard Zach, *University of Calgary*

### **Editorial Board**

Aldo Antonelli,<sup>†</sup> *University of California, Davis*

Andrew Arana, *Université de Lorraine*

Jeremy Avigad, *Carnegie Mellon University*

Tim Button, *University College London*

Walter Dean, *University of Warwick*

Gillian Russell, *Dianoia Institute of Philosophy*

Nicole Wyatt, *University of Calgary*

Audrey Yap, *University of Victoria*

### **Contributors**

Samara Burns, *Columbia University*

Dana Hägg, *University of Calgary*

Zesen Qian, *Carnegie Mellon University*

# Intermediate Logic

*An Open Introduction*

Remixed by Michael Hallett  
Richard Zach

FALL 2021 $\alpha$

The Open Logic Project would like to acknowledge the generous support of the **Taylor Institute of Teaching and Learning** of the University of Calgary, and the Alberta Open Educational Resources (ABOER) Initiative, which is made possible through an investment from the Alberta government.



**UNIVERSITY OF CALGARY**

Taylor Institute for Teaching and Learning



Cover illustrations by **Matthew Leadbeater**, used under a **Creative Commons Attribution-NonCommercial 4.0 International License**.

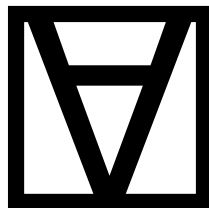
Typeset in Baskervald X and Nimbus Sans by L<sup>A</sup>T<sub>E</sub>X.

This version of *Intermediate Logic* is revision 596b290 (2024-01-02), with content generated from *Open Logic Text* revision 6c541de (2024-02-28). Free download at:

<https://builds.openlogicproject.org/courses/intermediate-logic/>



*Intermediate Logic* by **Michael Hallett Richard Zach** is licensed under a **Creative Commons Attribution 4.0 International License**. It is based on *The Open Logic Text* by the **Open Logic Project**, used under a **Creative Commons Attribution 4.0 International License**.



# Contents

<b>About this Book</b>	<b>xiv</b>
<b>I Sets, Relations, Functions</b>	<b>1</b>
<b>1 Sets</b>	<b>2</b>
1.1 Extensionality . . . . .	2
1.2 Subsets and Power Sets . . . . .	4
1.3 Some Important Sets . . . . .	5
1.4 Unions and Intersections . . . . .	6
1.5 Pairs, Tuples, Cartesian Products . . . . .	10
1.6 Russell's Paradox . . . . .	12
Summary . . . . .	14
Problems . . . . .	14
<b>2 Relations</b>	<b>16</b>
2.1 Relations as Sets . . . . .	16
2.2 Special Properties of Relations . . . . .	18
2.3 Equivalence Relations . . . . .	20
2.4 Orders . . . . .	21
2.5 Graphs . . . . .	24
2.6 Operations on Relations . . . . .	26
Summary . . . . .	27
Problems . . . . .	27

<b>3</b>	<b>Functions</b>	<b>29</b>
3.1	Basics . . . . .	29
3.2	Kinds of Functions . . . . .	32
3.3	Functions as Relations . . . . .	34
3.4	Inverses of Functions . . . . .	36
3.5	Composition of Functions . . . . .	39
3.6	Partial Functions . . . . .	40
	Summary . . . . .	41
	Problems . . . . .	42
<b>4</b>	<b>The Size of Sets</b>	<b>43</b>
4.1	Introduction . . . . .	43
4.2	Enumerations and Countable Sets . . . . .	43
4.3	Cantor's Zig-Zag Method . . . . .	48
4.4	Pairing Functions and Codes . . . . .	50
4.5	An Alternative Pairing Function . . . . .	52
4.6	Uncountable Sets . . . . .	54
4.7	Reduction . . . . .	58
4.8	Equinumerosity . . . . .	59
4.9	Sets of Different Sizes, and Cantor's Theorem . . . . .	61
4.10	The Notion of Size, and Schröder-Bernstein . . . . .	63
	Summary . . . . .	64
	Problems . . . . .	65
<b>II</b>	<b>First-order Logic</b>	<b>68</b>
<b>5</b>	<b>Introduction to First-Order Logic</b>	<b>69</b>
5.1	First-Order Logic . . . . .	69
5.2	Syntax . . . . .	71
5.3	Formulas . . . . .	72
5.4	Satisfaction . . . . .	74
5.5	Sentences . . . . .	76
5.6	Semantic Notions . . . . .	77
5.7	Substitution . . . . .	78
5.8	Models and Theories . . . . .	79
5.9	Soundness and Completeness . . . . .	81

<b>6</b>	<b>Syntax of First-Order Logic</b>	<b>83</b>
6.1	Introduction . . . . .	83
6.2	First-Order Languages . . . . .	84
6.3	Terms and Formulas . . . . .	86
6.4	Unique Readability . . . . .	90
6.5	Main operator of a Formula . . . . .	94
6.6	Subformulas . . . . .	95
6.7	Formation Sequences . . . . .	97
6.8	Free Variables and Sentences . . . . .	101
6.9	Substitution . . . . .	103
	Summary . . . . .	105
	Problems . . . . .	106
<b>7</b>	<b>Semantics of First-Order Logic</b>	<b>107</b>
7.1	Introduction . . . . .	107
7.2	Structures for First-order Languages . . . . .	108
7.3	Covered Structures for First-order Languages . . . . .	110
7.4	Satisfaction of a Formula in a Structure . . . . .	111
7.5	Variable Assignments . . . . .	118
7.6	Extensionality . . . . .	121
7.7	Semantic Notions . . . . .	123
	Summary . . . . .	126
	Problems . . . . .	126
<b>8</b>	<b>Theories and Their Models</b>	<b>129</b>
8.1	Introduction . . . . .	129
8.2	Expressing Properties of Structures . . . . .	131
8.3	Examples of First-Order Theories . . . . .	133
8.4	Expressing Relations in a Structure . . . . .	136
8.5	The Theory of Sets . . . . .	137
8.6	Expressing the Size of Structures . . . . .	141
	Summary . . . . .	142
	Problems . . . . .	143
<b>9</b>	<b>Derivation Systems</b>	<b>145</b>
9.1	Introduction . . . . .	145



9.2	The Sequent Calculus . . . . .	147
9.3	Natural Deduction . . . . .	148
9.4	Tableaux . . . . .	150
9.5	Axiomatic Derivations . . . . .	152
<b>10</b>	<b>The Sequent Calculus</b>	<b>155</b>
10.1	Rules and Derivations . . . . .	155
10.2	Propositional Rules . . . . .	156
10.3	Quantifier Rules . . . . .	157
10.4	Structural Rules . . . . .	159
10.5	Derivations . . . . .	160
10.6	Examples of Derivations . . . . .	162
10.7	Derivations with Quantifiers . . . . .	166
10.8	Proof-Theoretic Notions . . . . .	168
10.9	Derivability and Consistency . . . . .	171
10.10	Derivability and the Propositional Connectives . . . . .	172
10.11	Derivability and the Quantifiers . . . . .	174
10.12	Soundness . . . . .	175
10.13	Derivations with Identity predicate . . . . .	182
10.14	Soundness with Identity predicate . . . . .	183
	Summary . . . . .	184
	Problems . . . . .	184
<b>11</b>	<b>Natural Deduction</b>	<b>187</b>
11.1	Rules and Derivations . . . . .	187
11.2	Propositional Rules . . . . .	188
11.3	Quantifier Rules . . . . .	189
11.4	Derivations . . . . .	191
11.5	Examples of Derivations . . . . .	193
11.6	Derivations with Quantifiers . . . . .	198
11.7	Proof-Theoretic Notions . . . . .	202
11.8	Derivability and Consistency . . . . .	204
11.9	Derivability and the Propositional Connectives . . . . .	206
11.10	Derivability and the Quantifiers . . . . .	208
11.11	Soundness . . . . .	209
11.12	Derivations with Identity predicate . . . . .	214

11.13	Soundness with Identity predicate . . . . .	216
	Summary . . . . .	217
	Problems . . . . .	217
<b>12</b>	<b>The Completeness Theorem</b>	<b>221</b>
12.1	Introduction . . . . .	221
12.2	Outline of the Proof . . . . .	223
12.3	Complete Consistent Sets of Sentences . . . . .	226
12.4	Henkin Expansion . . . . .	227
12.5	Lindenbaum's Lemma . . . . .	230
12.6	Construction of a Model . . . . .	231
12.7	Identity . . . . .	234
12.8	The Completeness Theorem . . . . .	238
12.9	The Compactness Theorem . . . . .	239
12.10	A Direct Proof of the Compactness Theorem . . . . .	242
12.11	The Löwenheim–Skolem Theorem . . . . .	243
	Summary . . . . .	244
	Problems . . . . .	245
<b>13</b>	<b>Beyond First-order Logic</b>	<b>247</b>
13.1	Overview . . . . .	247
13.2	Many-Sorted Logic . . . . .	248
13.3	Second-Order logic . . . . .	250
13.4	Higher-Order logic . . . . .	255
13.5	Intuitionistic Logic . . . . .	258
13.6	Modal Logics . . . . .	264
13.7	Other Logics . . . . .	266
<b>III</b>	<b>Incompleteness</b>	<b>268</b>
<b>14</b>	<b>Introduction to Incompleteness</b>	<b>269</b>
14.1	Historical Background . . . . .	269
14.2	Definitions . . . . .	275
14.3	Overview of Incompleteness Results . . . . .	281
14.4	Undecidability and Incompleteness . . . . .	284
	Summary . . . . .	286

Problems . . . . .	287
<b>15 Recursive Functions</b>	<b>288</b>
15.1 Introduction . . . . .	288
15.2 Primitive Recursion . . . . .	289
15.3 Composition . . . . .	292
15.4 Primitive Recursion Functions . . . . .	294
15.5 Primitive Recursion Notations . . . . .	298
15.6 Primitive Recursive Functions are Computable . . . . .	299
15.7 Examples of Primitive Recursive Functions . . . . .	300
15.8 Primitive Recursive Relations . . . . .	303
15.9 Bounded Minimization . . . . .	306
15.10 Primes . . . . .	308
15.11 Sequences . . . . .	309
15.12 Trees . . . . .	313
15.13 Other Recursions . . . . .	314
15.14 Non-Primitive Recursive Functions . . . . .	315
15.15 Partial Recursive Functions . . . . .	317
15.16 The Normal Form Theorem . . . . .	320
15.17 The Halting Problem . . . . .	321
15.18 General Recursive Functions . . . . .	323
Summary . . . . .	323
Problems . . . . .	325
<b>16 Arithmetization of Syntax</b>	<b>327</b>
16.1 Introduction . . . . .	327
16.2 Coding Symbols . . . . .	329
16.3 Coding Terms . . . . .	331
16.4 Coding Formulas . . . . .	333
16.5 Substitution . . . . .	335
16.6 Derivations in <b>LK</b> . . . . .	336
16.7 Derivations in Natural Deduction . . . . .	341
Summary . . . . .	347
Problems . . . . .	348
<b>17 Representability in <math>\mathcal{Q}</math></b>	<b>350</b>

17.1	Introduction . . . . .	350
17.2	Functions Representable in $\mathbf{Q}$ are Computable . . . . .	353
17.3	The Beta Function Lemma . . . . .	355
17.4	Simulating Primitive Recursion . . . . .	359
17.5	Basic Functions are Representable in $\mathbf{Q}$ . . . . .	360
17.6	Composition is Representable in $\mathbf{Q}$ . . . . .	364
17.7	Regular Minimization is Representable in $\mathbf{Q}$ . . . . .	366
17.8	Computable Functions are Representable in $\mathbf{Q}$ . . . . .	370
17.9	Representing Relations . . . . .	371
17.10	Undecidability . . . . .	372
	Summary . . . . .	374
	Problems . . . . .	374
<b>18</b>	<b>Incompleteness and Provability</b>	<b>376</b>
18.1	Introduction . . . . .	376
18.2	The Fixed-Point Lemma . . . . .	378
18.3	The First Incompleteness Theorem . . . . .	381
18.4	Rosser's Theorem . . . . .	383
18.5	Comparison with Gödel's Original Paper . . . . .	385
18.6	The Derivability Conditions for $\mathbf{PA}$ . . . . .	386
18.7	The Second Incompleteness Theorem . . . . .	387
18.8	Löb's Theorem . . . . .	390
18.9	The Undefinability of Truth . . . . .	393
18.10	Tarski's Theorem and Löb's Theorem . . . . .	395
	Summary . . . . .	404
	Problems . . . . .	405
<b>19</b>	<b>Models of Arithmetic</b>	<b>407</b>
19.1	Introduction . . . . .	407
19.2	Reducts and Expansions . . . . .	408
19.3	Isomorphic Structures . . . . .	409
19.4	The Theory of a Structure . . . . .	412
19.5	Standard Models of Arithmetic . . . . .	413
19.6	Non-Standard Models . . . . .	416
19.7	Models of $\mathbf{Q}$ . . . . .	417
19.8	Models of $\mathbf{PA}$ . . . . .	420

19.9	Computable Models of Arithmetic	424
	Summary	426
	Problems	428
<b>A</b>	<b>Derivations in Arithmetic Theories</b>	<b>430</b>
<b>B</b>	<b>Proofs</b>	<b>438</b>
B.1	Introduction	438
B.2	Starting a Proof	440
B.3	Using Definitions	440
B.4	Inference Patterns	443
B.5	An Example	451
B.6	Another Example	455
B.7	Proof by Contradiction	457
B.8	Reading Proofs	462
B.9	I Can't Do It!	464
B.10	Other Resources	466
	Problems	467
<b>C</b>	<b>Induction</b>	<b>468</b>
C.1	Introduction	468
C.2	Induction on $\mathbb{N}$	469
C.3	Strong Induction	472
C.4	Inductive Definitions	473
C.5	Structural Induction	476
C.6	Relations and Functions	478
	Problems	482
<b>D</b>	<b>Biographies</b>	<b>483</b>
D.1	Georg Cantor	483
D.2	Alonzo Church	484
D.3	Gerhard Gentzen	485
D.4	Kurt Gödel	487
D.5	Emmy Noether	489
D.6	Rózsa Péter	491
D.7	Julia Robinson	493
D.8	Bertrand Russell	495

D.9 Alfred Tarski . . . . .	497
D.10 Alan Turing . . . . .	498
D.11 Ernst Zermelo . . . . .	500
<b>Photo Credits</b>	<b>503</b>
<b>Bibliography</b>	<b>505</b>
<b>About the Open Logic Project</b>	<b>514</b>

# *About this Book*

This book is an introduction to metalogic, aimed especially at students of computer science and philosophy. “Metalogic” is so-called because it is the discipline that studies logic itself. Logic proper is concerned with canons of valid inference, and its symbolic or formal version presents these canons using formal languages, such as those of propositional and predicate, a.k.a., first-order logic. Meta-logic investigates the properties of these language, and of the canons of correct inference that use them. It studies topics such as how to give precise meaning to the expressions of these formal languages, how to justify the canons of valid inference, what the properties of various derivation systems are, including their computational properties. These questions are important and interesting in their own right, because the languages and proof systems investigated are applied in many different areas—in mathematics, philosophy, computer science, and linguistics, especially—but they also serve as examples of how to study formal systems in general. The logical languages we study here are not the only ones people are interested in. For instance, linguists and philosophers are interested in languages that are much more complicated than those of propositional and first-order logic, and computer scientists are interested in other *kinds* of languages altogether, such as programming languages. And the methods we discuss here—how to give semantics for formal languages, how to prove results about formal languages, how to investigate the properties of formal languages—are applicable

in those cases as well.

Like any discipline, metalogic both has a set of results or facts, and a store of methods and techniques, and this text covers both. Some students won't need to know some of the results we discuss outside of this course, but they will need and use the methods we use to establish them. The Löwenheim-Skolem theorem, say, does not often make an appearance in computer science, but the methods we use to prove it do. On the other hand, many of the results we discuss do have relevance for certain debates, say, in the philosophy of science and in metaphysics. Philosophy students may not need to be able to prove these results outside this course, but they do need to understand what the results are—and you really only *understand* these results if you have thought through the definitions and proofs needed to establish them. These are, in part, the reasons for why the results and the methods covered in this text are recommended study—in some cases even required—for students of computer science and philosophy.

The material is divided into three parts. **Part I** concerns itself with the theory of sets. Logic and metalogic is historically connected very closely to what's called the “foundations of mathematics.” Mathematical foundations deal with how ultimately mathematical objects such as integers, rational, and real numbers, functions, spaces, etc., should be understood. Set theory provides one answer (there are others), and so set theory and logic have long been studied side-by-side. Sets, relations, and functions are also ubiquitous in any sort of formal investigation, not just in mathematics but also in computer science and in some of the more technical corners of philosophy. Certainly for the purposes of formulating and proving results about the semantics and proof theory of logic and the foundation of computability it is essential to have a language in which to do this. For instance, we will talk about sets of expressions, relations of consequence and provability, interpretations of predicate symbols (which turn out to be relations), computable functions, and various relations between and constructions using these. It will be good to have shorthand symbols for these, and think through the general prop-



erties of sets, relations, and functions in order to do that. If you are not used to thinking mathematically and to formulating mathematical proofs, then think of the first part on set theory as a training ground: all the basic definitions will be given, and we'll give increasingly complicated proofs using them. Note that understanding these proofs—and being able to find and formulate them yourself—is perhaps more important than understanding the results, and especially in the first part, and especially if you are new to mathematical thinking, it is important that you think through the examples and problems.

In the first part we will establish one important result, however. This result—Cantor's theorem—relies on one of the most striking examples of conceptual analysis to be found anywhere in the sciences, namely, Cantor's analysis of infinity. Infinity has puzzled mathematicians and philosophers alike for centuries. No one knew how to properly think about it. Many people even thought it was a mistake to think about it at all, that the notion of an infinite object or infinite collection itself was incoherent. Cantor made infinity into a subject we can coherently work with, and developed an entire theory of infinite collections—and infinite numbers with which we can measure the sizes of infinite collections—and showed that there are different levels of infinity. This theory of “transfinite” numbers is beautiful and intricate, and we won't get very far into it; but we will be able to show that there are different levels of infinity, specifically, that there are “countable” and “uncountable” levels of infinity. This result has important applications, but it is also really the kind of result that any self-respecting mathematician, computer scientist, or philosopher should know.

In **part II** we turn to first-order logic. We will define the language of first-order logic and its semantics, i.e., what first-order structures are and when a sentence of first-order logic is true in a structure. This will enable us to do two important things: (1) We can define, with mathematical precision, when a sentence is a logical consequence of another. (2) We can also consider how the relations that make up a first-order structure are described—

characterized—by the sentences that are true in them. This in particular leads us to a discussion of the axiomatic method, in which sentences of first-order languages are used to characterize certain kinds of structures. Proof theory will occupy us next, and we will consider the original version of the sequent calculus and natural deduction as defined in the 1930s by Gerhard Gentzen. (Your instructor may choose to cover only one, then any reference to “derivations” and “provability” will mean whatever system they chose.) The semantic notion of consequence and the syntactic notion of provability give us two completely different ways to make precise the idea that a sentence may follow from some others. The soundness and completeness theorems link these two characterizations. In particular, we will prove Gödel’s completeness theorem, which states that whenever a sentence is a semantic consequence of some others, there it is also provable from them. An equivalent formulation is: if a collection of sentences is consistent—in the sense that nothing contradictory can be proved from them—then there is a structure that makes all of them true.

The second formulation of the completeness theorem is perhaps the more surprising. Around the time Gödel proved this result (in 1929), the German mathematician David Hilbert famously held the view (which builds on a similar view of the leading French mathematician Henri Poincaré) that consistency (i.e., freedom from contradiction) is all that mathematical existence requires. In other words, whenever a mathematician can coherently describe a structure or class of structures, then they should be entitled to believe in the existence of such structures. At the time, many found this idea preposterous: just because you can describe a structure without contradicting yourself, it surely does not follow that such a structure actually exists. But that is exactly what Gödel’s completeness theorem says. In addition to this paradoxical—and certainly philosophically intriguing—aspect, the completeness theorem also has two important applications which allow us to prove further results about the existence of structures which make given sentences true. These are

the compactness and the Löwenheim-Skolem theorems.

In **part III**, we will present (in **chapter 18**) what have come to be called the *Gödel Incompleteness Theorems*, two results which constitute Gödel's second great contribution to mathematical logic.<sup>1</sup> We mentioned Hilbert above. Another of his extremely important contributions (from the 1890s) was the insistence on the axiomatic method, above all for studying independence proofs and relative consistency, and through this a much more refined notion of logical dependency than had been possible hitherto.<sup>2</sup> Part of the idea behind this was that theories (like that for natural number arithmetic or for sets) must be formulated by finitely presented axiom systems which are 'self-standing' in the sense that all the major arithmetical or set-theoretical facts can be derived in the axiom systems. This leads to two questions in particular:

1. How can we show that the axiom system we give are *consistent*, i.e., do not lead to the proof of contradictions? (Logic itself doesn't, as we will show, but special axioms for numbers or sets added to the logical system might.)
2. Are there truths, about numbers say, which cannot be derived from a natural axiom system for arithmetic? We can derive some truths, for example

$$2 + 2 = 4$$

or

$$\forall m, n (m + n = n + m),$$

but do we know that we have captured *all* of them?

Gödel's work relates directly to these questions.

---

<sup>1</sup>Believe it or not, there was a third, the proof of the relative consistency of the Axiom of Choice and the Generalised Continuum Hypothesis with the standard axioms of set theory. Any one of these contributions would have made Gödel a giant of the subject!

<sup>2</sup>Pursuit of the axiomatic method was what led to the "consistency  $\Rightarrow$  existence" thesis mentioned in the previous paragraph.

*The First Incompleteness Theorem* addresses the second question. What Gödel invented was a very powerful general method for the construction of *self-referential sentences* within the languages concerned (provided they have certain mechanisms available), for instance, sentences of the restricted language of arithmetic (see [chapter 16](#)). To be more specific, we will be able to form sentences in the special first-order arithmetical language which talk about themselves, and can describe fundamental properties that they possess. This is analogous to the following sentence of English: ‘I am a sentence which contains four commas, consists of 25 words, and, when written and set on Richard Zach’s computer, appears in black type’. If we’ve counted the number of commas and words correctly, this sentence is quite straightforwardly *true*. Gödel’s method allowed for the formulation of a perfectly ordinary (logically fairly simple, but rather long) sentence of arithmetic which says something like ‘I am not deducible from the system of axioms for arithmetic’, and it is easy to show that (given that the axiom system for arithmetic is consistent — and we will be in rather serious trouble if that isn’t correct) this sentence must be *true*, and that therefore what it says must be correct, from which it follows that it is indeed not deducible. Note this: *true* but *not deducible!* Hence, arithmetic is *incomplete*. Gödel’s procedure gives us a general recipe for the production of true but undervivable sentences for a very wide range of theories and their languages. For instance, suppose, out of frustration, we add the particular, undervivable sentence just described as a new axiom to get a new theory of arithmetic, then we can use the very same procedure to produce a *new* sentence which is true and cannot be derived from that expanded axiom system. It is the *general* nature of the procedure (what we now usually mean when we refer to ‘Gödel’s First Incompleteness Theorem’) which gives the Gödel result its enormous power, and ensures that it does not present a mere curio, but a genuine philosophical dilemma.

Closely related to Gödel’s First Incompleteness Theorem is *Tarski’s Theorem* on the undefinability of truth for reasonable axioms systems for arithmetic using the usual language, and this

we will also look at. We mentioned above the crucial importance of the assumption of the consistency of arithmetic. What Tarski's saw was that if truth *were* definable within the language of arithmetic, then we could show that arithmetic would in fact be *inconsistent*, since we could in effect produce in the language of arithmetic a version of a so-called 'Liar sentence', a sentence which is contradictory since it must be simultaneously both true and false, just like the sentence of ordinary English 'This sentence is not true'. (The parallel with the standard Gödel sentence which says something like 'This sentence is not derivable' is conscious and deliberate on Gödel's part. The Gödel sentence is a sentence of the language of arithmetic, the Liar sentence is not.)

Gödel produced a fairly straightforward (but difficult to execute) modification of the proof for the First Theorem which he used to show the following: the consistency of arithmetic can only be shown by a theory which has *stronger* theoretical resources than those available to arithmetic itself. This (expressed in the *Second Incompleteness Theorem*) addresses the *first* question posed above, for in a sense it says that we cannot really prove the consistency of arithmetic at all, at least, not if we are doing so to guarantee the 'health and safety' of formal arithmetic. (Of course, there are other reasons for studying the inferential structure of theories, and thus consistency, and these have been widely pursued.) Again, Gödel's Second Incompleteness Theorem has extraordinary range, and applies to a very wide swathe of ordinary, working theories, the theory of sets among them.

Closely related to the both the First and the Second Theorems is a result called *Löb's Theorem*, which we will also look at briefly. We will also examine another result which in a way presages Gödel's First Incompleteness Theorem, often called *Skolem's Theorem*, which concerns the austere and strange world of non-standard models of arithmetic. These contain the usual natural numbers, which behave exactly as we expect them to behave, but much (very much!) more, namely a vast multiplicity of what we call *non-standard numbers* with very odd structure (see [chapter 19](#).)

In sum, we will see that these results (some of the most important theorems of twentieth-century logic) tells us a good deal about the power and also the limitations of first-order logic. (But what's the alternative to first-order logic? This is touched on in [chapter 13](#).) From these central results, highly technical in nature, wider *philosophical* consequences start to flow, certainly important consequences for the philosophy of mathematics, for instance concerning the role of proof and provability, and the connections between proof and truth, and for Hilbert's programme, some of which we touched on above, but also consequences for philosophy more generally, for instance for the computational theory of mind, for the theory of truth (and thus philosophy of language), and consequences concerning the nature of mathematical and scientific theories.

The material here is very much cumulative, and the results emerge slowly. Hence, it's important both to keep up, and to be patient. Believe us, it's worth it! And remember all the time, that what we're doing is really proving things in an *informal* way about *formal* languages and derivation systems, even though it helps to be familiar with doing *formal* proofs (following rules of proof) and with the semantics of first-order languages.

Logic at this level is a difficult but very beautiful subject, and is very different in nature and approach from the subject of logic as presented in your first logic course. Being good at that does not, by any means, ensure you will be good at this. What this requires is, above all, not mathematical knowledge, or even mathematical ability, but rather something which we might, rather obscurely, refer to as *mathematical aptitude*.

Let us begin.



**PART I**

*Sets,  
Relations,  
Functions*



## CHAPTER 1

# Sets

### 1.1 Extensionality

A *set* is a collection of objects, considered as a single object. The objects making up the set are called *elements* or *members* of the set. If  $x$  is an element of a set  $a$ , we write  $x \in a$ ; if not, we write  $x \notin a$ . The set which has no elements is called the *empty set* and denoted “ $\emptyset$ ”.

It does not matter how we *specify* the set, or how we *order* its elements, or indeed how *many times* we count its elements. All that matters are what its elements are. We codify this in the following principle.

**Definition 1.1 (Extensionality).** If  $A$  and  $B$  are sets, then  $A = B$  iff every element of  $A$  is also an element of  $B$ , and vice versa.

Extensionality licenses some notation. In general, when we have some objects  $a_1, \dots, a_n$ , then  $\{a_1, \dots, a_n\}$  is *the* set whose elements are  $a_1, \dots, a_n$ . We emphasise the word “*the*”, since extensionality tells us that there can be only *one* such set. Indeed, extensionality also licenses the following:

$$\{a, a, b\} = \{a, b\} = \{b, a\}.$$

This delivers on the point that, when we consider sets, we don't care about the order of their elements, or how many times they are specified.

**Example 1.2.** Whenever you have a bunch of objects, you can collect them together in a set. The set of Richard's siblings, for instance, is a set that contains one person, and we could write it as  $S = \{\text{Ruth}\}$ . The set of positive integers less than 4 is  $\{1, 2, 3\}$ , but it can also be written as  $\{3, 2, 1\}$  or even as  $\{1, 2, 1, 2, 3\}$ . These are all the same set, by extensionality. For every element of  $\{1, 2, 3\}$  is also an element of  $\{3, 2, 1\}$  (and of  $\{1, 2, 1, 2, 3\}$ ), and vice versa.

Frequently we'll specify a set by some property that its elements share. We'll use the following shorthand notation for that:  $\{x : \varphi(x)\}$ , where the  $\varphi(x)$  stands for the property that  $x$  has to have in order to be counted among the elements of the set.

**Example 1.3.** In our example, we could have specified  $S$  also as

$$S = \{x : x \text{ is a sibling of Richard}\}.$$

**Example 1.4.** A number is called *perfect* iff it is equal to the sum of its proper divisors (i.e., numbers that evenly divide it but aren't identical to the number). For instance, 6 is perfect because its proper divisors are 1, 2, and 3, and  $6 = 1 + 2 + 3$ . In fact, 6 is the only positive integer less than 10 that is perfect. So, using extensionality, we can say:

$$\{6\} = \{x : x \text{ is perfect and } 0 \leq x \leq 10\}$$

We read the notation on the right as “the set of  $x$ 's such that  $x$  is perfect and  $0 \leq x \leq 10$ ”. The identity here confirms that, when we consider sets, we don't care about how they are specified. And, more generally, extensionality guarantees that there is always only one set of  $x$ 's such that  $\varphi(x)$ . So, extensionality justifies calling  $\{x : \varphi(x)\}$  *the* set of  $x$ 's such that  $\varphi(x)$ .

Extensionality gives us a way for showing that sets are identical: to show that  $A = B$ , show that whenever  $x \in A$  then also  $x \in B$ , and whenever  $y \in B$  then also  $y \in A$ .

## 1.2 Subsets and Power Sets

We will often want to compare sets. And one obvious kind of comparison one might make is as follows: *everything in one set is in the other too*. This situation is sufficiently important for us to introduce some new notation.

**Definition 1.5 (Subset).** If every element of a set  $A$  is also an element of  $B$ , then we say that  $A$  is a *subset* of  $B$ , and write  $A \subseteq B$ . If  $A$  is not a subset of  $B$  we write  $A \not\subseteq B$ . If  $A \subseteq B$  but  $A \neq B$ , we write  $A \subsetneq B$  and say that  $A$  is a *proper subset* of  $B$ .

**Example 1.6.** Every set is a subset of itself, and  $\emptyset$  is a subset of every set. The set of even numbers is a subset of the set of natural numbers. Also,  $\{a, b\} \subseteq \{a, b, c\}$ . But  $\{a, b, e\}$  is not a subset of  $\{a, b, c\}$ .

**Example 1.7.** The number 2 is an element of the set of integers, whereas the set of even numbers is a subset of the set of integers. However, a set may happen to *both* be an element and a subset of some other set, e.g.,  $\{0\} \in \{0, \{0\}\}$  and also  $\{0\} \subseteq \{0, \{0\}\}$ .

Extensionality gives a criterion of identity for sets:  $A = B$  iff every element of  $A$  is also an element of  $B$  and vice versa. The definition of “subset” defines  $A \subseteq B$  precisely as the first half of this criterion: every element of  $A$  is also an element of  $B$ . Of course the definition also applies if we switch  $A$  and  $B$ : that is,  $B \subseteq A$  iff every element of  $B$  is also an element of  $A$ . And that, in turn, is exactly the “vice versa” part of extensionality. In other words, extensionality entails that sets are equal iff they are subsets of one another.

**Proposition 1.8.**  $A = B$  iff both  $A \subseteq B$  and  $B \subseteq A$ .

Now is also a good opportunity to introduce some further bits of helpful notation. In defining when  $A$  is a subset of  $B$  we said that “every element of  $A$  is ...,” and filled the “...” with

“an element of  $B$ ”. But this is such a common *shape* of expression that it will be helpful to introduce some formal notation for it.

**Definition 1.9.**  $(\forall x \in A)\varphi$  abbreviates  $\forall x(x \in A \rightarrow \varphi)$ . Similarly,  $(\exists x \in A)\varphi$  abbreviates  $\exists x(x \in A \wedge \varphi)$ .

Using this notation, we can say that  $A \subseteq B$  iff  $(\forall x \in A)x \in B$ .

Now we move on to considering a certain kind of set: the set of all subsets of a given set.

**Definition 1.10 (Power Set).** The set consisting of all subsets of a set  $A$  is called the *power set* of  $A$ , written  $\wp(A)$ .

$$\wp(A) = \{B : B \subseteq A\}$$

**Example 1.11.** What are all the possible subsets of  $\{a, b, c\}$ ? They are:  $\emptyset$ ,  $\{a\}$ ,  $\{b\}$ ,  $\{c\}$ ,  $\{a, b\}$ ,  $\{a, c\}$ ,  $\{b, c\}$ ,  $\{a, b, c\}$ . The set of all these subsets is  $\wp(\{a, b, c\})$ :

$$\wp(\{a, b, c\}) = \{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{b, c\}, \{a, c\}, \{a, b, c\}\}$$

### 1.3 Some Important Sets

**Example 1.12.** We will mostly be dealing with sets whose elements are mathematical objects. Four such sets are important enough to have specific names:

$$\mathbb{N} = \{0, 1, 2, 3, \dots\}$$

the set of natural numbers

$$\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$$

the set of integers

$$\mathbb{Q} = \{m/n : m, n \in \mathbb{Z} \text{ and } n \neq 0\}$$

the set of rationals

$$\mathbb{R} = (-\infty, \infty)$$

the set of real numbers (the continuum)

These are all *infinite* sets, that is, they each have infinitely many elements.

As we move through these sets, we are adding *more* numbers to our stock. Indeed, it should be clear that  $\mathbb{N} \subseteq \mathbb{Z} \subseteq \mathbb{Q} \subseteq \mathbb{R}$ : after all, every natural number is an integer; every integer is a rational; and every rational is a real. Equally, it should be clear that  $\mathbb{N} \subsetneq \mathbb{Z} \subsetneq \mathbb{Q}$ , since  $-1$  is an integer but not a natural number, and  $1/2$  is rational but not integer. It is less obvious that  $\mathbb{Q} \subsetneq \mathbb{R}$ , i.e., that there are some real numbers which are not rational.

We'll sometimes also use the set of positive integers  $\mathbb{Z}^+ = \{1, 2, 3, \dots\}$  and the set containing just the first two natural numbers  $\mathbb{B} = \{0, 1\}$ .

**Example 1.13 (Strings).** Another interesting example is the set  $A^*$  of *finite strings* over an alphabet  $A$ : any finite sequence of elements of  $A$  is a string over  $A$ . We include the *empty string*  $\Lambda$  among the strings over  $A$ , for every alphabet  $A$ . For instance,

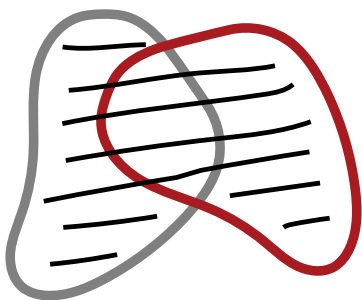
$$\mathbb{B}^* = \{\Lambda, 0, 1, 00, 01, 10, 11, \\ 000, 001, 010, 011, 100, 101, 110, 111, 0000, \dots\}.$$

If  $x = x_1 \dots x_n \in A^*$  is a string consisting of  $n$  “letters” from  $A$ , then we say *length* of the string is  $n$  and write  $\text{len}(x) = n$ .

**Example 1.14 (Infinite sequences).** For any set  $A$  we may also consider the set  $A^\omega$  of infinite sequences of elements of  $A$ . An infinite sequence  $a_1 a_2 a_3 a_4 \dots$  consists of a one-way infinite list of objects, each one of which is an element of  $A$ .

## 1.4 Unions and Intersections

In [section 1.1](#), we introduced definitions of sets by abstraction, i.e., definitions of the form  $\{x : \varphi(x)\}$ . Here, we invoke some property  $\varphi$ , and this property can mention sets we've already



*Figure 1.1:* The union  $A \cup B$  of two sets is set of elements of  $A$  together with those of  $B$ .

defined. So for instance, if  $A$  and  $B$  are sets, the set  $\{x : x \in A \vee x \in B\}$  consists of all those objects which are elements of either  $A$  or  $B$ , i.e., it's the set that combines the elements of  $A$  and  $B$ . We can visualize this as in [Figure 1.1](#), where the highlighted area indicates the elements of the two sets  $A$  and  $B$  together.

This operation on sets—combining them—is very useful and common, and so we give it a formal name and a symbol.

**Definition 1.15 (Union).** The *union* of two sets  $A$  and  $B$ , written  $A \cup B$ , is the set of all things which are elements of  $A$ ,  $B$ , or both.

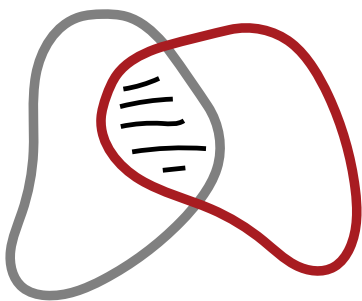
$$A \cup B = \{x : x \in A \vee x \in B\}$$

**Example 1.16.** Since the multiplicity of elements doesn't matter, the union of two sets which have an element in common contains that element only once, e.g.,  $\{a, b, c\} \cup \{a, 0, 1\} = \{a, b, c, 0, 1\}$ .

The union of a set and one of its subsets is just the bigger set:  $\{a, b, c\} \cup \{a\} = \{a, b, c\}$ .

The union of a set with the empty set is identical to the set:  $\{a, b, c\} \cup \emptyset = \{a, b, c\}$ .

We can also consider a “dual” operation to union. This is the operation that forms the set of all elements that are elements of  $A$



*Figure 1.2:* The intersection  $A \cap B$  of two sets is the set of elements they have in common.

and are also elements of  $B$ . This operation is called *intersection*, and can be depicted as in [Figure 1.2](#).

**Definition 1.17 (Intersection).** The *intersection* of two sets  $A$  and  $B$ , written  $A \cap B$ , is the set of all things which are elements of both  $A$  and  $B$ .

$$A \cap B = \{x : x \in A \wedge x \in B\}$$

Two sets are called *disjoint* if their intersection is empty. This means they have no elements in common.

**Example 1.18.** If two sets have no elements in common, their intersection is empty:  $\{a, b, c\} \cap \{0, 1\} = \emptyset$ .

If two sets do have elements in common, their intersection is the set of all those:  $\{a, b, c\} \cap \{a, b, d\} = \{a, b\}$ .

The intersection of a set with one of its subsets is just the smaller set:  $\{a, b, c\} \cap \{a, b\} = \{a, b\}$ .

The intersection of any set with the empty set is empty:  $\{a, b, c\} \cap \emptyset = \emptyset$ .

We can also form the union or intersection of more than two sets. An elegant way of dealing with this in general is the following: suppose you collect all the sets you want to form the union

(or intersection) of into a single set. Then we can define the union of all our original sets as the set of all objects which belong to at least one element of the set, and the intersection as the set of all objects which belong to every element of the set.

**Definition 1.19.** If  $A$  is a set of sets, then  $\bigcup A$  is the set of elements of elements of  $A$ :

$$\begin{aligned}\bigcup A &= \{x : x \text{ belongs to an element of } A\}, \text{ i.e.,} \\ &= \{x : \text{there is a } B \in A \text{ so that } x \in B\}\end{aligned}$$

**Definition 1.20.** If  $A$  is a set of sets, then  $\bigcap A$  is the set of objects which all elements of  $A$  have in common:

$$\begin{aligned}\bigcap A &= \{x : x \text{ belongs to every element of } A\}, \text{ i.e.,} \\ &= \{x : \text{for all } B \in A, x \in B\}\end{aligned}$$

**Example 1.21.** Suppose  $A = \{\{a, b\}, \{a, d, e\}, \{a, d\}\}$ . Then  $\bigcup A = \{a, b, d, e\}$  and  $\bigcap A = \{a\}$ .

We could also do the same for a sequence of sets  $A_1, A_2, \dots$

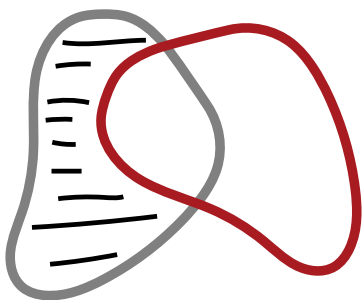
$$\begin{aligned}\bigcup_i A_i &= \{x : x \text{ belongs to one of the } A_i\} \\ \bigcap_i A_i &= \{x : x \text{ belongs to every } A_i\}.\end{aligned}$$

When we have an *index* of sets, i.e., some set  $I$  such that we are considering  $A_i$  for each  $i \in I$ , we may also use these abbreviations:

$$\begin{aligned}\bigcup_{i \in I} A_i &= \bigcup \{A_i : i \in I\} \\ \bigcap_{i \in I} A_i &= \bigcap \{A_i : i \in I\}\end{aligned}$$

Finally, we may want to think about the set of all elements in  $A$  which are not in  $B$ . We can depict this as in [Figure 1.3](#).





*Figure 1.3:* The difference  $A \setminus B$  of two sets is the set of those elements of  $A$  which are not also elements of  $B$ .

**Definition 1.22 (Difference).** The *set difference*  $A \setminus B$  is the set of all elements of  $A$  which are not also elements of  $B$ , i.e.,

$$A \setminus B = \{x : x \in A \text{ and } x \notin B\}.$$

## 1.5 Pairs, Tuples, Cartesian Products

It follows from extensionality that sets have no order to their elements. So if we want to represent order, we use *ordered pairs*  $\langle x, y \rangle$ . In an unordered pair  $\{x, y\}$ , the order does not matter:  $\{x, y\} = \{y, x\}$ . In an ordered pair, it does: if  $x \neq y$ , then  $\langle x, y \rangle \neq \langle y, x \rangle$ .

How should we think about ordered pairs in set theory? Crucially, we want to preserve the idea that ordered pairs are identical iff they share the same first element and share the same second element, i.e.:

$$\langle a, b \rangle = \langle c, d \rangle \text{ iff both } a = c \text{ and } b = d.$$

We can define ordered pairs in set theory using the Wiener–Kuratowski definition.

**Definition 1.23 (Ordered pair).**  $\langle a, b \rangle = \{\{a\}, \{a, b\}\}$ .

Having fixed a definition of an ordered pair, we can use it to define further sets. For example, sometimes we also want ordered sequences of more than two objects, e.g., *triples*  $\langle x, y, z \rangle$ , *quadruples*  $\langle x, y, z, u \rangle$ , and so on. We can think of triples as special ordered pairs, where the first element is itself an ordered pair:  $\langle x, y, z \rangle$  is  $\langle \langle x, y \rangle, z \rangle$ . The same is true for quadruples:  $\langle x, y, z, u \rangle$  is  $\langle \langle \langle x, y \rangle, z \rangle, u \rangle$ , and so on. In general, we talk of *ordered  $n$ -tuples*  $\langle x_1, \dots, x_n \rangle$ .

Certain sets of ordered pairs, or other ordered  $n$ -tuples, will be useful.

**Definition 1.24 (Cartesian product).** Given sets  $A$  and  $B$ , their *Cartesian product*  $A \times B$  is defined by

$$A \times B = \{\langle x, y \rangle : x \in A \text{ and } y \in B\}.$$

**Example 1.25.** If  $A = \{0, 1\}$ , and  $B = \{1, a, b\}$ , then their product is

$$A \times B = \{\langle 0, 1 \rangle, \langle 0, a \rangle, \langle 0, b \rangle, \langle 1, 1 \rangle, \langle 1, a \rangle, \langle 1, b \rangle\}.$$

**Example 1.26.** If  $A$  is a set, the product of  $A$  with itself,  $A \times A$ , is also written  $A^2$ . It is the set of *all* pairs  $\langle x, y \rangle$  with  $x, y \in A$ . The set of all triples  $\langle x, y, z \rangle$  is  $A^3$ , and so on. We can give a recursive definition:

$$\begin{aligned} A^1 &= A \\ A^{k+1} &= A^k \times A \end{aligned}$$

**Proposition 1.27.** *If  $A$  has  $n$  elements and  $B$  has  $m$  elements, then  $A \times B$  has  $n \cdot m$  elements.*

*Proof.* For every element  $x$  in  $A$ , there are  $m$  elements of the form  $\langle x, y \rangle \in A \times B$ . Let  $B_x = \{\langle x, y \rangle : y \in B\}$ . Since whenever  $x_1 \neq x_2$ ,  $\langle x_1, y \rangle \neq \langle x_2, y \rangle$ ,  $B_{x_1} \cap B_{x_2} = \emptyset$ . But if  $A = \{x_1, \dots, x_n\}$ , then  $A \times B = B_{x_1} \cup \dots \cup B_{x_n}$ , and so has  $n \cdot m$  elements.

To visualize this, arrange the elements of  $A \times B$  in a grid:

$$\begin{array}{llll} B_{x_1} = & \{\langle x_1, y_1 \rangle & \langle x_1, y_2 \rangle & \dots & \langle x_1, y_m \rangle\} \\ B_{x_2} = & \{\langle x_2, y_1 \rangle & \langle x_2, y_2 \rangle & \dots & \langle x_2, y_m \rangle\} \\ & \vdots & & \vdots & \\ B_{x_n} = & \{\langle x_n, y_1 \rangle & \langle x_n, y_2 \rangle & \dots & \langle x_n, y_m \rangle\} \end{array}$$

Since the  $x_i$  are all different, and the  $y_j$  are all different, no two of the pairs in this grid are the same, and there are  $n \cdot m$  of them.  $\square$

**Example 1.28.** If  $A$  is a set, a *word* over  $A$  is any sequence of elements of  $A$ . A sequence can be thought of as an  $n$ -tuple of elements of  $A$ . For instance, if  $A = \{a, b, c\}$ , then the sequence “ $bac$ ” can be thought of as the triple  $\langle b, a, c \rangle$ . Words, i.e., sequences of symbols, are of crucial importance in computer science. By convention, we count elements of  $A$  as sequences of length 1, and  $\emptyset$  as the sequence of length 0. The set of *all* words over  $A$  then is

$$A^* = \{\emptyset\} \cup A \cup A^2 \cup A^3 \cup \dots$$

## 1.6 Russell’s Paradox

Extensionality licenses the notation  $\{x : \varphi(x)\}$ , for *the* set of  $x$ ’s such that  $\varphi(x)$ . However, all that extensionality *really* licenses is the following thought. *If* there is a set whose members are all and only the  $\varphi$ ’s, *then* there is only one such set. Otherwise put: having fixed some  $\varphi$ , the set  $\{x : \varphi(x)\}$  is unique, *if it exists*.

But this conditional is important! Crucially, not every property lends itself to *comprehension*. That is, some properties do *not* define sets. If they all did, then we would run into outright contradictions. The most famous example of this is Russell’s Paradox.

Sets may be elements of other sets—for instance, the power set of a set  $A$  is made up of sets. And so it makes sense to ask or investigate whether a set is an element of another set. Can a set be a member of itself? Nothing about the idea of a set seems to rule this out. For instance, if *all* sets form a collection of objects, one might think that they can be collected into a single set—the set of all sets. And it, being a set, would be an element of the set of all sets.

Russell's Paradox arises when we consider the property of not having itself as an element, of being *non-self-membered*. What if we suppose that there is a set of all sets that do not have themselves as an element? Does

$$R = \{x : x \notin x\}$$

exist? It turns out that we can prove that it does not.

**Theorem 1.29 (Russell's Paradox).** *There is no set  $R = \{x : x \notin x\}$ .*

*Proof.* If  $R = \{x : x \notin x\}$  exists, then  $R \in R$  iff  $R \notin R$ , which is a contradiction.  $\square$

Let's run through this proof more slowly. If  $R$  exists, it makes sense to ask whether  $R \in R$  or not. Suppose that indeed  $R \in R$ . Now,  $R$  was defined as the set of all sets that are not elements of themselves. So, if  $R \in R$ , then  $R$  does not itself have  $R$ 's defining property. But only sets that have this property are in  $R$ , hence,  $R$  cannot be an element of  $R$ , i.e.,  $R \notin R$ . But  $R$  can't both be and not be an element of  $R$ , so we have a contradiction.

Since the assumption that  $R \in R$  leads to a contradiction, we have  $R \notin R$ . But this also leads to a contradiction! For if  $R \notin R$ , then  $R$  itself does have  $R$ 's defining property, and so  $R$  would be an element of  $R$  just like all the other non-self-membered sets. And again, it can't both not be and be an element of  $R$ .

How do we set up a set theory which avoids falling into Russell's Paradox, i.e., which avoids making the *inconsistent* claim that

$R = \{x : x \notin x\}$  exists? Well, we would need to lay down axioms which give us very precise conditions for stating when sets exist (and when they don't).

The set theory sketched in this chapter doesn't do this. It's *genuinely naïve*. It tells you only that sets obey extensionality and that, if you have some sets, you can form their union, intersection, etc. It is possible to develop set theory more rigorously than this.

## Summary

A **set** is a collection of objects, the elements of the set. We write  $x \in A$  if  $x$  is an element of  $A$ . Sets are **extensional**—they are completely determined by their elements. Sets are specified by **listing** the elements explicitly or by giving a property the elements share (**abstraction**). Extensionality means that the order or way of listing or specifying the elements of a set doesn't matter. To prove that  $A$  and  $B$  are the same set ( $A = B$ ) one has to prove that every element of  $X$  is an element of  $Y$  and vice versa.

Important sets include the natural ( $\mathbb{N}$ ), integer ( $\mathbb{Z}$ ), rational ( $\mathbb{Q}$ ), and real ( $\mathbb{R}$ ) numbers, but also **strings** ( $X^*$ ) and infinite **sequences** ( $X^\omega$ ) of objects.  $A$  is a **subset** of  $B$ ,  $A \subseteq B$ , if every element of  $A$  is also one of  $B$ . The collection of all subsets of a set  $B$  is itself a set, the **power set**  $\wp(B)$  of  $B$ . We can form the **union**  $A \cup B$  and **intersection**  $A \cap B$  of sets. An **ordered pair**  $\langle x, y \rangle$  consists of two objects  $x$  and  $y$ , but in that specific order. The pairs  $\langle x, y \rangle$  and  $\langle y, x \rangle$  are different pairs (unless  $x = y$ ). The set of all pairs  $\langle x, y \rangle$  where  $x \in A$  and  $y \in B$  is called the **Cartesian product**  $A \times B$  of  $A$  and  $B$ . We write  $A^2$  for  $A \times A$ ; so for instance  $\mathbb{N}^2$  is the set of pairs of natural numbers.

## Problems

**Problem 1.1.** Prove that there is at most one empty set, i.e., show that if  $A$  and  $B$  are sets without elements, then  $A = B$ .

**Problem 1.2.** List all subsets of  $\{a, b, c, d\}$ .

**Problem 1.3.** Show that if  $A$  has  $n$  elements, then  $\wp(A)$  has  $2^n$  elements.

**Problem 1.4.** Prove that if  $A \subseteq B$ , then  $A \cup B = B$ .

**Problem 1.5.** Prove rigorously that if  $A \subseteq B$ , then  $A \cap B = A$ .

**Problem 1.6.** Show that if  $A$  is a set and  $A \in B$ , then  $A \subseteq \bigcup B$ .

**Problem 1.7.** Prove that if  $A \subsetneq B$ , then  $B \setminus A \neq \emptyset$ .

**Problem 1.8.** Using **Definition 1.23**, prove that  $\langle a, b \rangle = \langle c, d \rangle$  iff both  $a = c$  and  $b = d$ .

**Problem 1.9.** List all elements of  $\{1, 2, 3\}^3$ .

**Problem 1.10.** Show, by induction on  $k$ , that for all  $k \geq 1$ , if  $A$  has  $n$  elements, then  $A^k$  has  $n^k$  elements.

## CHAPTER 2

# Relations

### 2.1 Relations as Sets

In [section 1.3](#), we mentioned some important sets:  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{Q}$ ,  $\mathbb{R}$ . You will no doubt remember some interesting relations between the elements of some of these sets. For instance, each of these sets has a completely standard *order relation* on it. There is also the relation *is identical with* that every object bears to itself and to no other thing. There are many more interesting relations that we'll encounter, and even more possible relations. Before we review them, though, we will start by pointing out that we can look at relations as a special sort of set.

For this, recall two things from [section 1.5](#). First, recall the notion of a *ordered pair*: given  $a$  and  $b$ , we can form  $\langle a, b \rangle$ . Importantly, the order of elements *does* matter here. So if  $a \neq b$  then  $\langle a, b \rangle \neq \langle b, a \rangle$ . (Contrast this with unordered pairs, i.e., 2-element sets, where  $\{a, b\} = \{b, a\}$ .) Second, recall the notion of a *Cartesian product*: if  $A$  and  $B$  are sets, then we can form  $A \times B$ , the set of all pairs  $\langle x, y \rangle$  with  $x \in A$  and  $y \in B$ . In particular,  $A^2 = A \times A$  is the set of all ordered pairs from  $A$ .

Now we will consider a particular relation on a set: the  $<$ -relation on the set  $\mathbb{N}$  of natural numbers. Consider the set of all pairs of numbers  $\langle n, m \rangle$  where  $n < m$ , i.e.,

$$R = \{\langle n, m \rangle : n, m \in \mathbb{N} \text{ and } n < m\}.$$

There is a close connection between  $n$  being less than  $m$ , and the pair  $\langle n, m \rangle$  being a member of  $R$ , namely:

$$n < m \text{ iff } \langle n, m \rangle \in R.$$

Indeed, without any loss of information, we can consider the set  $R$  to be the  $<$ -relation on  $\mathbb{N}$ .

In the same way we can construct a subset of  $\mathbb{N}^2$  for any relation between numbers. Conversely, given any set of pairs of numbers  $S \subseteq \mathbb{N}^2$ , there is a corresponding relation between numbers, namely, the relationship  $n$  bears to  $m$  if and only if  $\langle n, m \rangle \in S$ . This justifies the following definition:

**Definition 2.1 (Binary relation).** A *binary relation* on a set  $A$  is a subset of  $A^2$ . If  $R \subseteq A^2$  is a binary relation on  $A$  and  $x, y \in A$ , we sometimes write  $Rxy$  (or  $xRy$ ) for  $\langle x, y \rangle \in R$ .

**Example 2.2.** The set  $\mathbb{N}^2$  of pairs of natural numbers can be listed in a 2-dimensional matrix like this:

$$\begin{array}{ccccccc} \langle \mathbf{0}, \mathbf{0} \rangle & \langle 0, 1 \rangle & \langle 0, 2 \rangle & \langle 0, 3 \rangle & \dots & & \\ \langle 1, 0 \rangle & \langle \mathbf{1}, \mathbf{1} \rangle & \langle 1, 2 \rangle & \langle 1, 3 \rangle & \dots & & \\ \langle 2, 0 \rangle & \langle 2, 1 \rangle & \langle \mathbf{2}, \mathbf{2} \rangle & \langle 2, 3 \rangle & \dots & & \\ \langle 3, 0 \rangle & \langle 3, 1 \rangle & \langle 3, 2 \rangle & \langle \mathbf{3}, \mathbf{3} \rangle & \dots & & \\ \vdots & \vdots & \vdots & \vdots & \ddots & & \end{array}$$

We have put the diagonal, here, in bold, since the subset of  $\mathbb{N}^2$  consisting of the pairs lying on the diagonal, i.e.,

$$\{\langle 0, 0 \rangle, \langle 1, 1 \rangle, \langle 2, 2 \rangle, \dots\},$$

is the *identity relation* on  $\mathbb{N}$ . (Since the identity relation is popular, let's define  $\text{Id}_A = \{\langle x, x \rangle : x \in A\}$  for any set  $A$ .) The subset of all pairs lying above the diagonal, i.e.,

$$L = \{\langle 0, 1 \rangle, \langle 0, 2 \rangle, \dots, \langle 1, 2 \rangle, \langle 1, 3 \rangle, \dots, \langle 2, 3 \rangle, \langle 2, 4 \rangle, \dots\},$$



is the *less than* relation, i.e.,  $Lnm$  iff  $n < m$ . The subset of pairs below the diagonal, i.e.,

$$G = \{\langle 1, 0 \rangle, \langle 2, 0 \rangle, \langle 2, 1 \rangle, \langle 3, 0 \rangle, \langle 3, 1 \rangle, \langle 3, 2 \rangle, \dots\},$$

is the *greater than* relation, i.e.,  $Gnm$  iff  $n > m$ . The union of  $L$  with  $I$ , which we might call  $K = L \cup I$ , is the *less than or equal to* relation:  $Knm$  iff  $n \leq m$ . Similarly,  $H = G \cup I$  is the *greater than or equal to relation*. These relations  $L$ ,  $G$ ,  $K$ , and  $H$  are special kinds of relations called *orders*.  $L$  and  $G$  have the property that no number bears  $L$  or  $G$  to itself (i.e., for all  $n$ , neither  $Lnn$  nor  $Gnn$ ). Relations with this property are called *irreflexive*, and, if they also happen to be orders, they are called *strict orders*.

Although orders and identity are important and natural relations, it should be emphasized that according to our definition *any* subset of  $A^2$  is a relation on  $A$ , regardless of how unnatural or contrived it seems. In particular,  $\emptyset$  is a relation on any set (the *empty relation*, which no pair of elements bears), and  $A^2$  itself is a relation on  $A$  as well (one which every pair bears), called the *universal relation*. But also something like  $E = \{\langle n, m \rangle : n > 5 \text{ or } m \times n \geq 34\}$  counts as a relation.

## 2.2 Special Properties of Relations

Some kinds of relations turn out to be so common that they have been given special names. For instance,  $\leq$  and  $\subseteq$  both relate their respective domains (say,  $\mathbb{N}$  in the case of  $\leq$  and  $\wp(A)$  in the case of  $\subseteq$ ) in similar ways. To get at exactly how these relations are similar, and how they differ, we categorize them according to some special properties that relations can have. It turns out that (combinations of) some of these special properties are especially important: orders and equivalence relations.

**Definition 2.3 (Reflexivity).** A relation  $R \subseteq A^2$  is *reflexive* iff, for every  $x \in A$ ,  $Rxx$ .

**Definition 2.4 (Transitivity).** A relation  $R \subseteq A^2$  is *transitive* iff, whenever  $Rxy$  and  $Ryz$ , then also  $Rxz$ .

**Definition 2.5 (Symmetry).** A relation  $R \subseteq A^2$  is *symmetric* iff, whenever  $Rxy$ , then also  $Ryx$ .

**Definition 2.6 (Anti-symmetry).** A relation  $R \subseteq A^2$  is *anti-symmetric* iff, whenever both  $Rxy$  and  $Ryx$ , then  $x = y$  (or, in other words: if  $x \neq y$  then either  $\neg Rxy$  or  $\neg Ryx$ ).

In a symmetric relation,  $Rxy$  and  $Ryx$  always hold together, or neither holds. In an anti-symmetric relation, the only way for  $Rxy$  and  $Ryx$  to hold together is if  $x = y$ . Note that this does not *require* that  $Rxy$  and  $Ryx$  holds when  $x = y$ , only that it isn't ruled out. So an anti-symmetric relation can be reflexive, but it is not the case that every anti-symmetric relation is reflexive. Also note that being anti-symmetric and merely not being symmetric are different conditions. In fact, a relation can be both symmetric and anti-symmetric at the same time (e.g., the identity relation is).

**Definition 2.7 (Connectivity).** A relation  $R \subseteq A^2$  is *connected* if for all  $x, y \in A$ , if  $x \neq y$ , then either  $Rxy$  or  $Ryx$ .

**Definition 2.8 (Irreflexivity).** A relation  $R \subseteq A^2$  is called *irreflexive* if, for all  $x \in A$ , not  $Rxx$ .

**Definition 2.9 (Asymmetry).** A relation  $R \subseteq A^2$  is called *asymmetric* if for no pair  $x, y \in A$  we have both  $Rxy$  and  $Ryx$ .

Note that if  $A \neq \emptyset$ , then no irreflexive relation on  $A$  is reflexive and every asymmetric relation on  $A$  is also anti-symmetric. However, there are  $R \subseteq A^2$  that are not reflexive and also not irreflexive, and there are anti-symmetric relations that are not asymmetric.

## 2.3 Equivalence Relations

The identity relation on a set is reflexive, symmetric, and transitive. Relations  $R$  that have all three of these properties are very common.

**Definition 2.10 (Equivalence relation).** A relation  $R \subseteq A^2$  that is reflexive, symmetric, and transitive is called an *equivalence relation*. Elements  $x$  and  $y$  of  $A$  are said to be *R-equivalent* if  $Rxy$ .

Equivalence relations give rise to the notion of an *equivalence class*. An equivalence relation “chunks up” the domain into different partitions. Within each partition, all the objects are related to one another; and no objects from different partitions relate to one another. Sometimes, it’s helpful just to talk about these partitions *directly*. To that end, we introduce a definition:

**Definition 2.11.** Let  $R \subseteq A^2$  be an equivalence relation. For each  $x \in A$ , the *equivalence class* of  $x$  in  $A$  is the set  $[x]_R = \{y \in A : Rxy\}$ . The *quotient* of  $A$  under  $R$  is  $A/R = \{[x]_R : x \in A\}$ , i.e., the set of these equivalence classes.

The next result vindicates the definition of an equivalence class, in proving that the equivalence classes are indeed the partitions of  $A$ :

**Proposition 2.12.** *If  $R \subseteq A^2$  is an equivalence relation, then  $Rxy$  iff  $[x]_R = [y]_R$ .*

*Proof.* For the left-to-right direction, suppose  $Rxy$ , and let  $z \in [x]_R$ . By definition, then,  $Rxz$ . Since  $R$  is an equivalence relation,  $Ryz$ . (Spelling this out: as  $Rxy$  and  $R$  is symmetric we have  $Ryx$ , and as  $Rxz$  and  $R$  is transitive we have  $Ryz$ .) So  $z \in [y]_R$ . Generalising,  $[x]_R \subseteq [y]_R$ . But exactly similarly,  $[y]_R \subseteq [x]_R$ . So  $[x]_R = [y]_R$ , by extensionality.

For the right-to-left direction, suppose  $[x]_R = [y]_R$ . Since  $R$  is reflexive,  $Ryy$ , so  $y \in [y]_R$ . Thus also  $y \in [x]_R$  by the assumption that  $[x]_R = [y]_R$ . So  $Rxy$ .  $\square$

**Example 2.13.** A nice example of equivalence relations comes from modular arithmetic. For any  $a, b$ , and  $n \in \mathbb{Z}^+$ , say that  $a \equiv_n b$  iff dividing  $a$  by  $n$  gives the same remainder as dividing  $b$  by  $n$ . (Somewhat more symbolically:  $a \equiv_n b$  iff, for some  $k \in \mathbb{Z}$ ,  $a - b = kn$ .) Now,  $\equiv_n$  is an equivalence relation, for any  $n$ . And there are exactly  $n$  distinct equivalence classes generated by  $\equiv_n$ ; that is,  $\mathbb{N}/\equiv_n$  has  $n$  elements. These are: the set of numbers divisible by  $n$  without remainder, i.e.,  $[0]_{\equiv_n}$ ; the set of numbers divisible by  $n$  with remainder 1, i.e.,  $[1]_{\equiv_n}$ ; ...; and the set of numbers divisible by  $n$  with remainder  $n - 1$ , i.e.,  $[n - 1]_{\equiv_n}$ .

## 2.4 Orders

Many of our comparisons involve describing some objects as being “less than”, “equal to”, or “greater than” other objects, in a certain respect. These involve *order* relations. But there are different kinds of order relations. For instance, some require that any two objects be comparable, others don’t. Some include identity (like  $\leq$ ) and some exclude it (like  $<$ ). It will help us to have a taxonomy here.

**Definition 2.14 (Preorder).** A relation which is both reflexive and transitive is called a *preorder*.

**Definition 2.15 (Partial order).** A preorder which is also anti-symmetric is called a *partial order*.

**Definition 2.16 (Linear order).** A partial order which is also connected is called a *total order* or *linear order*.

**Example 2.17.** Every linear order is also a partial order, and every partial order is also a preorder, but the converses don't hold. The universal relation on  $A$  is a preorder, since it is reflexive and transitive. But, if  $A$  has more than one element, the universal relation is not anti-symmetric, and so not a partial order.

**Example 2.18.** Consider the *no longer than* relation  $\preceq$  on  $\mathbb{B}^*$ :  $x \preceq y$  iff  $\text{len}(x) \leq \text{len}(y)$ . This is a preorder (reflexive and transitive), and even connected, but not a partial order, since it is not anti-symmetric. For instance,  $01 \preceq 10$  and  $10 \preceq 01$ , but  $01 \neq 10$ .

**Example 2.19.** An important partial order is the relation  $\subseteq$  on a set of sets. This is not in general a linear order, since if  $a \neq b$  and we consider  $\wp(\{a, b\}) = \{\emptyset, \{a\}, \{b\}, \{a, b\}\}$ , we see that  $\{a\} \not\subseteq \{b\}$  and  $\{a\} \neq \{b\}$  and  $\{b\} \not\subseteq \{a\}$ .

**Example 2.20.** The relation of *divisibility without remainder* gives us a partial order which isn't a linear order. For integers  $n, m$ , we write  $n \mid m$  to mean  $n$  (evenly) divides  $m$ , i.e., iff there is some integer  $k$  so that  $m = kn$ . On  $\mathbb{N}$ , this is a partial order, but not a linear order: for instance,  $2 \nmid 3$  and also  $3 \nmid 2$ . Considered as a relation on  $\mathbb{Z}$ , divisibility is only a preorder since it is not anti-symmetric:  $1 \mid -1$  and  $-1 \mid 1$  but  $1 \neq -1$ .

**Definition 2.21 (Strict order).** A *strict order* is a relation which is irreflexive, asymmetric, and transitive.

**Definition 2.22 (Strict linear order).** A strict order which is also connected is called a *strict total order* or *strict linear order*.

**Example 2.23.**  $\leq$  is the linear order corresponding to the strict linear order  $<$ .  $\subseteq$  is the partial order corresponding to the strict order  $\subsetneq$ .

Any strict order  $R$  on  $A$  can be turned into a partial order by adding the diagonal  $\text{Id}_A$ , i.e., adding all the pairs  $\langle x, x \rangle$ . (This is called the *reflexive closure* of  $R$ .) Conversely, starting from a partial order, one can get a strict order by removing  $\text{Id}_A$ . These next two results make this precise.

**Proposition 2.24.** *If  $R$  is a strict order on  $A$ , then  $R^+ = R \cup \text{Id}_A$  is a partial order. Moreover, if  $R$  is a strict linear order, then  $R^+$  is a linear order.*

*Proof.* Suppose  $R$  is a strict order, i.e.,  $R \subseteq A^2$  and  $R$  is irreflexive, asymmetric, and transitive. Let  $R^+ = R \cup \text{Id}_A$ . We have to show that  $R^+$  is reflexive, anti-symmetric, and transitive.

$R^+$  is clearly reflexive, since  $\langle x, x \rangle \in \text{Id}_A \subseteq R^+$  for all  $x \in A$ .

To show  $R^+$  is anti-symmetric, suppose for reductio that  $R^+xy$  and  $R^+yx$  but  $x \neq y$ . Since  $\langle x, y \rangle \in R \cup \text{Id}_A$ , but  $\langle x, y \rangle \notin \text{Id}_A$ , we must have  $\langle x, y \rangle \in R$ , i.e.,  $Rxy$ . Similarly,  $Ryx$ . But this contradicts the assumption that  $R$  is asymmetric.

To establish transitivity, suppose that  $R^+xy$  and  $R^+yz$ . If both  $\langle x, y \rangle \in R$  and  $\langle y, z \rangle \in R$ , then  $\langle x, z \rangle \in R$  since  $R$  is transitive. Otherwise, either  $\langle x, y \rangle \in \text{Id}_A$ , i.e.,  $x = y$ , or  $\langle y, z \rangle \in \text{Id}_A$ , i.e.,  $y = z$ . In the first case, we have that  $R^+yz$  by assumption,  $x = y$ , hence  $R^+xz$ . Similarly in the second case. In either case,  $R^+xz$ , thus,  $R^+$  is also transitive.

Concerning the “moreover” clause, suppose that  $R$  is also connected. So for all  $x \neq y$ , either  $Rxy$  or  $Ryx$ , i.e., either  $\langle x, y \rangle \in R$  or  $\langle y, x \rangle \in R$ . Since  $R \subseteq R^+$ , this remains true of  $R^+$ , so  $R^+$  is connected as well.  $\square$

**Proposition 2.25.** *If  $R$  is a partial order on  $A$ , then  $R^- = R \setminus \text{Id}_A$  is a strict order. Moreover, if  $R$  is a linear order, then  $R^-$  is a strict linear order.*

*Proof.* This is left as an exercise.  $\square$

The following simple result establishes that strict linear orders satisfy an extensionality-like property:

**Proposition 2.26.** *If  $<$  is a strict linear order on  $A$ , then:*

$$(\forall a, b \in A)((\forall x \in A)(x < a \leftrightarrow x < b) \rightarrow a = b).$$

*Proof.* Suppose  $(\forall x \in A)(x < a \leftrightarrow x < b)$ . If  $a < b$ , then  $a < a$ , contradicting the fact that  $<$  is irreflexive; so  $a \not< b$ . Exactly similarly,  $b \not< a$ . So  $a = b$ , as  $<$  is connected.  $\square$

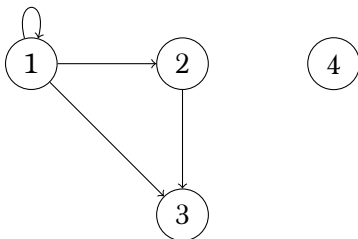
## 2.5 Graphs

A *graph* is a diagram in which points—called “nodes” or “vertices” (plural of “vertex”)—are connected by edges. Graphs are a ubiquitous tool in discrete mathematics and in computer science. They are incredibly useful for representing, and visualizing, relationships and structures, from concrete things like networks of various kinds to abstract structures such as the possible outcomes of decisions. There are many different kinds of graphs in the literature which differ, e.g., according to whether the edges are directed or not, have labels or not, whether there can be edges from a node to the same node, multiple edges between the same nodes, etc. *Directed graphs* have a special connection to relations.

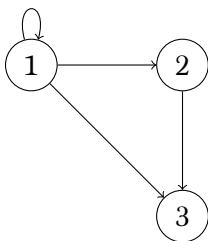
**Definition 2.27 (Directed graph).** A *directed graph*  $G = \langle V, E \rangle$  is a set of *vertices*  $V$  and a set of *edges*  $E \subseteq V^2$ .

According to our definition, a graph just is a set together with a relation on that set. Of course, when talking about graphs, it's only natural to expect that they are graphically represented: we can draw a graph by connecting two vertices  $v_1$  and  $v_2$  by an arrow iff  $\langle v_1, v_2 \rangle \in E$ . The only difference between a relation by itself and a graph is that a graph specifies the set of vertices, i.e., a graph may have isolated vertices. The important point, however, is that every relation  $R$  on a set  $X$  can be seen as a directed graph  $\langle X, R \rangle$ , and conversely, a directed graph  $\langle V, E \rangle$  can be seen as a relation  $E \subseteq V^2$  with the set  $V$  explicitly specified.

**Example 2.28.** The graph  $\langle V, E \rangle$  with  $V = \{1, 2, 3, 4\}$  and  $E = \{\langle 1, 1 \rangle, \langle 1, 2 \rangle, \langle 1, 3 \rangle, \langle 2, 3 \rangle\}$  looks like this:



This is a different graph than  $\langle V', E \rangle$  with  $V' = \{1, 2, 3\}$ , which looks like this:





## 2.6 Operations on Relations

It is often useful to modify or combine relations. In **Proposition 2.24**, we considered the *union* of relations, which is just the union of two relations considered as sets of pairs. Similarly, in **Proposition 2.25**, we considered the relative difference of relations. Here are some other operations we can perform on relations.

**Definition 2.29.** Let  $R, S$  be relations, and  $A$  be any set.

The *inverse* of  $R$  is  $R^{-1} = \{\langle y, x \rangle : \langle x, y \rangle \in R\}$ .

The *relative product* of  $R$  and  $S$  is  $(R \mid S) = \{\langle x, z \rangle : \exists y(Rxy \wedge Syz)\}$ .

The *restriction* of  $R$  to  $A$  is  $R \upharpoonright_A = R \cap A^2$ .

The *application* of  $R$  to  $A$  is  $R[A] = \{y : (\exists x \in A)Rxy\}$

**Example 2.30.** Let  $S \subseteq \mathbb{Z}^2$  be the successor relation on  $\mathbb{Z}$ , i.e.,  $S = \{\langle x, y \rangle \in \mathbb{Z}^2 : x + 1 = y\}$ , so that  $Sxy$  iff  $x + 1 = y$ .

$S^{-1}$  is the predecessor relation on  $\mathbb{Z}$ , i.e.,  $\{\langle x, y \rangle \in \mathbb{Z}^2 : x - 1 = y\}$ .

$S \mid S$  is  $\{\langle x, y \rangle \in \mathbb{Z}^2 : x + 2 = y\}$

$S \upharpoonright_{\mathbb{N}}$  is the successor relation on  $\mathbb{N}$ .

$S[\{1, 2, 3\}]$  is  $\{2, 3, 4\}$ .

**Definition 2.31 (Transitive closure).** Let  $R \subseteq A^2$  be a binary relation.

The *transitive closure* of  $R$  is  $R^+ = \bigcup_{0 < n \in \mathbb{N}} R^n$ , where we recursively define  $R^1 = R$  and  $R^{n+1} = R^n \mid R$ .

The *reflexive transitive closure* of  $R$  is  $R^* = R^+ \cup \text{Id}_A$ .

**Example 2.32.** Take the successor relation  $S \subseteq \mathbb{Z}^2$ .  $S^2xy$  iff  $x + 2 = y$ ,  $S^3xy$  iff  $x + 3 = y$ , etc. So  $S^+xy$  iff  $x + n = y$  for some  $n \geq 1$ . In other words,  $S^+xy$  iff  $x < y$ , and  $S^*xy$  iff  $x \leq y$ .

## Summary

A **relation**  $R$  on a set  $A$  is a way of relating elements of  $A$ . We write  $Rxy$  if the relation **holds** between  $x$  and  $y$ . Formally, we can consider  $R$  as the sets of pairs  $\langle x, y \rangle \in A^2$  such that  $Rxy$ . Being less than, greater than, equal to, evenly dividing, being the same length as, a subset of, and the same size as are all important examples of relations (on sets of numbers, strings, or of sets). **Graphs** are a general way of visually representing relations. But a graph can also be seen as a binary relation (the **edge** relation) together with the underlying set of **vertices**.

Some relations share certain features which makes them especially interesting or useful. A relation  $R$  is **reflexive** if everything is  $R$ -related to itself; **symmetric**, if with  $Rxy$  also  $Ryx$  holds for any  $x$  and  $y$ ; and **transitive** if  $Rxy$  and  $Ryz$  guarantees  $Rxz$ . Relations that have all three of these properties are **equivalence relations**. A relation is **anti-symmetric** if  $Rxy$  and  $Ryx$  guarantees  $x = y$ . **Partial orders** are those relations that are reflexive, anti-symmetric, and transitive. A **linear order** is any partial order which satisfies that for any  $x$  and  $y$ , either  $Rxy$  or  $x = y$  or  $Ryx$ . (Generally, a relation with this property is **connected**).

Since relations are sets (of pairs), they can be operated on as sets (e.g., we can form the union and intersection of relations). We can also chain them together (**relative product**  $R | S$ ). If we form the relative product of  $R$  with itself arbitrarily many times we get the **transitive closure**  $R^+$  of  $R$ .

## Problems

**Problem 2.1.** List the elements of the relation  $\subseteq$  on the set  $\wp(\{a, b, c\})$ .

**Problem 2.2.** Give examples of relations that are (a) reflexive and symmetric but not transitive, (b) reflexive and anti-symmetric, (c) anti-symmetric, transitive, but not reflexive, and

(d) reflexive, symmetric, and transitive. Do not use relations on numbers or sets.

**Problem 2.3.** Show that  $\equiv_n$  is an equivalence relation, for any  $n \in \mathbb{Z}^+$ , and that  $\mathbb{N}/\equiv_n$  has exactly  $n$  members.

**Problem 2.4.** Give a proof of [Proposition 2.25](#).

**Problem 2.5.** Consider the less-than-or-equal-to relation  $\leq$  on the set  $\{1, 2, 3, 4\}$  as a graph and draw the corresponding diagram.

**Problem 2.6.** Show that the transitive closure of  $R$  is in fact transitive.

## CHAPTER 3

# Functions

### 3.1 Basics

A *function* is a map which sends each element of a given set to a specific element in some (other) given set. For instance, the operation of adding 1 defines a function: each number  $n$  is mapped to a unique number  $n + 1$ .

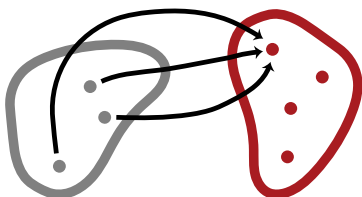
More generally, functions may take pairs, triples, etc., as inputs and return some kind of output. Many functions are familiar to us from basic arithmetic. For instance, addition and multiplication are functions. They take in two numbers and return a third.

In this mathematical, abstract sense, a function is a *black box*: what matters is only what output is paired with what input, not the method for calculating the output.

**Definition 3.1 (Function).** A *function*  $f: A \rightarrow B$  is a mapping of each element of  $A$  to an element of  $B$ .

We call  $A$  the *domain* of  $f$  and  $B$  the *codomain* of  $f$ . The elements of  $A$  are called inputs or *arguments* of  $f$ , and the element of  $B$  that is paired with an argument  $x$  by  $f$  is called the *value of  $f$*  for argument  $x$ , written  $f(x)$ .

The *range*  $\text{ran}(f)$  of  $f$  is the subset of the codomain consisting of the values of  $f$  for some argument;  $\text{ran}(f) = \{f(x) : x \in A\}$ .



*Figure 3.1:* A function is a mapping of each element of one set to an element of another. An arrow points from an argument in the domain to the corresponding value in the codomain.

The diagram in [Figure 3.1](#) may help to think about functions. The ellipse on the left represents the function’s *domain*; the ellipse on the right represents the function’s *codomain*; and an arrow points from an *argument* in the domain to the corresponding *value* in the codomain.

**Example 3.2.** Multiplication takes pairs of natural numbers as inputs and maps them to natural numbers as outputs, so goes from  $\mathbb{N} \times \mathbb{N}$  (the domain) to  $\mathbb{N}$  (the codomain). As it turns out, the range is also  $\mathbb{N}$ , since every  $n \in \mathbb{N}$  is  $n \times 1$ .

**Example 3.3.** Multiplication is a function because it pairs each input—each pair of natural numbers—with a single output:  $\times: \mathbb{N}^2 \rightarrow \mathbb{N}$ . By contrast, the square root operation applied to the domain  $\mathbb{N}$  is not functional, since each positive integer  $n$  has two square roots:  $\sqrt{n}$  and  $-\sqrt{n}$ . We can make it functional by only returning the positive square root:  $\sqrt{\phantom{x}}: \mathbb{N} \rightarrow \mathbb{R}$ .

**Example 3.4.** The relation that pairs each student in a class with their final grade is a function—no student can get two different final grades in the same class. The relation that pairs each student in a class with their parents is not a function: students can have zero, or two, or more parents.

We can define functions by specifying in some precise way what the value of the function is for every possible argument.

Different ways of doing this are by giving a formula, describing a method for computing the value, or listing the values for each argument. However functions are defined, we must make sure that for each argument we specify one, and only one, value.

**Example 3.5.** Let  $f: \mathbb{N} \rightarrow \mathbb{N}$  be defined such that  $f(x) = x + 1$ . This is a definition that specifies  $f$  as a function which takes in natural numbers and outputs natural numbers. It tells us that, given a natural number  $x$ ,  $f$  will output its successor  $x + 1$ . In this case, the codomain  $\mathbb{N}$  is not the range of  $f$ , since the natural number 0 is not the successor of any natural number. The range of  $f$  is the set of all positive integers,  $\mathbb{Z}^+$ .

**Example 3.6.** Let  $g: \mathbb{N} \rightarrow \mathbb{N}$  be defined such that  $g(x) = x + 2 - 1$ . This tells us that  $g$  is a function which takes in natural numbers and outputs natural numbers. Given a natural number  $n$ ,  $g$  will output the predecessor of the successor of the successor of  $x$ , i.e.,  $x + 1$ .

We just considered two functions,  $f$  and  $g$ , with different *definitions*. However, these are the *same function*. After all, for any natural number  $n$ , we have that  $f(n) = n + 1 = n + 2 - 1 = g(n)$ . Otherwise put: our definitions for  $f$  and  $g$  specify the same mapping by means of different equations. Implicitly, then, we are relying upon a principle of extensionality for functions,

$$\text{if } \forall x f(x) = g(x), \text{ then } f = g$$

provided that  $f$  and  $g$  share the same domain and codomain.

**Example 3.7.** We can also define functions by cases. For instance, we could define  $h: \mathbb{N} \rightarrow \mathbb{N}$  by

$$h(x) = \begin{cases} \frac{x}{2} & \text{if } x \text{ is even} \\ \frac{x+1}{2} & \text{if } x \text{ is odd.} \end{cases}$$

Since every natural number is either even or odd, the output of this function will always be a natural number. Just remember that

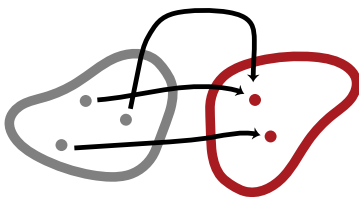


Figure 3.2: A surjective function has every element of the codomain as a value.

if you define a function by cases, every possible input must fall into exactly one case. In some cases, this will require a proof that the cases are exhaustive and exclusive.

## 3.2 Kinds of Functions

It will be useful to introduce a kind of taxonomy for some of the kinds of functions which we encounter most frequently.

To start, we might want to consider functions which have the property that every member of the codomain is a value of the function. Such functions are called surjective, and can be pictured as in Figure 3.2.

**Definition 3.8 (Surjective function).** A function  $f: A \rightarrow B$  is *surjective* iff  $B$  is also the range of  $f$ , i.e., for every  $y \in B$  there is at least one  $x \in A$  such that  $f(x) = y$ , or in symbols:

$$(\forall y \in B)(\exists x \in A)f(x) = y.$$

We call such a function a surjection from  $A$  to  $B$ .

If you want to show that  $f$  is a surjection, then you need to show that every object in  $f$ 's codomain is the value of  $f(x)$  for some input  $x$ .

Note that any function *induces* a surjection. After all, given a function  $f: A \rightarrow B$ , let  $f': A \rightarrow \text{ran}(f)$  be defined by  $f'(x) =$



*Figure 3.3:* An injective function never maps two different arguments to the same value.

$f(x)$ . Since  $\text{ran}(f)$  is *defined* as  $\{f(x) \in B : x \in A\}$ , this function  $f'$  is guaranteed to be a surjection

Now, any function maps each possible input to a unique output. But there are also functions which never map different inputs to the same outputs. Such functions are called injective, and can be pictured as in **Figure 3.3**.

**Definition 3.9 (Injective function).** A function  $f: A \rightarrow B$  is *injective* iff for each  $y \in B$  there is at most one  $x \in A$  such that  $f(x) = y$ . We call such a function an injection from  $A$  to  $B$ .

If you want to show that  $f$  is an injection, you need to show that for any elements  $x$  and  $y$  of  $f$ 's domain, if  $f(x) = f(y)$ , then  $x = y$ .

**Example 3.10.** The constant function  $f: \mathbb{N} \rightarrow \mathbb{N}$  given by  $f(x) = 1$  is neither injective, nor surjective.

The identity function  $f: \mathbb{N} \rightarrow \mathbb{N}$  given by  $f(x) = x$  is both injective and surjective.

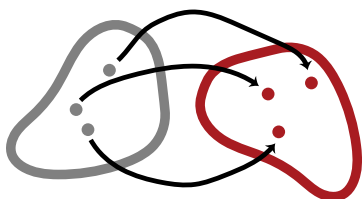
The successor function  $f: \mathbb{N} \rightarrow \mathbb{N}$  given by  $f(x) = x + 1$  is injective but not surjective.

The function  $f: \mathbb{N} \rightarrow \mathbb{N}$  defined by:

$$f(x) = \begin{cases} \frac{x}{2} & \text{if } x \text{ is even} \\ \frac{x+1}{2} & \text{if } x \text{ is odd.} \end{cases}$$

is surjective, but not injective.





*Figure 3.4:* A bijective function uniquely pairs the elements of the codomain with those of the domain.

Often enough, we want to consider functions which are both injective and surjective. We call such functions bijective. They look like the function pictured in [Figure 3.4](#). Bijections are also sometimes called *one-to-one correspondences*, since they uniquely pair elements of the codomain with elements of the domain.

**Definition 3.11 (Bijection).** A function  $f: A \rightarrow B$  is *bijective* iff it is both surjective and injective. We call such a function a bijection from  $A$  to  $B$  (or between  $A$  and  $B$ ).

### 3.3 Functions as Relations

A function which maps elements of  $A$  to elements of  $B$  obviously defines a relation between  $A$  and  $B$ , namely the relation which holds between  $x$  and  $y$  iff  $f(x) = y$ . In fact, we might even—if we are interested in reducing the building blocks of mathematics for instance—*identify* the function  $f$  with this relation, i.e., with a set of pairs. This then raises the question: which relations define functions in this way?

**Definition 3.12 (Graph of a function).** Let  $f: A \rightarrow B$  be a function. The *graph* of  $f$  is the relation  $R_f \subseteq A \times B$  defined by

$$R_f = \{\langle x, y \rangle : f(x) = y\}.$$

The graph of a function is uniquely determined, by extensionality. Moreover, extensionality (on sets) will immediately vindicate the implicit principle of extensionality for functions, whereby if  $f$  and  $g$  share a domain and codomain then they are identical if they agree on all values.

Similarly, if a relation is “functional”, then it is the graph of a function.

**Proposition 3.13.** *Let  $R \subseteq A \times B$  be such that:*

1. *If  $Rxy$  and  $Rxz$  then  $y = z$ ; and*
2. *for every  $x \in A$  there is some  $y \in B$  such that  $\langle x, y \rangle \in R$ .*

*Then  $R$  is the graph of the function  $f: A \rightarrow B$  defined by  $f(x) = y$  iff  $Rxy$ .*

*Proof.* Suppose there is a  $y$  such that  $Rxy$ . If there were another  $z \neq y$  such that  $Rxz$ , the condition on  $R$  would be violated. Hence, if there is a  $y$  such that  $Rxy$ , this  $y$  is unique, and so  $f$  is well-defined. Obviously,  $R_f = R$ .  $\square$

Every function  $f: A \rightarrow B$  has a graph, i.e., a relation on  $A \times B$  defined by  $f(x) = y$ . On the other hand, every relation  $R \subseteq A \times B$  with the properties given in **Proposition 3.13** is the graph of a function  $f: A \rightarrow B$ . Because of this close connection between functions and their graphs, we can think of a function simply as its graph. In other words, functions can be identified with certain relations, i.e., with certain sets of tuples. We can now consider performing similar operations on functions as we performed on relations (see **section 2.6**). In particular:

**Definition 3.14.** Let  $f: A \rightarrow B$  be a function with  $C \subseteq A$ .

The *restriction* of  $f$  to  $C$  is the function  $f \upharpoonright_C: C \rightarrow B$  defined by  $(f \upharpoonright_C)(x) = f(x)$  for all  $x \in C$ . In other words,  $f \upharpoonright_C = \{\langle x, y \rangle \in R_f : x \in C\}$ .

The *application* of  $f$  to  $C$  is  $f[C] = \{f(x) : x \in C\}$ . We also

call this the *image* of  $C$  under  $f$ .

It follows from these definitions that  $\text{ran}(f) = f[\text{dom}(f)]$ , for any function  $f$ . These notions are exactly as one would expect, given the definitions in [section 2.6](#) and our identification of functions with relations. But two other operations—inverses and relative products—require a little more detail. We will provide that in [section 3.4](#) and [section 3.5](#).

### 3.4 Inverses of Functions

We think of functions as maps. An obvious question to ask about functions, then, is whether the mapping can be “reversed.” For instance, the successor function  $f(x) = x + 1$  can be reversed, in the sense that the function  $g(y) = y - 1$  “undoes” what  $f$  does.

But we must be careful. Although the definition of  $g$  defines a function  $\mathbb{Z} \rightarrow \mathbb{Z}$ , it does not define a *function*  $\mathbb{N} \rightarrow \mathbb{N}$ , since  $g(0) \notin \mathbb{N}$ . So even in simple cases, it is not quite obvious whether a function can be reversed; it may depend on the domain and codomain.

This is made more precise by the notion of an inverse of a function.

**Definition 3.15.** A function  $g: B \rightarrow A$  is an *inverse* of a function  $f: A \rightarrow B$  if  $f(g(y)) = y$  and  $g(f(x)) = x$  for all  $x \in A$  and  $y \in B$ .

If  $f$  has an inverse  $g$ , we often write  $f^{-1}$  instead of  $g$ .

Now we will determine when functions have inverses. A good candidate for an inverse of  $f: A \rightarrow B$  is  $g: B \rightarrow A$  “defined by”

$$g(y) = \text{“the” } x \text{ such that } f(x) = y.$$

But the scare quotes around “defined by” (and “the”) suggest that this is not a definition. At least, it will not always work, with complete generality. For, in order for this definition to specify a

function, there has to be one and only one  $x$  such that  $f(x) = y$ —the output of  $g$  has to be uniquely specified. Moreover, it has to be specified for every  $y \in B$ . If there are  $x_1$  and  $x_2 \in A$  with  $x_1 \neq x_2$  but  $f(x_1) = f(x_2)$ , then  $g(y)$  would not be uniquely specified for  $y = f(x_1) = f(x_2)$ . And if there is no  $x$  at all such that  $f(x) = y$ , then  $g(y)$  is not specified at all. In other words, for  $g$  to be defined,  $f$  must be both injective and surjective.

Let's go slowly. We'll divide the question into two: Given a function  $f: A \rightarrow B$ , when is there a function  $g: B \rightarrow A$  so that  $g(f(x)) = x$ ? Such a  $g$  “undoes” what  $f$  does, and is called a *left inverse* of  $f$ . Secondly, when is there a function  $h: B \rightarrow A$  so that  $f(h(y)) = y$ ? Such an  $h$  is called a *right inverse* of  $f$ — $f$  “undoes” what  $h$  does.

**Proposition 3.16.** *If  $f: A \rightarrow B$  is injective, then there is a left inverse  $g: B \rightarrow A$  of  $f$  so that  $g(f(x)) = x$  for all  $x \in A$ .*

*Proof.* Suppose that  $f: A \rightarrow B$  is injective. Consider a  $y \in B$ . If  $y \in \text{ran}(f)$ , there is an  $x \in A$  so that  $f(x) = y$ . Because  $f$  is injective, there is only one such  $x \in A$ . Then we can define:  $g(y) = x$ , i.e.,  $g(y)$  is “the”  $x \in A$  such that  $f(x) = y$ . If  $y \notin \text{ran}(f)$ , we can map it to any  $a \in A$ . So, we can pick an  $a \in A$  and define  $g: B \rightarrow A$  by:

$$g(y) = \begin{cases} x & \text{if } f(x) = y \\ a & \text{if } y \notin \text{ran}(f). \end{cases}$$

It is defined for all  $y \in B$ , since for each such  $y \in \text{ran}(f)$  there is exactly one  $x \in A$  such that  $f(x) = y$ . By definition, if  $y = f(x)$ , then  $g(y) = x$ , i.e.,  $g(f(x)) = x$ .  $\square$

**Proposition 3.17.** *If  $f: A \rightarrow B$  is surjective, then there is a right inverse  $h: B \rightarrow A$  of  $f$  so that  $f(h(y)) = y$  for all  $y \in B$ .*

*Proof.* Suppose that  $f: A \rightarrow B$  is surjective. Consider a  $y \in B$ . Since  $f$  is surjective, there is an  $x_y \in A$  with  $f(x_y) = y$ . Then we can define:  $h(y) = x_y$ , i.e., for each  $y \in B$  we choose some  $x \in A$

so that  $f(x) = y$ ; since  $f$  is surjective there is always at least one to choose from.<sup>1</sup> By definition, if  $x = h(y)$ , then  $f(x) = y$ , i.e., for any  $y \in B$ ,  $f(h(y)) = y$ .  $\square$

By combining the ideas in the previous proof, we now get that every bijection has an inverse, i.e., there is a single function which is both a left and right inverse of  $f$ .

**Proposition 3.18.** *If  $f: A \rightarrow B$  is bijective, there is a function  $f^{-1}: B \rightarrow A$  so that for all  $x \in A$ ,  $f^{-1}(f(x)) = x$  and for all  $y \in B$ ,  $f(f^{-1}(y)) = y$ .*

*Proof.* Exercise.  $\square$

There is a slightly more general way to extract inverses. We saw in [section 3.2](#) that every function  $f$  induces a surjection  $f': A \rightarrow \text{ran}(f)$  by letting  $f'(x) = f(x)$  for all  $x \in A$ . Clearly, if  $f$  is injective, then  $f'$  is bijective, so that it has a unique inverse by [Proposition 3.18](#). By a very minor abuse of notation, we sometimes call the inverse of  $f'$  simply “the inverse of  $f$ .”

**Proposition 3.19.** *Show that if  $f: A \rightarrow B$  has a left inverse  $g$  and a right inverse  $h$ , then  $h = g$ .*

*Proof.* Exercise.  $\square$

---

<sup>1</sup>Since  $f$  is surjective, for every  $y \in B$  the set  $\{x : f(x) = y\}$  is nonempty. Our definition of  $h$  requires that we choose a single  $x$  from each of these sets. That this is always possible is actually not obvious—the possibility of making these choices is simply assumed as an axiom. In other words, this proposition assumes the so-called Axiom of Choice, an issue we will gloss over. However, in many specific cases, e.g., when  $A = \mathbb{N}$  or is finite, or when  $f$  is bijective, the Axiom of Choice is not required. (In the particular case when  $f$  is bijective, for each  $y \in B$  the set  $\{x : f(x) = y\}$  has exactly one element, so that there is no choice to make.)

**Proposition 3.20.** *Every function  $f$  has at most one inverse.*

*Proof.* Suppose  $g$  and  $h$  are both inverses of  $f$ . Then in particular  $g$  is a left inverse of  $f$  and  $h$  is a right inverse. By **Proposition 3.19**,  $g = h$ .  $\square$

## 3.5 Composition of Functions

We saw in **section 3.4** that the inverse  $f^{-1}$  of a bijection  $f$  is itself a function. Another operation on functions is composition: we can define a new function by composing two functions,  $f$  and  $g$ , i.e., by first applying  $f$  and then  $g$ . Of course, this is only possible if the ranges and domains match, i.e., the range of  $f$  must be a subset of the domain of  $g$ . This operation on functions is the analogue of the operation of relative product on relations from **section 2.6**.

A diagram might help to explain the idea of composition. In **Figure 3.5**, we depict two functions  $f: A \rightarrow B$  and  $g: B \rightarrow C$  and their composition  $(g \circ f)$ . The function  $(g \circ f): A \rightarrow C$  pairs each element of  $A$  with an element of  $C$ . We specify which element of  $C$  an element of  $A$  is paired with as follows: given an input  $x \in A$ , first apply the function  $f$  to  $x$ , which will output some  $f(x) = y \in B$ , then apply the function  $g$  to  $y$ , which will output some  $g(f(x)) = g(y) = z \in C$ .

**Definition 3.21 (Composition).** Let  $f: A \rightarrow B$  and  $g: B \rightarrow C$  be functions. The *composition* of  $f$  with  $g$  is  $g \circ f: A \rightarrow C$ , where  $(g \circ f)(x) = g(f(x))$ .

**Example 3.22.** Consider the functions  $f(x) = x + 1$ , and  $g(x) = 2x$ . Since  $(g \circ f)(x) = g(f(x))$ , for each input  $x$  you must first take its successor, then multiply the result by two. So their composition is given by  $(g \circ f)(x) = 2(x + 1)$ .

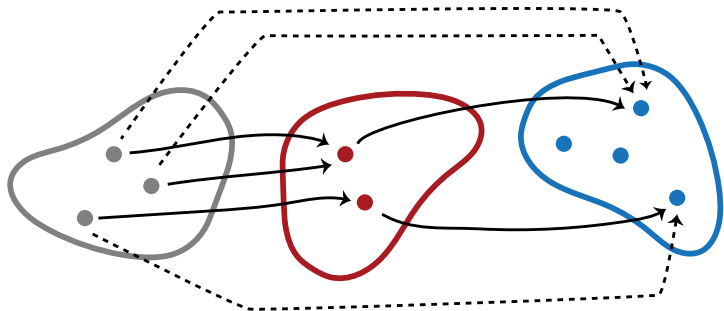


Figure 3.5: The composition  $g \circ f$  of two functions  $f$  and  $g$ .

## 3.6 Partial Functions

It is sometimes useful to relax the definition of function so that it is not required that the output of the function is defined for all possible inputs. Such mappings are called *partial functions*.

**Definition 3.23.** A *partial function*  $f: A \rightarrow B$  is a mapping which assigns to every element of  $A$  at most one element of  $B$ . If  $f$  assigns an element of  $B$  to  $x \in A$ , we say  $f(x)$  is *defined*, and otherwise *undefined*. If  $f(x)$  is defined, we write  $f(x) \downarrow$ , otherwise  $f(x) \uparrow$ . The *domain* of a partial function  $f$  is the subset of  $A$  where it is defined, i.e.,  $\text{dom}(f) = \{x \in A : f(x) \downarrow\}$ .

**Example 3.24.** Every function  $f: A \rightarrow B$  is also a partial function. Partial functions that are defined everywhere on  $A$ —i.e., what we so far have simply called a function—are also called *total functions*.

**Example 3.25.** The partial function  $f: \mathbb{R} \rightarrow \mathbb{R}$  given by  $f(x) = 1/x$  is undefined for  $x = 0$ , and defined everywhere else.

**Definition 3.26 (Graph of a partial function).** Let  $f: A \rightarrow B$  be a partial function. The *graph* of  $f$  is the relation  $R_f \subseteq A \times B$  defined by

$$R_f = \{\langle x, y \rangle : f(x) = y\}.$$

**Proposition 3.27.** *Suppose  $R \subseteq A \times B$  has the property that whenever  $Rxy$  and  $Rxy'$  then  $y = y'$ . Then  $R$  is the graph of the partial function  $f: X \rightarrow Y$  defined by: if there is a  $y$  such that  $Rxy$ , then  $f(x) = y$ , otherwise  $f(x) \uparrow$ . If  $R$  is also serial, i.e., for each  $x \in X$  there is a  $y \in Y$  such that  $Rxy$ , then  $f$  is total.*

*Proof.* Suppose there is a  $y$  such that  $Rxy$ . If there were another  $y' \neq y$  such that  $Rxy'$ , the condition on  $R$  would be violated. Hence, if there is a  $y$  such that  $Rxy$ , that  $y$  is unique, and so  $f$  is well-defined. Obviously,  $R_f = R$  and  $f$  is total if  $R$  is serial.  $\square$

## Summary

A **function**  $f: A \rightarrow B$  maps every element of the **domain**  $A$  to a unique element of the **codomain**  $B$ . If  $x \in A$ , we call the  $y$  that  $f$  maps  $x$  to the **value**  $f(x)$  of  $f$  for **argument**  $x$ . If  $A$  is a set of pairs, we can think of the function  $f$  as taking two arguments. The **range**  $\text{ran}(f)$  of  $f$  is the subset of  $B$  that consists of all the values of  $f$ .

If  $\text{ran}(f) = B$  then  $f$  is called **surjective**. The value  $f(x)$  is unique in that  $f$  maps  $x$  to only one  $f(x)$ , never more than one. If  $f(x)$  is also unique in the sense that no two different arguments are mapped to the same value,  $f$  is called **injective**. Functions which are both injective and surjective are called **bijective**.

Bijective functions have a unique **inverse function**  $f^{-1}$ . Functions can also be chained together: the function  $(g \circ f)$  is the **composition** of  $f$  with  $g$ . Compositions of injective functions are injective, and of surjective functions are surjective, and  $(f^{-1} \circ f)$  is the identity function.



If we relax the requirement that  $f$  must have a value for every  $x \in A$ , we get the notion of a **partial functions**. If  $f: A \rightarrow B$  is partial, we say  $f(x)$  is **defined**,  $f(x) \downarrow$  if  $f$  has a value for argument  $x$ , and otherwise we say that  $f(x)$  is **undefined**,  $f(x) \uparrow$ . Any (partial) function  $f$  is associated with the **graph**  $R_f$  of  $f$ , the relation that holds iff  $f(x) = y$ .

## Problems

**Problem 3.1.** Show that if  $f: A \rightarrow B$  has a left inverse  $g$ , then  $f$  is injective.

**Problem 3.2.** Show that if  $f: A \rightarrow B$  has a right inverse  $h$ , then  $f$  is surjective.

**Problem 3.3.** Prove **Proposition 3.18**. You have to define  $f^{-1}$ , show that it is a function, and show that it is an inverse of  $f$ , i.e.,  $f^{-1}(f(x)) = x$  and  $f(f^{-1}(y)) = y$  for all  $x \in A$  and  $y \in B$ .

**Problem 3.4.** Prove **Proposition 3.19**.

**Problem 3.5.** Show that if  $f: A \rightarrow B$  and  $g: B \rightarrow C$  are both injective, then  $g \circ f: A \rightarrow C$  is injective.

**Problem 3.6.** Show that if  $f: A \rightarrow B$  and  $g: B \rightarrow C$  are both surjective, then  $g \circ f: A \rightarrow C$  is surjective.

**Problem 3.7.** Suppose  $f: A \rightarrow B$  and  $g: B \rightarrow C$ . Show that the graph of  $g \circ f$  is  $R_f \mid R_g$ .

**Problem 3.8.** Given  $f: A \rightarrow B$ , define the partial function  $g: B \rightarrow A$  by: for any  $y \in B$ , if there is a unique  $x \in A$  such that  $f(x) = y$ , then  $g(y) = x$ ; otherwise  $g(y) \uparrow$ . Show that if  $f$  is injective, then  $g(f(x)) = x$  for all  $x \in \text{dom}(f)$ , and  $f(g(y)) = y$  for all  $y \in \text{ran}(f)$ .

## CHAPTER 4

# *The Size of Sets*

### 4.1 Introduction

When Georg Cantor developed set theory in the 1870s, one of his aims was to make palatable the idea of an infinite collection—an actual infinity, as the medievals would say. A key part of this was his treatment of the *size* of different sets. If  $a$ ,  $b$  and  $c$  are all distinct, then the set  $\{a, b, c\}$  is intuitively *larger* than  $\{a, b\}$ . But what about infinite sets? Are they all as large as each other? It turns out that they are not.

The first important idea here is that of an enumeration. We can list every finite set by listing all its elements. For some infinite sets, we can also list all their elements if we allow the list itself to be infinite. Such sets are called countable. Cantor's surprising result, which we will fully understand by the end of this chapter, was that some infinite sets are not countable.

### 4.2 Enumerations and Countable Sets

We've already given examples of sets by listing their elements. Let's discuss in more general terms how and when we can list the elements of a set, even if that set is infinite.

**Definition 4.1 (Enumeration, informally).** Informally, an *enumeration* of a set  $A$  is a list (possibly infinite) of elements of  $A$  such that every element of  $A$  appears on the list at some finite position. If  $A$  has an enumeration, then  $A$  is said to be *countable*.

A couple of points about enumerations:

1. We count as enumerations only lists which have a beginning and in which every element other than the first has a single element immediately preceding it. In other words, there are only finitely many elements between the first element of the list and any other element. In particular, this means that every element of an enumeration has a finite position: the first element has position 1, the second position 2, etc.
2. We can have different enumerations of the same set  $A$  which differ by the order in which the elements appear: 4, 1, 25, 16, 9 enumerates the (set of the) first five square numbers just as well as 1, 4, 9, 16, 25 does.
3. Redundant enumerations are still enumerations: 1, 1, 2, 2, 3, 3, ... enumerates the same set as 1, 2, 3, ... does.
4. Order and redundancy *do* matter when we specify an enumeration: we can enumerate the positive integers beginning with 1, 2, 3, 1, ..., but the pattern is easier to see when enumerated in the standard way as 1, 2, 3, 4, ...
5. Enumerations must have a beginning: ..., 3, 2, 1 is not an enumeration of the positive integers because it has no first element. To see how this follows from the informal definition, ask yourself, “at what position in the list does the number 76 appear?”
6. The following is not an enumeration of the positive integers: 1, 3, 5, ..., 2, 4, 6, ... The problem is that the even

numbers occur at places  $\infty + 1$ ,  $\infty + 2$ ,  $\infty + 3$ , rather than at finite positions.

7. The empty set is enumerable: it is enumerated by the empty list!

**Proposition 4.2.** *If  $A$  has an enumeration, it has an enumeration without repetitions.*

*Proof.* Suppose  $A$  has an enumeration  $x_1, x_2, \dots$  in which each  $x_i$  is an element of  $A$ . We can remove repetitions from an enumeration by removing repeated elements. For instance, we can turn the enumeration into a new one in which we list  $x_i$  if it is an element of  $A$  that is not among  $x_1, \dots, x_{i-1}$  or remove  $x_i$  from the list if it already appears among  $x_1, \dots, x_{i-1}$ .  $\square$

The last argument shows that in order to get a good handle on enumerations and countable sets and to prove things about them, we need a more precise definition. The following provides it.

**Definition 4.3 (Enumeration, formally).** *An enumeration of a set  $A \neq \emptyset$  is any surjective function  $f: \mathbb{Z}^+ \rightarrow A$ .*

Let's convince ourselves that the formal definition and the informal definition using a possibly infinite list are equivalent. First, any surjective function from  $\mathbb{Z}^+$  to a set  $A$  enumerates  $A$ . Such a function determines an enumeration as defined informally above: the list  $f(1), f(2), f(3), \dots$ . Since  $f$  is surjective, every element of  $A$  is guaranteed to be the value of  $f(n)$  for some  $n \in \mathbb{Z}^+$ . Hence, every element of  $A$  appears at some finite position in the list. Since the function may not be injective, the list may be redundant, but that is acceptable (as noted above).

On the other hand, given a list that enumerates all elements of  $A$ , we can define a surjective function  $f: \mathbb{Z}^+ \rightarrow A$  by letting  $f(n)$  be the  $n$ th element of the list, or the final element of the

list if there is no  $n$ th element. The only case where this does not produce a surjective function is when  $A$  is empty, and hence the list is empty. So, every non-empty list determines a surjective function  $f: \mathbb{Z}^+ \rightarrow A$ .

**Definition 4.4.** A set  $A$  is countable iff it is empty or has an enumeration.

**Example 4.5.** A function enumerating the positive integers ( $\mathbb{Z}^+$ ) is simply the identity function given by  $f(n) = n$ . A function enumerating the natural numbers  $\mathbb{N}$  is the function  $g(n) = n - 1$ .

**Example 4.6.** The functions  $f: \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$  and  $g: \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$  given by

$$\begin{aligned} f(n) &= 2n \text{ and} \\ g(n) &= 2n - 1 \end{aligned}$$

enumerate the even positive integers and the odd positive integers, respectively. However, neither function is an enumeration of  $\mathbb{Z}^+$ , since neither is surjective.

**Example 4.7.** The function  $f(n) = (-1)^n \lceil \frac{(n-1)}{2} \rceil$  (where  $\lceil x \rceil$  denotes the *ceiling* function, which rounds  $x$  up to the nearest integer) enumerates the set of integers  $\mathbb{Z}$ . Notice how  $f$  generates the values of  $\mathbb{Z}$  by “hopping” back and forth between positive and negative integers:

$$\begin{array}{cccccccc} f(1) & f(2) & f(3) & f(4) & f(5) & f(6) & f(7) & \dots \\ -\lceil \frac{0}{2} \rceil & \lceil \frac{1}{2} \rceil & -\lceil \frac{2}{2} \rceil & \lceil \frac{3}{2} \rceil & -\lceil \frac{4}{2} \rceil & \lceil \frac{5}{2} \rceil & -\lceil \frac{6}{2} \rceil & \dots \\ 0 & 1 & -1 & 2 & -2 & 3 & \dots & \end{array}$$

You can also think of  $f$  as defined by cases as follows:

$$f(n) = \begin{cases} 0 & \text{if } n = 1 \\ n/2 & \text{if } n \text{ is even} \\ -(n-1)/2 & \text{if } n \text{ is odd and } > 1 \end{cases}$$

Although it is perhaps more natural when listing the elements of a set to start counting from the 1st element, mathematicians like to use the natural numbers  $\mathbb{N}$  for counting things. They talk about the 0th, 1st, 2nd, and so on, elements of a list. Correspondingly, we can define an enumeration as a surjective function from  $\mathbb{N}$  to  $A$ . Of course, the two definitions are equivalent.

**Proposition 4.8.** *There is a surjection  $f: \mathbb{Z}^+ \rightarrow A$  iff there is a surjection  $g: \mathbb{N} \rightarrow A$ .*

*Proof.* Given a surjection  $f: \mathbb{Z}^+ \rightarrow A$ , we can define  $g(n) = f(n+1)$  for all  $n \in \mathbb{N}$ . It is easy to see that  $g: \mathbb{N} \rightarrow A$  is surjective. Conversely, given a surjection  $g: \mathbb{N} \rightarrow A$ , define  $f(n) = g(n-1)$ .  $\square$

This gives us the following result:

**Corollary 4.9.** *A set  $A$  is countable iff it is empty or there is a surjective function  $f: \mathbb{N} \rightarrow A$ .*

We discussed above that an list of elements of a set  $A$  can be turned into a list without repetitions. This is also true for enumerations, but a bit harder to formulate and prove rigorously. Any function  $f: \mathbb{Z}^+ \rightarrow A$  must be defined for all  $n \in \mathbb{Z}^+$ . If there are only finitely many elements in  $A$  then we clearly cannot have a function defined on the infinitely many elements of  $\mathbb{Z}^+$  that takes as values all the elements of  $A$  but never takes the same value twice. In that case, i.e., in the case where the list without repetitions is finite, we must choose a different domain for  $f$ , one with only finitely many elements. Not having repetitions means that  $f$  must be injective. Since it is also surjective, we are looking for a bijection between some finite set  $\{1, \dots, n\}$  or  $\mathbb{Z}^+$  and  $A$ .

**Proposition 4.10.** *If  $f: \mathbb{Z}^+ \rightarrow A$  is surjective (i.e., an enumeration of  $A$ ), there is a bijection  $g: Z \rightarrow A$  where  $Z$  is either  $\mathbb{Z}^+$  or  $\{1, \dots, n\}$  for some  $n \in \mathbb{Z}^+$ .*

*Proof.* We define the function  $g$  recursively: Let  $g(1) = f(1)$ . If  $g(i)$  has already been defined, let  $g(i+1)$  be the first value of  $f(1), f(2), \dots$  not already among  $g(1), \dots, g(i)$ , if there is one. If  $A$  has just  $n$  elements, then  $g(1), \dots, g(n)$  are all defined, and so we have defined a function  $g: \{1, \dots, n\} \rightarrow A$ . If  $A$  has infinitely many elements, then for any  $i$  there must be an element of  $A$  in the enumeration  $f(1), f(2), \dots$ , which is not already among  $g(1), \dots, g(i)$ . In this case we have defined a function  $g: \mathbb{Z}^+ \rightarrow A$ .

The function  $g$  is surjective, since any element of  $A$  is among  $f(1), f(2), \dots$  (since  $f$  is surjective) and so will eventually be a value of  $g(i)$  for some  $i$ . It is also injective, since if there were  $j < i$  such that  $g(j) = g(i)$ , then  $g(i)$  would already be among  $g(1), \dots, g(i-1)$ , contrary to how we defined  $g$ .  $\square$

**Corollary 4.11.** *A set  $A$  is countable iff it is empty or there is a bijection  $f: \mathbb{N} \rightarrow A$  where either  $\mathbb{N} = \mathbb{N}$  or  $\mathbb{N} = \{0, \dots, n\}$  for some  $n \in \mathbb{N}$ .*

*Proof.*  $A$  is countable iff  $A$  is empty or there is a surjective  $f: \mathbb{Z}^+ \rightarrow A$ . By **Proposition 4.10**, the latter holds iff there is a bijective function  $f: Z \rightarrow A$  where  $Z = \mathbb{Z}^+$  or  $Z = \{1, \dots, n\}$  for some  $n \in \mathbb{Z}^+$ . By the same argument as in the proof of **Proposition 4.8**, that in turn is the case iff there is a bijection  $g: \mathbb{N} \rightarrow A$  where either  $\mathbb{N} = \mathbb{N}$  or  $\mathbb{N} = \{0, \dots, n-1\}$ .  $\square$

### 4.3 Cantor's Zig-Zag Method

We've already considered some "easy" enumerations. Now we will consider something a bit harder. Consider the set of pairs of natural numbers, which we defined in **section 1.5** thus:

$$\mathbb{N} \times \mathbb{N} = \{\langle n, m \rangle : n, m \in \mathbb{N}\}$$

We can organize these ordered pairs into an *array*, like so:

	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	...
<b>0</b>	$\langle 0, 0 \rangle$	$\langle 0, 1 \rangle$	$\langle 0, 2 \rangle$	$\langle 0, 3 \rangle$	...
<b>1</b>	$\langle 1, 0 \rangle$	$\langle 1, 1 \rangle$	$\langle 1, 2 \rangle$	$\langle 1, 3 \rangle$	...
<b>2</b>	$\langle 2, 0 \rangle$	$\langle 2, 1 \rangle$	$\langle 2, 2 \rangle$	$\langle 2, 3 \rangle$	...
<b>3</b>	$\langle 3, 0 \rangle$	$\langle 3, 1 \rangle$	$\langle 3, 2 \rangle$	$\langle 3, 3 \rangle$	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Clearly, every ordered pair in  $\mathbb{N} \times \mathbb{N}$  will appear exactly once in the array. In particular,  $\langle n, m \rangle$  will appear in the  $n$ th row and  $m$ th column. But how do we organize the elements of such an array into a “one-dimensional” list? The pattern in the array below demonstrates one way to do this (although of course there are many other options):

	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	...
<b>0</b>	0	1	3	6	10	...
<b>1</b>	2	4	7	11	...	...
<b>2</b>	5	8	12	...	...	...
<b>3</b>	9	13	...	...	...	...
<b>4</b>	14	...	...	...	...	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	...	$\ddots$

This pattern is called *Cantor’s zig-zag method*. It enumerates  $\mathbb{N} \times \mathbb{N}$  as follows:

$$\langle 0, 0 \rangle, \langle 0, 1 \rangle, \langle 1, 0 \rangle, \langle 0, 2 \rangle, \langle 1, 1 \rangle, \langle 2, 0 \rangle, \langle 0, 3 \rangle, \langle 1, 2 \rangle, \langle 2, 1 \rangle, \langle 3, 0 \rangle, \dots$$

And this establishes the following:

**Proposition 4.12.**  $\mathbb{N} \times \mathbb{N}$  is countable.

*Proof.* Let  $f: \mathbb{N} \rightarrow \mathbb{N} \times \mathbb{N}$  take each  $k \in \mathbb{N}$  to the tuple  $\langle n, m \rangle \in \mathbb{N} \times \mathbb{N}$  such that  $k$  is the value of the  $n$ th row and  $m$ th column in Cantor’s zig-zag array.  $\square$



This technique also generalises rather nicely. For example, we can use it to enumerate the set of ordered triples of natural numbers, i.e.:

$$\mathbb{N} \times \mathbb{N} \times \mathbb{N} = \{\langle n, m, k \rangle : n, m, k \in \mathbb{N}\}$$

We think of  $\mathbb{N} \times \mathbb{N} \times \mathbb{N}$  as the Cartesian product of  $\mathbb{N} \times \mathbb{N}$  with  $\mathbb{N}$ , that is,

$$\mathbb{N}^3 = (\mathbb{N} \times \mathbb{N}) \times \mathbb{N} = \{\langle \langle n, m \rangle, k \rangle : n, m, k \in \mathbb{N}\}$$

and thus we can enumerate  $\mathbb{N}^3$  with an array by labelling one axis with the enumeration of  $\mathbb{N}$ , and the other axis with the enumeration of  $\mathbb{N}^2$ :

	0	1	2	3	...
$\langle 0, 0 \rangle$	$\langle 0, 0, 0 \rangle$	$\langle 0, 0, 1 \rangle$	$\langle 0, 0, 2 \rangle$	$\langle 0, 0, 3 \rangle$	...
$\langle 0, 1 \rangle$	$\langle 0, 1, 0 \rangle$	$\langle 0, 1, 1 \rangle$	$\langle 0, 1, 2 \rangle$	$\langle 0, 1, 3 \rangle$	...
$\langle 1, 0 \rangle$	$\langle 1, 0, 0 \rangle$	$\langle 1, 0, 1 \rangle$	$\langle 1, 0, 2 \rangle$	$\langle 1, 0, 3 \rangle$	...
$\langle 0, 2 \rangle$	$\langle 0, 2, 0 \rangle$	$\langle 0, 2, 1 \rangle$	$\langle 0, 2, 2 \rangle$	$\langle 0, 2, 3 \rangle$	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Thus, by using a method like Cantor's zig-zag method, we may similarly obtain an enumeration of  $\mathbb{N}^3$ . And we can keep going, obtaining enumerations of  $\mathbb{N}^n$  for any natural number  $n$ . So, we have:

**Proposition 4.13.**  $\mathbb{N}^n$  is countable, for every  $n \in \mathbb{N}$ .

## 4.4 Pairing Functions and Codes

Cantor's zig-zag method makes the enumerability of  $\mathbb{N}^n$  visually evident. But let us focus on our array depicting  $\mathbb{N}^2$ . Following the zig-zag line in the array and counting the places, we can check that  $\langle 1, 2 \rangle$  is associated with the number 7. However, it would be nice if we could compute this more directly. That is, it would

be nice to have to hand the *inverse* of the zig-zag enumeration,  $g: \mathbb{N}^2 \rightarrow \mathbb{N}$ , such that

$$g(\langle 0,0 \rangle) = 0, \quad g(\langle 0,1 \rangle) = 1, \quad g(\langle 1,0 \rangle) = 2, \quad \dots, \quad g(\langle 1,2 \rangle) = 7, \quad \dots$$

This would enable us to calculate exactly where  $\langle n, m \rangle$  will occur in our enumeration.

In fact, we can define  $g$  directly by making two observations. First: if the  $n$ th row and  $m$ th column contains value  $v$ , then the  $(n+1)$ st row and  $(m-1)$ st column contains value  $v+1$ . Second: the first row of our enumeration consists of the triangular numbers, starting with 0, 1, 3, 6, etc. The  $k$ th triangular number is the sum of the natural numbers  $< k$ , which can be computed as  $k(k+1)/2$ . Putting these two observations together, consider this function:

$$g(n, m) = \frac{(n+m+1)(n+m)}{2} + n$$

We often just write  $g(n, m)$  rather than  $g(\langle n, m \rangle)$ , since it is easier on the eyes. This tells you first to determine the  $(n+m)$ <sup>th</sup> triangle number, and then add  $n$  to it. And it populates the array in exactly the way we would like. So in particular, the pair  $\langle 1, 2 \rangle$  is sent to  $\frac{4 \times 3}{2} + 1 = 7$ .

This function  $g$  is the *inverse* of an enumeration of a set of pairs. Such functions are called *pairing functions*.

**Definition 4.14 (Pairing function).** A function  $f: A \times B \rightarrow \mathbb{N}$  is an arithmetical *pairing function* if  $f$  is injective. We also say that  $f$  *encodes*  $A \times B$ , and that  $f(x, y)$  is the *code* for  $\langle x, y \rangle$ .

We can use pairing functions to encode, e.g., pairs of natural numbers; or, in other words, we can represent each *pair* of elements using a *single* number. Using the inverse of the pairing function, we can *decode* the number, i.e., find out which pair it represents.

## 4.5 An Alternative Pairing Function

There are other enumerations of  $\mathbb{N}^2$  that make it easier to figure out what their inverses are. Here is one. Instead of visualizing the enumeration in an array, start with the list of positive integers associated with (initially) empty spaces. Imagine filling these spaces successively with pairs  $\langle n, m \rangle$  as follows. Starting with the pairs that have 0 in the first place (i.e., pairs  $\langle 0, m \rangle$ ), put the first (i.e.,  $\langle 0, 0 \rangle$ ) in the first empty place, then skip an empty space, put the second (i.e.,  $\langle 0, 2 \rangle$ ) in the next empty place, skip one again, and so forth. The (incomplete) beginning of our enumeration now looks like this

<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	...
$\langle 0, 1 \rangle$	$\langle 0, 2 \rangle$	$\langle 0, 3 \rangle$	$\langle 0, 4 \rangle$	$\langle 0, 5 \rangle$						...

Repeat this with pairs  $\langle 1, m \rangle$  for the place that still remain empty, again skipping every other empty place:

<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	...
$\langle 0, 0 \rangle$	$\langle 1, 0 \rangle$	$\langle 0, 1 \rangle$	$\langle 0, 2 \rangle$	$\langle 1, 1 \rangle$	$\langle 0, 3 \rangle$	$\langle 0, 4 \rangle$	$\langle 1, 2 \rangle$			...

Enter pairs  $\langle 2, m \rangle$ ,  $\langle 2, m \rangle$ , etc., in the same way. Our completed enumeration thus starts like this:

<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	...
$\langle 0, 0 \rangle$	$\langle 1, 0 \rangle$	$\langle 0, 1 \rangle$	$\langle 2, 0 \rangle$	$\langle 0, 2 \rangle$	$\langle 1, 1 \rangle$	$\langle 0, 3 \rangle$	$\langle 3, 0 \rangle$	$\langle 0, 4 \rangle$	$\langle 1, 2 \rangle$	...

If we number the cells in the array above according to this enumeration, we will not find a neat zig-zag line, but this arrange-

ment:

	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	...
<b>0</b>	1	3	5	7	9	11	...
<b>1</b>	2	6	10	14	18	...	...
<b>2</b>	4	12	20	28	...	...	...
<b>3</b>	8	24	40	...	...	...	...
<b>4</b>	16	48	...	...	...	...	...
<b>5</b>	32	...	...	...	...	...	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

We can see that the pairs in row 0 are in the odd numbered places of our enumeration, i.e., pair  $\langle 0, m \rangle$  is in place  $2m + 1$ ; pairs in the second row,  $\langle 1, m \rangle$ , are in places whose number is the double of an odd number, specifically,  $2 \cdot (2m + 1)$ ; pairs in the third row,  $\langle 2, m \rangle$ , are in places whose number is four times an odd number,  $4 \cdot (2m + 1)$ ; and so on. The factors of  $(2m + 1)$  for each row, 1, 2, 4, 8, ..., are exactly the powers of 2:  $1 = 2^0$ ,  $2 = 2^1$ ,  $4 = 2^2$ ,  $8 = 2^3$ , ... In fact, the relevant exponent is always the first member of the pair in question. Thus, for pair  $\langle n, m \rangle$  the factor is  $2^n$ . This gives us the general formula:  $2^n \cdot (2m + 1)$ . However, this is a mapping of pairs to *positive* integers, i.e.,  $\langle 0, 0 \rangle$  has position 1. If we want to begin at position 0 we must subtract 1 from the result. This gives us:

**Example 4.15.** The function  $h: \mathbb{N}^2 \rightarrow \mathbb{N}$  given by

$$h(n, m) = 2^n(2m + 1) - 1$$

is a pairing function for the set of pairs of natural numbers  $\mathbb{N}^2$ .

Accordingly, in our second enumeration of  $\mathbb{N}^2$ , the pair  $\langle 0, 0 \rangle$  has code  $h(0, 0) = 2^0(2 \cdot 0 + 1) - 1 = 0$ ;  $\langle 1, 2 \rangle$  has code  $2^1 \cdot (2 \cdot 2 + 1) - 1 = 2 \cdot 5 - 1 = 9$ ;  $\langle 2, 6 \rangle$  has code  $2^2 \cdot (2 \cdot 6 + 1) - 1 = 51$ .

Sometimes it is enough to encode pairs of natural numbers  $\mathbb{N}^2$  without requiring that the encoding is surjective. Such encodings have inverses that are only partial functions.

**Example 4.16.** The function  $j: \mathbb{N}^2 \rightarrow \mathbb{N}^+$  given by

$$j(n, m) = 2^n 3^m$$

is an injective function  $\mathbb{N}^2 \rightarrow \mathbb{N}$ .

## 4.6 Uncountable Sets

Some sets, such as the set  $\mathbb{Z}^+$  of positive integers, are infinite. So far we've seen examples of infinite sets which were all countable. However, there are also infinite sets which do not have this property. Such sets are called *uncountable*.

First of all, it is perhaps already surprising that there are uncountable sets. For any countable set  $A$  there is a surjective function  $f: \mathbb{Z}^+ \rightarrow A$ . If a set is uncountable there is no such function. That is, no function mapping the infinitely many elements of  $\mathbb{Z}^+$  to  $A$  can exhaust all of  $A$ . So there are “more” elements of  $A$  than the infinitely many positive integers.

How would one prove that a set is uncountable? You have to show that no such surjective function can exist. Equivalently, you have to show that the elements of  $A$  cannot be enumerated in a one way infinite list. The best way to do this is to show that every list of elements of  $A$  must leave at least one element out; or that no function  $f: \mathbb{Z}^+ \rightarrow A$  can be surjective. We can do this using Cantor's *diagonal method*. Given a list of elements of  $A$ , say,  $x_1, x_2, \dots$ , we construct another element of  $A$  which, by its construction, cannot possibly be on that list.

Our first example is the set  $\mathbb{B}^\omega$  of all infinite, non-gappy sequences of 0's and 1's.

**Theorem 4.17.**  $\mathbb{B}^\omega$  is uncountable.

*Proof.* Suppose, by way of contradiction, that  $\mathbb{B}^\omega$  is countable, i.e., suppose that there is a list  $s_1, s_2, s_3, s_4, \dots$  of all elements of  $\mathbb{B}^\omega$ . Each of these  $s_i$  is itself an infinite sequence of 0's and 1's.

Let's call the  $j$ -th element of the  $i$ -th sequence in this list  $s_i(j)$ . Then the  $i$ -th sequence  $s_i$  is

$$s_i(1), s_i(2), s_i(3), \dots$$

We may arrange this list, and the elements of each sequence  $s_i$  in it, in an array:

	1	2	3	4	...
1	<b><math>s_1(1)</math></b>	$s_1(2)$	$s_1(3)$	$s_1(4)$	...
2	$s_2(1)$	<b><math>s_2(2)</math></b>	$s_2(3)$	$s_2(4)$	...
3	$s_3(1)$	$s_3(2)$	<b><math>s_3(3)</math></b>	$s_3(4)$	...
4	$s_4(1)$	$s_4(2)$	$s_4(3)$	<b><math>s_4(4)</math></b>	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

The labels down the side give the number of the sequence in the list  $s_1, s_2, \dots$ ; the numbers across the top label the elements of the individual sequences. For instance,  $s_1(1)$  is a name for whatever number, a 0 or a 1, is the first element in the sequence  $s_1$ , and so on.

Now we construct an infinite sequence,  $\bar{s}$ , of 0's and 1's which cannot possibly be on this list. The definition of  $\bar{s}$  will depend on the list  $s_1, s_2, \dots$ . Any infinite list of infinite sequences of 0's and 1's gives rise to an infinite sequence  $\bar{s}$  which is guaranteed to not appear on the list.

To define  $\bar{s}$ , we specify what all its elements are, i.e., we specify  $\bar{s}(n)$  for all  $n \in \mathbb{Z}^+$ . We do this by reading down the diagonal of the array above (hence the name "diagonal method") and then changing every 1 to a 0 and every 0 to a 1. More abstractly, we define  $\bar{s}(n)$  to be 0 or 1 according to whether the  $n$ -th element of the diagonal,  $s_n(n)$ , is 1 or 0.

$$\bar{s}(n) = \begin{cases} 1 & \text{if } s_n(n) = 0 \\ 0 & \text{if } s_n(n) = 1. \end{cases}$$

If you like formulas better than definitions by cases, you could also define  $\bar{s}(n) = 1 - s_n(n)$ .

Clearly  $\bar{s}$  is an infinite sequence of 0's and 1's, since it is just the mirror sequence to the sequence of 0's and 1's that appear on the diagonal of our array. So  $\bar{s}$  is an element of  $\mathbb{B}^\omega$ . But it cannot be on the list  $s_1, s_2, \dots$ . Why not?

It can't be the first sequence in the list,  $s_1$ , because it differs from  $s_1$  in the first element. Whatever  $s_1(1)$  is, we defined  $\bar{s}(1)$  to be the opposite. It can't be the second sequence in the list, because  $\bar{s}$  differs from  $s_2$  in the second element: if  $s_2(2)$  is 0,  $\bar{s}(2)$  is 1, and vice versa. And so on.

More precisely: if  $\bar{s}$  were on the list, there would be some  $k$  so that  $\bar{s} = s_k$ . Two sequences are identical iff they agree at every place, i.e., for any  $n$ ,  $\bar{s}(n) = s_k(n)$ . So in particular, taking  $n = k$  as a special case,  $\bar{s}(k) = s_k(k)$  would have to hold.  $s_k(k)$  is either 0 or 1. If it is 0 then  $\bar{s}(k)$  must be 1—that's how we defined  $\bar{s}$ . But if  $s_k(k) = 1$  then, again because of the way we defined  $\bar{s}$ ,  $\bar{s}(k) = 0$ . In either case  $\bar{s}(k) \neq s_k(k)$ .

We started by assuming that there is a list of elements of  $\mathbb{B}^\omega$ ,  $s_1, s_2, \dots$ . From this list we constructed a sequence  $\bar{s}$  which we proved cannot be on the list. But it definitely is a sequence of 0's and 1's if all the  $s_i$  are sequences of 0's and 1's, i.e.,  $\bar{s} \in \mathbb{B}^\omega$ . This shows in particular that there can be no list of *all* elements of  $\mathbb{B}^\omega$ , since for any such list we could also construct a sequence  $\bar{s}$  guaranteed to not be on the list, so the assumption that there is a list of all sequences in  $\mathbb{B}^\omega$  leads to a contradiction.  $\square$

This proof method is called “diagonalization” because it uses the diagonal of the array to define  $\bar{s}$ . Diagonalization need not involve the presence of an array: we can show that sets are not countable by using a similar idea even when no array and no actual diagonal is involved.

**Theorem 4.18.**  $\wp(\mathbb{Z}^+)$  is not countable.

*Proof.* We proceed in the same way, by showing that for every list of subsets of  $\mathbb{Z}^+$  there is a subset of  $\mathbb{Z}^+$  which cannot be on the

list. Suppose the following is a given list of subsets of  $\mathbb{Z}^+$ :

$$Z_1, Z_2, Z_3, \dots$$

We now define a set  $\overline{Z}$  such that for any  $n \in \mathbb{Z}^+$ ,  $n \in \overline{Z}$  iff  $n \notin Z_n$ :

$$\overline{Z} = \{n \in \mathbb{Z}^+ : n \notin Z_n\} \quad \square$$

$\overline{Z}$  is clearly a set of positive integers, since by assumption each  $Z_n$  is, and thus  $\overline{Z} \in \wp(\mathbb{Z}^+)$ . But  $\overline{Z}$  cannot be on the list. To show this, we'll establish that for each  $k \in \mathbb{Z}^+$ ,  $\overline{Z} \neq Z_k$ .

So let  $k \in \mathbb{Z}^+$  be arbitrary. We've defined  $\overline{Z}$  so that for any  $n \in \mathbb{Z}^+$ ,  $n \in \overline{Z}$  iff  $n \notin Z_n$ . In particular, taking  $n = k$ ,  $k \in \overline{Z}$  iff  $k \notin Z_k$ . But this shows that  $\overline{Z} \neq Z_k$ , since  $k$  is an element of one but not the other, and so  $\overline{Z}$  and  $Z_k$  have different elements. Since  $k$  was arbitrary,  $\overline{Z}$  is not on the list  $Z_1, Z_2, \dots$

The preceding proof did not mention a diagonal, but you can think of it as involving a diagonal if you picture it this way: Imagine the sets  $Z_1, Z_2, \dots$ , written in an array, where each element  $j \in Z_i$  is listed in the  $j$ -th column. Say the first four sets on that list are  $\{1, 2, 3, \dots\}$ ,  $\{2, 4, 6, \dots\}$ ,  $\{1, 2, 5\}$ , and  $\{3, 4, 5, \dots\}$ . Then the array would begin with

$$\begin{array}{l} Z_1 = \{\mathbf{1}, 2, 3, 4, 5, 6, \dots\} \\ Z_2 = \{ \mathbf{2}, 4, 6, \dots\} \\ Z_3 = \{1, 2, \mathbf{5} \} \\ Z_4 = \{ \mathbf{3}, \mathbf{4}, 5, 6, \dots\} \\ \vdots \end{array} \quad \begin{array}{l} \\ \\ \\ \ddots \end{array}$$

Then  $\overline{Z}$  is the set obtained by going down the diagonal, leaving out any numbers that appear along the diagonal and include those  $j$  where the array has a gap in the  $j$ -th row/column. In the above case, we would leave out 1 and 2, include 3, leave out 4, etc.



## 4.7 Reduction

We showed  $\wp(\mathbb{Z}^+)$  to be uncountable by a diagonalization argument. We already had a proof that  $\mathbb{B}^\omega$ , the set of all infinite sequences of 0s and 1s, is uncountable. Here's another way we can prove that  $\wp(\mathbb{Z}^+)$  is uncountable: Show that *if  $\wp(\mathbb{Z}^+)$  is countable then  $\mathbb{B}^\omega$  is also countable*. Since we know  $\mathbb{B}^\omega$  is not countable,  $\wp(\mathbb{Z}^+)$  can't be either. This is called *reducing* one problem to another—in this case, we reduce the problem of enumerating  $\mathbb{B}^\omega$  to the problem of enumerating  $\wp(\mathbb{Z}^+)$ . A solution to the latter—an enumeration of  $\wp(\mathbb{Z}^+)$ —would yield a solution to the former—an enumeration of  $\mathbb{B}^\omega$ .

How do we reduce the problem of enumerating a set  $B$  to that of enumerating a set  $A$ ? We provide a way of turning an enumeration of  $A$  into an enumeration of  $B$ . The easiest way to do that is to define a surjective function  $f: A \rightarrow B$ . If  $x_1, x_2, \dots$  enumerates  $A$ , then  $f(x_1), f(x_2), \dots$  would enumerate  $B$ . In our case, we are looking for a surjective function  $f: \wp(\mathbb{Z}^+) \rightarrow \mathbb{B}^\omega$ .

*Proof of Theorem 4.18 by reduction.* Suppose that  $\wp(\mathbb{Z}^+)$  were countable, and thus that there is an enumeration of it,  $Z_1, Z_2, Z_3, \dots$

Define the function  $f: \wp(\mathbb{Z}^+) \rightarrow \mathbb{B}^\omega$  by letting  $f(Z)$  be the sequence  $s_k$  such that  $s_k(n) = 1$  iff  $n \in Z$ , and  $s_k(n) = 0$  otherwise. This clearly defines a function, since whenever  $Z \subseteq \mathbb{Z}^+$ , any  $n \in \mathbb{Z}^+$  either is an element of  $Z$  or isn't. For instance, the set  $2\mathbb{Z}^+ = \{2, 4, 6, \dots\}$  of positive even numbers gets mapped to the sequence 010101..., the empty set gets mapped to 0000... and the set  $\mathbb{Z}^+$  itself to 1111....

It also is surjective: Every sequence of 0s and 1s corresponds to some set of positive integers, namely the one which has as its members those integers corresponding to the places where the sequence has 1s. More precisely, suppose  $s \in \mathbb{B}^\omega$ . Define  $Z \subseteq \mathbb{Z}^+$  by:

$$Z = \{n \in \mathbb{Z}^+ : s(n) = 1\}$$

Then  $f(Z) = s$ , as can be verified by consulting the definition of  $f$ .

Now consider the list

$$f(Z_1), f(Z_2), f(Z_3), \dots$$

Since  $f$  is surjective, every member of  $\mathbb{B}^\omega$  must appear as a value of  $f$  for some argument, and so must appear on the list. This list must therefore enumerate all of  $\mathbb{B}^\omega$ .

So if  $\wp(\mathbb{Z}^+)$  were countable,  $\mathbb{B}^\omega$  would be countable. But  $\mathbb{B}^\omega$  is uncountable ([Theorem 4.17](#)). Hence  $\wp(\mathbb{Z}^+)$  is uncountable.  $\square$

It is easy to be confused about the direction the reduction goes in. For instance, a surjective function  $g: \mathbb{B}^\omega \rightarrow B$  does *not* establish that  $B$  is uncountable. (Consider  $g: \mathbb{B}^\omega \rightarrow \mathbb{B}$  defined by  $g(s) = s(1)$ , the function that maps a sequence of 0's and 1's to its first element. It is surjective, because some sequences start with 0 and some start with 1. But  $\mathbb{B}$  is finite.) Note also that the function  $f$  must be surjective, or otherwise the argument does not go through:  $f(x_1), f(x_2), \dots$  would then not be guaranteed to include all the elements of  $B$ . For instance,

$$h(n) = \underbrace{000 \dots 0}_{n \text{ 0's}}$$

defines a function  $h: \mathbb{Z}^+ \rightarrow \mathbb{B}^\omega$ , but  $\mathbb{Z}^+$  is countable.

## 4.8 Equinumerosity

We have an intuitive notion of “size” of sets, which works fine for finite sets. But what about infinite sets? If we want to come up with a formal way of comparing the sizes of two sets of *any* size, it is a good idea to start by defining when sets are the same size. Here is Frege:

If a waiter wants to be sure that he has laid exactly as many knives as plates on the table, he does not need

to count either of them, if he simply lays a knife to the right of each plate, so that every knife on the table lies to the right of some plate. The plates and knives are thus uniquely correlated to each other, and indeed through that same spatial relationship. (Frege, 1884, §70)

The insight of this passage can be brought out through a formal definition:

**Definition 4.19.**  $A$  is *equinumerous* with  $B$ , written  $A \approx B$ , iff there is a bijection  $f: A \rightarrow B$ .

**Proposition 4.20.** *Equinumerosity is an equivalence relation.*

*Proof.* We must show that equinumerosity is reflexive, symmetric, and transitive. Let  $A, B$ , and  $C$  be sets.

*Reflexivity.* The identity map  $\text{Id}_A: A \rightarrow A$ , where  $\text{Id}_A(x) = x$  for all  $x \in A$ , is a bijection. So  $A \approx A$ .

*Symmetry.* Suppose  $A \approx B$ , i.e., there is a bijection  $f: A \rightarrow B$ . Since  $f$  is bijective, its inverse  $f^{-1}$  exists and is also bijective. Hence,  $f^{-1}: B \rightarrow A$  is a bijection, so  $B \approx A$ .

*Transitivity.* Suppose that  $A \approx B$  and  $B \approx C$ , i.e., there are bijections  $f: A \rightarrow B$  and  $g: B \rightarrow C$ . Then the composition  $g \circ f: A \rightarrow C$  is bijective, so that  $A \approx C$ .  $\square$

**Proposition 4.21.** *If  $A \approx B$ , then  $A$  is countable if and only if  $B$  is.*

*Proof.* Suppose  $A \approx B$ , so there is some bijection  $f: A \rightarrow B$ , and suppose that  $A$  is countable. Then either  $A = \emptyset$  or there is a surjective function  $g: \mathbb{Z}^+ \rightarrow A$ . If  $A = \emptyset$ , then  $B = \emptyset$  also (otherwise there would be an element  $y \in B$  but no  $x \in A$  with  $g(x) = y$ ). If, on the other hand,  $g: \mathbb{Z}^+ \rightarrow A$  is surjective, then  $f \circ g: \mathbb{Z}^+ \rightarrow B$  is surjective. To see this, let  $y \in B$ . Since  $f$

is surjective, there is an  $x \in A$  such that  $f(x) = y$ . Since  $g$  is surjective, there is an  $n \in \mathbb{Z}^+$  such that  $g(n) = x$ . Hence,

$$(f \circ g)(n) = f(g(n)) = f(x) = y$$

and thus  $f \circ g$  is surjective. We have that  $f \circ g$  is an enumeration of  $B$ , and so  $B$  is countable.

If  $B$  is countable, we obtain that  $A$  is countable by repeating the argument with the bijection  $f^{-1}: B \rightarrow A$  instead of  $f$ .  $\square$

## 4.9 Sets of Different Sizes, and Cantor's Theorem

We have offered a precise statement of the idea that two sets have the same size. We can also offer a precise statement of the idea that one set is smaller than another. Our definition of “is smaller than (or equinumerous)” will require, instead of a bijection between the sets, an injection from the first set to the second. If such a function exists, the size of the first set is less than or equal to the size of the second. Intuitively, an injection from one set to another guarantees that the range of the function has at least as many elements as the domain, since no two elements of the domain map to the same element of the range.

**Definition 4.22.**  $A$  is *no larger than*  $B$ , written  $A \preceq B$ , iff there is an injection  $f: A \rightarrow B$ .

It is clear that this is a reflexive and transitive relation, but that it is not symmetric (this is left as an exercise). We can also introduce a notion, which states that one set is (strictly) smaller than another.

**Definition 4.23.**  $A$  is *smaller than*  $B$ , written  $A \prec B$ , iff there is an injection  $f: A \rightarrow B$  but no bijection  $g: A \rightarrow B$ , i.e.,  $A \preceq B$  and  $A \not\approx B$ .

It is clear that this relation is irreflexive and transitive. (This is left as an exercise.) Using this notation, we can say that a set  $A$  is countable iff  $A \preceq \mathbb{N}$ , and that  $A$  is uncountable iff  $\mathbb{N} \prec A$ . This allows us to restate **Theorem 4.18** as the observation that  $\mathbb{Z}^+ \prec \wp(\mathbb{Z}^+)$ . In fact, **Cantor (1892)** proved that this last point is *perfectly general*:

**Theorem 4.24 (Cantor).**  $A \prec \wp(A)$ , for any set  $A$ .

*Proof.* The map  $f(x) = \{x\}$  is an injection  $f: A \rightarrow \wp(A)$ , since if  $x \neq y$ , then also  $\{x\} \neq \{y\}$  by extensionality, and so  $f(x) \neq f(y)$ . So we have that  $A \preceq \wp(A)$ .

We will now show that there cannot be a surjective function  $g: A \rightarrow \wp(A)$ , let alone a bijective one, and hence that  $A \not\approx \wp(A)$ . For suppose that  $g: A \rightarrow \wp(A)$ . Since  $g$  is total, every  $x \in A$  is mapped to a subset  $g(x) \subseteq A$ . We can show that  $g$  cannot be surjective. To do this, we define a subset  $\bar{A} \subseteq A$  which by definition cannot be in the range of  $g$ . Let

$$\bar{A} = \{x \in A : x \notin g(x)\}.$$

Since  $g(x)$  is defined for all  $x \in A$ ,  $\bar{A}$  is clearly a well-defined subset of  $A$ . But, it cannot be in the range of  $g$ . Let  $x \in A$  be arbitrary, we will show that  $\bar{A} \neq g(x)$ . If  $x \in g(x)$ , then it does not satisfy  $x \notin g(x)$ , and so by the definition of  $\bar{A}$ , we have  $x \notin \bar{A}$ . If  $x \in \bar{A}$ , it must satisfy the defining property of  $\bar{A}$ , i.e.,  $x \in A$  and  $x \notin g(x)$ . Since  $x$  was arbitrary, this shows that for each  $x \in \bar{A}$ ,  $x \in g(x)$  iff  $x \notin \bar{A}$ , and so  $g(x) \neq \bar{A}$ . In other words,  $\bar{A}$  cannot be in the range of  $g$ , contradicting the assumption that  $g$  is surjective.  $\square$

It's instructive to compare the proof of **Theorem 4.24** to that of **Theorem 4.18**. There we showed that for any list  $Z_1, Z_2, \dots$ , of subsets of  $\mathbb{Z}^+$  one can construct a set  $\bar{Z}$  of numbers guaranteed not to be on the list. It was guaranteed not to be on the list because, for every  $n \in \mathbb{Z}^+$ ,  $n \in Z_n$  iff  $n \notin \bar{Z}$ . This way, there is always some number that is an element of one of  $Z_n$  or  $\bar{Z}$  but not

the other. We follow the same idea here, except the indices  $n$  are now elements of  $A$  instead of  $\mathbb{Z}^+$ . The set  $\overline{B}$  is defined so that it is different from  $g(x)$  for each  $x \in A$ , because  $x \in g(x)$  iff  $x \notin \overline{B}$ . Again, there is always an element of  $A$  which is an element of one of  $g(x)$  and  $\overline{B}$  but not the other. And just as  $\overline{Z}$  therefore cannot be on the list  $Z_1, Z_2, \dots$ ,  $\overline{B}$  cannot be in the range of  $g$ .

The proof is also worth comparing with the proof of Russell's Paradox, [Theorem 1.29](#). Indeed, Cantor's Theorem was the inspiration for Russell's own paradox.

## 4.10 The Notion of Size, and Schröder-Bernstein

Here is an intuitive thought: if  $A$  is no larger than  $B$  and  $B$  is no larger than  $A$ , then  $A$  and  $B$  are equinumerous. To be honest, if this thought were *wrong*, then we could scarcely justify the thought that our defined notion of equinumerosity has anything to do with comparisons of “sizes” between sets! Fortunately, though, the intuitive thought is correct. This is justified by the Schröder-Bernstein Theorem.

**Theorem 4.25 (Schröder-Bernstein).** *If  $A \preceq B$  and  $B \preceq A$ , then  $A \approx B$ .*

In other words, if there is an injection from  $A$  to  $B$ , and an injection from  $B$  to  $A$ , then there is a bijection from  $A$  to  $B$ .

This result, however, is really rather *difficult* to prove. Indeed, although Cantor stated the result, others proved it.<sup>1</sup> For now, you can (and must) take it on trust.

Fortunately, Schröder-Bernstein is *correct*, and it vindicates our thinking of the relations we defined, i.e.,  $A \approx B$  and  $A \preceq B$ , as having something to do with “size”. Moreover, Schröder-Bernstein is very *useful*. It can be difficult to think of a bijection between two equinumerous sets. The Schröder-Bernstein Theorem allows us

<sup>1</sup>For more on the history, see e.g., [Potter \(2004, pp. 165–6\)](#).

to break the comparison down into cases so we only have to think of an injection from the first to the second, and vice-versa.

## Summary

The size of a set  $A$  can be measured by a natural number if the set is finite, and sizes can be compared by comparing these numbers. If sets are infinite, things are more complicated. The first level of infinity is that of **countably infinite** sets. A set  $A$  is countable if its elements can be arranged in an **enumeration**, a one-way infinite list, i.e., when there is a surjective function  $f: \mathbb{Z}^+ \rightarrow A$ . It is countably infinite if it is countable but not finite. Cantor's **zig-zag method** shows that the sets of pairs of elements of countably infinite sets is also countable; and this can be used to show that even the set of rational numbers  $\mathbb{Q}$  is countable.

There are, however, infinite sets that are not countable: these sets are called **uncountable**. There are two ways of showing that a set is uncountable: directly, using a **diagonal argument**, or by **reduction**. To give a diagonal argument, we assume that the set  $A$  in question is countable, and use a hypothetical enumeration to define an element of  $A$  which, by the very way we define it, is guaranteed to be different from every element in the enumeration. So the enumeration can't be an enumeration of all of  $A$  after all, and we've shown that no enumeration of  $A$  can exist. A reduction shows that  $A$  is uncountable by associating every element of  $A$  with an element of some known uncountable set  $B$  in a surjective way. If this is possible, then a hypothetical enumeration of  $A$  would yield an enumeration of  $B$ . Since  $B$  is uncountable, no enumeration of  $A$  can exist.

In general, infinite sets can be compared sizewise:  $A$  and  $B$  are the same size, or **equinumerous**, if there is a bijection between them. We can also define that  $A$  is no larger than  $B$  ( $A \preceq B$ ) if there is an injective function from  $A$  to  $B$ . By the Schröder-Bernstein Theorem, this in fact provides a sizewise order of infinite sets. Finally, **Cantor's theorem** says that for any

$A$ ,  $A \prec \wp(A)$ . This is a generalization of our result that  $\wp(\mathbb{Z}^+)$  is uncountable, and shows that there are not just two, but infinitely many levels of infinity.

## Problems

**Problem 4.1.** Define an enumeration of the positive squares 1, 4, 9, 16, ...

**Problem 4.2.** Show that if  $A$  and  $B$  are countable, so is  $A \cup B$ . To do this, suppose there are surjective functions  $f: \mathbb{Z}^+ \rightarrow A$  and  $g: \mathbb{Z}^+ \rightarrow B$ , and define a surjective function  $h: \mathbb{Z}^+ \rightarrow A \cup B$  and prove that it is surjective. Also consider the cases where  $A$  or  $B = \emptyset$ .

**Problem 4.3.** Show that if  $B \subseteq A$  and  $A$  is countable, so is  $B$ . To do this, suppose there is a surjective function  $f: \mathbb{Z}^+ \rightarrow A$ . Define a surjective function  $g: \mathbb{Z}^+ \rightarrow B$  and prove that it is surjective. What happens if  $B = \emptyset$ ?

**Problem 4.4.** Show by induction on  $n$  that if  $A_1, A_2, \dots, A_n$  are all countable, so is  $A_1 \cup \dots \cup A_n$ . You may assume the fact that if two sets  $A$  and  $B$  are countable, so is  $A \cup B$ .

**Problem 4.5.** According to **Definition 4.4**, a set  $A$  is enumerable iff  $A = \emptyset$  or there is a surjective  $f: \mathbb{Z}^+ \rightarrow A$ . It is also possible to define “countable set” precisely by: a set is enumerable iff there is an injective function  $g: A \rightarrow \mathbb{Z}^+$ . Show that the definitions are equivalent, i.e., show that there is an injective function  $g: A \rightarrow \mathbb{Z}^+$  iff either  $A = \emptyset$  or there is a surjective  $f: \mathbb{Z}^+ \rightarrow A$ .

**Problem 4.6.** Show that  $(\mathbb{Z}^+)^n$  is countable, for every  $n \in \mathbb{N}$ .

**Problem 4.7.** Show that  $(\mathbb{Z}^+)^*$  is countable. You may assume **problem 4.6**.



**Problem 4.8.** Give an enumeration of the set of all non-negative rational numbers.

**Problem 4.9.** Show that  $\mathbb{Q}$  is countable. Recall that any rational number can be written as a fraction  $z/m$  with  $z \in \mathbb{Z}$ ,  $m \in \mathbb{N}^+$ .

**Problem 4.10.** Define an enumeration of  $\mathbb{B}^*$ .

**Problem 4.11.** Recall from your introductory logic course that each possible truth table expresses a truth function. In other words, the truth functions are all functions from  $\mathbb{B}^k \rightarrow \mathbb{B}$  for some  $k$ . Prove that the set of all truth functions is enumerable.

**Problem 4.12.** Show that the set of all finite subsets of an arbitrary infinite countable set is countable.

**Problem 4.13.** A subset of  $\mathbb{N}$  is said to be *cofinite* iff it is the complement of a finite set  $\mathbb{N}$ ; that is,  $A \subseteq \mathbb{N}$  is cofinite iff  $\mathbb{N} \setminus A$  is finite. Let  $I$  be the set whose elements are exactly the finite and cofinite subsets of  $\mathbb{N}$ . Show that  $I$  is countable.

**Problem 4.14.** Show that the countable union of countable sets is countable. That is, whenever  $A_1, A_2, \dots$  are sets, and each  $A_i$  is countable, then the union  $\bigcup_{i=1}^{\infty} A_i$  of all of them is also countable. [NB: this is hard!]

**Problem 4.15.** Let  $f: A \times B \rightarrow \mathbb{N}$  be an arbitrary pairing function. Show that the inverse of  $f$  is an enumeration of  $A \times B$ .

**Problem 4.16.** Specify a function that encodes  $\mathbb{N}^3$ .

**Problem 4.17.** Show that  $\wp(\mathbb{N})$  is uncountable by a diagonal argument.

**Problem 4.18.** Show that the set of functions  $f: \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$  is uncountable by an explicit diagonal argument. That is, show that if  $f_1, f_2, \dots$ , is a list of functions and each  $f_i: \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$ , then there is some  $\bar{f}: \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$  not on this list.

**Problem 4.19.** Show that if there is an injective function  $g: B \rightarrow A$ , and  $B$  is uncountable, then so is  $A$ . Do this by showing how you can use  $g$  to turn an enumeration of  $A$  into one of  $B$ .

**Problem 4.20.** Show that the set of all *sets of* pairs of positive integers is uncountable by a reduction argument.

**Problem 4.21.** Show that the set  $X$  of all functions  $f: \mathbb{N} \rightarrow \mathbb{N}$  is uncountable by a reduction argument (Hint: give a surjective function from  $X$  to  $\mathbb{B}^\omega$ .)

**Problem 4.22.** Show that  $\mathbb{N}^\omega$ , the set of infinite sequences of natural numbers, is uncountable by a reduction argument.

**Problem 4.23.** Let  $P$  be the set of functions from the set of positive integers to the set  $\{0\}$ , and let  $Q$  be the set of *partial* functions from the set of positive integers to the set  $\{0\}$ . Show that  $P$  is countable and  $Q$  is not. (Hint: reduce the problem of enumerating  $\mathbb{B}^\omega$  to enumerating  $Q$ ).

**Problem 4.24.** Let  $S$  be the set of all surjective functions from the set of positive integers to the set  $\{0,1\}$ , i.e.,  $S$  consists of all surjective  $f: \mathbb{Z}^+ \rightarrow \mathbb{B}$ . Show that  $S$  is uncountable.

**Problem 4.25.** Show that the set  $\mathbb{R}$  of all real numbers is uncountable.

**Problem 4.26.** Show that if  $A \approx C$  and  $B \approx D$ , and  $A \cap B = C \cap D = \emptyset$ , then  $A \cup B \approx C \cup D$ .

**Problem 4.27.** Show that if  $A$  is infinite and countable, then  $A \approx \mathbb{N}$ .

**Problem 4.28.** Show that there cannot be an injection  $g: \wp(A) \rightarrow A$ , for any set  $A$ . Hint: Suppose  $g: \wp(A) \rightarrow A$  is injective. Consider  $D = \{g(B) : B \subseteq A \text{ and } g(B) \notin B\}$ . Let  $x = g(D)$ . Use the fact that  $g$  is injective to derive a contradiction.

**PART II**

*First-order  
Logic*

## CHAPTER 5

# *Introduction to First-Order Logic*

### 5.1 First-Order Logic

You are probably familiar with first-order logic from your first introduction to formal logic.<sup>1</sup> You may know it as “quantificational logic” or “predicate logic.” First-order logic, first of all, is a formal language. That means, it has a certain vocabulary, and its expressions are strings from this vocabulary. But not every string is permitted. There are different kinds of permitted expressions: terms, formulas, and sentences. We are mainly interested in sentences of first-order logic: they provide us with a formal analogue of sentences of English, and about them we can ask the questions a logician typically is interested in. For instance:

- Does  $B$  follow from  $A$  logically?
- Is  $A$  logically true, logically false, or contingent?

---

<sup>1</sup>In fact, we more or less assume you are! If you’re not, you could review a more elementary textbook, such as *forall x* (Magnus et al., 2021).

- Are  $A$  and  $B$  equivalent?

These questions are primarily questions about the “meaning” of sentences of first-order logic. For instance, a philosopher would analyze the question of whether  $B$  follows logically from  $A$  as asking: is there a case where  $A$  is true but  $B$  is false ( $B$  doesn’t follow from  $A$ ), or does every case that makes  $A$  true also make  $B$  true ( $B$  does follow from  $A$ )? But we haven’t been told yet what a “case” is—that is the job of *semantics*. The semantics of first-order logic provides a mathematically precise model of the philosopher’s intuitive idea of “case,” and also—and this is important—of what it is for a sentence  $A$  to be *true in* a case. We call the mathematically precise model that we will develop a *structure*. The relation which makes “true in” precise, is called the relation of *satisfaction*. So what we will define is “ $A$  is satisfied in  $M$ ” (in symbols:  $M \models A$ ) for sentences  $A$  and structures  $M$ . Once this is done, we can also give precise definitions of the other semantical terms such as “follows from” or “is logically true.” These definitions will make it possible to settle, again with mathematical precision, whether, e.g.,  $\forall x (A(x) \rightarrow B(x)), \exists x A(x) \models \exists x B(x)$ . The answer will, of course, be “yes.” If you’ve already been trained to symbolize sentences of English in first-order logic, you will recognize this as, e.g., the symbolizations of, say, “All ants are insects, there are ants, therefore there are insects.” That is obviously a valid argument, and so our mathematical model of “follows from” for our formal language should give the same answer.

Another topic you probably remember from your first introduction to formal logic is that there are *derivations*. If you have taken a first formal logic course, your instructor will have made you practice finding such derivations, perhaps even a derivation that shows that the above entailment holds. There are many different ways to give derivations: you may have done something called “natural deduction” or “truth trees,” but there are many others. The purpose of derivation systems is to provide tools using which the logicians’ questions above can be answered: e.g., a natural deduction derivation in which  $\forall x (A(x) \rightarrow B(x))$  and

$\exists x A(x)$  are premises and  $\exists x B(x)$  is the conclusion (last line) *verifies* that  $\exists x B(x)$  logically follows from  $\forall x (A(x) \rightarrow B(x))$  and  $\exists x A(x)$ .

But why is that? On the face of it, derivation systems have nothing to do with semantics: giving a formal derivation merely involves arranging symbols in certain rule-governed ways; they don't mention "cases" or "true in" at all. The connection between derivation systems and semantics has to be established by a meta-logical investigation. What's needed is a mathematical proof, e.g., that a formal derivation of  $\exists x B(x)$  from premises  $\forall x (A(x) \rightarrow B(x))$  and  $\exists x A(x)$  is possible, if, and only if,  $\forall x (A(x) \rightarrow B(x))$  and  $\exists x A(x)$  together entail  $\exists x B(x)$ . Before this can be done, however, a lot of painstaking work has to be carried out to get the definitions of syntax and semantics correct.

## 5.2 Syntax

We first must make precise what strings of symbols count as sentences of first-order logic. We'll do this later; for now we'll just proceed by example. The basic building blocks—the vocabulary—of first-order logic divides into two parts. The first part is the symbols we use to say specific things or to pick out specific things. We pick out things using constant symbols, and we say stuff about the things we pick out using predicate symbols. E.g, we might use  $a$  as a constant symbol to pick out a single thing, and then say something about it using the sentence  $P(a)$ . If you have meanings for " $a$ " and " $P$ " in mind, you can read  $P(a)$  as a sentence of English (and you probably have done so when you first learned formal logic). Once you have such simple sentences of first-order logic, you can build more complex ones using the second part of the vocabulary: the logical symbols (connectives and quantifiers). So, for instance, we can form expressions like  $(P(a) \wedge Q(b))$  or  $\exists x P(x)$ .

In order to provide the precise definitions of semantics and the rules of our derivation systems required for rigorous meta-

logical study, we first of all have to give a precise definition of what counts as a sentence of first-order logic. The basic idea is easy enough to understand: there are some simple sentences we can form from just predicate symbols and constant symbols, such as  $P(a)$ . And then from these we form more complex ones using the connectives and quantifiers. But what exactly are the rules by which we are allowed to form more complex sentences? These must be specified, otherwise we have not defined “sentence of first-order logic” precisely enough. There are a few issues. The first one is to get the right strings to count as sentences. The second one is to do this in such a way that we can give mathematical proofs about *all* sentences. Finally, we’ll have to also give precise definitions of some rudimentary operations with sentences, such as “replace every  $x$  in  $A$  by  $b$ .” The trouble is that the quantifiers and variables we have in first-order logic make it not entirely obvious how this should be done. E.g., should  $\exists x P(a)$  count as a sentence? What about  $\exists x \exists x P(x)$ ? What should the result of “replace  $x$  by  $b$  in  $(P(x) \wedge \exists x P(x))$ ” be?

### 5.3 Formulas

Here is the approach we will use to rigorously specify sentences of first-order logic and to deal with the issues arising from the use of variables. We first define a *different* set of expressions: formulas. Once we’ve done that, we can consider the role variables play in them—and on the basis of some other ideas, namely those of “free” and “bound” variables, we can define what a sentence is (namely, a formula without free variables). We do this not just because it makes the definition of “sentence” more manageable, but also because it will be crucial to the way we define the semantic notion of satisfaction.

Let’s define “formula” for a simple first-order language, one containing only a single predicate symbol  $P$  and a single constant symbol  $a$ , and only the logical symbols  $\neg$ ,  $\wedge$ , and  $\exists$ . Our full definitions will be much more general: we’ll allow infinitely

many predicate symbols and constant symbols. In fact, we will also consider function symbols which can be combined with constant symbols and variables to form “terms.” For now,  $a$  and the variables will be our only terms. We do need infinitely many variables. We’ll officially use the symbols  $v_0, v_1, \dots$ , as variables.

**Definition 5.1.** The set of *formulas*  $\text{Frm}$  is defined as follows:

1.  $P(a)$  and  $P(v_i)$  are formulas ( $i \in \mathbb{N}$ ).
2. If  $A$  is a formula, then  $\neg A$  is formula.
3. If  $A$  and  $B$  are formulas, then  $(A \wedge B)$  is a formula.
4. If  $A$  is a formula and  $x$  is a variable, then  $\exists x A$  is a formula.
5. Nothing else is a formula.

(1) tells us that  $P(a)$  and  $P(v_i)$  are formulas, for any  $i \in \mathbb{N}$ . These are the so-called *atomic* formulas. They give us something to start from. The other clauses give us ways of forming new formulas from ones we have already formed. So for instance, by (2), we get that  $\neg P(v_2)$  is a formula, since  $P(v_2)$  is already a formula by (1). Then, by (4), we get that  $\exists v_2 \neg P(v_2)$  is another formula, and so on. (5) tells us that *only* strings we can form in this way count as formulas. In particular,  $\exists v_0 P(a)$  and  $\exists v_0 \exists v_0 P(a)$  *do* count as formulas, and  $(\neg P(a))$  does not, because of the extraneous outer parentheses.

This way of defining formulas is called an *inductive definition*, and it allows us to prove things about formulas using a version of proof by induction called *structural induction*. These are discussed in a general way in [appendix C.4](#) and [appendix C.5](#), which you should review before delving into the proofs later on. Basically, the idea is that if you want to give a proof that something is true for all formulas, you show first that it is true for the atomic formulas, and then that *if* it’s true for any formula  $A$  (and  $B$ ), it’s *also* true for  $\neg A$ ,  $(A \wedge B)$ , and  $\exists x A$ . For instance, this proves



that it's true for  $\exists v_2 \neg P(v_2)$ : from the first part you know that it's true for the atomic formula  $P(v_2)$ . Then you get that it's true for  $\neg P(v_2)$  by the second part, and then again that it's true for  $\exists v_2 \neg P(v_2)$  itself. Since all formulas are inductively generated from atomic formulas, this works for any of them.

## 5.4 Satisfaction

We can already skip ahead to the semantics of first-order logic once we know what formulas are: here, the basic definition is that of a structure. For our simple language, a structure  $M$  has just three components: a non-empty set  $|M|$  called the *domain*, what  $a$  picks out in  $M$ , and what  $P$  is true of in  $M$ . The object picked out by  $a$  is denoted  $a^M$  and the set of things  $P$  is true of by  $P^M$ . A structure  $M$  consists of just these three things:  $|M|$ ,  $a^M \in |M|$  and  $P^M \subseteq |M|$ . The general case will be more complicated, since there will be many predicate symbols and constant symbols, the constant symbols can have more than one place, and there will also be function symbols.

This is enough to give a definition of satisfaction for formulas that don't contain variables. The idea is to give an inductive definition that mirrors the way we have defined formulas. We specify when an atomic formula is satisfied in  $M$ , and then when, e.g.,  $\neg A$  is satisfied in  $M$  on the basis of whether or not  $A$  is satisfied in  $M$ . E.g., we could define:

1.  $P(a)$  is satisfied in  $M$  iff  $a^M \in P^M$ .
2.  $\neg A$  is satisfied in  $M$  iff  $A$  is not satisfied in  $M$ .
3.  $(A \wedge B)$  is satisfied in  $M$  iff  $A$  is satisfied in  $M$ , and  $B$  is satisfied in  $M$  as well.

Let's say that  $|M| = \{0,1,2\}$ ,  $a^M = 1$ , and  $P^M = \{1,2\}$ . This definition would tell us that  $P(a)$  is satisfied in  $M$  (since  $a^M = 1 \in \{1,2\} = P^M$ ). It tells us further that  $\neg P(a)$  is not satisfied

in  $M$ , and that in turn  $\neg\neg P(a)$  is and  $(\neg P(a) \wedge P(a))$  is not satisfied, and so on.

The trouble comes when we want to give a definition for the quantifiers: we'd like to say something like, " $\exists v_0 P(v_0)$  is satisfied iff  $P(v_0)$  is satisfied." But the structure  $M$  doesn't tell us what to do about variables. What we actually want to say is that  $P(v_0)$  is satisfied *for some value of*  $v_0$ . To make this precise we need a way to assign elements of  $|M|$  not just to  $a$  but also to  $v_0$ . To this end, we introduce variable *assignments*. A variable assignment is simply a function  $s$  that maps variables to elements of  $|M|$  (in our example, to one of 1, 2, or 3). Since we don't know beforehand which variables might appear in a formula we can't limit which variables  $s$  assigns values to. The simple solution is to require that  $s$  assigns values to *all* variables  $v_0, v_1, \dots$ . We'll just use only the ones we need.

Instead of defining satisfaction of formulas just relative to a structure, we'll define it relative to a structure  $M$  and a variable assignment  $s$ , and write  $M, s \models A$  for short. Our definition will now include an additional clause to deal with atomic formulas containing variables:

1.  $M, s \models P(a)$  iff  $a^M \in P^M$ .
2.  $M, s \models P(v_i)$  iff  $s(v_i) \in P^M$ .
3.  $M, s \models \neg A$  iff not  $M, s \models A$ .
4.  $M, s \models (A \wedge B)$  iff  $M, s \models A$  and  $M, s \models B$ .

Ok, this solves one problem: we can now say when  $M$  satisfies  $P(v_0)$  for the value  $s(v_0)$ . To get the definition right for  $\exists v_0 P(v_0)$  we have to do one more thing: We want to have that  $M, s \models \exists v_0 P(v_0)$  iff  $M, s' \models P(v_0)$  for *some* way  $s'$  of assigning a value to  $v_0$ . But the value assigned to  $v_0$  does not necessarily have to be the value that  $s(v_0)$  picks out. We'll introduce a notation for that: if  $m \in |M|$ , then we let  $s[m/v_0]$  be the assignment that is just like  $s$  (for all variables other than  $v_0$ ), except to  $v_0$  it assigns  $m$ . Now our definition can be:

5.  $M, s \models \exists v_i A$  iff  $M, s[m/v_i] \models A$  for some  $m \in |M|$ .

Does it work out? Let's say we let  $s(v_i) = 0$  for all  $i \in \mathbb{N}$ .  $M, s \models \exists v_0 P(v_0)$  iff there is an  $m \in |M|$  so that  $M, s[m/v_0] \models P(v_0)$ . And there is: we can choose  $m = 1$  or  $m = 2$ . Note that this is true even if the value  $s(v_0)$  assigned to  $v_0$  by  $s$  itself—in this case, 0—doesn't do the job. We have  $M, s[1/v_0] \models P(v_0)$  but not  $M, s \models P(v_0)$ .

If this looks confusing and cumbersome: it is. But the added complexity is required to give a precise, inductive definition of satisfaction for all formulas, and we need something like it to precisely define the semantic notions. There are other ways of doing it, but they are all equally (in)elegant.

## 5.5 Sentences

Ok, now we have a (sketch of a) definition of satisfaction (“true in”) for structures and formulas. But it needs this additional bit—a variable assignment—and what we wanted is a definition of sentences. How do we get rid of assignments, and what are sentences?

You probably remember a discussion in your first introduction to formal logic about the relation between variables and quantifiers. A quantifier is always followed by a variable, and then in the part of the sentence to which that quantifier applies (its “scope”), we understand that the variable is “bound” by that quantifier. In formulas it was not required that every variable has a matching quantifier, and variables without matching quantifiers are “free” or “unbound.” We will take sentences to be all those formulas that have no free variables.

Again, the intuitive idea of when an occurrence of a variable in a formula  $A$  is bound, which quantifier binds it, and when it is free, is not difficult to get. You may have learned a method for testing this, perhaps involving counting parentheses. We have to insist on a precise definition—and because we have defined formulas by induction, we can give a definition of the free and bound

occurrences of a variable  $x$  in a formula  $A$  also by induction. E.g., it might look like this for our simplified language:

1. If  $A$  is atomic, all occurrences of  $x$  in it are free (that is, the occurrence of  $x$  in  $P(x)$  is free).
2. If  $A$  is of the form  $\neg B$ , then an occurrence of  $x$  in  $\neg B$  is free iff the corresponding occurrence of  $x$  is free in  $B$  (that is, the free occurrences of variables in  $B$  are exactly the corresponding occurrences in  $\neg B$ ).
3. If  $A$  is of the form  $(B \wedge C)$ , then an occurrence of  $x$  in  $(B \wedge C)$  is free iff the corresponding occurrence of  $x$  is free in  $B$  or in  $C$ .
4. If  $A$  is of the form  $\exists x B$ , then no occurrence of  $x$  in  $A$  is free; if it is of the form  $\exists y B$  where  $y$  is a different variable than  $x$ , then an occurrence of  $x$  in  $\exists y B$  is free iff the corresponding occurrence of  $x$  is free in  $B$ .

Once we have a precise definition of free and bound occurrences of variables, we can simply say: a sentence is any formula without free occurrences of variables.

## 5.6 Semantic Notions

We mentioned above that when we consider whether  $M, s \models A$  holds, we (for convenience) let  $s$  assign values to all variables, but only the values it assigns to variables in  $A$  are used. In fact, it's only the values of *free* variables in  $A$  that matter. Of course, because we're careful, we are going to prove this fact. Since sentences have no free variables,  $s$  doesn't matter at all when it comes to whether or not they are satisfied in a structure. So, when  $A$  is a sentence we can define  $M \models A$  to mean " $M, s \models A$  for all  $s$ ," which as it happens is true iff  $M, s \models A$  for at least one  $s$ . We need to introduce variable assignments to get a working definition of satisfaction for formulas, but for sentences, satisfaction is independent of the variable assignments.

Once we have a definition of “ $M \models A$ ,” we know what “case” and “true in” mean as far as sentences of first-order logic are concerned. On the basis of the definition of  $M \models A$  for sentences we can then define the basic semantic notions of validity, entailment, and satisfiability. A sentence is valid,  $\models A$ , if every structure satisfies it. It is entailed by a set of sentences,  $\Gamma \models A$ , if every structure that satisfies all the sentences in  $\Gamma$  also satisfies  $A$ . And a set of sentences is satisfiable if some structure satisfies all sentences in it at the same time.

Because formulas are inductively defined, and satisfaction is in turn defined by induction on the structure of formulas, we can use induction to prove properties of our semantics and to relate the semantic notions defined. We’ll collect and prove some of these properties, partly because they are individually interesting, but mainly because many of them will come in handy when we go on to investigate the relation between semantics and derivation systems. In order to do so, we’ll also have to define (precisely, i.e., by induction) some syntactic notions and operations we haven’t mentioned yet.

## 5.7 Substitution

We’ll discuss an example to illustrate how things hang together, and how the development of syntax and semantics lays the foundation for our more advanced investigations later. Our derivation systems should let us derive  $P(a)$  from  $\forall v_0 P(v_0)$ . Maybe we even want to state this as a rule of inference. However, to do so, we must be able to state it in the most general terms: not just for  $P$ ,  $a$ , and  $v_0$ , but for any formula  $A$ , and term  $t$ , and variable  $x$ . (Recall that constant symbols are terms, but we’ll consider also more complicated terms built from constant symbols and function symbols.) So we want to be able to say something like, “whenever you have derived  $\forall x A(x)$  you are justified in inferring  $A(t)$ —the result of removing  $\forall x$  and replacing  $x$  by  $t$ .” But what exactly does “replacing  $x$  by  $t$ ” mean? What is the relation between  $A(x)$

and  $A(t)$ ? Does this always work?

To make this precise, we define the operation of *substitution*. Substitution is actually tricky, because we can't just replace all  $x$ 's in  $A$  by  $t$ , and not every  $t$  can be substituted for any  $x$ . We'll deal with this, again, using inductive definitions. But once this is done, specifying an inference rule as “infer  $A(t)$  from  $\forall x A(x)$ ” becomes a precise definition. Moreover, we'll be able to show that this is a good inference rule in the sense that  $\forall x A(x)$  entails  $A(t)$ . But to prove this, we have to again prove something that may at first glance prompt you to ask “why are we doing this?” That  $\forall x A(x)$  entails  $A(t)$  relies on the fact that whether or not  $M \models A(t)$  holds depends only on the value of the term  $t$ , i.e., if we let  $m$  be whatever element of  $|M|$  is picked out by  $t$ , then  $M, s \models A(t)$  iff  $M, s[m/x] \models A(x)$ . This holds even when  $t$  contains variables, but we'll have to be careful with how exactly we state the result.

## 5.8 Models and Theories

Once we've defined the syntax and semantics of first-order logic, we can get to work investigating the properties of structures and the semantic notions. We can also define derivation systems, and investigate those. For a set of sentences, we can ask: what structures make all the sentences in that set true? Given a set of sentences  $\Gamma$ , a structure  $M$  that satisfies them is called a *model of  $\Gamma$* . We might start from  $\Gamma$  and try to find its models—what do they look like? How big or small do they have to be? But we might also start with a single structure or collection of structures and ask: what sentences are true in them? Are there sentences that *characterize* these structures in the sense that they, and only they, are true in them? These kinds of questions are the domain of *model theory*. They also underlie the *axiomatic method*: describing a collection of structures by a set of sentences, the axioms of a theory. This is made possible by the observation that exactly those sentences entailed in first-order logic by the axioms are true in all models of the axioms.

As a very simple example, consider preorders. A preorder is a relation  $R$  on some set  $A$  which is both reflexive and transitive. A set  $A$  with a two-place relation  $R \subseteq A \times A$  on it is exactly what we would need to give a structure for a first-order language with a single two-place relation symbol  $P$ : we would set  $|M| = A$  and  $P^M = R$ . Since  $R$  is a preorder, it is reflexive and transitive, and we can find a set  $\Gamma$  of sentences of first-order logic that say this:

$$\begin{aligned} & \forall v_0 P(v_0, v_0) \\ & \forall v_0 \forall v_1 \forall v_2 ((P(v_0, v_1) \wedge P(v_1, v_2)) \rightarrow P(v_0, v_2)) \end{aligned}$$

These sentences are just the symbolizations of “for any  $x$ ,  $Rxx$ ” ( $R$  is reflexive) and “whenever  $Rxy$  and  $Ryz$  then also  $Rxz$ ” ( $R$  is transitive). We see that a structure  $M$  is a model of these two sentences  $\Gamma$  iff  $R$  (i.e.,  $P^M$ ), is a preorder on  $A$  (i.e.,  $|M|$ ). In other words, the models of  $\Gamma$  are exactly the preorders. Any property of all preorders that can be expressed in the first-order language with just  $P$  as predicate symbol (like reflexivity and transitivity above), is entailed by the two sentences in  $\Gamma$  and vice versa. So anything we can prove about models of  $\Gamma$  we have proved about all preorders.

For any particular theory and class of models (such as  $\Gamma$  and all preorders), there will be interesting questions about what can be expressed in the corresponding first-order language, and what cannot be expressed. There are some properties of structures that are interesting for all languages and classes of models, namely those concerning the size of the domain. One can always express, for instance, that the domain contains exactly  $n$  elements, for any  $n \in \mathbb{Z}^+$ . One can also express, using a set of infinitely many sentences, that the domain is infinite. But one cannot express that the domain is finite, or that the domain is uncountable. These results about the limitations of first-order languages are consequences of the compactness and Löwenheim–Skolem theorems.

## 5.9 Soundness and Completeness

We'll also introduce derivation systems for first-order logic. There are many derivation systems that logicians have developed, but they all define the same derivability relation between sentences. We say that  $\Gamma$  *derives*  $A$ ,  $\Gamma \vdash A$ , if there is a derivation of a certain precisely defined sort. Derivations are always finite arrangements of symbols—perhaps a list of sentences, or some more complicated structure. The purpose of derivation systems is to provide a tool to determine if a sentence is entailed by some set  $\Gamma$ . In order to serve that purpose, it must be true that  $\Gamma \vDash A$  if, and only if,  $\Gamma \vdash A$ .

If  $\Gamma \vdash A$  but not  $\Gamma \vDash A$ , our derivation system would be too strong, prove too much. The property that if  $\Gamma \vdash A$  then  $\Gamma \vDash A$  is called *soundness*, and it is a minimal requirement on any good derivation system. On the other hand, if  $\Gamma \vDash A$  but not  $\Gamma \vdash A$ , then our derivation system is too weak, it doesn't prove enough. The property that if  $\Gamma \vDash A$  then  $\Gamma \vdash A$  is called *completeness*. Soundness is usually relatively easy to prove (by induction on the structure of derivations, which are inductively defined). Completeness is harder to prove.

Soundness and completeness have a number of important consequences. If a set of sentences  $\Gamma$  derives a contradiction (such as  $A \wedge \neg A$ ) it is called *inconsistent*. Inconsistent  $\Gamma$ s cannot have any models, they are unsatisfiable. From completeness the converse follows: any  $\Gamma$  that is not inconsistent—or, as we will say, *consistent*—has a model. In fact, this is equivalent to completeness, and is the form of completeness we will actually prove. It is a deep and perhaps surprising result: just because you cannot prove  $A \wedge \neg A$  from  $\Gamma$  guarantees that there is a structure that is as  $\Gamma$  describes it. So completeness gives an answer to the question: which sets of sentences have models? Answer: all and only consistent sets do.

The soundness and completeness theorems have two important consequences: the compactness and the Löwenheim–Skolem theorem. These are important results in the theory of models,



and can be used to establish many interesting results. We've already mentioned two: first-order logic cannot express that the domain of a structure is finite or that it is uncountable.

Historically, all of this—how to define syntax and semantics of first-order logic, how to define good derivation systems, how to prove that they are sound and complete, getting clear about what can and cannot be expressed in first-order languages—took a long time to figure out and get right. We now know how to do it, but going through all the details can still be confusing and tedious. But it's also important, because the methods developed here for the formal language of first-order logic are applied all over the place in logic, computer science, and linguistics. So working through the details pays off in the long run.

## CHAPTER 6

# *Syntax of First-Order Logic*

### 6.1 Introduction

In order to develop the theory and metatheory of first-order logic, we must first define the syntax and semantics of its expressions. The expressions of first-order logic are terms and formulas. Terms are formed from variables, constant symbols, and function symbols. Formulas, in turn, are formed from predicate symbols together with terms (these form the smallest, “atomic” formulas), and then from atomic formulas we can form more complex ones using logical connectives and quantifiers. There are many different ways to set down the formation rules; we give just one possible one. Other systems will chose different symbols, will select different sets of connectives as primitive, will use parentheses differently (or even not at all, as in the case of so-called Polish notation). What all approaches have in common, though, is that the formation rules define the set of terms and formulas *inductively*. If done properly, every expression can result essentially

in only one way according to the formation rules. The inductive definition resulting in expressions that are *uniquely readable* means we can give meanings to these expressions using the same method—inductive definition.

## 6.2 First-Order Languages

Expressions of first-order logic are built up from a basic vocabulary containing *variables*, *constant symbols*, *predicate symbols* and sometimes *function symbols*. From them, together with logical connectives, quantifiers, and punctuation symbols such as parentheses and commas, *terms* and *formulas* are formed.

Informally, predicate symbols are names for properties and relations, constant symbols are names for individual objects, and function symbols are names for mappings. These, except for the identity predicate  $=$ , are the *non-logical symbols* and together make up a language. Any first-order language  $\mathcal{L}$  is determined by its non-logical symbols. In the most general case,  $\mathcal{L}$  contains infinitely many symbols of each kind.

In the general case, we make use of the following symbols in first-order logic:

1. Logical symbols
  - a) Logical connectives:  $\neg$  (negation),  $\wedge$  (conjunction),  $\vee$  (disjunction),  $\rightarrow$  (conditional),  $\forall$  (universal quantifier),  $\exists$  (existential quantifier).
  - b) The propositional constant for falsity  $\perp$ .
  - c) The two-place identity predicate  $=$ .
  - d) A countably infinite set of variables:  $v_0, v_1, v_2, \dots$
2. Non-logical symbols, making up the *standard language* of first-order logic
  - a) A countably infinite set of  $n$ -place predicate symbols for each  $n > 0$ :  $A_0^n, A_1^n, A_2^n, \dots$

- b) A countably infinite set of constant symbols:  $c_0, c_1, c_2, \dots$
  - c) A countably infinite set of  $n$ -place function symbols for each  $n > 0$ :  $f_0^n, f_1^n, f_2^n, \dots$
3. Punctuation marks: (, ), and the comma.

Most of our definitions and results will be formulated for the full standard language of first-order logic. However, depending on the application, we may also restrict the language to only a few predicate symbols, constant symbols, and function symbols.

**Example 6.1.** The language  $\mathcal{L}_A$  of arithmetic contains a single two-place predicate symbol  $<$ , a single constant symbol  $0$ , one one-place function symbol  $\iota$ , and two two-place function symbols  $+$  and  $\times$ .

**Example 6.2.** The language of set theory  $\mathcal{L}_Z$  contains only the single two-place predicate symbol  $\in$ .

**Example 6.3.** The language of orders  $\mathcal{L}_{\leq}$  contains only the two-place predicate symbol  $\leq$ .

Again, these are conventions: officially, these are just aliases, e.g.,  $<$ ,  $\in$ , and  $\leq$  are aliases for  $A_0^2$ ,  $0$  for  $c_0$ ,  $\iota$  for  $f_0^1$ ,  $+$  for  $f_0^2$ ,  $\times$  for  $f_1^2$ .

In addition to the primitive connectives and quantifiers introduced above, we also use the following *defined* symbols:  $\leftrightarrow$  (biconditional), truth  $\top$

A defined symbol is not officially part of the language, but is introduced as an informal abbreviation: it allows us to abbreviate formulas which would, if we only used primitive symbols, get quite long. This is obviously an advantage. The bigger advantage, however, is that proofs become shorter. If a symbol is primitive, it has to be treated separately in proofs. The more primitive symbols, therefore, the longer our proofs.

You may be familiar with different terminology and symbols than the ones we use above. Logic texts (and teachers) commonly use  $\sim$ ,  $\neg$ , or  $!$  for “negation”,  $\wedge$ ,  $\cdot$ , or  $\&$  for “conjunction”. Commonly used symbols for the “conditional” or “implication” are  $\rightarrow$ ,  $\Rightarrow$ , and  $\supset$ . Symbols for “biconditional,” “bi-implication,” or “(material) equivalence” are  $\leftrightarrow$ ,  $\Leftrightarrow$ , and  $\equiv$ . The  $\perp$  symbol is variously called “falsity,” “falsum,” “absurdity,” or “bottom.” The  $\top$  symbol is variously called “truth,” “verum,” or “top.”

It is conventional to use lower case letters (e.g.,  $a$ ,  $b$ ,  $c$ ) from the beginning of the Latin alphabet for constant symbols (sometimes called names), and lower case letters from the end (e.g.,  $x$ ,  $y$ ,  $z$ ) for variables. Quantifiers combine with variables, e.g.,  $x$ ; notational variations include  $\forall x$ ,  $(\forall x)$ ,  $(x)$ ,  $\Pi x$ ,  $\bigwedge_x$  for the universal quantifier and  $\exists x$ ,  $(\exists x)$ ,  $(Ex)$ ,  $\Sigma x$ ,  $\bigvee_x$  for the existential quantifier.

We might treat all the propositional operators and both quantifiers as primitive symbols of the language. We might instead choose a smaller stock of primitive symbols and treat the other logical operators as defined. “Truth functionally complete” sets of Boolean operators include  $\{\neg, \vee\}$ ,  $\{\neg, \wedge\}$ , and  $\{\neg, \rightarrow\}$ —these can be combined with either quantifier for an expressively complete first-order language.

You may be familiar with two other logical operators: the Sheffer stroke  $|$  (named after Henry Sheffer), and Peirce’s arrow  $\downarrow$ , also known as Quine’s dagger. When given their usual readings of “nand” and “nor” (respectively), these operators are truth functionally complete by themselves.

## 6.3 Terms and Formulas

Once a first-order language  $\mathcal{L}$  is given, we can define expressions built up from the basic vocabulary of  $\mathcal{L}$ . These include in particular *terms* and *formulas*.

**Definition 6.4 (Terms).** The set of *terms*  $\text{Trm}(\mathcal{L})$  of  $\mathcal{L}$  is defined inductively by:

1. Every variable is a term.
2. Every constant symbol of  $\mathcal{L}$  is a term.
3. If  $f$  is an  $n$ -place function symbol and  $t_1, \dots, t_n$  are terms, then  $f(t_1, \dots, t_n)$  is a term.
4. Nothing else is a term.

A term containing no variables is a *closed term*.

The constant symbols appear in our specification of the language and the terms as a separate category of symbols, but they could instead have been included as zero-place function symbols. We could then do without the second clause in the definition of terms. We just have to understand  $f(t_1, \dots, t_n)$  as just  $f$  by itself if  $n = 0$ .

**Definition 6.5 (Formulas).** The set of *formulas*  $\text{Frm}(\mathcal{L})$  of the language  $\mathcal{L}$  is defined inductively as follows:

1.  $\perp$  is an atomic formula.
2. If  $R$  is an  $n$ -place predicate symbol of  $\mathcal{L}$  and  $t_1, \dots, t_n$  are terms of  $\mathcal{L}$ , then  $R(t_1, \dots, t_n)$  is an atomic formula.
3. If  $t_1$  and  $t_2$  are terms of  $\mathcal{L}$ , then  $=(t_1, t_2)$  is an atomic formula.
4. If  $A$  is a formula, then  $\neg A$  is formula.
5. If  $A$  and  $B$  are formulas, then  $(A \wedge B)$  is a formula.
6. If  $A$  and  $B$  are formulas, then  $(A \vee B)$  is a formula.
7. If  $A$  and  $B$  are formulas, then  $(A \rightarrow B)$  is a formula.

8. If  $A$  is a formula and  $x$  is a variable, then  $\forall x A$  is a formula.
9. If  $A$  is a formula and  $x$  is a variable, then  $\exists x A$  is a formula.
10. Nothing else is a formula.

The definitions of the set of terms and that of formulas are *inductive definitions*. Essentially, we construct the set of formulas in infinitely many stages. In the initial stage, we pronounce all atomic formulas to be formulas; this corresponds to the first few cases of the definition, i.e., the cases for  $\perp$ ,  $R(t_1, \dots, t_n)$  and  $=(t_1, t_2)$ . “Atomic formula” thus means any formula of this form.

The other cases of the definition give rules for constructing new formulas out of formulas already constructed. At the second stage, we can use them to construct formulas out of atomic formulas. At the third stage, we construct new formulas from the atomic formulas and those obtained in the second stage, and so on. A formula is anything that is eventually constructed at such a stage, and nothing else.

By convention, we write  $=$  between its arguments and leave out the parentheses:  $t_1 = t_2$  is an abbreviation for  $=(t_1, t_2)$ . Moreover,  $\neg=(t_1, t_2)$  is abbreviated as  $t_1 \neq t_2$ . When writing a formula  $(B * C)$  constructed from  $B, C$  using a two-place connective  $*$ , we will often leave out the outermost pair of parentheses and write simply  $B * C$ .

Some logic texts require that the variable  $x$  must occur in  $A$  in order for  $\exists x A$  and  $\forall x A$  to count as formulas. Nothing bad happens if you don't require this, and it makes things easier.

**Definition 6.6.** Formulas constructed using the defined operators are to be understood as follows:

1.  $\top$  abbreviates  $\neg\perp$ .
2.  $A \leftrightarrow B$  abbreviates  $(A \rightarrow B) \wedge (B \rightarrow A)$ .

If we work in a language for a specific application, we will often write two-place predicate symbols and function symbols between the respective terms, e.g.,  $t_1 < t_2$  and  $(t_1 + t_2)$  in the language of arithmetic and  $t_1 \in t_2$  in the language of set theory. The successor function in the language of arithmetic is even written conventionally *after* its argument:  $t'$ . Officially, however, these are just conventional abbreviations for  $A_0^2(t_1, t_2)$ ,  $f_0^2(t_1, t_2)$ ,  $A_0^2(t_1, t_2)$  and  $f_0^1(t)$ , respectively.

**Definition 6.7 (Syntactic identity).** The symbol  $\equiv$  expresses syntactic identity between strings of symbols, i.e.,  $A \equiv B$  iff  $A$  and  $B$  are strings of symbols of the same length and which contain the same symbol in each place.

The  $\equiv$  symbol may be flanked by strings obtained by concatenation, e.g.,  $A \equiv (B \vee C)$  means: the string of symbols  $A$  is the same string as the one obtained by concatenating an opening parenthesis, the string  $B$ , the  $\vee$  symbol, the string  $C$ , and a closing parenthesis, in this order. If this is the case, then we know that the first symbol of  $A$  is an opening parenthesis,  $A$  contains  $B$  as a substring (starting at the second symbol), that substring is followed by  $\vee$ , etc.

As terms and formulas are built up from basic elements via inductive definitions, we can use the following induction principles to prove things about them.

**Lemma 6.8 (Principle of induction on terms).** Let  $\mathcal{L}$  be a first-order language. If some property  $P$  is such that

1. it holds for every variable  $v$ ,
2. it holds for every constant symbol  $a$  of  $\mathcal{L}$ , and
3. it holds for  $f(t_1, \dots, t_n)$  whenever it holds for  $t_1, \dots, t_n$  and  $f$  is an  $n$ -place function symbol of  $\mathcal{L}$



(assuming  $t_1, \dots, t_n$  are terms of  $\mathcal{L}$ ), then  $P$  holds for every term in  $\text{Trm}(\mathcal{L})$ .

**Lemma 6.9 (Principle of induction on formulas).** Let  $\mathcal{L}$  be a first-order language. If some property  $P$  holds for all the atomic formulas and is such that

1. it holds for  $\neg A$  whenever it holds for  $A$ ;
2. it holds for  $(A \wedge B)$  whenever it holds for  $A$  and  $B$ ;
3. it holds for  $(A \vee B)$  whenever it holds for  $A$  and  $B$ ;
4. it holds for  $(A \rightarrow B)$  whenever it holds for  $A$  and  $B$ ;
5. it holds for  $\exists x A$  whenever it holds for  $A$ ;
6. it holds for  $\forall x A$  whenever it holds for  $A$ ;

(assuming  $A$  and  $B$  are formulas of  $\mathcal{L}$ ), then  $P$  holds for all formulas in  $\text{Frm}(\mathcal{L})$ .

## 6.4 Unique Readability

The way we defined formulas guarantees that every formula has a *unique reading*, i.e., there is essentially only one way of constructing it according to our formation rules for formulas and only one way of “interpreting” it. If this were not so, we would have ambiguous formulas, i.e., formulas that have more than one reading or interpretation—and that is clearly something we want to avoid. But more importantly, without this property, most of the definitions and proofs we are going to give will not go through.

Perhaps the best way to make this clear is to see what would happen if we had given bad rules for forming formulas that would not guarantee unique readability. For instance, we could have

forgotten the parentheses in the formation rules for connectives, e.g., we might have allowed this:

If  $A$  and  $B$  are formulas, then so is  $A \rightarrow B$ .

Starting from an atomic formula  $D$ , this would allow us to form  $D \rightarrow D$ . From this, together with  $D$ , we would get  $D \rightarrow D \rightarrow D$ . But there are two ways to do this:

1. We take  $D$  to be  $A$  and  $D \rightarrow D$  to be  $B$ .
2. We take  $A$  to be  $D \rightarrow D$  and  $B$  is  $D$ .

Correspondingly, there are two ways to “read” the formula  $D \rightarrow D \rightarrow D$ . It is of the form  $B \rightarrow C$  where  $B$  is  $D$  and  $C$  is  $D \rightarrow D$ , but *it is also* of the form  $B \rightarrow C$  with  $B$  being  $D \rightarrow D$  and  $C$  being  $D$ .

If this happens, our definitions will not always work. For instance, when we define the main operator of a formula, we say: in a formula of the form  $B \rightarrow C$ , the main operator is the indicated occurrence of  $\rightarrow$ . But if we can match the formula  $D \rightarrow D \rightarrow D$  with  $B \rightarrow C$  in the two different ways mentioned above, then in one case we get the first occurrence of  $\rightarrow$  as the main operator, and in the second case the second occurrence. But we intend the main operator to be a *function* of the formula, i.e., every formula must have exactly one main operator occurrence.

**Lemma 6.10.** *The number of left and right parentheses in a formula  $A$  are equal.*

*Proof.* We prove this by induction on the way  $A$  is constructed. This requires two things: (a) We have to prove first that all atomic formulas have the property in question (the induction basis). (b) Then we have to prove that when we construct new formulas out of given formulas, the new formulas have the property provided the old ones do.

Let  $l(A)$  be the number of left parentheses, and  $r(A)$  the number of right parentheses in  $A$ , and  $l(t)$  and  $r(t)$  similarly the number of left and right parentheses in a term  $t$ .

1.  $A \equiv \perp$ :  $A$  has 0 left and 0 right parentheses.
2.  $A \equiv R(t_1, \dots, t_n)$ :  $l(A) = 1 + l(t_1) + \dots + l(t_n) = 1 + r(t_1) + \dots + r(t_n) = r(A)$ . Here we make use of the fact, left as an exercise, that  $l(t) = r(t)$  for any term  $t$ .
3.  $A \equiv t_1 = t_2$ :  $l(A) = l(t_1) + l(t_2) = r(t_1) + r(t_2) = r(A)$ .
4.  $A \equiv \neg B$ : By induction hypothesis,  $l(B) = r(B)$ . Thus  $l(A) = l(B) = r(B) = r(A)$ .
5.  $A \equiv (B * C)$ : By induction hypothesis,  $l(B) = r(B)$  and  $l(C) = r(C)$ . Thus  $l(A) = 1 + l(B) + l(C) = 1 + r(B) + r(C) = r(A)$ .
6.  $A \equiv \forall x B$ : By induction hypothesis,  $l(B) = r(B)$ . Thus,  $l(A) = l(B) = r(B) = r(A)$ .
7.  $A \equiv \exists x B$ : Similarly. □

**Definition 6.11 (Proper prefix).** A string of symbols  $B$  is a *proper prefix* of a string of symbols  $A$  if concatenating  $B$  and a non-empty string of symbols yields  $A$ .

**Lemma 6.12.** *If  $A$  is a formula, and  $B$  is a proper prefix of  $A$ , then  $B$  is not a formula.*

*Proof.* Exercise. □

**Proposition 6.13.** *If  $A$  is an atomic formula, then it satisfies one, and only one of the following conditions.*

1.  $A \equiv \perp$ .
2.  $A \equiv R(t_1, \dots, t_n)$  where  $R$  is an  $n$ -place predicate symbol,  $t_1, \dots, t_n$  are terms, and each of  $R, t_1, \dots, t_n$  is uniquely determined.

3.  $A \equiv t_1 = t_2$  where  $t_1$  and  $t_2$  are uniquely determined terms.

*Proof.* Exercise. □

**Proposition 6.14 (Unique Readability).** *Every formula satisfies one, and only one of the following conditions.*

1.  $A$  is atomic.
2.  $A$  is of the form  $\neg B$ .
3.  $A$  is of the form  $(B \wedge C)$ .
4.  $A$  is of the form  $(B \vee C)$ .
5.  $A$  is of the form  $(B \rightarrow C)$ .
6.  $A$  is of the form  $\forall x B$ .
7.  $A$  is of the form  $\exists x B$ .

*Moreover, in each case  $B$ , or  $B$  and  $C$ , are uniquely determined. This means that, e.g., there are no different pairs  $B, C$  and  $B', C'$  so that  $A$  is both of the form  $(B \rightarrow C)$  and  $(B' \rightarrow C')$ .*

*Proof.* The formation rules require that if a formula is not atomic, it must start with an opening parenthesis  $($ ,  $\neg$ , or a quantifier. On the other hand, every formula that starts with one of the following symbols must be atomic: a predicate symbol, a function symbol, a constant symbol,  $\perp$ .

So we really only have to show that if  $A$  is of the form  $(B * C)$  and also of the form  $(B' *' C')$ , then  $B \equiv B'$ ,  $C \equiv C'$ , and  $* = *'$ .

So suppose both  $A \equiv (B * C)$  and  $A \equiv (B' *' C')$ . Then either  $B \equiv B'$  or not. If it is, clearly  $* = *'$  and  $C \equiv C'$ , since they then are substrings of  $A$  that begin in the same place and are of the same length. The other case is  $B \not\equiv B'$ . Since  $B$  and  $B'$  are both substrings of  $A$  that begin at the same place, one must be a proper prefix of the other. But this is impossible by [Lemma 6.12](#). □

## 6.5 Main operator of a Formula

It is often useful to talk about the last operator used in constructing a formula  $A$ . This operator is called the *main operator* of  $A$ . Intuitively, it is the “outermost” operator of  $A$ . For example, the main operator of  $\neg A$  is  $\neg$ , the main operator of  $(A \vee B)$  is  $\vee$ , etc.

**Definition 6.15 (Main operator).** The *main operator* of a formula  $A$  is defined as follows:

1.  $A$  is atomic:  $A$  has no main operator.
2.  $A \equiv \neg B$ : the main operator of  $A$  is  $\neg$ .
3.  $A \equiv (B \wedge C)$ : the main operator of  $A$  is  $\wedge$ .
4.  $A \equiv (B \vee C)$ : the main operator of  $A$  is  $\vee$ .
5.  $A \equiv (B \rightarrow C)$ : the main operator of  $A$  is  $\rightarrow$ .
6.  $A \equiv \forall x B$ : the main operator of  $A$  is  $\forall$ .
7.  $A \equiv \exists x B$ : the main operator of  $A$  is  $\exists$ .

In each case, we intend the specific indicated *occurrence* of the main operator in the formula. For instance, since the formula  $((D \rightarrow E) \rightarrow (E \rightarrow D))$  is of the form  $(B \rightarrow C)$  where  $B$  is  $(D \rightarrow E)$  and  $C$  is  $(E \rightarrow D)$ , the second occurrence of  $\rightarrow$  is the main operator.

This is a *recursive* definition of a function which maps all non-atomic formulas to their main operator occurrence. Because of the way formulas are defined inductively, every formula  $A$  satisfies one of the cases in **Definition 6.15**. This guarantees that for each non-atomic formula  $A$  a main operator exists. Because each formula satisfies only one of these conditions, and because the smaller formulas from which  $A$  is constructed are uniquely determined in each case, the main operator occurrence of  $A$  is unique, and so we have defined a function.

We call formulas by the names in Table 6.1 depending on which symbol their main operator is. Recall, however, that defined operators do not officially appear in formulas. They are just abbreviations, so officially they cannot be the main operator of a formula. In proofs about all formulas they therefore do not have to be treated separately.

Main operator	Type of formula	Example
none	atomic (formula)	$\perp, R(t_1, \dots, t_n), t_1 = t_2$
$\neg$	negation	$\neg A$
$\wedge$	conjunction	$(A \wedge B)$
$\vee$	disjunction	$(A \vee B)$
$\rightarrow$	conditional	$(A \rightarrow B)$
$\leftrightarrow$	biconditional	$(A \leftrightarrow B)$
$\forall$	universal (formula)	$\forall x A$
$\exists$	existential (formula)	$\exists x A$

Table 6.1: Main operator and names of formulas

## 6.6 Subformulas

It is often useful to talk about the formulas that “make up” a given formula. We call these its *subformulas*. Any formula counts as a subformula of itself; a subformula of  $A$  other than  $A$  itself is a *proper subformula*.

**Definition 6.16 (Immediate Subformula).** If  $A$  is a formula, the *immediate subformulas* of  $A$  are defined inductively as follows:

1. Atomic formulas have no immediate subformulas.
2.  $A \equiv \neg B$ : The only immediate subformula of  $A$  is  $B$ .
3.  $A \equiv (B * C)$ : The immediate subformulas of  $A$  are  $B$  and  $C$  ( $*$  is any one of the two-place connectives).
4.  $A \equiv \forall x B$ : The only immediate subformula of  $A$  is  $B$ .

5.  $A \equiv \exists x B$ : The only immediate subformula of  $A$  is  $B$ .

**Definition 6.17 (Proper Subformula).** If  $A$  is a formula, the *proper subformulas* of  $A$  are defined recursively as follows:

1. Atomic formulas have no proper subformulas.
2.  $A \equiv \neg B$ : The proper subformulas of  $A$  are  $B$  together with all proper subformulas of  $B$ .
3.  $A \equiv (B * C)$ : The proper subformulas of  $A$  are  $B$ ,  $C$ , together with all proper subformulas of  $B$  and those of  $C$ .
4.  $A \equiv \forall x B$ : The proper subformulas of  $A$  are  $B$  together with all proper subformulas of  $B$ .
5.  $A \equiv \exists x B$ : The proper subformulas of  $A$  are  $B$  together with all proper subformulas of  $B$ .

**Definition 6.18 (Subformula).** The subformulas of  $A$  are  $A$  itself together with all its proper subformulas.

Note the subtle difference in how we have defined immediate subformulas and proper subformulas. In the first case, we have directly defined the immediate subformulas of a formula  $A$  for each possible form of  $A$ . It is an explicit definition by cases, and the cases mirror the inductive definition of the set of formulas. In the second case, we have also mirrored the way the set of all formulas is defined, but in each case we have also included the proper subformulas of the smaller formulas  $B$ ,  $C$  in addition to these formulas themselves. This makes the definition *recursive*. In general, a definition of a function on an inductively defined set (in our case, formulas) is recursive if the cases in the definition of the function make use of the function itself. To be well defined, we must make sure, however, that we only ever use the values

of the function for arguments that come “before” the one we are defining—in our case, when defining “proper subformula” for  $(B * C)$  we only use the proper subformulas of the “earlier” formulas  $B$  and  $C$ .

**Proposition 6.19.** *Suppose  $B$  is a subformula of  $A$  and  $C$  is a subformula of  $B$ . Then  $C$  is a subformula of  $A$ . In other words, the subformula relation is transitive.*

**Proposition 6.20.** *Suppose  $A$  is a formula with  $n$  connectives and quantifiers. Then  $A$  has at most  $2n + 1$  subformulas.*

## 6.7 Formation Sequences

Defining formulas via an inductive definition, and the complementary technique of proving properties of formulas via induction, is an elegant and efficient approach. However, it can also be useful to consider a more bottom-up, step-by-step approach to the construction of formulas, which we do here using the notion of a *formation sequence*. To show how terms and formulas can be introduced in this way without needing to refer to their inductive definitions, we first introduce the notion of an arbitrary string of symbols drawn from some language  $\mathcal{L}$ .

**Definition 6.21 (Strings).** Suppose  $\mathcal{L}$  is a first-order language. An  $\mathcal{L}$ -string is a finite sequence of symbols of  $\mathcal{L}$ . Where the language  $\mathcal{L}$  is clearly fixed by the context, we will often refer to a  $\mathcal{L}$ -string simply as a *string*.

**Example 6.22.** For any first-order language  $\mathcal{L}$ , all  $\mathcal{L}$ -formulas are  $\mathcal{L}$ -strings, but not conversely. For example,

$$)(v_0 \rightarrow \exists$$

is an  $\mathcal{L}$ -string but not an  $\mathcal{L}$ -formula.



**Definition 6.23 (Formation sequences for terms).** A finite sequence of  $\mathcal{L}$ -strings  $\langle t_0, \dots, t_n \rangle$  is a *formation sequence* for a term  $t$  if  $t \equiv t_n$  and for all  $i \leq n$ , either  $t_i$  is a variable or a constant symbol, or  $\mathcal{L}$  contains a  $k$ -ary function symbol  $f$  and there exist  $m_0, \dots, m_k < i$  such that  $t_i \equiv f(t_{m_0}, \dots, t_{m_k})$ .

**Example 6.24.** The sequence

$$\langle c_0, v_0, f_0^2(c_0, v_0), f_0^1(f_0^2(c_0, v_0)) \rangle$$

is a formation sequence for the term  $f_0^1(f_0^2(c_0, v_0))$ , as is

$$\langle v_0, c_0, f_0^2(c_0, v_0), f_0^1(f_0^2(c_0, v_0)) \rangle.$$

**Definition 6.25 (Formation sequences for formulas).** A finite sequence of  $\mathcal{L}$ -strings  $\langle A_0, \dots, A_n \rangle$  is a *formation sequence* for  $A$  if  $A \equiv A_n$  and for all  $i \leq n$ , either  $A_i$  is an atomic formula or there exist  $j, k < i$  and a variable  $x$  such that one of the following holds:

1.  $A_i \equiv \neg A_j$ .
2.  $A_i \equiv (A_j \wedge A_k)$ .
3.  $A_i \equiv (A_j \vee A_k)$ .
4.  $A_i \equiv (A_j \rightarrow A_k)$ .
5.  $A_i \equiv \forall x A_j$ .
6.  $A_i \equiv \exists x A_j$ .

**Example 6.26.**

$$\langle A_0^1(v_0), A_1^1(c_1), (A_1^1(c_1) \wedge A_0^1(v_0)), \exists v_0 (A_1^1(c_1) \wedge A_0^1(v_0)) \rangle$$

is a formation sequence of  $\exists v_0 (A_1^1(c_1) \wedge A_0^1(v_0))$ , as is

$$\langle A_0^1(v_0), A_1^1(c_1), (A_1^1(c_1) \wedge A_0^1(v_0)), A_1^1(c_1), \\ \forall v_1 A_0^1(v_0), \exists v_0 (A_1^1(c_1) \wedge A_0^1(v_0)) \rangle.$$

As can be seen from the second example, formation sequences may contain “junk”: formulas which are redundant or do not contribute to the construction.

**Proposition 6.27.** *Every formula  $A$  in  $\text{Frm}(\mathcal{L})$  has a formation sequence.*

*Proof.* Suppose  $A$  is atomic. Then the sequence  $\langle A \rangle$  is a formation sequence for  $A$ . Now suppose that  $B$  and  $C$  have formation sequences  $\langle B_0, \dots, B_n \rangle$  and  $\langle C_0, \dots, C_m \rangle$  respectively.

1. If  $A \equiv \neg B$ , then  $\langle B_0, \dots, B_n, \neg B_n \rangle$  is a formation sequence for  $A$ .
2. If  $A \equiv (B \wedge C)$ , then  $\langle B_0, \dots, B_n, C_0, \dots, C_m, (B_n \wedge C_m) \rangle$  is a formation sequence for  $A$ .
3. If  $A \equiv (B \vee C)$ , then  $\langle B_0, \dots, B_n, C_0, \dots, C_m, (B_n \vee C_m) \rangle$  is a formation sequence for  $A$ .
4. If  $A \equiv (B \rightarrow C)$ , then  $\langle B_0, \dots, B_n, C_0, \dots, C_m, (B_n \rightarrow C_m) \rangle$  is a formation sequence for  $A$ .
5. If  $A \equiv \forall x B$ , then  $\langle B_0, \dots, B_n, \forall x B_n \rangle$  is a formation sequence for  $A$ .
6. If  $A \equiv \exists x B$ , then  $\langle B_0, \dots, B_n, \exists x B_n \rangle$  is a formation sequence for  $A$ .

By the principle of induction on formulas, every formula has a formation sequence.  $\square$

We can also prove the converse. This is important because it shows that our two ways of defining formulas are equivalent: they give the same results. It also means that we can prove theorems about formulas by using ordinary induction on the length of formation sequences.

**Lemma 6.28.** *Suppose that  $\langle A_0, \dots, A_n \rangle$  is a formation sequence for  $A_n$ , and that  $k \leq n$ . Then  $\langle A_0, \dots, A_k \rangle$  is a formation sequence for  $A_k$ .*

*Proof.* Exercise. □

**Theorem 6.29.**  *$\text{Frm}(\mathcal{L})$  is the set of all expressions (strings of symbols) in the language  $\mathcal{L}$  with a formation sequence.*

*Proof.* Let  $F$  be the set of all strings of symbols in the language  $\mathcal{L}$  that have a formation sequence. We have seen in [Proposition 6.27](#) that  $\text{Frm}(\mathcal{L}) \subseteq F$ , so now we prove the converse.

Suppose  $A$  has a formation sequence  $\langle A_0, \dots, A_n \rangle$ . We prove that  $A \in \text{Frm}(\mathcal{L})$  by strong induction on  $n$ . Our induction hypothesis is that every string of symbols with a formation sequence of length  $m < n$  is in  $\text{Frm}(\mathcal{L})$ . By the definition of a formation sequence, either  $A \equiv A_n$  is atomic or there must exist  $j, k < n$  such that one of the following is the case:

1.  $A \equiv \neg A_j$ .
2.  $A \equiv (A_j \wedge A_k)$ .
3.  $A \equiv (A_j \vee A_k)$ .
4.  $A \equiv (A_j \rightarrow A_k)$ .
5.  $A \equiv \forall x A_j$ .
6.  $A \equiv \exists x A_j$ .

Now we reason by cases. If  $A$  is atomic then  $A_n \in \text{Frm}(\mathcal{L}_0)$ . Suppose instead that  $A \equiv (A_j \wedge A_k)$ . By [Lemma 6.28](#),  $\langle A_0, \dots, A_j \rangle$  and  $\langle A_0, \dots, A_k \rangle$  are formation sequences for  $A_j$  and  $A_k$ , respectively. Since these are proper initial subsequences of the formation sequence for  $A$ , they both have length less than  $n$ . Therefore by the induction hypothesis,  $A_j$  and  $A_k$  are in  $\text{Frm}(\mathcal{L}_0)$ , and by the definition of a formula, so is  $(A_j \wedge A_k)$ . The other cases follow by parallel reasoning. □

Formation sequences for terms have similar properties to those for formulas.

**Proposition 6.30.**  $\text{Trm}(\mathcal{L})$  is the set of all expressions  $t$  in the language  $\mathcal{L}$  such that there exists a (term) formation sequence for  $t$ .

*Proof.* Exercise. □

There are two types of “junk” that can appear in formation sequences: repeated elements, and elements that are irrelevant to the construction of the formation or term. We can eliminate both by looking at minimal formation sequences.

**Definition 6.31 (Minimal formation sequences).** A formation sequence  $\langle A_0, \dots, A_n \rangle$  for  $A$  is a *minimal formation sequence* for  $A$  if for every other formation sequence  $s$  for  $A$ , the length of  $s$  is greater than or equal to  $n + 1$ .

**Proposition 6.32.** *The following are equivalent:*

1.  $B$  is a sub-formula of  $A$ .
2.  $B$  occurs in every formation sequence of  $A$ .
3.  $B$  occurs in a minimal formation sequence of  $A$ .

*Proof.* Exercise. □

## 6.8 Free Variables and Sentences

**Definition 6.33 (Free occurrences of a variable).** The *free* occurrences of a variable in a formula are defined inductively as follows:

1.  $A$  is atomic: all variable occurrences in  $A$  are free.
2.  $A \equiv \neg B$ : the free variable occurrences of  $A$  are exactly

those of  $B$ .

3.  $A \equiv (B * C)$ : the free variable occurrences of  $A$  are those in  $B$  together with those in  $C$ .
4.  $A \equiv \forall x B$ : the free variable occurrences in  $A$  are all of those in  $B$  except for occurrences of  $x$ .
5.  $A \equiv \exists x B$ : the free variable occurrences in  $A$  are all of those in  $B$  except for occurrences of  $x$ .

**Definition 6.34 (Bound Variables).** An occurrence of a variable in a formula  $A$  is *bound* if it is not free.

**Definition 6.35 (Scope).** If  $\forall x B$  is an occurrence of a subformula in a formula  $A$ , then the corresponding occurrence of  $B$  in  $A$  is called the *scope* of the corresponding occurrence of  $\forall x$ . Similarly for  $\exists x$ .

If  $B$  is the scope of a quantifier occurrence  $\forall x$  or  $\exists x$  in  $A$ , then the free occurrences of  $x$  in  $B$  are bound in  $\forall x B$  and  $\exists x B$ . We say that these occurrences are *bound by* the mentioned quantifier occurrence.

**Example 6.36.** Consider the following formula:

$$\exists v_0 \underbrace{A_0^2(v_0, v_1)}_B$$

$B$  represents the scope of  $\exists v_0$ . The quantifier binds the occurrence of  $v_0$  in  $B$ , but does not bind the occurrence of  $v_1$ . So  $v_1$  is a free variable in this case.

We can now see how this might work in a more complicated formula  $A$ :

$$\forall v_0 \underbrace{(A_0^1(v_0) \rightarrow A_0^2(v_0, v_1))}_B \rightarrow \exists v_1 \underbrace{(A_1^2(v_0, v_1) \vee \forall v_0 \overbrace{\neg A_1^1(v_0)}^D)}_C$$

$B$  is the scope of the first  $\forall v_0$ ,  $C$  is the scope of  $\exists v_1$ , and  $D$  is the scope of the second  $\forall v_0$ . The first  $\forall v_0$  binds the occurrences of  $v_0$  in  $B$ ,  $\exists v_1$  binds the occurrence of  $v_1$  in  $C$ , and the second  $\forall v_0$  binds the occurrence of  $v_0$  in  $D$ . The first occurrence of  $v_1$  and the fourth occurrence of  $v_0$  are free in  $A$ . The last occurrence of  $v_0$  is free in  $D$ , but bound in  $C$  and  $A$ .

**Definition 6.37 (Sentence).** A formula  $A$  is a *sentence* iff it contains no free occurrences of variables.

## 6.9 Substitution

**Definition 6.38 (Substitution in a term).** We define  $s[t/x]$ , the result of *substituting*  $t$  for every occurrence of  $x$  in  $s$ , recursively:

1.  $s \equiv c$ :  $s[t/x]$  is just  $s$ .
2.  $s \equiv y$ :  $s[t/x]$  is also just  $s$ , provided  $y$  is a variable and  $y \neq x$ .
3.  $s \equiv x$ :  $s[t/x]$  is  $t$ .
4.  $s \equiv f(t_1, \dots, t_n)$ :  $s[t/x]$  is  $f(t_1[t/x], \dots, t_n[t/x])$ .

**Definition 6.39.** A term  $t$  is *free for*  $x$  in  $A$  if none of the free occurrences of  $x$  in  $A$  occur in the scope of a quantifier that binds a variable in  $t$ .

**Example 6.40.**

1.  $v_8$  is free for  $v_1$  in  $\exists v_3 A_4^2(v_3, v_1)$
2.  $f_1^2(v_1, v_2)$  is *not* free for  $v_0$  in  $\forall v_2 A_4^2(v_0, v_2)$

**Definition 6.41 (Substitution in a formula).** If  $A$  is a formula,  $x$  is a variable, and  $t$  is a term free for  $x$  in  $A$ , then  $A[t/x]$  is the result of substituting  $t$  for all free occurrences of  $x$  in  $A$ .

1.  $A \equiv \perp$ :  $A[t/x]$  is  $\perp$ .
2.  $A \equiv P(t_1, \dots, t_n)$ :  $A[t/x]$  is  $P(t_1[t/x], \dots, t_n[t/x])$ .
3.  $A \equiv t_1 = t_2$ :  $A[t/x]$  is  $t_1[t/x] = t_2[t/x]$ .
4.  $A \equiv \neg B$ :  $A[t/x]$  is  $\neg B[t/x]$ .
5.  $A \equiv (B \wedge C)$ :  $A[t/x]$  is  $(B[t/x] \wedge C[t/x])$ .
6.  $A \equiv (B \vee C)$ :  $A[t/x]$  is  $(B[t/x] \vee C[t/x])$ .
7.  $A \equiv (B \rightarrow C)$ :  $A[t/x]$  is  $(B[t/x] \rightarrow C[t/x])$ .
8.  $A \equiv \forall y B$ :  $A[t/x]$  is  $\forall y B[t/x]$ , provided  $y$  is a variable other than  $x$ ; otherwise  $A[t/x]$  is just  $A$ .
9.  $A \equiv \exists y B$ :  $A[t/x]$  is  $\exists y B[t/x]$ , provided  $y$  is a variable other than  $x$ ; otherwise  $A[t/x]$  is just  $A$ .

Note that substitution may be vacuous: If  $x$  does not occur in  $A$  at all, then  $A[t/x]$  is just  $A$ .

The restriction that  $t$  must be free for  $x$  in  $A$  is necessary to exclude cases like the following. If  $A \equiv \exists y x < y$  and  $t \equiv y$ ,

then  $A[t/x]$  would be  $\exists y y < y$ . In this case the free variable  $y$  is “captured” by the quantifier  $\exists y$  upon substitution, and that is undesirable. For instance, we would like it to be the case that whenever  $\forall x B$  holds, so does  $B[t/x]$ . But consider  $\forall x \exists y x < y$  (here  $B$  is  $\exists y x < y$ ). It is a sentence that is true about, e.g., the natural numbers: for every number  $x$  there is a number  $y$  greater than it. If we allowed  $y$  as a possible substitution for  $x$ , we would end up with  $B[y/x] \equiv \exists y y < y$ , which is false. We prevent this by requiring that none of the free variables in  $t$  would end up being bound by a quantifier in  $A$ .

We often use the following convention to avoid cumbersome notation: If  $A$  is a formula which may contain the variable  $x$  free, we also write  $A(x)$  to indicate this. When it is clear which  $A$  and  $x$  we have in mind, and  $t$  is a term (assumed to be free for  $x$  in  $A(x)$ ), then we write  $A(t)$  as short for  $A[t/x]$ . So for instance, we might say, “we call  $A(t)$  an instance of  $\forall x A(x)$ .” By this we mean that if  $A$  is any formula,  $x$  a variable, and  $t$  a term that’s free for  $x$  in  $A$ , then  $A[t/x]$  is an instance of  $\forall x A$ .

## Summary

A **first-order language** consists of **constant**, **function**, and **predicate** symbols. Function and constant symbols take a specified number of arguments. In the **language of arithmetic**, e.g., we have a single constant symbol  $0$ , one 1-place function symbol  $\iota$ , two 2-place function symbols  $+$  and  $\times$ , and one 2-place predicate symbol  $<$ . From **variables** and constant and function symbols we form the **terms** of a language. From the terms of a language together with its predicate symbols, as well as the **identity symbol**  $=$ , we form the **atomic formulas**. And in turn from them, using the logical connectives  $\neg$ ,  $\vee$ ,  $\wedge$ ,  $\rightarrow$ ,  $\leftrightarrow$  and the quantifiers  $\forall$  and  $\exists$  we form its formulas. Since we are careful to always include necessary parentheses in the process of forming terms and formulas, there is always exactly one way of reading a formula. This makes it possible to define things by induction on



the structure of formulas.

Occurrences of variables in formulas are sometimes governed by a corresponding quantifier: if a variable occurs in the **scope** of a quantifier it is considered **bound**, otherwise **free**. These concepts all have inductive definitions, and we also inductively define the operation of **substitution** of a term for a variable in a formula. Formulas without free variable occurrences are called **sentences**.

## Problems

**Problem 6.1.** Prove [Lemma 6.8](#).

**Problem 6.2.** Prove that for any term  $t$ ,  $l(t) = r(t)$ .

**Problem 6.3.** Prove [Lemma 6.12](#).

**Problem 6.4.** Prove [Proposition 6.13](#) (Hint: Formulate and prove a version of [Lemma 6.12](#) for terms.)

**Problem 6.5.** Prove [Proposition 6.19](#).

**Problem 6.6.** Prove [Proposition 6.20](#).

**Problem 6.7.** Prove [Lemma 6.28](#).

**Problem 6.8.** Prove [Proposition 6.30](#). Hint: use a similar strategy to that used in the proof of [Theorem 6.29](#).

**Problem 6.9.** Prove [Proposition 6.32](#).

**Problem 6.10.** Give an inductive definition of the bound variable occurrences along the lines of [Definition 6.33](#).

## CHAPTER 7

# *Semantics of First-Order Logic*

### 7.1 Introduction

Giving the meaning of expressions is the domain of semantics. The central concept in semantics is that of satisfaction in a structure. A structure gives meaning to the building blocks of the language: a domain is a non-empty set of objects. The quantifiers are interpreted as ranging over this domain, constant symbols are assigned elements in the domain, function symbols are assigned functions from the domain to itself, and predicate symbols are assigned relations on the domain. The domain together with assignments to the basic vocabulary constitutes a structure. Variables may appear in formulas, and in order to give a semantics, we also have to assign elements of the domain to them—this is a variable assignment. The satisfaction relation, finally, brings these together. A formula may be satisfied in a structure  $M$  relative to a variable assignment  $s$ , written as  $M, s \models A$ . This relation is also defined by induction on the structure of  $A$ , using the truth

tables for the logical connectives to define, say, satisfaction of  $(A \wedge B)$  in terms of satisfaction (or not) of  $A$  and  $B$ . It then turns out that the variable assignment is irrelevant if the formula  $A$  is a sentence, i.e., has no free variables, and so we can talk of sentences being simply satisfied (or not) in structures.

On the basis of the satisfaction relation  $M \models A$  for sentences we can then define the basic semantic notions of validity, entailment, and satisfiability. A sentence is valid,  $\models A$ , if every structure satisfies it. It is entailed by a set of sentences,  $\Gamma \models A$ , if every structure that satisfies all the sentences in  $\Gamma$  also satisfies  $A$ . And a set of sentences is satisfiable if some structure satisfies all sentences in it at the same time. Because formulas are inductively defined, and satisfaction is in turn defined by induction on the structure of formulas, we can use induction to prove properties of our semantics and to relate the semantic notions defined.

## 7.2 Structures for First-order Languages

First-order languages are, by themselves, *uninterpreted*: the constant symbols, function symbols, and predicate symbols have no specific meaning attached to them. Meanings are given by specifying a *structure*. It specifies the *domain*, i.e., the objects which the constant symbols pick out, the function symbols operate on, and the quantifiers range over. In addition, it specifies which constant symbols pick out which objects, how a function symbol maps objects to objects, and which objects the predicate symbols apply to. Structures are the basis for *semantic* notions in logic, e.g., the notion of consequence, validity, satisfiability. They are variously called “structures,” “interpretations,” or “models” in the literature.

**Definition 7.1 (Structures).** A structure  $M$ , for a language  $\mathcal{L}$  of first-order logic consists of the following elements:

1. *Domain*: a non-empty set,  $|M|$

2. *Interpretation of constant symbols:* for each constant symbol  $c$  of  $\mathcal{L}$ , an element  $c^M \in |M|$
3. *Interpretation of predicate symbols:* for each  $n$ -place predicate symbol  $R$  of  $\mathcal{L}$  (other than  $=$ ), an  $n$ -place relation  $R^M \subseteq |M|^n$
4. *Interpretation of function symbols:* for each  $n$ -place function symbol  $f$  of  $\mathcal{L}$ , an  $n$ -place function  $f^M: |M|^n \rightarrow |M|$

**Example 7.2.** A structure  $M$  for the language of arithmetic consists of a set, an element of  $|M|$ ,  $0^M$ , as interpretation of the constant symbol  $0$ , a one-place function  $\iota^M: |M| \rightarrow |M|$ , two two-place functions  $+^M$  and  $\times^M$ , both  $|M|^2 \rightarrow |M|$ , and a two-place relation  $<^M \subseteq |M|^2$ .

An obvious example of such a structure is the following:

1.  $|N| = \mathbb{N}$
2.  $0^N = 0$
3.  $\iota^N(n) = n + 1$  for all  $n \in \mathbb{N}$
4.  $+^N(n, m) = n + m$  for all  $n, m \in \mathbb{N}$
5.  $\times^N(n, m) = n \cdot m$  for all  $n, m \in \mathbb{N}$
6.  $<^N = \{\langle n, m \rangle : n \in \mathbb{N}, m \in \mathbb{N}, n < m\}$

The structure  $N$  for  $\mathcal{L}_A$  so defined is called the *standard model of arithmetic*, because it interprets the non-logical constants of  $\mathcal{L}_A$  exactly how you would expect.

However, there are many other possible structures for  $\mathcal{L}_A$ . For instance, we might take as the domain the set  $\mathbb{Z}$  of integers instead of  $\mathbb{N}$ , and define the interpretations of  $0$ ,  $\iota$ ,  $+$ ,  $\times$ ,  $<$  accordingly. But we can also define structures for  $\mathcal{L}_A$  which have nothing even remotely to do with numbers.

**Example 7.3.** A structure  $M$  for the language  $\mathcal{L}_Z$  of set theory requires just a set and a single-two place relation. So technically, e.g., the set of people plus the relation “ $x$  is older than  $y$ ” could be used as a structure for  $\mathcal{L}_Z$ , as well as  $\mathbb{N}$  together with  $n \geq m$  for  $n, m \in \mathbb{N}$ .

A particularly interesting structure for  $\mathcal{L}_Z$  in which the elements of the domain are actually sets, and the interpretation of  $\in$  actually is the relation “ $x$  is an element of  $y$ ” is the structure  $HF$  of *hereditarily finite sets*:

1.  $|HF| = \emptyset \cup \wp(\emptyset) \cup \wp(\wp(\emptyset)) \cup \wp(\wp(\wp(\emptyset))) \cup \dots$ ;
2.  $\in^{HF} = \{\langle x, y \rangle : x, y \in |HF|, x \in y\}$ .

The stipulations we make as to what counts as a structure impact our logic. For example, the choice to prevent empty domains ensures, given the usual account of satisfaction (or truth) for quantified sentences, that  $\exists x (A(x) \vee \neg A(x))$  is valid—that is, a logical truth. And the stipulation that all constant symbols must refer to an object in the domain ensures that the existential generalization is a sound pattern of inference:  $A(a)$ , therefore  $\exists x A(x)$ . If we allowed names to refer outside the domain, or to not refer, then we would be on our way to a *free logic*, in which existential generalization requires an additional premise:  $A(a)$  and  $\exists x x = a$ , therefore  $\exists x A(x)$ .

### 7.3 Covered Structures for First-order Languages

Recall that a term is *closed* if it contains no variables.

**Definition 7.4 (Value of closed terms).** If  $t$  is a closed term of the language  $\mathcal{L}$  and  $M$  is a structure for  $\mathcal{L}$ , the *value*  $\text{Val}^M(t)$  is defined as follows:

1. If  $t$  is just the constant symbol  $c$ , then  $\text{Val}^M(c) = c^M$ .

2. If  $t$  is of the form  $f(t_1, \dots, t_n)$ , then

$$\text{Val}^M(t) = f^M(\text{Val}^M(t_1), \dots, \text{Val}^M(t_n)).$$

**Definition 7.5 (Covered structure).** A structure is *covered* if every element of the domain is the value of some closed term.

**Example 7.6.** Let  $\mathcal{L}$  be the language with constant symbols *zero*, *one*, *two*, ..., the binary predicate symbol  $<$ , and the binary function symbols  $+$  and  $\times$ . Then a structure  $M$  for  $\mathcal{L}$  is the one with domain  $|M| = \{0, 1, 2, \dots\}$  and assignments  $\text{zero}^M = 0$ ,  $\text{one}^M = 1$ ,  $\text{two}^M = 2$ , and so forth. For the binary relation symbol  $<$ , the set  $<^M$  is the set of all pairs  $\langle c_1, c_2 \rangle \in |M|^2$  such that  $c_1$  is less than  $c_2$ : for example,  $\langle 1, 3 \rangle \in <^M$  but  $\langle 2, 2 \rangle \notin <^M$ . For the binary function symbol  $+$ , define  $+^M$  in the usual way—for example,  $+^M(2, 3)$  maps to 5, and similarly for the binary function symbol  $\times$ . Hence, the value of *four* is just 4, and the value of  $\times(\text{two}, +(\text{three}, \text{zero}))$  (or in infix notation,  $\text{two} \times (\text{three} + \text{zero})$ ) is

$$\begin{aligned} \text{Val}^M(\times(\text{two}, +(\text{three}, \text{zero}))) &= \\ &= \times^M(\text{Val}^M(\text{two}), \text{Val}^M(+(\text{three}, \text{zero}))) \\ &= \times^M(\text{Val}^M(\text{two}), +^M(\text{Val}^M(\text{three}), \text{Val}^M(\text{zero}))) \\ &= \times^M(\text{two}^M, +^M(\text{three}^M, \text{zero}^M)) \\ &= \times^M(2, +^M(3, 0)) \\ &= \times^M(2, 3) \\ &= 6 \end{aligned}$$

## 7.4 Satisfaction of a Formula in a Structure

The basic notion that relates expressions such as terms and formulas, on the one hand, and structures on the other, are those of *value* of a term and *satisfaction* of a formula. Informally, the

value of a term is an element of a structure—if the term is just a constant, its value is the object assigned to the constant by the structure, and if it is built up using function symbols, the value is computed from the values of constants and the functions assigned to the functions in the term. A formula is *satisfied* in a structure if the interpretation given to the predicates makes the formula true in the domain of the structure. This notion of satisfaction is specified inductively: the specification of the structure directly states when atomic formulas are satisfied, and we define when a complex formula is satisfied depending on the main connective or quantifier and whether or not the immediate subformulas are satisfied.

The case of the quantifiers here is a bit tricky, as the immediate subformula of a quantified formula has a free variable, and structures don't specify the values of variables. In order to deal with this difficulty, we also introduce *variable assignments* and define satisfaction not with respect to a structure alone, but with respect to a structure plus a variable assignment.

**Definition 7.7 (Variable Assignment).** A *variable assignment*  $s$  for a structure  $M$  is a function which maps each variable to an element of  $|M|$ , i.e.,  $s: \text{Var} \rightarrow |M|$ .

A structure assigns a value to each constant symbol, and a variable assignment to each variable. But we want to use terms built up from them to also name elements of the domain. For this we define the value of terms inductively. For constant symbols and variables the value is just as the structure or the variable assignment specifies it; for more complex terms it is computed recursively using the functions the structure assigns to the function symbols.

**Definition 7.8 (Value of Terms).** If  $t$  is a term of the language  $\mathcal{L}$ ,  $M$  is a structure for  $\mathcal{L}$ , and  $s$  is a variable assignment for  $M$ , the *value*  $\text{Val}_s^M(t)$  is defined as follows:

$$1. t \equiv c: \text{Val}_s^M(t) = c^M.$$

$$2. t \equiv x: \text{Val}_s^M(t) = s(x).$$

$$3. t \equiv f(t_1, \dots, t_n):$$

$$\text{Val}_s^M(t) = f^M(\text{Val}_s^M(t_1), \dots, \text{Val}_s^M(t_n)).$$

**Definition 7.9 (*x*-Variant).** If  $s$  is a variable assignment for a structure  $M$ , then any variable assignment  $s'$  for  $M$  which differs from  $s$  at most in what it assigns to  $x$  is called an *x*-variant of  $s$ . If  $s'$  is an *x*-variant of  $s$  we write  $s' \sim_x s$ .

Note that an *x*-variant of an assignment  $s$  does not *have* to assign something different to  $x$ . In fact, every assignment counts as an *x*-variant of itself.

**Definition 7.10.** If  $s$  is a variable assignment for a structure  $M$  and  $m \in |M|$ , then the assignment  $s[m/x]$  is the variable assignment defined by

$$s[m/x](y) = \begin{cases} m & \text{if } y \equiv x \\ s(y) & \text{otherwise.} \end{cases}$$

In other words,  $s[m/x]$  is the particular *x*-variant of  $s$  which assigns the domain element  $m$  to  $x$ , and assigns the same things to variables other than  $x$  that  $s$  does.

**Definition 7.11 (Satisfaction).** Satisfaction of a formula  $A$  in a structure  $M$  relative to a variable assignment  $s$ , in symbols:  $M, s \models A$ , is defined recursively as follows. (We write  $M, s \not\models A$  to mean “not  $M, s \models A$ .”)

$$1. A \equiv \perp: M, s \not\models A.$$

$$2. A \equiv R(t_1, \dots, t_n): M, s \models A \text{ iff } \langle \text{Val}_s^M(t_1), \dots, \text{Val}_s^M(t_n) \rangle \in$$



$R^M.$ 

3.  $A \equiv t_1 = t_2$ :  $M, s \models A$  iff  $\text{Val}_s^M(t_1) = \text{Val}_s^M(t_2)$ .
4.  $A \equiv \neg B$ :  $M, s \models A$  iff  $M, s \not\models B$ .
5.  $A \equiv (B \wedge C)$ :  $M, s \models A$  iff  $M, s \models B$  and  $M, s \models C$ .
6.  $A \equiv (B \vee C)$ :  $M, s \models A$  iff  $M, s \models B$  or  $M, s \models C$  (or both).
7.  $A \equiv (B \rightarrow C)$ :  $M, s \models A$  iff  $M, s \not\models B$  or  $M, s \models C$  (or both).
8.  $A \equiv \forall x B$ :  $M, s \models A$  iff for every element  $m \in |M|$ ,  $M, s[m/x] \models B$ .
9.  $A \equiv \exists x B$ :  $M, s \models A$  iff for at least one element  $m \in |M|$ ,  $M, s[m/x] \models B$ .

The variable assignments are important in the last two clauses. We cannot define satisfaction of  $\forall x B(x)$  by “for all  $m \in |M|$ ,  $M \models B(m)$ .” We cannot define satisfaction of  $\exists x B(x)$  by “for at least one  $m \in |M|$ ,  $M \models B(m)$ .” The reason is that if  $m \in |M|$ , it is not a symbol of the language, and so  $B(m)$  is not a formula (that is,  $B[m/x]$  is undefined). We also cannot assume that we have constant symbols or terms available that name every element of  $M$ , since there is nothing in the definition of structures that requires it. In the standard language, the set of constant symbols is countably infinite, so if  $|M|$  is not countable there aren’t even enough constant symbols to name every object.

We solve this problem by introducing variable assignments, which allow us to link variables directly with elements of the domain. Then instead of saying that, e.g.,  $\exists x B(x)$  is satisfied in  $M$  iff for at least one  $m \in |M|$ , we say it is satisfied in  $M$  *relative to*  $s$  iff  $B(x)$  is satisfied relative to  $s[m/x]$  for at least one  $m \in |M|$ .

**Example 7.12.** Let  $\mathcal{L} = \{a, b, f, R\}$  where  $a$  and  $b$  are constant symbols,  $f$  is a two-place function symbol, and  $R$  is a two-place predicate symbol. Consider the structure  $M$  defined by:

1.  $|M| = \{1, 2, 3, 4\}$
2.  $a^M = 1$
3.  $b^M = 2$
4.  $f^M(x, y) = x + y$  if  $x + y \leq 3$  and  $= 3$  otherwise.
5.  $R^M = \{\langle 1, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, \langle 2, 4 \rangle\}$

The function  $s(x) = 1$  that assigns  $1 \in |M|$  to every variable is a variable assignment for  $M$ .

Then

$$\text{Val}_s^M(f(a, b)) = f^M(\text{Val}_s^M(a), \text{Val}_s^M(b)).$$

Since  $a$  and  $b$  are constant symbols,  $\text{Val}_s^M(a) = a^M = 1$  and  $\text{Val}_s^M(b) = b^M = 2$ . So

$$\text{Val}_s^M(f(a, b)) = f^M(1, 2) = 1 + 2 = 3.$$

To compute the value of  $f(f(a, b), a)$  we have to consider

$$\text{Val}_s^M(f(f(a, b), a)) = f^M(\text{Val}_s^M(f(a, b)), \text{Val}_s^M(a)) = f^M(3, 1) = 3,$$

since  $3 + 1 > 3$ . Since  $s(x) = 1$  and  $\text{Val}_s^M(x) = s(x)$ , we also have

$$\text{Val}_s^M(f(f(a, b), x)) = f^M(\text{Val}_s^M(f(a, b)), \text{Val}_s^M(x)) = f^M(3, 1) = 3,$$

An atomic formula  $R(t_1, t_2)$  is satisfied if the tuple of values of its arguments, i.e.,  $\langle \text{Val}_s^M(t_1), \text{Val}_s^M(t_2) \rangle$ , is an element of  $R^M$ . So, e.g., we have  $M, s \models R(b, f(a, b))$  since  $\langle \text{Val}_s^M(b), \text{Val}_s^M(f(a, b)) \rangle = \langle 2, 3 \rangle \in R^M$ , but  $M, s \not\models R(x, f(a, b))$  since  $\langle 1, 3 \rangle \notin R^M[s]$ .

To determine if a non-atomic formula  $A$  is satisfied, you apply the clauses in the inductive definition that applies to the main connective. For instance, the main connective in  $R(a, a) \rightarrow (R(b, x) \vee R(x, b))$  is the  $\rightarrow$ , and

$$M, s \models R(a, a) \rightarrow (R(b, x) \vee R(x, b)) \text{ iff}$$

$$M, s \not\models R(a, a) \text{ or } M, s \models R(b, x) \vee R(x, b)$$

Since  $M, s \models R(a, a)$  (because  $\langle 1, 1 \rangle \in R^M$ ) we can't yet determine the answer and must first figure out if  $M, s \models R(b, x) \vee R(x, b)$ :

$$\begin{aligned} M, s \models R(b, x) \vee R(x, b) &\text{ iff} \\ M, s \models R(b, x) \text{ or } M, s \models R(x, b) \end{aligned}$$

And this is the case, since  $M, s \models R(x, b)$  (because  $\langle 1, 2 \rangle \in R^M$ ).

Recall that an  $x$ -variant of  $s$  is a variable assignment that differs from  $s$  at most in what it assigns to  $x$ . For every element of  $|M|$ , there is an  $x$ -variant of  $s$ :

$$\begin{aligned} s_1 &= s[1/x], & s_2 &= s[2/x], \\ s_3 &= s[3/x], & s_4 &= s[4/x]. \end{aligned}$$

So, e.g.,  $s_2(x) = 2$  and  $s_2(y) = s(y) = 1$  for all variables  $y$  other than  $x$ . These are all the  $x$ -variants of  $s$  for the structure  $M$ , since  $|M| = \{1, 2, 3, 4\}$ . Note, in particular, that  $s_1 = s$  ( $s$  is always an  $x$ -variant of itself).

To determine if an existentially quantified formula  $\exists x A(x)$  is satisfied, we have to determine if  $M, s[m/x] \models A(x)$  for at least one  $m \in |M|$ . So,

$$M, s \models \exists x (R(b, x) \vee R(x, b)),$$

since  $M, s[1/x] \models R(b, x) \vee R(x, b)$  ( $s[3/x]$  would also fit the bill). But,

$$M, s \not\models \exists x (R(b, x) \wedge R(x, b))$$

since, whichever  $m \in |M|$  we pick,  $M, s[m/x] \not\models R(b, x) \wedge R(x, b)$ .

To determine if a universally quantified formula  $\forall x A(x)$  is satisfied, we have to determine if  $M, s[m/x] \models A(x)$  for all  $m \in |M|$ . So,

$$M, s \models \forall x (R(x, a) \rightarrow R(a, x)),$$

since  $M, s[m/x] \models R(x, a) \rightarrow R(a, x)$  for all  $m \in |M|$ . For  $m = 1$ , we have  $M, s[1/x] \models R(a, x)$  so the consequent is true; for  $m = 2, 3$ , and  $4$ , we have  $M, s[m/x] \not\models R(x, a)$ , so the antecedent is false. But,

$$M, s \not\models \forall x (R(a, x) \rightarrow R(x, a))$$

since  $M, s[2/x] \not\models R(a, x) \rightarrow R(x, a)$  (because  $M, s[2/x] \models R(a, x)$  and  $M, s[2/x] \not\models R(x, a)$ ).

For a more complicated case, consider

$$\forall x (R(a, x) \rightarrow \exists y R(x, y)).$$

Since  $M, s[3/x] \not\models R(a, x)$  and  $M, s[4/x] \not\models R(a, x)$ , the interesting cases where we have to worry about the consequent of the conditional are only  $m = 1$  and  $m = 2$ . Does  $M, s[1/x] \models \exists y R(x, y)$  hold? It does if there is at least one  $n \in |M|$  so that  $M, s[1/x][n/y] \models R(x, y)$ . In fact, if we take  $n = 1$ , we have  $s[1/x][n/y] = s[1/y] = s$ . Since  $s(x) = 1$ ,  $s(y) = 1$ , and  $\langle 1, 1 \rangle \in R^M$ , the answer is yes.

To determine if  $M, s[2/x] \models \exists y R(x, y)$ , we have to look at the variable assignments  $s[2/x][n/y]$ . Here, for  $n = 1$ , this assignment is  $s_2 = s[2/x]$ , which does not satisfy  $R(x, y)$  ( $s_2(x) = 2$ ,  $s_2(y) = 1$ , and  $\langle 2, 1 \rangle \notin R^M$ ). However, consider  $s[2/x][3/y] = s_2[3/y]$ .  $M, s_2[3/y] \models R(x, y)$  since  $\langle 2, 3 \rangle \in R^M$ , and so  $M, s_2 \models \exists y R(x, y)$ .

So, for all  $n \in |M|$ , either  $M, s[m/x] \not\models R(a, x)$  (if  $m = 3, 4$ ) or  $M, s[m/x] \models \exists y R(x, y)$  (if  $m = 1, 2$ ), and so

$$M, s \models \forall x (R(a, x) \rightarrow \exists y R(x, y)).$$

On the other hand,

$$M, s \not\models \exists x (R(a, x) \wedge \forall y R(x, y)).$$

We have  $M, s[m/x] \models R(a, x)$  only for  $m = 1$  and  $m = 2$ . But for both of these values of  $m$ , there is in turn an  $n \in |M|$ , namely  $n = 4$ , so that  $M, s[m/x][n/y] \not\models R(x, y)$  and so  $M, s[m/x] \not\models \forall y R(x, y)$  for  $m = 1$  and  $m = 2$ . In sum, there is no  $m \in |M|$  such that  $M, s[m/x] \models R(a, x) \wedge \forall y R(x, y)$ .

## 7.5 Variable Assignments

A variable assignment  $s$  provides a value for *every* variable—and there are infinitely many of them. This is of course not necessary. We require variable assignments to assign values to all variables simply because it makes things a lot easier. The value of a term  $t$ , and whether or not a formula  $A$  is satisfied in a structure with respect to  $s$ , only depend on the assignments  $s$  makes to the variables in  $t$  and the free variables of  $A$ . This is the content of the next two propositions. To make the idea of “depends on” precise, we show that any two variable assignments that agree on all the variables in  $t$  give the same value, and that  $A$  is satisfied relative to one iff it is satisfied relative to the other if two variable assignments agree on all free variables of  $A$ .

**Proposition 7.13.** *If the variables in a term  $t$  are among  $x_1, \dots, x_n$ , and  $s_1(x_i) = s_2(x_i)$  for  $i = 1, \dots, n$ , then  $\text{Val}_{s_1}^M(t) = \text{Val}_{s_2}^M(t)$ .*

*Proof.* By induction on the complexity of  $t$ . For the base case,  $t$  can be a constant symbol or one of the variables  $x_1, \dots, x_n$ . If  $t = c$ , then  $\text{Val}_{s_1}^M(t) = c^M = \text{Val}_{s_2}^M(t)$ . If  $t = x_i$ ,  $s_1(x_i) = s_2(x_i)$  by the hypothesis of the proposition, and so  $\text{Val}_{s_1}^M(t) = s_1(x_i) = s_2(x_i) = \text{Val}_{s_2}^M(t)$ .

For the inductive step, assume that  $t = f(t_1, \dots, t_k)$  and that the claim holds for  $t_1, \dots, t_k$ . Then

$$\begin{aligned} \text{Val}_{s_1}^M(t) &= \text{Val}_{s_1}^M(f(t_1, \dots, t_k)) = \\ &= f^M(\text{Val}_{s_1}^M(t_1), \dots, \text{Val}_{s_1}^M(t_k)) \end{aligned}$$

For  $j = 1, \dots, k$ , the variables of  $t_j$  are among  $x_1, \dots, x_n$ . By induction hypothesis,  $\text{Val}_{s_1}^M(t_j) = \text{Val}_{s_2}^M(t_j)$ . So,

$$\begin{aligned} \text{Val}_{s_1}^M(t) &= \text{Val}_{s_1}^M(f(t_1, \dots, t_k)) = \\ &= f^M(\text{Val}_{s_1}^M(t_1), \dots, \text{Val}_{s_1}^M(t_k)) = \\ &= f^M(\text{Val}_{s_2}^M(t_1), \dots, \text{Val}_{s_2}^M(t_k)) = \\ &= \text{Val}_{s_2}^M(f(t_1, \dots, t_k)) = \text{Val}_{s_2}^M(t). \end{aligned}$$

□

**Proposition 7.14.** *If the free variables in  $A$  are among  $x_1, \dots, x_n$ , and  $s_1(x_i) = s_2(x_i)$  for  $i = 1, \dots, n$ , then  $M, s_1 \vDash A$  iff  $M, s_2 \vDash A$ .*

*Proof.* We use induction on the complexity of  $A$ . For the base case, where  $A$  is atomic,  $A$  can be:  $\perp$ ,  $R(t_1, \dots, t_k)$  for a  $k$ -place predicate  $R$  and terms  $t_1, \dots, t_k$ , or  $t_1 = t_2$  for terms  $t_1$  and  $t_2$ .

1.  $A \equiv \perp$ : both  $M, s_1 \not\vDash A$  and  $M, s_2 \not\vDash A$ .

2.  $A \equiv R(t_1, \dots, t_k)$ : let  $M, s_1 \vDash A$ . Then

$$\langle \text{Val}_{s_1}^M(t_1), \dots, \text{Val}_{s_1}^M(t_k) \rangle \in R^M.$$

For  $i = 1, \dots, k$ ,  $\text{Val}_{s_1}^M(t_i) = \text{Val}_{s_2}^M(t_i)$  by **Proposition 7.13**.

So we also have  $\langle \text{Val}_{s_2}^M(t_1), \dots, \text{Val}_{s_2}^M(t_k) \rangle \in R^M$ .

3.  $A \equiv t_1 = t_2$ : suppose  $M, s_1 \vDash A$ . Then  $\text{Val}_{s_1}^M(t_1) = \text{Val}_{s_1}^M(t_2)$ . So,

$$\begin{aligned} \text{Val}_{s_2}^M(t_1) &= \text{Val}_{s_1}^M(t_1) && \text{(by Proposition 7.13)} \\ &= \text{Val}_{s_1}^M(t_2) && \text{(since } M, s_1 \vDash t_1 = t_2\text{)} \\ &= \text{Val}_{s_2}^M(t_2) && \text{(by Proposition 7.13),} \end{aligned}$$

so  $M, s_2 \vDash t_1 = t_2$ .

Now assume  $M, s_1 \vDash B$  iff  $M, s_2 \vDash B$  for all formulas  $B$  less complex than  $A$ . The induction step proceeds by cases determined by the main operator of  $A$ . In each case, we only demonstrate the forward direction of the biconditional; the proof of the reverse direction is symmetrical. In all cases except those for the quantifiers, we apply the induction hypothesis to sub-formulas  $B$  of  $A$ . The free variables of  $B$  are among those of  $A$ . Thus, if  $s_1$  and  $s_2$  agree on the free variables of  $A$ , they also agree on those of  $B$ , and the induction hypothesis applies to  $B$ .

1.  $A \equiv \neg B$ : if  $M, s_1 \vDash A$ , then  $M, s_1 \not\vDash B$ , so by the induction hypothesis,  $M, s_2 \not\vDash B$ , hence  $M, s_2 \vDash A$ .

2.  $A \equiv B \wedge C$ : exercise.
3.  $A \equiv B \vee C$ : if  $M, s_1 \models A$ , then  $M, s_1 \models B$  or  $M, s_1 \models C$ . By induction hypothesis,  $M, s_2 \models B$  or  $M, s_2 \models C$ , so  $M, s_2 \models A$ .
4.  $A \equiv B \rightarrow C$ : exercise.
5.  $A \equiv \exists x B$ : if  $M, s_1 \models A$ , there is an  $m \in |M|$  so that  $M, s_1[m/x] \models B$ . Let  $s'_1 = s_1[m/x]$  and  $s'_2 = s_2[m/x]$ . The free variables of  $B$  are among  $x_1, \dots, x_n$ , and  $x$ .  $s'_1(x_i) = s'_2(x_i)$ , since  $s'_1$  and  $s'_2$  are  $x$ -variants of  $s_1$  and  $s_2$ , respectively, and by hypothesis  $s_1(x_i) = s_2(x_i)$ .  $s'_1(x) = s'_2(x) = m$  by the way we have defined  $s'_1$  and  $s'_2$ . Then the induction hypothesis applies to  $B$  and  $s'_1, s'_2$ , so  $M, s'_2 \models B$ . Hence, since  $s'_2 = s_2[m/x]$ , there is an  $m \in |M|$  such that  $M, s_2[m/x] \models B$ , and so  $M, s_2 \models A$ .
6.  $A \equiv \forall x B$ : exercise.

By induction, we get that  $M, s_1 \models A$  iff  $M, s_2 \models A$  whenever the free variables in  $A$  are among  $x_1, \dots, x_n$  and  $s_1(x_i) = s_2(x_i)$  for  $i = 1, \dots, n$ .  $\square$

Sentences have no free variables, so any two variable assignments assign the same things to all the (zero) free variables of any sentence. The proposition just proved then means that whether or not a sentence is satisfied in a structure relative to a variable assignment is completely independent of the assignment. We'll record this fact. It justifies the definition of satisfaction of a sentence in a structure (without mentioning a variable assignment) that follows.

**Corollary 7.15.** *If  $A$  is a sentence and  $s$  a variable assignment, then  $M, s \models A$  iff  $M, s' \models A$  for every variable assignment  $s'$ .*

*Proof.* Let  $s'$  be any variable assignment. Since  $A$  is a sentence, it has no free variables, and so every variable assignment  $s'$  trivially assigns the same things to all free variables of  $A$  as does  $s$ . So the

condition of **Proposition 7.14** is satisfied, and we have  $M, s \models A$  iff  $M, s' \models A$ .  $\square$

**Definition 7.16.** If  $A$  is a sentence, we say that a structure  $M$  satisfies  $A$ ,  $M \models A$ , iff  $M, s \models A$  for all variable assignments  $s$ .

If  $M \models A$ , we also simply say that  $A$  is true in  $M$ .

**Proposition 7.17.** Let  $M$  be a structure,  $A$  be a sentence, and  $s$  a variable assignment.  $M \models A$  iff  $M, s \models A$ .

*Proof.* Exercise.  $\square$

**Proposition 7.18.** Suppose  $A(x)$  only contains  $x$  free, and  $M$  is a structure. Then:

1.  $M \models \exists x A(x)$  iff  $M, s \models A(x)$  for at least one variable assignment  $s$ .
2.  $M \models \forall x A(x)$  iff  $M, s \models A(x)$  for all variable assignments  $s$ .

*Proof.* Exercise.  $\square$

## 7.6 Extensionality

Extensionality, sometimes called relevance, can be expressed informally as follows: the only factors that bear upon the satisfaction of formula  $A$  in a structure  $M$  relative to a variable assignment  $s$ , are the size of the domain and the assignments made by  $M$  and  $s$  to the elements of the language that actually appear in  $A$ .

One immediate consequence of extensionality is that where two structures  $M$  and  $M'$  agree on all the elements of the language appearing in a sentence  $A$  and have the same domain,  $M$  and  $M'$  must also agree on whether or not  $A$  itself is true.



**Proposition 7.19 (Extensionality).** *Let  $A$  be a formula, and  $M_1$  and  $M_2$  be structures with  $|M_1| = |M_2|$ , and  $s$  a variable assignment on  $|M_1| = |M_2|$ . If  $c^{M_1} = c^{M_2}$ ,  $R^{M_1} = R^{M_2}$ , and  $f^{M_1} = f^{M_2}$  for every constant symbol  $c$ , relation symbol  $R$ , and function symbol  $f$  occurring in  $A$ , then  $M_1, s \models A$  iff  $M_2, s \models A$ .*

*Proof.* First prove (by induction on  $t$ ) that for every term,  $\text{Val}_s^{M_1}(t) = \text{Val}_s^{M_2}(t)$ . Then prove the proposition by induction on  $A$ , making use of the claim just proved for the induction basis (where  $A$  is atomic).  $\square$

**Corollary 7.20 (Extensionality for Sentences).** *Let  $A$  be a sentence and  $M_1, M_2$  as in Proposition 7.19. Then  $M_1 \models A$  iff  $M_2 \models A$ .*

*Proof.* Follows from Proposition 7.19 by Corollary 7.15.  $\square$

Moreover, the value of a term, and whether or not a structure satisfies a formula, only depend on the values of its subterms.

**Proposition 7.21.** *Let  $M$  be a structure,  $t$  and  $t'$  terms, and  $s$  a variable assignment. Then  $\text{Val}_s^M(t[t'/x]) = \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(t)$ .*

*Proof.* By induction on  $t$ .

1. If  $t$  is a constant, say,  $t \equiv c$ , then  $t[t'/x] = c$ , and  $\text{Val}_s^M(c) = c^M = \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(c)$ .
2. If  $t$  is a variable other than  $x$ , say,  $t \equiv y$ , then  $t[t'/x] = y$ , and  $\text{Val}_s^M(y) = \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(y)$  since  $s \sim_x s[\text{Val}_s^M(t')/x]$ .
3. If  $t \equiv x$ , then  $t[t'/x] = t'$ . But  $\text{Val}_{s[\text{Val}_s^M(t')/x]}^M(x) = \text{Val}_s^M(t')$  by definition of  $s[\text{Val}_s^M(t')/x]$ .
4. If  $t \equiv f(t_1, \dots, t_n)$  then we have:

$$\text{Val}_s^M(t[t'/x]) =$$

$$\begin{aligned}
&= \text{Val}_s^M(f(t_1[t'/x], \dots, t_n[t'/x])) \\
&\quad \text{by definition of } t[t'/x] \\
&= f^M(\text{Val}_s^M(t_1[t'/x]), \dots, \text{Val}_s^M(t_n[t'/x])) \\
&\quad \text{by definition of } \text{Val}_s^M(f(\dots)) \\
&= f^M(\text{Val}_{s[\text{Val}_s^M(t')/x]}^M(t_1), \dots, \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(t_n)) \\
&\quad \text{by induction hypothesis} \\
&= \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(t) \text{ by definition of } \text{Val}_{s[\text{Val}_s^M(t')/x]}^M(f(\dots)) \quad \square
\end{aligned}$$

**Proposition 7.22.** *Let  $M$  be a structure,  $A$  a formula,  $t'$  a term, and  $s$  a variable assignment. Then  $M, s \models A[t'/x]$  iff  $M, s[\text{Val}_s^M(t')/x] \models A$ .*

*Proof.* Exercise. □

The point of **Propositions 7.21** and **7.22** is the following. Suppose we have a term  $t$  or a formula  $A$  and some term  $t'$ , and we want to know the value of  $t[t'/x]$  or whether or not  $A[t'/x]$  is satisfied in a structure  $M$  relative to a variable assignment  $s$ . Then we can either perform the substitution first and then consider the value or satisfaction relative to  $M$  and  $s$ , or we can first determine the value  $m = \text{Val}_s^M(t')$  of  $t'$  in  $M$  relative to  $s$ , change the variable assignment to  $s[m/x]$  and then consider the value of  $t$  in  $M$  and  $s[m/x]$ , or whether  $M, s[m/x] \models A$ . **Propositions 7.21** and **7.22** guarantee that the answer will be the same, whichever way we do it.

## 7.7 Semantic Notions

Given the definition of structures for first-order languages, we can define some basic semantic properties of and relationships between sentences. The simplest of these is the notion of *validity* of a sentence. A sentence is valid if it is satisfied in every structure. Valid sentences are those that are satisfied regardless of how

the non-logical symbols in it are interpreted. Valid sentences are therefore also called *logical truths*—they are true, i.e., satisfied, in any structure and hence their truth depends only on the logical symbols occurring in them and their syntactic structure, but not on the non-logical symbols or their interpretation.

**Definition 7.23 (Validity).** A sentence  $A$  is *valid*,  $\vDash A$ , iff  $M \vDash A$  for every structure  $M$ .

**Definition 7.24 (Entailment).** A set of sentences  $\Gamma$  *entails* a sentence  $A$ ,  $\Gamma \vDash A$ , iff for every structure  $M$  with  $M \vDash \Gamma$ ,  $M \vDash A$ .

**Definition 7.25 (Satisfiability).** A set of sentences  $\Gamma$  is *satisfiable* if  $M \vDash \Gamma$  for some structure  $M$ . If  $\Gamma$  is not satisfiable it is called *unsatisfiable*.

**Proposition 7.26.** A sentence  $A$  is valid iff  $\Gamma \vDash A$  for every set of sentences  $\Gamma$ .

*Proof.* For the forward direction, let  $A$  be valid, and let  $\Gamma$  be a set of sentences. Let  $M$  be a structure so that  $M \vDash \Gamma$ . Since  $A$  is valid,  $M \vDash A$ , hence  $\Gamma \vDash A$ .

For the contrapositive of the reverse direction, let  $A$  be invalid, so there is a structure  $M$  with  $M \not\vDash A$ . When  $\Gamma = \{\top\}$ , since  $\top$  is valid,  $M \vDash \Gamma$ . Hence, there is a structure  $M$  so that  $M \vDash \Gamma$  but  $M \not\vDash A$ , hence  $\Gamma$  does not entail  $A$ .  $\square$

**Proposition 7.27.**  $\Gamma \vDash A$  iff  $\Gamma \cup \{\neg A\}$  is unsatisfiable.

*Proof.* For the forward direction, suppose  $\Gamma \vDash A$  and suppose to the contrary that there is a structure  $M$  so that  $M \vDash \Gamma \cup \{\neg A\}$ . Since  $M \vDash \Gamma$  and  $\Gamma \vDash A$ ,  $M \vDash A$ . Also, since  $M \vDash \Gamma \cup \{\neg A\}$ ,  $M \vDash$

$\neg A$ , so we have both  $M \models A$  and  $M \not\models A$ , a contradiction. Hence, there can be no such structure  $M$ , so  $\Gamma \cup \{\neg A\}$  is unsatisfiable.

For the reverse direction, suppose  $\Gamma \cup \{\neg A\}$  is unsatisfiable. So for every structure  $M$ , either  $M \not\models \Gamma$  or  $M \models A$ . Hence, for every structure  $M$  with  $M \models \Gamma$ ,  $M \models A$ , so  $\Gamma \models A$ .  $\square$

**Proposition 7.28.** *If  $\Gamma \subseteq \Gamma'$  and  $\Gamma \models A$ , then  $\Gamma' \models A$ .*

*Proof.* Suppose that  $\Gamma \subseteq \Gamma'$  and  $\Gamma \models A$ . Let  $M$  be a structure such that  $M \models \Gamma'$ ; then  $M \models \Gamma$ , and since  $\Gamma \models A$ , we get that  $M \models A$ . Hence, whenever  $M \models \Gamma'$ ,  $M \models A$ , so  $\Gamma' \models A$ .  $\square$

**Theorem 7.29 (Semantic Deduction Theorem).**  *$\Gamma \cup \{A\} \models B$  iff  $\Gamma \models A \rightarrow B$ .*

*Proof.* For the forward direction, let  $\Gamma \cup \{A\} \models B$  and let  $M$  be a structure so that  $M \models \Gamma$ . If  $M \models A$ , then  $M \models \Gamma \cup \{A\}$ , so since  $\Gamma \cup \{A\}$  entails  $B$ , we get  $M \models B$ . Therefore,  $M \models A \rightarrow B$ , so  $\Gamma \models A \rightarrow B$ .

For the reverse direction, let  $\Gamma \models A \rightarrow B$  and  $M$  be a structure so that  $M \models \Gamma \cup \{A\}$ . Then  $M \models \Gamma$ , so  $M \models A \rightarrow B$ , and since  $M \models A$ ,  $M \models B$ . Hence, whenever  $M \models \Gamma \cup \{A\}$ ,  $M \models B$ , so  $\Gamma \cup \{A\} \models B$ .  $\square$

**Proposition 7.30.** *Let  $M$  be a structure, and  $A(x)$  a formula with one free variable  $x$ , and  $t$  a closed term. Then:*

1.  $A(t) \models \exists x A(x)$
2.  $\forall x A(x) \models A(t)$

*Proof.* 1. Suppose  $M \models A(t)$ . Let  $s$  be a variable assignment with  $s(x) = \text{Val}^M(t)$ . Then  $M, s \models A(t)$  since  $A(t)$  is a sentence. By **Proposition 7.22**,  $M, s \models A(x)$ . By **Proposition 7.18**,  $M \models \exists x A(x)$ .

2. Exercise.  $\square$

## Summary

The **semantics** for a first-order language is given by a **structure** for that language. It consists of a **domain** and elements of that domain are assigned to each constant symbol. Function symbols are interpreted by functions and relation symbols by relation on the domain. A function from the set of variables to the domain is a **variable assignment**. The relation of **satisfaction** relates structures, variable assignments and formulas;  $M, s \models A$  is defined by induction on the structure of  $A$ .  $M, s \models A$  only depends on the interpretation of the symbols actually occurring in  $A$ , and in particular does not depend on  $s$  if  $A$  contains no free variables. So if  $A$  is a sentence,  $M \models A$  if  $M, s \models A$  for any (or all)  $s$ .

The satisfaction relation is the basis for all semantic notions. A sentence is **valid**,  $\models A$ , if it is satisfied in every structure. A sentence  $A$  is **entailed** by set of sentences  $\Gamma$ ,  $\Gamma \models A$ , iff  $M \models A$  for all  $M$  which satisfy every sentence in  $\Gamma$ . A set  $\Gamma$  is **satisfiable** iff there is some structure that satisfies every sentence in  $\Gamma$ , otherwise **unsatisfiable**. These notions are interrelated, e.g.,  $\Gamma \models A$  iff  $\Gamma \cup \{\neg A\}$  is unsatisfiable.

## Problems

**Problem 7.1.** Is  $N$ , the standard model of arithmetic, covered? Explain.

**Problem 7.2.** Let  $\mathcal{L} = \{c, f, A\}$  with one constant symbol, one one-place function symbol and one two-place predicate symbol, and let the structure  $M$  be given by

1.  $|M| = \{1, 2, 3\}$
2.  $c^M = 3$
3.  $f^M(1) = 2, f^M(2) = 3, f^M(3) = 2$
4.  $A^M = \{\langle 1, 2 \rangle, \langle 2, 3 \rangle, \langle 3, 3 \rangle\}$

(a) Let  $s(v) = 1$  for all variables  $v$ . Find out whether

$$M, s \models \exists x (A(f(z), c) \rightarrow \forall y (A(y, x) \vee A(f(y), x)))$$

Explain why or why not.

(b) Give a different structure and variable assignment in which the formula is not satisfied.

**Problem 7.3.** Complete the proof of [Proposition 7.14](#).

**Problem 7.4.** Prove [Proposition 7.17](#)

**Problem 7.5.** Prove [Proposition 7.18](#).

**Problem 7.6.** Suppose  $\mathcal{L}$  is a language without function symbols. Given a structure  $M$ ,  $c$  a constant symbol and  $a \in |M|$ , define  $M[a/c]$  to be the structure that is just like  $M$ , except that  $c^{M[a/c]} = a$ . Define  $M \models A$  for sentences  $A$  by:

1.  $A \equiv \perp$ :  $\text{not } M \models A$ .
2.  $A \equiv R(d_1, \dots, d_n)$ :  $M \models A$  iff  $\langle d_1^M, \dots, d_n^M \rangle \in R^M$ .
3.  $A \equiv d_1 = d_2$ :  $M \models A$  iff  $d_1^M = d_2^M$ .
4.  $A \equiv \neg B$ :  $M \models A$  iff  $\text{not } M \models B$ .
5.  $A \equiv (B \wedge C)$ :  $M \models A$  iff  $M \models B$  and  $M \models C$ .
6.  $A \equiv (B \vee C)$ :  $M \models A$  iff  $M \models B$  or  $M \models C$  (or both).
7.  $A \equiv (B \rightarrow C)$ :  $M \models A$  iff  $\text{not } M \models B$  or  $M \models C$  (or both).
8.  $A \equiv \forall x B$ :  $M \models A$  iff for all  $a \in |M|$ ,  $M[a/c] \models B[c/x]$ , if  $c$  does not occur in  $B$ .
9.  $A \equiv \exists x B$ :  $M \models A$  iff there is an  $a \in |M|$  such that  $M[a/c] \models B[c/x]$ , if  $c$  does not occur in  $B$ .

Let  $x_1, \dots, x_n$  be all free variables in  $A$ ,  $c_1, \dots, c_n$  constant symbols not in  $A$ ,  $a_1, \dots, a_n \in |M|$ , and  $s(x_i) = a_i$ .

Show that  $M, s \models A$  iff  $M[a_1/c_1, \dots, a_n/c_n] \models A[c_1/x_1] \dots [c_n/x_n]$ .

(This problem shows that it is possible to give a semantics for first-order logic that makes do without variable assignments.)

**Problem 7.7.** Suppose that  $f$  is a function symbol not in  $A(x, y)$ . Show that there is a structure  $M$  such that  $M \models \forall x \exists y A(x, y)$  iff there is an  $M'$  such that  $M' \models \forall x A(x, f(x))$ .

(This problem is a special case of what's known as Skolem's Theorem;  $\forall x A(x, f(x))$  is called a *Skolem normal form* of  $\forall x \exists y A(x, y)$ .)

**Problem 7.8.** Carry out the proof of [Proposition 7.19](#) in detail.

**Problem 7.9.** Prove [Proposition 7.22](#)

**Problem 7.10.** 1. Show that  $\Gamma \models \perp$  iff  $\Gamma$  is unsatisfiable.

2. Show that  $\Gamma \cup \{A\} \models \perp$  iff  $\Gamma \models \neg A$ .

3. Suppose  $c$  does not occur in  $A$  or  $\Gamma$ . Show that  $\Gamma \models \forall x A$  iff  $\Gamma \models A[c/x]$ .

**Problem 7.11.** Complete the proof of [Proposition 7.30](#).

## CHAPTER 8

# *Theories and Their Models*

### 8.1 Introduction

The development of the axiomatic method is a significant achievement in the history of science, and is of special importance in the history of mathematics. An axiomatic development of a field involves the clarification of many questions: What is the field about? What are the most fundamental concepts? How are they related? Can all the concepts of the field be defined in terms of these fundamental concepts? What laws do, and must, these concepts obey?

The axiomatic method and logic were made for each other. Formal logic provides the tools for formulating axiomatic theories, for proving theorems from the axioms of the theory in a precisely specified way, for studying the properties of all systems satisfying the axioms in a systematic way.

**Definition 8.1.** A set of sentences  $\Gamma$  is *closed* iff, whenever  $\Gamma \vDash A$  then  $A \in \Gamma$ . The *closure* of a set of sentences  $\Gamma$  is  $\{A : \Gamma \vDash A\}$ .

We say that  $\Gamma$  is *axiomatized by* a set of sentences  $\Delta$  if  $\Gamma$  is the closure of  $\Delta$ .



We can think of an axiomatic theory as the set of sentences that is axiomatized by its set of axioms  $\mathcal{A}$ . In other words, when we have a first-order language which contains non-logical symbols for the primitives of the axiomatically developed science we wish to study, together with a set of sentences that express the fundamental laws of the science, we can think of the theory as represented by all the sentences in this language that are entailed by the axioms. This ranges from simple examples with only a single primitive and simple axioms, such as the theory of partial orders, to complex theories such as Newtonian mechanics.

The important logical facts that make this formal approach to the axiomatic method so important are the following. Suppose  $\Gamma$  is an axiom system for a theory, i.e., a set of sentences.

1. We can state precisely when an axiom system captures an intended class of structures. That is, if we are interested in a certain class of structures, we will successfully capture that class by an axiom system  $\Gamma$  iff the structures are exactly those  $M$  such that  $M \models \Gamma$ .
2. We may fail in this respect because there are  $M$  such that  $M \models \Gamma$ , but  $M$  is not one of the structures we intend. This may lead us to add axioms which are not true in  $M$ .
3. If we are successful at least in the respect that  $\Gamma$  is true in all the intended structures, then a sentence  $A$  is true in all intended structures whenever  $\Gamma \models A$ . Thus we can use logical tools (such as derivation methods) to show that sentences are true in all intended structures simply by showing that they are entailed by the axioms.
4. Sometimes we don't have intended structures in mind, but instead start from the axioms themselves: we begin with some primitives that we want to satisfy certain laws which we codify in an axiom system. One thing that we would like to verify right away is that the axioms do not contradict each other: if they do, there can be no concepts that obey

these laws, and we have tried to set up an incoherent theory. We can verify that this doesn't happen by finding a model of  $\Gamma$ . And if there are models of our theory, we can use logical methods to investigate them, and we can also use logical methods to construct models.

5. The independence of the axioms is likewise an important question. It may happen that one of the axioms is actually a consequence of the others, and so is redundant. We can prove that an axiom  $A$  in  $\Gamma$  is redundant by proving  $\Gamma \setminus \{A\} \vDash A$ . We can also prove that an axiom is not redundant by showing that  $(\Gamma \setminus \{A\}) \cup \{\neg A\}$  is satisfiable. For instance, this is how it was shown that the parallel postulate is independent of the other axioms of geometry.
6. Another important question is that of definability of concepts in a theory: The choice of the language determines what the models of a theory consist of. But not every aspect of a theory must be represented separately in its models. For instance, every ordering  $\leq$  determines a corresponding strict ordering  $<$ —given one, we can define the other. So it is not necessary that a model of a theory involving such an order must *also* contain the corresponding strict ordering. When is it the case, in general, that one relation can be defined in terms of others? When is it impossible to define a relation in terms of others (and hence must add it to the primitives of the language)?

## 8.2 Expressing Properties of Structures

It is often useful and important to express conditions on functions and relations, or more generally, that the functions and relations in a structure satisfy these conditions. For instance, we would like to have ways of distinguishing those structures for a language which “capture” what we want the predicate symbols to “mean” from those that do not. Of course we're completely

free to specify which structures we “intend,” e.g., we can specify that the interpretation of the predicate symbol  $\leq$  must be an ordering, or that we are only interested in interpretations of  $\mathcal{L}$  in which the domain consists of sets and  $\in$  is interpreted by the “is an element of” relation. But can we do this with sentences of the language? In other words, which conditions on a structure  $M$  can we express by a sentence (or perhaps a set of sentences) in the language of  $M$ ? There are some conditions that we will not be able to express. For instance, there is no sentence of  $\mathcal{L}_A$  which is only true in a structure  $M$  if  $|M| = \mathbb{N}$ . We cannot express “the domain contains only natural numbers.” But there are “structural properties” of structures that we perhaps can express. Which properties of structures can we express by sentences? Or, to put it another way, which collections of structures can we describe as those making a sentence (or set of sentences) true?

**Definition 8.2 (Model of a set).** Let  $\Gamma$  be a set of sentences in a language  $\mathcal{L}$ . We say that a structure  $M$  is a *model of  $\Gamma$*  if  $M \models A$  for all  $A \in \Gamma$ .

**Example 8.3.** The sentence  $\forall x x \leq x$  is true in  $M$  iff  $\leq^M$  is a reflexive relation. The sentence  $\forall x \forall y ((x \leq y \wedge y \leq x) \rightarrow x = y)$  is true in  $M$  iff  $\leq^M$  is anti-symmetric. The sentence  $\forall x \forall y \forall z ((x \leq y \wedge y \leq z) \rightarrow x \leq z)$  is true in  $M$  iff  $\leq^M$  is transitive. Thus, the models of

$$\left\{ \begin{array}{l} \forall x x \leq x, \\ \forall x \forall y ((x \leq y \wedge y \leq x) \rightarrow x = y), \\ \forall x \forall y \forall z ((x \leq y \wedge y \leq z) \rightarrow x \leq z) \end{array} \right\}$$

are exactly those structures in which  $\leq^M$  is reflexive, anti-symmetric, and transitive, i.e., a partial order. Hence, we can take them as axioms for the *first-order theory of partial orders*.

### 8.3 Examples of First-Order Theories

**Example 8.4.** The theory of strict linear orders in the language  $\mathcal{L}_<$  is axiomatized by the set

$$\left\{ \begin{array}{l} \forall x \neg x < x, \\ \forall x \forall y ((x < y \vee y < x) \vee x = y), \\ \forall x \forall y \forall z ((x < y \wedge y < z) \rightarrow x < z) \end{array} \right\}$$

It completely captures the intended structures: every strict linear order is a model of this axiom system, and vice versa, if  $R$  is a linear order on a set  $X$ , then the structure  $M$  with  $|M| = X$  and  $<^M = R$  is a model of this theory.

**Example 8.5.** The theory of groups in the language  $\mathcal{L}_1$  (constant symbol),  $\cdot$  (two-place function symbol) is axiomatized by

$$\begin{array}{l} \forall x (x \cdot 1) = x \\ \forall x \forall y \forall z (x \cdot (y \cdot z)) = ((x \cdot y) \cdot z) \\ \forall x \exists y (x \cdot y) = 1 \end{array}$$

**Example 8.6.** The theory of Peano arithmetic is axiomatized by the following sentences in the language of arithmetic  $\mathcal{L}_A$ .

$$\begin{array}{l} \forall x \forall y (x' = y' \rightarrow x = y) \\ \forall x 0 \neq x' \\ \forall x (x + 0) = x \\ \forall x \forall y (x + y') = (x + y)' \\ \forall x (x \times 0) = 0 \\ \forall x \forall y (x \times y') = ((x \times y) + x) \\ \forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y) \end{array}$$

plus all sentences of the form

$$(A(0) \wedge \forall x (A(x) \rightarrow A(x')))) \rightarrow \forall x A(x)$$

Since there are infinitely many sentences of the latter form, this axiom system is infinite. The latter form is called the *induction schema*. (Actually, the induction schema is a bit more complicated than we let on here.)

The last axiom is an *explicit definition* of  $<$ .

**Example 8.7.** The theory of pure sets plays an important role in the foundations (and in the philosophy) of mathematics. A set is pure if all its elements are also pure sets. The empty set counts therefore as pure, but a set that has something as an element that is not a set would not be pure. So the pure sets are those that are formed just from the empty set and no “urelements,” i.e., objects that are not themselves sets.

The following might be considered as an axiom system for a theory of pure sets:

$$\begin{aligned} & \exists x \neg \exists y y \in x \\ & \forall x \forall y (\forall z (z \in x \leftrightarrow z \in y) \rightarrow x = y) \\ & \forall x \forall y \exists z \forall u (u \in z \leftrightarrow (u = x \vee u = y)) \\ & \forall x \exists y \forall z (z \in y \leftrightarrow \exists u (z \in u \wedge u \in x)) \end{aligned}$$

plus all sentences of the form

$$\exists x \forall y (y \in x \leftrightarrow A(y))$$

The first axiom says that there is a set with no elements (i.e.,  $\emptyset$  exists); the second says that sets are extensional; the third that for any sets  $X$  and  $Y$ , the set  $\{X, Y\}$  exists; the fourth that for any set  $X$ , the set  $\cup X$  exists, where  $\cup X$  is the union of all the elements of  $X$ .

The sentences mentioned last are collectively called the *naive comprehension scheme*. It essentially says that for every  $A(x)$ , the set  $\{x : A(x)\}$  exists—so at first glance a true, useful, and perhaps even necessary axiom. It is called “naive” because, as it turns out, it makes this theory unsatisfiable: if you take  $A(y)$  to be  $\neg y \in y$ , you get the sentence

$$\exists x \forall y (y \in x \leftrightarrow \neg y \in y)$$

and this sentence is not satisfied in any structure.

**Example 8.8.** In the area of *mereology*, the relation of *parthood* is a fundamental relation. Just like theories of sets, there are theories of parthood that axiomatize various conceptions (sometimes conflicting) of this relation.

The language of mereology contains a single two-place predicate symbol  $P$ , and  $P(x, y)$  “means” that  $x$  is a part of  $y$ . When we have this interpretation in mind, a structure for this language is called a *parthood structure*. Of course, not every structure for a single two-place predicate will really deserve this name. To have a chance of capturing “parthood,”  $P^M$  must satisfy some conditions, which we can lay down as axioms for a theory of parthood. For instance, parthood is a partial order on objects: every object is a part (albeit an *improper* part) of itself; no two different objects can be parts of each other; a part of a part of an object is itself part of that object. Note that in this sense “is a part of” resembles “is a subset of,” but does not resemble “is an element of” which is neither reflexive nor transitive.

$$\forall x P(x, x)$$

$$\forall x \forall y ((P(x, y) \wedge P(y, x)) \rightarrow x = y)$$

$$\forall x \forall y \forall z ((P(x, y) \wedge P(y, z)) \rightarrow P(x, z))$$

Moreover, any two objects have a mereological sum (an object that has these two objects as parts, and is minimal in this respect).

$$\forall x \forall y \exists z \forall u (P(z, u) \leftrightarrow (P(x, u) \wedge P(y, u)))$$

These are only some of the basic principles of parthood considered by metaphysicians. Further principles, however, quickly become hard to formulate or write down without first introducing some defined relations. For instance, most metaphysicians interested in mereology also view the following as a valid principle: whenever an object  $x$  has a proper part  $y$ , it also has a part  $z$  that has no parts in common with  $y$ , and so that the fusion of  $y$  and  $z$  is  $x$ .

## 8.4 Expressing Relations in a Structure

One main use formulas can be put to is to express properties and relations in a structure  $M$  in terms of the primitives of the language  $\mathcal{L}$  of  $M$ . By this we mean the following: the domain of  $M$  is a set of objects. The constant symbols, function symbols, and predicate symbols are interpreted in  $M$  by some objects in  $|M|$ , functions on  $|M|$ , and relations on  $|M|$ . For instance, if  $A_0^2$  is in  $\mathcal{L}$ , then  $M$  assigns to it a relation  $R = A_0^2{}^M$ . Then the formula  $A_0^2(v_1, v_2)$  expresses that very relation, in the following sense: if a variable assignment  $s$  maps  $v_1$  to  $a \in |M|$  and  $v_2$  to  $b \in |M|$ , then

$$Rab \quad \text{iff} \quad M, s \models A_0^2(v_1, v_2).$$

Note that we have to involve variable assignments here: we can't just say " $Rab$  iff  $M \models A_0^2(a, b)$ " because  $a$  and  $b$  are not symbols of our language: they are elements of  $|M|$ .

Since we don't just have atomic formulas, but can combine them using the logical connectives and the quantifiers, more complex formulas can define other relations which aren't directly built into  $M$ . We're interested in how to do that, and specifically, which relations we can define in a structure.

**Definition 8.9.** Let  $A(v_1, \dots, v_n)$  be a formula of  $\mathcal{L}$  in which only  $v_1, \dots, v_n$  occur free, and let  $M$  be a structure for  $\mathcal{L}$ .  $A(v_1, \dots, v_n)$  expresses the relation  $R \subseteq |M|^n$  iff

$$Ra_1 \dots a_n \quad \text{iff} \quad M, s \models A(v_1, \dots, v_n)$$

for any variable assignment  $s$  with  $s(v_i) = a_i$  ( $i = 1, \dots, n$ ).

**Example 8.10.** In the standard model of arithmetic  $\mathbb{N}$ , the formula  $v_1 < v_2 \vee v_1 = v_2$  expresses the  $\leq$  relation on  $\mathbb{N}$ . The formula  $v_2 = v_1'$  expresses the successor relation, i.e., the relation  $R \subseteq \mathbb{N}^2$  where  $Rnm$  holds if  $m$  is the successor of  $n$ . The formula  $v_1 = v_2'$  expresses the predecessor relation. The formulas  $\exists v_3 (v_3 \neq 0 \wedge v_2 = (v_1 + v_3))$  and  $\exists v_3 (v_1 + v_3') = v_2$  both express

the  $<$  relation. This means that the predicate symbol  $<$  is actually superfluous in the language of arithmetic; it can be defined.

This idea is not just interesting in specific structures, but generally whenever we use a language to describe an intended model or models, i.e., when we consider theories. These theories often only contain a few predicate symbols as basic symbols, but in the domain they are used to describe often many other relations play an important role. If these other relations can be systematically expressed by the relations that interpret the basic predicate symbols of the language, we say we can *define* them in the language.

## 8.5 The Theory of Sets

Almost all of mathematics can be developed in the theory of sets. Developing mathematics in this theory involves a number of things. First, it requires a set of axioms for the relation  $\in$ . A number of different axiom systems have been developed, sometimes with conflicting properties of  $\in$ . The axiom system known as **ZFC**, Zermelo–Fraenkel set theory with the axiom of choice stands out: it is by far the most widely used and studied, because it turns out that its axioms suffice to prove almost all the things mathematicians expect to be able to prove. But before that can be established, it first is necessary to make clear how we can even *express* all the things mathematicians would like to express. For starters, the language contains no constant symbols or function symbols, so it seems at first glance unclear that we can talk about particular sets (such as  $\emptyset$  or  $\mathbb{N}$ ), can talk about operations on sets (such as  $X \cup Y$  and  $\wp(X)$ ), let alone other constructions which involve things other than sets, such as relations and functions.

To begin with, “is an element of” is not the only relation we are interested in: “is a subset of” seems almost as important. But we can *define* “is a subset of” in terms of “is an element of.” To do this, we have to find a formula  $A(x, y)$  in the language of set theory which is satisfied by a pair of sets  $\langle X, Y \rangle$  iff  $X \subseteq Y$ . But  $X$



is a subset of  $Y$  just in case all elements of  $X$  are also elements of  $Y$ . So we can define  $\subseteq$  by the formula

$$\forall z (z \in x \rightarrow z \in y)$$

Now, whenever we want to use the relation  $\subseteq$  in a formula, we could instead use that formula (with  $x$  and  $y$  suitably replaced, and the bound variable  $z$  renamed if necessary). For instance, extensionality of sets means that if any sets  $x$  and  $y$  are contained in each other, then  $x$  and  $y$  must be the same set. This can be expressed by  $\forall x \forall y ((x \subseteq y \wedge y \subseteq x) \rightarrow x = y)$ , or, if we replace  $\subseteq$  by the above definition, by

$$\forall x \forall y ((\forall z (z \in x \rightarrow z \in y) \wedge \forall z (z \in y \rightarrow z \in x)) \rightarrow x = y).$$

This is in fact one of the axioms of **ZFC**, the “axiom of extensionality.”

There is no constant symbol for  $\emptyset$ , but we can express “ $x$  is empty” by  $\neg \exists y y \in x$ . Then “ $\emptyset$  exists” becomes the sentence  $\exists x \neg \exists y y \in x$ . This is another axiom of **ZFC**. (Note that the axiom of extensionality implies that there is only one empty set.) Whenever we want to talk about  $\emptyset$  in the language of set theory, we would write this as “there is a set that’s empty and ...” As an example, to express the fact that  $\emptyset$  is a subset of every set, we could write

$$\exists x (\neg \exists y y \in x \wedge \forall z x \subseteq z)$$

where, of course,  $x \subseteq z$  would in turn have to be replaced by its definition.

To talk about operations on sets, such as  $X \cup Y$  and  $\wp(X)$ , we have to use a similar trick. There are no function symbols in the language of set theory, but we can express the functional relations  $X \cup Y = Z$  and  $\wp(X) = Y$  by

$$\forall u ((u \in x \vee u \in y) \leftrightarrow u \in z)$$

$$\forall u (u \subseteq x \leftrightarrow u \in y)$$

since the elements of  $X \cup Y$  are exactly the sets that are either elements of  $X$  or elements of  $Y$ , and the elements of  $\wp(X)$  are exactly the subsets of  $X$ . However, this doesn't allow us to use  $x \cup y$  or  $\wp(x)$  as if they were terms: we can only use the entire formulas that define the relations  $X \cup Y = Z$  and  $\wp(X) = Y$ . In fact, we do not know that these relations are ever satisfied, i.e., we do not know that unions and power sets always exist. For instance, the sentence  $\forall x \exists y \wp(x) = y$  is another axiom of **ZFC** (the power set axiom).

Now what about talk of ordered pairs or functions? Here we have to explain how we can think of ordered pairs and functions as special kinds of sets. One way to define the ordered pair  $\langle x, y \rangle$  is as the set  $\{\{x\}, \{x, y\}\}$ . But like before, we cannot introduce a function symbol that names this set; we can only define the relation  $\langle x, y \rangle = z$ , i.e.,  $\{\{x\}, \{x, y\}\} = z$ :

$$\forall u (u \in z \leftrightarrow (\forall v (v \in u \leftrightarrow v = x) \vee \forall v (v \in u \leftrightarrow (v = x \vee v = y))))$$

This says that the elements  $u$  of  $z$  are exactly those sets which either have  $x$  as its only element or have  $x$  and  $y$  as its only elements (in other words, those sets that are either identical to  $\{x\}$  or identical to  $\{x, y\}$ ). Once we have this, we can say further things, e.g., that  $X \times Y = Z$ :

$$\forall z (z \in Z \leftrightarrow \exists x \exists y (x \in X \wedge y \in Y \wedge \langle x, y \rangle = z))$$

A function  $f: X \rightarrow Y$  can be thought of as the relation  $f(x) = y$ , i.e., as the set of pairs  $\{\langle x, y \rangle : f(x) = y\}$ . We can then say that a set  $f$  is a function from  $X$  to  $Y$  if (a) it is a relation  $\subseteq X \times Y$ , (b) it is total, i.e., for all  $x \in X$  there is some  $y \in Y$  such that  $\langle x, y \rangle \in f$  and (c) it is functional, i.e., whenever  $\langle x, y \rangle, \langle x, y' \rangle \in f$ ,  $y = y'$  (because values of functions must be unique). So “ $f$  is a function from  $X$  to  $Y$ ” can be written as:

$$\begin{aligned} &\forall u (u \in f \rightarrow \exists x \exists y (x \in X \wedge y \in Y \wedge \langle x, y \rangle = u)) \wedge \\ &\forall x (x \in X \rightarrow (\exists y (y \in Y \wedge \text{maps}(f, x, y)) \wedge \\ &\quad (\forall y \forall y' ((\text{maps}(f, x, y) \wedge \text{maps}(f, x, y')) \rightarrow y = y')))) \end{aligned}$$

where  $\text{maps}(f, x, y)$  abbreviates  $\exists v (v \in f \wedge \langle x, y \rangle = v)$  (this formula expresses “ $f(x) = y$ ”).

It is now also not hard to express that  $f: X \rightarrow Y$  is injective, for instance:

$$f: X \rightarrow Y \wedge \forall x \forall x' ((x \in X \wedge x' \in X \wedge \exists y (\text{maps}(f, x, y) \wedge \text{maps}(f, x', y))) \rightarrow x = x')$$

A function  $f: X \rightarrow Y$  is injective iff, whenever  $f$  maps  $x, x' \in X$  to a single  $y$ ,  $x = x'$ . If we abbreviate this formula as  $\text{inj}(f, X, Y)$ , we’re already in a position to state in the language of set theory something as non-trivial as Cantor’s theorem: there is no injective function from  $\wp(X)$  to  $X$ :

$$\forall X \forall Y (\wp(X) = Y \rightarrow \neg \exists f \text{inj}(f, Y, X))$$

One might think that set theory requires another axiom that guarantees the existence of a set for every defining property. If  $A(x)$  is a formula of set theory with the variable  $x$  free, we can consider the sentence

$$\exists y \forall x (x \in y \leftrightarrow A(x)).$$

This sentence states that there is a set  $y$  whose elements are all and only those  $x$  that satisfy  $A(x)$ . This schema is called the “comprehension principle.” It looks very useful; unfortunately it is inconsistent. Take  $A(x) \equiv \neg x \in x$ , then the comprehension principle states

$$\exists y \forall x (x \in y \leftrightarrow x \notin x),$$

i.e., it states the existence of a set of all sets that are not elements of themselves. No such set can exist—this is Russell’s Paradox. **ZFC**, in fact, contains a restricted—and consistent—version of this principle, the separation principle:

$$\forall z \exists y \forall x (x \in y \leftrightarrow (x \in z \wedge A(x))).$$

## 8.6 Expressing the Size of Structures

There are some properties of structures we can express even without using the non-logical symbols of a language. For instance, there are sentences which are true in a structure iff the domain of the structure has at least, at most, or exactly a certain number  $n$  of elements.

**Proposition 8.11.** *The sentence*

$$\begin{aligned}
 A_{\geq n} \equiv & \exists x_1 \exists x_2 \dots \exists x_n \\
 & (x_1 \neq x_2 \wedge x_1 \neq x_3 \wedge x_1 \neq x_4 \wedge \dots \wedge x_1 \neq x_n \wedge \\
 & \quad x_2 \neq x_3 \wedge x_2 \neq x_4 \wedge \dots \wedge x_2 \neq x_n \wedge \\
 & \quad \vdots \\
 & \quad x_{n-1} \neq x_n)
 \end{aligned}$$

*is true in a structure  $M$  iff  $|M|$  contains at least  $n$  elements. Consequently,  $M \models \neg A_{\geq n+1}$  iff  $|M|$  contains at most  $n$  elements.*

**Proposition 8.12.** *The sentence*

$$\begin{aligned}
 A_{=n} \equiv & \exists x_1 \exists x_2 \dots \exists x_n \\
 & (x_1 \neq x_2 \wedge x_1 \neq x_3 \wedge x_1 \neq x_4 \wedge \dots \wedge x_1 \neq x_n \wedge \\
 & \quad x_2 \neq x_3 \wedge x_2 \neq x_4 \wedge \dots \wedge x_2 \neq x_n \wedge \\
 & \quad \vdots \\
 & \quad x_{n-1} \neq x_n \wedge \\
 & \quad \forall y (y = x_1 \vee \dots \vee y = x_n))
 \end{aligned}$$

*is true in a structure  $M$  iff  $|M|$  contains exactly  $n$  elements.*

**Proposition 8.13.** *A structure is infinite iff it is a model of*

$$\{A_{\geq 1}, A_{\geq 2}, A_{\geq 3}, \dots\}.$$

There is no single purely logical sentence which is true in  $M$  iff  $|M|$  is infinite. However, one can give sentences with non-logical predicate symbols which only have infinite models (although not every infinite structure is a model of them). The property of being a finite structure, and the property of being a uncountable structure cannot even be expressed with an infinite set of sentences. These facts follow from the compactness and Löwenheim–Skolem theorems.

## Summary

Sets of sentences in a sense describe the structures in which they are jointly true; these structures are their **models**. Conversely, if we start with a structure or set of structures, we might be interested in the set of sentences they are models of, this is the **theory** of the structure or set of structures. Any such set of sentences has the property that every sentence entailed by them is already in the set; they are **closed**. More generally, we call a set  $\Gamma$  a theory if it is closed under entailment, and say  $\Gamma$  is **axiomatized** by  $\Delta$  if  $\Gamma$  consists of all sentences entailed by  $\Delta$ .

Mathematics yields many examples of theories, e.g., the theories of linear orders, of groups, or theories of arithmetic, e.g., the theory axiomatized by Peano's axioms. But there are many examples of important theories in other disciplines as well, e.g., relational databases may be thought of as theories, and metaphysics concerns itself with theories of parthood which can be axiomatized.

One significant question when setting up a theory for study is whether its language is expressive enough to allow us to formulate everything we want the theory to talk about, and another is whether it is strong enough to prove what we want it to prove. To **express** a relation we need a formula with the requisite number

of free variables. In **set theory**, we only have  $\in$  as a relation symbol, but it allows us to express  $x \subseteq y$  using  $\forall u (u \in x \rightarrow u \in y)$ . **Zermelo-Fraenkel set theory ZFC**, in fact, is strong enough to both express (almost) every mathematical claim and to (almost) prove every mathematical theorem using a handful of axioms and a chain of increasingly complicated definitions such as that of  $\subseteq$ .

## Problems

**Problem 8.1.** Find formulas in  $\mathcal{L}_A$  which define the following relations:

1.  $n$  is between  $i$  and  $j$ ;
2.  $n$  evenly divides  $m$  (i.e.,  $m$  is a multiple of  $n$ );
3.  $n$  is a prime number (i.e., no number other than 1 and  $n$  evenly divides  $n$ ).

**Problem 8.2.** Suppose the formula  $A(v_1, v_2)$  expresses the relation  $R \subseteq |M|^2$  in a structure  $M$ . Find formulas that express the following relations:

1. the inverse  $R^{-1}$  of  $R$ ;
2. the relative product  $R \mid R$ ;

Can you find a way to express  $R^+$ , the transitive closure of  $R$ ?

**Problem 8.3.** Let  $\mathcal{L}$  be the language containing a 2-place predicate symbol  $<$  only (no other constant symbols, function symbols or predicate symbols— except of course  $=$ ). Let  $N$  be the structure such that  $|N| = \mathbb{N}$ , and  $<^N = \{\langle n, m \rangle : n < m\}$ . Prove the following:

1.  $\{0\}$  is definable in  $N$ ;
2.  $\{1\}$  is definable in  $N$ ;

3.  $\{2\}$  is definable in  $N$ ;
4. for each  $n \in \mathbb{N}$ , the set  $\{n\}$  is definable in  $N$ ;
5. every finite subset of  $|N|$  is definable in  $N$ ;
6. every co-finite subset of  $|N|$  is definable in  $N$  (where  $X \subseteq \mathbb{N}$  is co-finite iff  $\mathbb{N} \setminus X$  is finite).

**Problem 8.4.** Show that the comprehension principle is inconsistent by giving a derivation that shows

$$\exists y \forall x (x \in y \leftrightarrow x \notin x) \vdash \perp.$$

It may help to first show  $(A \rightarrow \neg A) \wedge (\neg A \rightarrow A) \vdash \perp$ .

## CHAPTER 9

# *Derivation Systems*

### 9.1 Introduction

Logics commonly have both a semantics and a derivation system. The semantics concerns concepts such as truth, satisfiability, validity, and entailment. The purpose of derivation systems is to provide a purely syntactic method of establishing entailment and validity. They are purely syntactic in the sense that a derivation in such a system is a finite syntactic object, usually a sequence (or other finite arrangement) of sentences or formulas. Good derivation systems have the property that any given sequence or arrangement of sentences or formulas can be verified mechanically to be “correct.”

The simplest (and historically first) derivation systems for first-order logic were *axiomatic*. A sequence of formulas counts as a derivation in such a system if each individual formula in it is either among a fixed set of “axioms” or follows from formulas coming before it in the sequence by one of a fixed number of “inference rules”—and it can be mechanically verified if a formula is an axiom and whether it follows correctly from other formulas by one of the inference rules. Axiomatic derivation systems are easy to describe—and also easy to handle meta-theoretically—



but derivations in them are hard to read and understand, and are also hard to produce.

Other derivation systems have been developed with the aim of making it easier to construct derivations or easier to understand derivations once they are complete. Examples are natural deduction, truth trees, also known as tableaux proofs, and the sequent calculus. Some derivation systems are designed especially with mechanization in mind, e.g., the resolution method is easy to implement in software (but its derivations are essentially impossible to understand). Most of these other derivation systems represent derivations as trees of formulas rather than sequences. This makes it easier to see which parts of a derivation depend on which other parts.

So for a given logic, such as first-order logic, the different derivation systems will give different explications of what it is for a sentence to be a *theorem* and what it means for a sentence to be derivable from some others. However that is done (via axiomatic derivations, natural deductions, sequent derivations, truth trees, resolution refutations), we want these relations to match the semantic notions of validity and entailment. Let's write  $\vdash A$  for “ $A$  is a theorem” and “ $\Gamma \vdash A$ ” for “ $A$  is derivable from  $\Gamma$ .” However  $\vdash$  is defined, we want it to match up with  $\vDash$ , that is:

1.  $\vdash A$  if and only if  $\vDash A$
2.  $\Gamma \vdash A$  if and only if  $\Gamma \vDash A$

The “only if” direction of the above is called *soundness*. A derivation system is sound if derivability guarantees entailment (or validity). Every decent derivation system has to be sound; unsound derivation systems are not useful at all. After all, the entire purpose of a derivation is to provide a syntactic guarantee of validity or entailment. We'll prove soundness for the derivation systems we present.

The converse “if” direction is also important: it is called *completeness*. A complete derivation system is strong enough to show

that  $A$  is a theorem whenever  $A$  is valid, and that  $\Gamma \vdash A$  whenever  $\Gamma \models A$ . Completeness is harder to establish, and some logics have no complete derivation systems. First-order logic does. Kurt Gödel was the first one to prove completeness for a derivation system of first-order logic in his 1929 dissertation.

Another concept that is connected to derivation systems is that of *consistency*. A set of sentences is called inconsistent if anything whatsoever can be derived from it, and consistent otherwise. Inconsistency is the syntactic counterpart to unsatisfiability: like unsatisfiable sets, inconsistent sets of sentences do not make good theories, they are defective in a fundamental way. Consistent sets of sentences may not be true or useful, but at least they pass that minimal threshold of logical usefulness. For different derivation systems the specific definition of consistency of sets of sentences might differ, but like  $\vdash$ , we want consistency to coincide with its semantic counterpart, satisfiability. We want it to always be the case that  $\Gamma$  is consistent if and only if it is satisfiable. Here, the “if” direction amounts to completeness (consistency guarantees satisfiability), and the “only if” direction amounts to soundness (satisfiability guarantees consistency). In fact, for classical first-order logic, the two versions of soundness and completeness are equivalent.

## 9.2 The Sequent Calculus

While many derivation systems operate with arrangements of sentences, the sequent calculus operates with *sequents*. A sequent is an expression of the form

$$A_1, \dots, A_m \Rightarrow B_1, \dots, B_n,$$

that is a pair of sequences of sentences, separated by the sequent symbol  $\Rightarrow$ . Either sequence may be empty. A derivation in the sequent calculus is a tree of sequents, where the topmost sequents are of a special form (they are called “initial sequents” or “axioms”) and every other sequent follows from the sequents imme-

diately above it by one of the rules of inference. The rules of inference either manipulate the sentences in the sequents (adding, removing, or rearranging them on either the left or the right), or they introduce a complex formula in the conclusion of the rule. For instance, the  $\wedge$ L rule allows the inference from  $A, \Gamma \Rightarrow \Delta$  to  $A \wedge B, \Gamma \Rightarrow \Delta$ , and the  $\rightarrow$ R allows the inference from  $A, \Gamma \Rightarrow \Delta, B$  to  $\Gamma \Rightarrow \Delta, A \rightarrow B$ , for any  $\Gamma, \Delta, A$ , and  $B$ . (In particular,  $\Gamma$  and  $\Delta$  may be empty.)

The  $\vdash$  relation based on the sequent calculus is defined as follows:  $\Gamma \vdash A$  iff there is some sequence  $\Gamma_0$  such that every  $A$  in  $\Gamma_0$  is in  $\Gamma$  and there is a derivation with the sequent  $\Gamma_0 \Rightarrow A$  at its root.  $A$  is a theorem in the sequent calculus if the sequent  $\Rightarrow A$  has a derivation. For instance, here is a derivation that shows that  $\vdash (A \wedge B) \rightarrow A$ :

$$\frac{\frac{A \Rightarrow A}{A \wedge B \Rightarrow A} \wedge\text{L}}{\Rightarrow (A \wedge B) \rightarrow A} \rightarrow\text{R}$$

A set  $\Gamma$  is inconsistent in the sequent calculus if there is a derivation of  $\Gamma_0 \Rightarrow$  (where every  $A \in \Gamma_0$  is in  $\Gamma$  and the right side of the sequent is empty). Using the rule WR, any sentence can be derived from an inconsistent set.

The sequent calculus was invented in the 1930s by Gerhard Gentzen. Because of its systematic and symmetric design, it is a very useful formalism for developing a theory of derivations. It is relatively easy to find derivations in the sequent calculus, but these derivations are often hard to read and their connection to proofs are sometimes not easy to see. It has proved to be a very elegant approach to derivation systems, however, and many logics have sequent calculus systems.

### 9.3 Natural Deduction

Natural deduction is a derivation system intended to mirror actual reasoning (especially the kind of regimented reasoning em-

ployed by mathematicians). Actual reasoning proceeds by a number of “natural” patterns. For instance, proof by cases allows us to establish a conclusion on the basis of a disjunctive premise, by establishing that the conclusion follows from either of the disjuncts. Indirect proof allows us to establish a conclusion by showing that its negation leads to a contradiction. Conditional proof establishes a conditional claim “if ... then ...” by showing that the consequent follows from the antecedent. Natural deduction is a formalization of some of these natural inferences. Each of the logical connectives and quantifiers comes with two rules, an introduction and an elimination rule, and they each correspond to one such natural inference pattern. For instance,  $\rightarrow$ Intro corresponds to conditional proof, and  $\vee$ Elim to proof by cases. A particularly simple rule is  $\wedge$ Elim which allows the inference from  $A \wedge B$  to  $A$  (or  $B$ ).

One feature that distinguishes natural deduction from other derivation systems is its use of assumptions. A derivation in natural deduction is a tree of formulas. A single formula stands at the root of the tree of formulas, and the “leaves” of the tree are formulas from which the conclusion is derived. In natural deduction, some leaf formulas play a role inside the derivation but are “used up” by the time the derivation reaches the conclusion. This corresponds to the practice, in actual reasoning, of introducing hypotheses which only remain in effect for a short while. For instance, in a proof by cases, we assume the truth of each of the disjuncts; in conditional proof, we assume the truth of the antecedent; in indirect proof, we assume the truth of the negation of the conclusion. This way of introducing hypothetical assumptions and then doing away with them in the service of establishing an intermediate step is a hallmark of natural deduction. The formulas at the leaves of a natural deduction derivation are called assumptions, and some of the rules of inference may “discharge” them. For instance, if we have a derivation of  $B$  from some assumptions which include  $A$ , then the  $\rightarrow$ Intro rule allows us to infer  $A \rightarrow B$  and discharge any assumption of the form  $A$ . (To keep track of which assumptions are discharged at which in-

ferences, we label the inference and the assumptions it discharges with a number.) The assumptions that remain undischarged at the end of the derivation are together sufficient for the truth of the conclusion, and so a derivation establishes that its undischarged assumptions entail its conclusion.

The relation  $\Gamma \vdash A$  based on natural deduction holds iff there is a derivation in which  $A$  is the last sentence in the tree, and every leaf which is undischarged is in  $\Gamma$ .  $A$  is a theorem in natural deduction iff there is a derivation in which  $A$  is the last sentence and all assumptions are discharged. For instance, here is a derivation that shows that  $\vdash (A \wedge B) \rightarrow A$ :

$$1 \frac{\frac{[A \wedge B]^1}{A} \wedge\text{Elim}}{(A \wedge B) \rightarrow A} \rightarrow\text{Intro}$$

The label 1 indicates that the assumption  $A \wedge B$  is discharged at the  $\rightarrow\text{Intro}$  inference.

A set  $\Gamma$  is inconsistent iff  $\Gamma \vdash \perp$  in natural deduction. The rule  $\perp_I$  makes it so that from an inconsistent set, any sentence can be derived.

Natural deduction systems were developed by Gerhard Gentzen and Stanisław Jaśkowski in the 1930s, and later developed by Dag Prawitz and Frederic Fitch. Because its inferences mirror natural methods of proof, it is favored by philosophers. The versions developed by Fitch are often used in introductory logic textbooks. In the philosophy of logic, the rules of natural deduction have sometimes been taken to give the meanings of the logical operators (“proof-theoretic semantics”).

## 9.4 Tableaux

While many derivation systems operate with arrangements of sentences, tableaux operate with signed formulas. A signed formula is a pair consisting of a truth value sign ( $\mathbb{T}$  or  $\mathbb{F}$ ) and a sentence

$$\mathbb{T} A \text{ or } \mathbb{F} A.$$

A tableau consists of signed formulas arranged in a downward-branching tree. It begins with a number of *assumptions* and continues with signed formulas which result from one of the signed formulas above it by applying one of the rules of inference. Each rule allows us to add one or more signed formulas to the end of a branch, or two signed formulas side by side—in this case a branch splits into two, with the two added signed formulas forming the ends of the two branches.

A rule applied to a complex signed formula results in the addition of signed formulas which are immediate sub-formulas. They come in pairs, one rule for each of the two signs. For instance, the  $\wedge\mathbb{T}$  rule applies to  $\mathbb{T} A \wedge B$ , and allows the addition of both the two signed formulas  $\mathbb{T} A$  and  $\mathbb{T} B$  to the end of any branch containing  $\mathbb{T} A \wedge B$ , and the rule  $A \wedge B\mathbb{F}$  allows a branch to be split by adding  $\mathbb{F} A$  and  $\mathbb{F} B$  side-by-side. A tableau is closed if every one of its branches contains a matching pair of signed formulas  $\mathbb{T} A$  and  $\mathbb{F} A$ .

The  $\vdash$  relation based on tableaux is defined as follows:  $\Gamma \vdash A$  iff there is some finite set  $\Gamma_0 = \{B_1, \dots, B_n\} \subseteq \Gamma$  such that there is a closed tableau for the assumptions

$$\{\mathbb{F} A, \mathbb{T} B_1, \dots, \mathbb{T} B_n\}$$

For instance, here is a closed tableau that shows that  $\vdash (A \wedge B) \rightarrow A$ :

1.	$\mathbb{F} (A \wedge B) \rightarrow A$	Assumption
2.	$\mathbb{T} A \wedge B$	$\rightarrow\mathbb{F}1$
3.	$\mathbb{F} A$	$\rightarrow\mathbb{F}1$
4.	$\mathbb{T} A$	$\rightarrow\mathbb{T}2$
5.	$\mathbb{T} B$	$\rightarrow\mathbb{T}2$
	$\otimes$	

A set  $\Gamma$  is inconsistent in the tableau calculus if there is a closed tableau for assumptions

$$\{\mathbb{T} B_1, \dots, \mathbb{T} B_n\}$$

for some  $B_i \in \Gamma$ .

Tableaux were invented in the 1950s independently by Evert Beth and Jaakko Hintikka, and simplified and popularized by Raymond Smullyan. They are very easy to use, since constructing a tableau is a very systematic procedure. Because of the systematic nature of tableaux, they also lend themselves to implementation by computer. However, a tableau is often hard to read and their connection to proofs are sometimes not easy to see. The approach is also quite general, and many different logics have tableau systems. Tableaux also help us to find structures that satisfy given (sets of) sentences: if the set is satisfiable, it won't have a closed tableau, i.e., any tableau will have an open branch. The satisfying structure can be "read off" an open branch, provided every rule it is possible to apply has been applied on that branch. There is also a very close connection to the sequent calculus: essentially, a closed tableau is a condensed derivation in the sequent calculus, written upside-down.

## 9.5 Axiomatic Derivations

Axiomatic derivations are the oldest and simplest logical derivation systems. Its derivations are simply sequences of sentences. A sequence of sentences counts as a correct derivation if every sentence  $A$  in it satisfies one of the following conditions:

1.  $A$  is an axiom, or
2.  $A$  is an element of a given set  $\Gamma$  of sentences, or
3.  $A$  is justified by a rule of inference.

To be an axiom,  $A$  has to have the form of one of a number of fixed sentence schemas. There are many sets of axiom schemas that provide a satisfactory (sound and complete) derivation system for first-order logic. Some are organized according to the connectives they govern, e.g., the schemas

$$A \rightarrow (B \rightarrow A) \quad B \rightarrow (B \vee C) \quad (B \wedge C) \rightarrow B$$

are common axioms that govern  $\rightarrow$ ,  $\vee$  and  $\wedge$ . Some axiom systems aim at a minimal number of axioms. Depending on the connectives that are taken as primitives, it is even possible to find axiom systems that consist of a single axiom.

A rule of inference is a conditional statement that gives a sufficient condition for a sentence in a derivation to be justified. Modus ponens is one very common such rule: it says that if  $A$  and  $A \rightarrow B$  are already justified, then  $B$  is justified. This means that a line in a derivation containing the sentence  $B$  is justified, provided that both  $A$  and  $A \rightarrow B$  (for some sentence  $A$ ) appear in the derivation before  $B$ .

The  $\vdash$  relation based on axiomatic derivations is defined as follows:  $\Gamma \vdash A$  iff there is a derivation with the sentence  $A$  as its last formula (and  $\Gamma$  is taken as the set of sentences in that derivation which are justified by (2) above).  $A$  is a theorem if  $A$  has a derivation where  $\Gamma$  is empty, i.e., every sentence in the derivation is justified either by (1) or (3). For instance, here is a derivation that shows that  $\vdash A \rightarrow (B \rightarrow (B \vee A))$ :

1.  $B \rightarrow (B \vee A)$
2.  $(B \rightarrow (B \vee A)) \rightarrow (A \rightarrow (B \rightarrow (B \vee A)))$
3.  $A \rightarrow (B \rightarrow (B \vee A))$

The sentence on line 1 is of the form of the axiom  $A \rightarrow (A \vee B)$  (with the roles of  $A$  and  $B$  reversed). The sentence on line 2 is of the form of the axiom  $A \rightarrow (B \rightarrow A)$ . Thus, both lines are justified. Line 3 is justified by modus ponens: if we abbreviate it as  $D$ , then line 2 has the form  $C \rightarrow D$ , where  $C$  is  $B \rightarrow (B \vee A)$ , i.e., line 1.

A set  $\Gamma$  is inconsistent if  $\Gamma \vdash \perp$ . A complete axiom system will also prove that  $\perp \rightarrow A$  for any  $A$ , and so if  $\Gamma$  is inconsistent, then  $\Gamma \vdash A$  for any  $A$ .

Systems of axiomatic derivations for logic were first given by Gottlob Frege in his 1879 *Begriffsschrift*, which for this reason is often considered the first work of modern logic. They were perfected in Alfred North Whitehead and Bertrand Russell's *Principia Mathematica* and by David Hilbert and his students in the



1920s. They are thus often called “Frege systems” or “Hilbert systems.” They are very versatile in that it is often easy to find an axiomatic system for a logic. Because derivations have a very simple structure and only one or two inference rules, it is also relatively easy to prove things *about* them. However, they are very hard to use in practice, i.e., it is difficult to find and write proofs.

## CHAPTER 10

# *The Sequent Calculus*

### 10.1 Rules and Derivations

For the following, let  $\Gamma, \Delta, \Pi, \Lambda$  represent finite sequences of sentences.

**Definition 10.1 (Sequent).** A *sequent* is an expression of the form

$$\Gamma \Rightarrow \Delta$$

where  $\Gamma$  and  $\Delta$  are finite (possibly empty) sequences of sentences of the language  $\mathcal{L}$ .  $\Gamma$  is called the *antecedent*, while  $\Delta$  is the *succedent*.

The intuitive idea behind a sequent is: if all of the sentences in the antecedent hold, then at least one of the sentences in the succedent holds. That is, if  $\Gamma = \langle A_1, \dots, A_m \rangle$  and  $\Delta = \langle B_1, \dots, B_n \rangle$ , then  $\Gamma \Rightarrow \Delta$  holds iff

$$(A_1 \wedge \dots \wedge A_m) \rightarrow (B_1 \vee \dots \vee B_n)$$

holds. There are two special cases: where  $\Gamma$  is empty and when  $\Delta$  is empty. When  $\Gamma$  is empty, i.e.,  $m = 0$ ,  $\Rightarrow \Delta$  holds iff  $B_1 \vee \dots \vee$

$B_n$  holds. When  $\Delta$  is empty, i.e.,  $n = 0$ ,  $\Gamma \Rightarrow$  holds iff  $\neg(A_1 \wedge \cdots \wedge A_m)$  does. We say a sequent is valid iff the corresponding sentence is valid.

If  $\Gamma$  is a sequence of sentences, we write  $\Gamma, A$  for the result of appending  $A$  to the right end of  $\Gamma$  (and  $A, \Gamma$  for the result of appending  $A$  to the left end of  $\Gamma$ ). If  $\Delta$  is a sequence of sentences also, then  $\Gamma, \Delta$  is the concatenation of the two sequences.

**Definition 10.2 (Initial Sequent).** An *initial sequent* is a sequent of one of the following forms:

1.  $A \Rightarrow A$
2.  $\perp \Rightarrow$

for any sentence  $A$  in the language.

Derivations in the sequent calculus are certain trees of sequents, where the topmost sequents are initial sequents, and if a sequent stands below one or two other sequents, it must follow correctly by a rule of inference. The rules for **LK** are divided into two main types: *logical* rules and *structural* rules. The logical rules are named for the main operator of the sentence containing  $A$  and/or  $B$  in the lower sequent. Each one comes in two versions, one for inferring a sequent with the sentence containing the logical operator on the left, and one with the sentence on the right.

## 10.2 Propositional Rules

### Rules for $\neg$

$$\frac{\Gamma \Rightarrow \Delta, A}{\neg A, \Gamma \Rightarrow \Delta} \neg\text{L}$$

$$\frac{A, \Gamma \Rightarrow \Delta}{\Gamma \Rightarrow \Delta, \neg A} \neg\text{R}$$

### Rules for $\wedge$

$$\frac{A, \Gamma \Rightarrow \Delta}{A \wedge B, \Gamma \Rightarrow \Delta} \wedge L$$

$$\frac{B, \Gamma \Rightarrow \Delta}{A \wedge B, \Gamma \Rightarrow \Delta} \wedge L$$

$$\frac{\Gamma \Rightarrow \Delta, A \quad \Gamma \Rightarrow \Delta, B}{\Gamma \Rightarrow \Delta, A \wedge B} \wedge R$$

**Rules for  $\vee$** 

$$\frac{A, \Gamma \Rightarrow \Delta \quad B, \Gamma \Rightarrow \Delta}{A \vee B, \Gamma \Rightarrow \Delta} \vee L$$

$$\frac{\Gamma \Rightarrow \Delta, A}{\Gamma \Rightarrow \Delta, A \vee B} \vee R$$

$$\frac{\Gamma \Rightarrow \Delta, B}{\Gamma \Rightarrow \Delta, A \vee B} \vee R$$

**Rules for  $\rightarrow$** 

$$\frac{\Gamma \Rightarrow \Delta, A \quad B, \Pi \Rightarrow \Lambda}{A \rightarrow B, \Gamma, \Pi \Rightarrow \Delta, \Lambda} \rightarrow L$$

$$\frac{A, \Gamma \Rightarrow \Delta, B}{\Gamma \Rightarrow \Delta, A \rightarrow B} \rightarrow R$$

**10.3 Quantifier Rules****Rules for  $\forall$** 

$$\frac{A(t), \Gamma \Rightarrow \Delta}{\forall x A(x), \Gamma \Rightarrow \Delta} \forall L$$

$$\frac{\Gamma \Rightarrow \Delta, A(a)}{\Gamma \Rightarrow \Delta, \forall x A(x)} \forall R$$

In  $\forall L$ ,  $t$  is a closed term (i.e., one without variables). In  $\forall R$ ,  $a$  is a constant symbol which must not occur anywhere in the lower sequent of the  $\forall R$  rule. We call  $a$  the *eigenvariable* of the  $\forall R$  inference.<sup>1</sup>

**Rules for  $\exists$** 

<sup>1</sup>We use the term “eigenvariable” even though  $a$  in the above rule is a constant symbol. This has historical reasons.

$$\frac{A(a), \Gamma \Rightarrow \Delta}{\exists x A(x), \Gamma \Rightarrow \Delta} \exists L \qquad \frac{\Gamma \Rightarrow \Delta, A(t)}{\Gamma \Rightarrow \Delta, \exists x A(x)} \exists R$$

Again,  $t$  is a closed term, and  $a$  is a constant symbol which does not occur in the lower sequent of the  $\exists L$  rule. We call  $a$  the *eigenvariable* of the  $\exists L$  inference.

The condition that an eigenvariable not occur in the lower sequent of the  $\forall R$  or  $\exists L$  inference is called the *eigenvariable condition*.

Recall the convention that when  $A$  is a formula with the variable  $x$  free, we indicate this by writing  $A(x)$ . In the same context,  $A(t)$  then is short for  $A[t/x]$ . So we could also write the  $\exists R$  rule as:

$$\frac{\Gamma \Rightarrow \Delta, A[t/x]}{\Gamma \Rightarrow \Delta, \exists x A} \exists R$$

Note that  $t$  may already occur in  $A$ , e.g.,  $A$  might be  $P(t, x)$ . Thus, inferring  $\Gamma \Rightarrow \Delta, \exists x P(t, x)$  from  $\Gamma \Rightarrow \Delta, P(t, t)$  is a correct application of  $\exists R$ —you may “replace” one or more, and not necessarily all, occurrences of  $t$  in the premise by the bound variable  $x$ . However, the eigenvariable conditions in  $\forall R$  and  $\exists L$  require that the constant symbol  $a$  does not occur in  $A$ . So, you cannot correctly infer  $\Gamma \Rightarrow \Delta, \forall x P(a, x)$  from  $\Gamma \Rightarrow \Delta, P(a, a)$  using  $\forall R$ .

In  $\exists R$  and  $\forall L$  there are no restrictions on the term  $t$ . On the other hand, in the  $\exists L$  and  $\forall R$  rules, the eigenvariable condition requires that the constant symbol  $a$  does not occur anywhere outside of  $A(a)$  in the upper sequent. It is necessary to ensure that the system is sound, i.e., only derives sequents that are valid. Without this condition, the following would be allowed:

$$\frac{A(a) \Rightarrow A(a)}{\exists x A(x) \Rightarrow A(a)} * \exists L \qquad \frac{A(a) \Rightarrow A(a)}{A(a) \Rightarrow \forall x A(x)} * \forall R$$

$$\frac{\exists x A(x) \Rightarrow A(a)}{\exists x A(x) \Rightarrow \forall x A(x)} \forall R \qquad \frac{A(a) \Rightarrow \forall x A(x)}{\exists x A(x) \Rightarrow \forall x A(x)} \exists L$$

However,  $\exists x A(x) \Rightarrow \forall x A(x)$  is not valid.

## 10.4 Structural Rules

We also need a few rules that allow us to rearrange sentences in the left and right side of a sequent. Since the logical rules require that the sentences in the premise which the rule acts upon stand either to the far left or to the far right, we need an “exchange” rule that allows us to move sentences to the right position. It’s also important sometimes to be able to combine two identical sentences into one, and to add a sentence on either side.

### Weakening

$$\frac{\Gamma \Rightarrow \Delta}{A, \Gamma \Rightarrow \Delta} \text{WL}$$

$$\frac{\Gamma \Rightarrow \Delta}{\Gamma \Rightarrow \Delta, A} \text{WR}$$

### Contraction

$$\frac{A, A, \Gamma \Rightarrow \Delta}{A, \Gamma \Rightarrow \Delta} \text{CL}$$

$$\frac{\Gamma \Rightarrow \Delta, A, A}{\Gamma \Rightarrow \Delta, A} \text{CR}$$

### Exchange

$$\frac{\Gamma, A, B, \Pi \Rightarrow \Delta}{\Gamma, B, A, \Pi \Rightarrow \Delta} \text{XL}$$

$$\frac{\Gamma \Rightarrow \Delta, A, B, A}{\Gamma \Rightarrow \Delta, B, A, A} \text{XR}$$

A series of weakening, contraction, and exchange inferences will often be indicated by double inference lines.

The following rule, called “cut,” is not strictly speaking necessary, but makes it a lot easier to reuse and combine derivations.

$$\frac{\Gamma \Rightarrow \Delta, A \quad A, \Pi \Rightarrow \Lambda}{\Gamma, \Pi \Rightarrow \Delta, \Lambda} \text{Cut}$$

## 10.5 Derivations

We've said what an initial sequent looks like, and we've given the rules of inference. Derivations in the sequent calculus are inductively generated from these: each derivation either is an initial sequent on its own, or consists of one or two derivations followed by an inference.

**Definition 10.3 (LK derivation).** An **LK-derivation** of a sequent  $S$  is a finite tree of sequents satisfying the following conditions:

1. The topmost sequents of the tree are initial sequents.
2. The bottommost sequent of the tree is  $S$ .
3. Every sequent in the tree except  $S$  is a premise of a correct application of an inference rule whose conclusion stands directly below that sequent in the tree.

We then say that  $S$  is the *end-sequent* of the derivation and that  $S$  is *derivable in LK* (or **LK-derivable**).

**Example 10.4.** Every initial sequent, e.g.,  $C \Rightarrow C$  is a derivation. We can obtain a new derivation from this by applying, say, the WL rule,

$$\frac{\Gamma \Rightarrow \Delta}{A, \Gamma \Rightarrow \Delta} \text{WL}$$

The rule, however, is meant to be general: we can replace the  $A$  in the rule with any sentence, e.g., also with  $D$ . If the premise matches our initial sequent  $C \Rightarrow C$ , that means that both  $\Gamma$  and  $\Delta$  are just  $C$ , and the conclusion would then be  $D, C \Rightarrow C$ . So, the following is a derivation:

$$\frac{C \Rightarrow C}{D, C \Rightarrow C} \text{WL}$$

We can now apply another rule, say XL, which allows us to switch two sentences on the left. So, the following is also a correct derivation:

$$\frac{\frac{C \Rightarrow C}{D, C \Rightarrow C} \text{WL}}{C, D \Rightarrow C} \text{XL}$$

In this application of the rule, which was given as

$$\frac{\Gamma, A, B, \Pi \Rightarrow \Delta}{\Gamma, B, A, \Pi \Rightarrow \Delta} \text{XL}$$

both  $\Gamma$  and  $\Pi$  were empty,  $\Delta$  is  $C$ , and the roles of  $A$  and  $B$  are played by  $D$  and  $C$ , respectively. In much the same way, we also see that

$$\frac{D \Rightarrow D}{C, D \Rightarrow D} \text{WL}$$

is a derivation. Now we can take these two derivations, and combine them using  $\wedge R$ . That rule was

$$\frac{\Gamma \Rightarrow \Delta, A \quad \Gamma \Rightarrow \Delta, B}{\Gamma \Rightarrow \Delta, A \wedge B} \wedge R$$

In our case, the premises must match the last sequents of the derivations ending in the premises. That means that  $\Gamma$  is  $C, D$ ,  $\Delta$  is empty,  $A$  is  $C$  and  $B$  is  $D$ . So the conclusion, if the inference should be correct, is  $C, D \Rightarrow C \wedge D$ .

$$\frac{\frac{\frac{C \Rightarrow C}{D, C \Rightarrow C} \text{WL}}{C, D \Rightarrow C} \text{XL} \quad \frac{D \Rightarrow D}{C, D \Rightarrow D} \text{WL}}{C, D \Rightarrow C \wedge D} \wedge R$$

Of course, we can also reverse the premises, then  $A$  would be  $D$  and  $B$  would be  $C$ .

$$\frac{\frac{D \Rightarrow D}{C, D \Rightarrow D} \text{WL} \quad \frac{\frac{C \Rightarrow C}{D, C \Rightarrow C} \text{WL}}{C, D \Rightarrow C} \text{XL}}{C, D \Rightarrow D \wedge C} \wedge R$$



## 10.6 Examples of Derivations

**Example 10.5.** Give an **LK**-derivation for the sequent  $A \wedge B \Rightarrow A$ .

We begin by writing the desired end-sequent at the bottom of the derivation.

$$\overline{A \wedge B \Rightarrow A}$$

Next, we need to figure out what kind of inference could have a lower sequent of this form. This could be a structural rule, but it is a good idea to start by looking for a logical rule. The only logical connective occurring in the lower sequent is  $\wedge$ , so we're looking for an  $\wedge$  rule, and since the  $\wedge$  symbol occurs in the antecedent, we're looking at the  $\wedge$ L rule.

$$\overline{A \wedge B \Rightarrow A} \wedge\text{L}$$

There are two options for what could have been the upper sequent of the  $\wedge$ L inference: we could have an upper sequent of  $A \Rightarrow A$ , or of  $B \Rightarrow A$ . Clearly,  $A \Rightarrow A$  is an initial sequent (which is a good thing), while  $B \Rightarrow A$  is not derivable in general. We fill in the upper sequent:

$$\frac{A \Rightarrow A}{A \wedge B \Rightarrow A} \wedge\text{L}$$

We now have a correct **LK**-derivation of the sequent  $A \wedge B \Rightarrow A$ .

**Example 10.6.** Give an **LK**-derivation for the sequent  $\neg A \vee B \Rightarrow A \rightarrow B$ .

Begin by writing the desired end-sequent at the bottom of the derivation.

$$\overline{\neg A \vee B \Rightarrow A \rightarrow B}$$

To find a logical rule that could give us this end-sequent, we look at the logical connectives in the end-sequent:  $\neg$ ,  $\vee$ , and  $\rightarrow$ . We only care at the moment about  $\vee$  and  $\rightarrow$  because they are main

operators of sentences in the end-sequent, while  $\neg$  is inside the scope of another connective, so we will take care of it later. Our options for logical rules for the final inference are therefore the  $\vee\text{L}$  rule and the  $\rightarrow\text{R}$  rule. We could pick either rule, really, but let's pick the  $\rightarrow\text{R}$  rule (if for no reason other than it allows us to put off splitting into two branches). According to the form of  $\rightarrow\text{R}$  inferences which can yield the lower sequent, this must look like:

$$\frac{\overline{A, \neg A \vee B \Rightarrow B}}{\neg A \vee B \Rightarrow A \rightarrow B} \rightarrow\text{R}$$

If we move  $\neg A \vee B$  to the outside of the antecedent, we can apply the  $\vee\text{L}$  rule. According to the schema, this must split into two upper sequents as follows:

$$\frac{\overline{\neg A, A \Rightarrow B} \quad \overline{B, A \Rightarrow B}}{\neg A \vee B, A \Rightarrow B} \vee\text{L}$$

$$\frac{\overline{A, \neg A \vee B \Rightarrow B}}{\neg A \vee B \Rightarrow A \rightarrow B} \text{XR}$$

$$\rightarrow\text{R}$$

Remember that we are trying to wind our way up to initial sequents; we seem to be pretty close! The right branch is just one weakening and one exchange away from an initial sequent and then it is done:

$$\frac{\overline{\neg A, A \Rightarrow B} \quad \overline{B \Rightarrow B}}{A, B \Rightarrow B} \text{WL}$$

$$\frac{\overline{A, B \Rightarrow B}}{B, A \Rightarrow B} \text{XL}$$

$$\frac{\overline{\neg A, A \Rightarrow B} \quad \overline{B, A \Rightarrow B}}{\neg A \vee B, A \Rightarrow B} \vee\text{L}$$

$$\frac{\overline{A, \neg A \vee B \Rightarrow B}}{\neg A \vee B \Rightarrow A \rightarrow B} \text{XR}$$

$$\rightarrow\text{R}$$

Now looking at the left branch, the only logical connective in any sentence is the  $\neg$  symbol in the antecedent sentences, so we're looking at an instance of the  $\neg\text{L}$  rule.

$$\begin{array}{c}
\frac{}{A \Rightarrow B, A} \quad \frac{B \Rightarrow B}{A, B \Rightarrow B} \text{WL} \\
\frac{}{\neg A, A \Rightarrow B} \neg\text{L} \quad \frac{A, B \Rightarrow B}{B, A \Rightarrow B} \text{XL} \\
\frac{}{\neg A \vee B, A \Rightarrow B} \vee\text{L} \\
\frac{}{A, \neg A \vee B \Rightarrow B} \text{XR} \\
\frac{}{\neg A \vee B \Rightarrow A \rightarrow B} \rightarrow\text{R}
\end{array}$$

Similarly to how we finished off the right branch, we are just one weakening and one exchange away from finishing off this left branch as well.

$$\begin{array}{c}
\frac{A \Rightarrow A}{A \Rightarrow A, B} \text{WR} \quad \frac{B \Rightarrow B}{A, B \Rightarrow B} \text{WL} \\
\frac{A \Rightarrow A, B}{A \Rightarrow B, A} \text{XR} \quad \frac{A, B \Rightarrow B}{B, A \Rightarrow B} \text{XL} \\
\frac{}{\neg A, A \Rightarrow B} \neg\text{L} \quad \frac{}{\neg A \vee B, A \Rightarrow B} \vee\text{L} \\
\frac{}{A, \neg A \vee B \Rightarrow B} \text{XR} \\
\frac{}{\neg A \vee B \Rightarrow A \rightarrow B} \rightarrow\text{R}
\end{array}$$

**Example 10.7.** Give an **LK**-derivation of the sequent  $\neg A \vee \neg B \Rightarrow \neg(A \wedge B)$

Using the techniques from above, we start by writing the desired end-sequent at the bottom.

$$\overline{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)}$$

The available main connectives of sentences in the end-sequent are the  $\vee$  symbol and the  $\neg$  symbol. It would work to apply either the  $\vee\text{L}$  or the  $\neg\text{R}$  rule here, but we start with the  $\neg\text{R}$  rule because it avoids splitting up into two branches for a moment:

$$\frac{A \wedge B, \neg A \vee \neg B \Rightarrow}{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)} \neg\text{R}$$

Now we have a choice of whether to look at the  $\wedge\text{L}$  or the  $\vee\text{L}$  rule. Let's see what happens when we apply the  $\wedge\text{L}$  rule: we have a choice to start with either the sequent  $A, \neg A \vee \neg B \Rightarrow$  or the sequent  $B, \neg A \vee \neg B \Rightarrow$ . Since the derivation is symmetric with regards to  $A$  and  $B$ , let's go with the former:

$$\frac{\frac{\frac{A, \neg A \vee \neg B \Rightarrow}{A \wedge B, \neg A \vee \neg B \Rightarrow} \wedge L}{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)} \neg R$$

Continuing to fill in the derivation, we see that we run into a problem:

$$\frac{\frac{\frac{A \Rightarrow A}{\neg A, A \Rightarrow} \neg L \quad \frac{\frac{\frac{A \Rightarrow B}{\neg B, A \Rightarrow} ?}{\neg A \vee \neg B, A \Rightarrow} \neg L}{\neg A \vee \neg B, A \Rightarrow} \vee L}{\frac{\frac{A, \neg A \vee \neg B \Rightarrow}{A, \neg A \vee \neg B \Rightarrow} XL}{A \wedge B, \neg A \vee \neg B \Rightarrow} \wedge L}{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)} \neg R$$

The top of the right branch cannot be reduced any further, and it cannot be brought by way of structural inferences to an initial sequent, so this is not the right path to take. So clearly, it was a mistake to apply the  $\wedge L$  rule above. Going back to what we had before and carrying out the  $\vee L$  rule instead, we get

$$\frac{\frac{\frac{\neg A, A \wedge B \Rightarrow}{\neg A \vee \neg B, A \wedge B \Rightarrow} \vee L}{A \wedge B, \neg A \vee \neg B \Rightarrow} XL}{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)} \neg R$$

Completing each branch as we've done before, we get

$$\frac{\frac{\frac{A \Rightarrow A}{A \wedge B \Rightarrow A} \wedge L}{\neg A, A \wedge B \Rightarrow} \neg L \quad \frac{\frac{B \Rightarrow B}{A \wedge B \Rightarrow B} \wedge L}{\neg B, A \wedge B \Rightarrow} \neg L}{\neg A \vee \neg B, A \wedge B \Rightarrow} \vee L}{\frac{A \wedge B, \neg A \vee \neg B \Rightarrow}{\neg A \vee \neg B \Rightarrow \neg(A \wedge B)} XL} \neg R$$

(We could have carried out the  $\wedge$  rules lower than the  $\neg$  rules in these steps and still obtained a correct derivation).

**Example 10.8.** So far we haven't used the contraction rule, but it is sometimes required. Here's an example where that happens. Suppose we want to prove  $\Rightarrow A \vee \neg A$ . Applying  $\vee R$  backwards would give us one of these two derivations:

$$\frac{\overline{\Rightarrow A}}{\Rightarrow A \vee \neg A} \vee R \qquad \frac{\overline{A \Rightarrow}}{\Rightarrow \neg A} \neg R \quad \frac{}{\Rightarrow A \vee \neg A} \vee R$$

Neither of these of course ends in an initial sequent. The trick is to realize that the contraction rule allows us to combine two copies of a sentence into one—and when we're searching for a proof, i.e., going from bottom to top, we can keep a copy of  $A \vee \neg A$  in the premise, e.g.,

$$\frac{\overline{\Rightarrow A \vee \neg A, A}}{\Rightarrow A \vee \neg A, A \vee \neg A} \vee R \quad \frac{}{\Rightarrow A \vee \neg A} CR$$

Now we can apply  $\vee R$  a second time, and also get  $\neg A$ , which leads to a complete derivation.

$$\frac{\frac{A \Rightarrow A}{\Rightarrow A, \neg A} \neg R}{\Rightarrow A, A \vee \neg A} \vee R \quad \frac{}{\Rightarrow A \vee \neg A, A} XR \quad \frac{}{\Rightarrow A \vee \neg A, A \vee \neg A} \vee R \quad \frac{}{\Rightarrow A \vee \neg A} CR$$

## 10.7 Derivations with Quantifiers

**Example 10.9.** Give an LK-derivation of the sequent  $\exists x \neg A(x) \Rightarrow \neg \forall x A(x)$ .

When dealing with quantifiers, we have to make sure not to violate the eigenvariable condition, and sometimes this requires us to play around with the order of carrying out certain inferences. In general, it helps to try and take care of rules subject

to the eigenvariable condition first (they will be lower down in the finished proof). Also, it is a good idea to try and look ahead and try to guess what the initial sequent might look like. In our case, it will have to be something like  $A(a) \Rightarrow A(a)$ . That means that when we are “reversing” the quantifier rules, we will have to pick the same term—what we will call  $a$ —for both the  $\forall$  and the  $\exists$  rule. If we picked different terms for each rule, we would end up with something like  $A(a) \Rightarrow A(b)$ , which, of course, is not derivable.

Starting as usual, we write

$$\overline{\exists x \neg A(x) \Rightarrow \neg \forall x A(x)}$$

We could either carry out the  $\exists\text{L}$  rule or the  $\neg\text{R}$  rule. Since the  $\exists\text{L}$  rule is subject to the eigenvariable condition, it’s a good idea to take care of it sooner rather than later, so we’ll do that one first.

$$\frac{\overline{\neg A(a) \Rightarrow \neg \forall x A(x)}}{\exists x \neg A(x) \Rightarrow \neg \forall x A(x)} \exists\text{L}$$

Applying the  $\neg\text{L}$  and  $\neg\text{R}$  rules backwards, we get

$$\frac{\frac{\frac{\overline{\forall x A(x) \Rightarrow A(a)}}{\neg A(a), \forall x A(x) \Rightarrow} \neg\text{L}}{\forall x A(x), \neg A(a) \Rightarrow} \text{XL}}{\neg A(a) \Rightarrow \neg \forall x A(x)} \neg\text{R}}{\exists x \neg A(x) \Rightarrow \neg \forall x A(x)} \exists\text{L}$$

At this point, our only option is to carry out the  $\forall\text{L}$  rule. Since this rule is not subject to the eigenvariable restriction, we’re in the clear. Remember, we want to try and obtain an initial sequent (of the form  $A(a) \Rightarrow A(a)$ ), so we should choose  $a$  as our argument for  $A$  when we apply the rule.

$$\begin{array}{c}
 \frac{A(a) \Rightarrow A(a)}{\forall x A(x) \Rightarrow A(a)} \forall L \\
 \frac{\quad}{\neg A(a), \forall x A(x) \Rightarrow} \neg L \\
 \frac{\quad}{\forall x A(x), \neg A(a) \Rightarrow} XL \\
 \frac{\quad}{\neg A(a) \Rightarrow \neg \forall x A(x)} \neg R \\
 \frac{\quad}{\exists x \neg A(x) \Rightarrow \neg \forall x A(x)} \exists L
 \end{array}$$

It is important, especially when dealing with quantifiers, to double check at this point that the eigenvariable condition has not been violated. Since the only rule we applied that is subject to the eigenvariable condition was  $\exists L$ , and the eigenvariable  $a$  does not occur in its lower sequent (the end-sequent), this is a correct derivation.

## 10.8 Proof-Theoretic Notions

Just as we've defined a number of important semantic notions (validity, entailment, satisfiability), we now define corresponding *proof-theoretic notions*. These are not defined by appeal to satisfaction of sentences in structures, but by appeal to the derivability or non-derivability of certain sequents. It was an important discovery that these notions coincide. That they do is the content of the *soundness and completeness theorem*.

**Definition 10.10 (Theorems).** A sentence  $A$  is a *theorem* if there is a derivation in **LK** of the sequent  $\Rightarrow A$ . We write  $\vdash A$  if  $A$  is a theorem and  $\not\vdash A$  if it is not.

**Definition 10.11 (Derivability).** A sentence  $A$  is *derivable from* a set of sentences  $\Gamma$ ,  $\Gamma \vdash A$ , iff there is a finite subset  $\Gamma_0 \subseteq \Gamma$  and a sequence  $\Gamma'_0$  of the sentences in  $\Gamma_0$  such that **LK** derives  $\Gamma'_0 \Rightarrow A$ . If  $A$  is not derivable from  $\Gamma$  we write  $\Gamma \not\vdash A$ .

Because of the contraction, weakening, and exchange rules, the order and number of sentences in  $\Gamma'_0$  does not matter: if a

sequent  $\Gamma'_0 \Rightarrow A$  is derivable, then so is  $\Gamma''_0 \Rightarrow A$  for any  $\Gamma''_0$  that contains the same sentences as  $\Gamma'_0$ . For instance, if  $\Gamma_0 = \{B, C\}$  then both  $\Gamma'_0 = \langle B, B, C \rangle$  and  $\Gamma''_0 = \langle C, C, B \rangle$  are sequences containing just the sentences in  $\Gamma_0$ . If a sequent containing one is derivable, so is the other, e.g.:

$$\begin{array}{c} \vdots \\ \vdots \\ \frac{B, B, C \Rightarrow A}{B, C \Rightarrow A} \text{CL} \\ \frac{B, C \Rightarrow A}{C, B \Rightarrow A} \text{XL} \\ \frac{C, B \Rightarrow A}{C, C, B \Rightarrow A} \text{WL} \end{array}$$

From now on we'll say that if  $\Gamma_0$  is a finite set of sentences then  $\Gamma_0 \Rightarrow A$  is any sequent where the antecedent is a sequence of sentences in  $\Gamma_0$  and tacitly include contractions, exchanges, and weakenings if necessary.

**Definition 10.12 (Consistency).** A set of sentences  $\Gamma$  is *inconsistent* iff there is a finite subset  $\Gamma_0 \subseteq \Gamma$  such that **LK** derives  $\Gamma_0 \Rightarrow$  . If  $\Gamma$  is not inconsistent, i.e., if for every finite  $\Gamma_0 \subseteq \Gamma$ , **LK** does not derive  $\Gamma_0 \Rightarrow$  , we say it is *consistent*.

**Proposition 10.13 (Reflexivity).** If  $A \in \Gamma$ , then  $\Gamma \vdash A$ .

*Proof.* The initial sequent  $A \Rightarrow A$  is derivable, and  $\{A\} \subseteq \Gamma$ .  $\square$

**Proposition 10.14 (Monotonicity).** If  $\Gamma \subseteq \Delta$  and  $\Gamma \vdash A$ , then  $\Delta \vdash A$ .

*Proof.* Suppose  $\Gamma \vdash A$ , i.e., there is a finite  $\Gamma_0 \subseteq \Gamma$  such that  $\Gamma_0 \Rightarrow A$  is derivable. Since  $\Gamma \subseteq \Delta$ , then  $\Gamma_0$  is also a finite subset of  $\Delta$ . The derivation of  $\Gamma_0 \Rightarrow A$  thus also shows  $\Delta \vdash A$ .  $\square$



**Proposition 10.15 (Transitivity).** *If  $\Gamma \vdash A$  and  $\{A\} \cup \Delta \vdash B$ , then  $\Gamma \cup \Delta \vdash B$ .*

*Proof.* If  $\Gamma \vdash A$ , there is a finite  $\Gamma_0 \subseteq \Gamma$  and a derivation  $\pi_0$  of  $\Gamma_0 \Rightarrow A$ . If  $\{A\} \cup \Delta \vdash B$ , then for some finite subset  $\Delta_0 \subseteq \Delta$ , there is a derivation  $\pi_1$  of  $A, \Delta_0 \Rightarrow B$ . Consider the following derivation:

$$\frac{\begin{array}{c} \vdots \\ \pi_0 \\ \vdots \\ \Gamma_0 \Rightarrow A \end{array} \quad \begin{array}{c} \vdots \\ \pi_1 \\ \vdots \\ A, \Delta_0 \Rightarrow B \end{array}}{\Gamma_0, \Delta_0 \Rightarrow B} \text{Cut}$$

Since  $\Gamma_0 \cup \Delta_0 \subseteq \Gamma \cup \Delta$ , this shows  $\Gamma \cup \Delta \vdash B$ .  $\square$

Note that this means that in particular if  $\Gamma \vdash A$  and  $A \vdash B$ , then  $\Gamma \vdash B$ . It follows also that if  $A_1, \dots, A_n \vdash B$  and  $\Gamma \vdash A_i$  for each  $i$ , then  $\Gamma \vdash B$ .

**Proposition 10.16.**  *$\Gamma$  is inconsistent iff  $\Gamma \vdash A$  for every sentence  $A$ .*

*Proof.* Exercise.  $\square$

**Proposition 10.17 (Compactness).** *1. If  $\Gamma \vdash A$  then there is a finite subset  $\Gamma_0 \subseteq \Gamma$  such that  $\Gamma_0 \vdash A$ .*

*2. If every finite subset of  $\Gamma$  is consistent, then  $\Gamma$  is consistent.*

*Proof.* 1. If  $\Gamma \vdash A$ , then there is a finite subset  $\Gamma_0 \subseteq \Gamma$  such that the sequent  $\Gamma_0 \Rightarrow A$  has a derivation. Consequently,  $\Gamma_0 \vdash A$ .

2. If  $\Gamma$  is inconsistent, there is a finite subset  $\Gamma_0 \subseteq \Gamma$  such that **LK** derives  $\Gamma_0 \Rightarrow \perp$ . But then  $\Gamma_0$  is a finite subset of  $\Gamma$  that is inconsistent.  $\square$

## 10.9 Derivability and Consistency

We will now establish a number of properties of the derivability relation. They are independently interesting, but each will play a role in the proof of the completeness theorem.

**Proposition 10.18.** *If  $\Gamma \vdash A$  and  $\Gamma \cup \{A\}$  is inconsistent, then  $\Gamma$  is inconsistent.*

*Proof.* There are finite  $\Gamma_0$  and  $\Gamma_1 \subseteq \Gamma$  such that **LK** derives  $\Gamma_0 \Rightarrow A$  and  $A, \Gamma_1 \Rightarrow$  . Let the **LK**-derivation of  $\Gamma_0 \Rightarrow A$  be  $\pi_0$  and the **LK**-derivation of  $\Gamma_1, A \Rightarrow$  be  $\pi_1$ . We can then derive

$$\frac{\begin{array}{c} \vdots \\ \vdots \pi_0 \\ \vdots \\ \Gamma_0 \Rightarrow A \end{array} \quad \begin{array}{c} \vdots \\ \vdots \pi_1 \\ \vdots \\ A, \Gamma_1 \Rightarrow \end{array}}{\Gamma_0, \Gamma_1 \Rightarrow} \text{Cut}$$

Since  $\Gamma_0 \subseteq \Gamma$  and  $\Gamma_1 \subseteq \Gamma$ ,  $\Gamma_0 \cup \Gamma_1 \subseteq \Gamma$ , hence  $\Gamma$  is inconsistent.  $\square$

**Proposition 10.19.**  *$\Gamma \vdash A$  iff  $\Gamma \cup \{\neg A\}$  is inconsistent.*

*Proof.* First suppose  $\Gamma \vdash A$ , i.e., there is a derivation  $\pi_0$  of  $\Gamma \Rightarrow A$ . By adding a  $\neg$ L rule, we obtain a derivation of  $\neg A, \Gamma \Rightarrow$  , i.e.,  $\Gamma \cup \{\neg A\}$  is inconsistent.

If  $\Gamma \cup \{\neg A\}$  is inconsistent, there is a derivation  $\pi_1$  of  $\neg A, \Gamma \Rightarrow$  . The following is a derivation of  $\Gamma \Rightarrow A$ :

$$\frac{\frac{A \Rightarrow A}{\Rightarrow A, \neg A} \neg\text{R} \quad \begin{array}{c} \vdots \\ \vdots \pi_1 \\ \vdots \\ \neg A, \Gamma \Rightarrow \end{array}}{\Gamma \Rightarrow A} \text{Cut} \quad \square$$

**Proposition 10.20.** *If  $\Gamma \vdash A$  and  $\neg A \in \Gamma$ , then  $\Gamma$  is inconsistent.*

*Proof.* Suppose  $\Gamma \vdash A$  and  $\neg A \in \Gamma$ . Then there is a derivation  $\pi$  of a sequent  $\Gamma_0 \Rightarrow A$ . The sequent  $\neg A, \Gamma_0 \Rightarrow$  is also derivable:

$$\frac{\frac{\frac{\vdots \pi}{\Gamma_0 \Rightarrow A} \quad \frac{A \Rightarrow A}{\neg A, A \Rightarrow} \neg\text{L}}{A, \neg A \Rightarrow} \text{XL}}{\Gamma_0, \neg A \Rightarrow} \text{Cut}$$

Since  $\neg A \in \Gamma$  and  $\Gamma_0 \subseteq \Gamma$ , this shows that  $\Gamma$  is inconsistent.  $\square$

**Proposition 10.21.** *If  $\Gamma \cup \{A\}$  and  $\Gamma \cup \{\neg A\}$  are both inconsistent, then  $\Gamma$  is inconsistent.*

*Proof.* There are finite sets  $\Gamma_0 \subseteq \Gamma$  and  $\Gamma_1 \subseteq \Gamma$  and **LK**-derivations  $\pi_0$  and  $\pi_1$  of  $A, \Gamma_0 \Rightarrow$  and  $\neg A, \Gamma_1 \Rightarrow$ , respectively. We can then derive

$$\frac{\frac{\frac{\vdots \pi_0}{A, \Gamma_0 \Rightarrow} \quad \frac{\Gamma_0 \Rightarrow \neg A}{\neg R}}{\Gamma_0 \Rightarrow \neg A} \quad \frac{\frac{\vdots \pi_1}{\neg A, \Gamma_1 \Rightarrow}}{\Gamma_0, \Gamma_1 \Rightarrow} \text{Cut}}$$

Since  $\Gamma_0 \subseteq \Gamma$  and  $\Gamma_1 \subseteq \Gamma$ ,  $\Gamma_0 \cup \Gamma_1 \subseteq \Gamma$ . Hence  $\Gamma$  is inconsistent.  $\square$

## 10.10 Derivability and the Propositional Connectives

We establish that the derivability relation  $\vdash$  of the sequent calculus is strong enough to establish some basic facts involving the propositional connectives, such as that  $A \wedge B \vdash A$  and  $A, A \rightarrow B \vdash B$  (modus ponens). These facts are needed for the proof of the completeness theorem.

**Proposition 10.22.** 1. Both  $A \wedge B \vdash A$  and  $A \wedge B \vdash B$ .

2.  $A, B \vdash A \wedge B$ .

*Proof.* 1. Both sequents  $A \wedge B \Rightarrow A$  and  $A \wedge B \Rightarrow B$  are derivable:

$$\frac{A \Rightarrow A}{A \wedge B \Rightarrow A} \wedge\text{L} \qquad \frac{B \Rightarrow B}{A \wedge B \Rightarrow B} \wedge\text{L}$$

2. Here is a derivation of the sequent  $A, B \Rightarrow A \wedge B$ :

$$\frac{A \Rightarrow A \quad B \Rightarrow B}{A, B \Rightarrow A \wedge B} \wedge\text{R} \qquad \square$$

**Proposition 10.23.** 1.  $A \vee B, \neg A, \neg B$  is inconsistent.

2. Both  $A \vdash A \vee B$  and  $B \vdash A \vee B$ .

*Proof.* 1. We give a derivation of the sequent  $A \vee B, \neg A, \neg B \Rightarrow$ :

$$\frac{\frac{\frac{A \Rightarrow A}{\neg A, A \Rightarrow} \neg\text{L}}{A, \neg A, \neg B \Rightarrow} \quad \frac{\frac{B \Rightarrow B}{\neg B, B \Rightarrow} \neg\text{L}}{B, \neg A, \neg B \Rightarrow}}{A \vee B, \neg A, \neg B \Rightarrow} \vee\text{L}$$

(Recall that double inference lines indicate several weakening, contraction, and exchange inferences.)

2. Both sequents  $A \Rightarrow A \vee B$  and  $B \Rightarrow A \vee B$  have derivations:

$$\frac{A \Rightarrow A}{A \Rightarrow A \vee B} \vee\text{R} \qquad \frac{B \Rightarrow B}{B \Rightarrow A \vee B} \vee\text{R} \qquad \square$$

**Proposition 10.24.** 1.  $A, A \rightarrow B \vdash B$ .

2. Both  $\neg A \vdash A \rightarrow B$  and  $B \vdash A \rightarrow B$ .

*Proof.* 1. The sequent  $A \rightarrow B, A \Rightarrow B$  is derivable:

$$\frac{A \Rightarrow A \quad B \Rightarrow B}{A \rightarrow B, A \Rightarrow B} \rightarrow\text{L}$$

2. Both sequents  $\neg A \Rightarrow A \rightarrow B$  and  $B \Rightarrow A \rightarrow B$  are derivable:

$$\frac{\frac{\frac{A \Rightarrow A}{\neg A, A \Rightarrow} \neg\text{L}}{A, \neg A \Rightarrow} \text{XL}}{A, \neg A \Rightarrow B} \text{WR} \rightarrow\text{R} \quad \frac{B \Rightarrow B}{A, B \Rightarrow B} \text{WL} \rightarrow\text{R} \quad \square$$

## 10.11 Derivability and the Quantifiers

The completeness theorem also requires that the sequent calculus rules yield the facts about  $\vdash$  established in this section.

**Theorem 10.25.** *If  $c$  is a constant not occurring in  $\Gamma$  or  $A(x)$  and  $\Gamma \vdash A(c)$ , then  $\Gamma \vdash \forall x A(x)$ .*

*Proof.* Let  $\pi_0$  be an **LK**-derivation of  $\Gamma_0 \Rightarrow A(c)$  for some finite  $\Gamma_0 \subseteq \Gamma$ . By adding a  $\forall\text{R}$  inference, we obtain a derivation of  $\Gamma_0 \Rightarrow \forall x A(x)$ , since  $c$  does not occur in  $\Gamma$  or  $A(x)$  and thus the eigenvariable condition is satisfied.  $\square$

**Proposition 10.26.** 1.  $A(t) \vdash \exists x A(x)$ .

2.  $\forall x A(x) \vdash A(t)$ .

*Proof.* 1. The sequent  $A(t) \Rightarrow \exists x A(x)$  is derivable:

$$\frac{A(t) \Rightarrow A(t)}{A(t) \Rightarrow \exists x A(x)} \exists R$$

2. The sequent  $\forall x A(x) \Rightarrow A(t)$  is derivable:

$$\frac{A(t) \Rightarrow A(t)}{\forall x A(x) \Rightarrow A(t)} \forall L$$

□

## 10.12 Soundness

A derivation system, such as the sequent calculus, is *sound* if it cannot derive things that do not actually hold. Soundness is thus a kind of guaranteed safety property for derivation systems. Depending on which proof theoretic property is in question, we would like to know for instance, that

1. every derivable  $A$  is valid;
2. if a sentence is derivable from some others, it is also a consequence of them;
3. if a set of sentences is inconsistent, it is unsatisfiable.

These are important properties of a derivation system. If any of them do not hold, the derivation system is deficient—it would derive too much. Consequently, establishing the soundness of a derivation system is of the utmost importance.

Because all these proof-theoretic properties are defined via derivability in the sequent calculus of certain sequents, proving (1)–(3) above requires proving something about the semantic properties of derivable sequents. We will first define what it means for a sequent to be *valid*, and then show that every derivable sequent is valid. (1)–(3) then follow as corollaries from this result.

**Definition 10.27.** A structure  $M$  satisfies a sequent  $\Gamma \Rightarrow \Delta$  iff either  $M \not\models A$  for some  $A \in \Gamma$  or  $M \models A$  for some  $A \in \Delta$ .

A sequent is *valid* iff every structure  $M$  satisfies it.

**Theorem 10.28 (Soundness).** If LK derives  $\Theta \Rightarrow \Xi$ , then  $\Theta \Rightarrow \Xi$  is valid.

*Proof.* Let  $\pi$  be a derivation of  $\Theta \Rightarrow \Xi$ . We proceed by induction on the number of inferences  $n$  in  $\pi$ .

If the number of inferences is 0, then  $\pi$  consists only of an initial sequent. Every initial sequent  $A \Rightarrow A$  is obviously valid, since for every  $M$ , either  $M \not\models A$  or  $M \models A$ .

If the number of inferences is greater than 0, we distinguish cases according to the type of the lowermost inference. By induction hypothesis, we can assume that the premises of that inference are valid, since the number of inferences in the derivation of any premise is smaller than  $n$ .

First, we consider the possible inferences with only one premise.

1. The last inference is a weakening. Then  $\Theta \Rightarrow \Xi$  is either  $A, \Gamma \Rightarrow \Delta$  (if the last inference is WL) or  $\Gamma \Rightarrow \Delta, A$  (if it's WR), and the derivation ends in one of

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ \Gamma \Rightarrow \Delta \end{array}}{A, \Gamma \Rightarrow \Delta} \text{WL} \qquad \frac{\begin{array}{c} \vdots \\ \vdots \\ \Gamma \Rightarrow \Delta \end{array}}{\Gamma \Rightarrow \Delta, A} \text{WR}$$

By induction hypothesis,  $\Gamma \Rightarrow \Delta$  is valid, i.e., for every structure  $M$ , either there is some  $C \in \Gamma$  such that  $M \not\models C$  or there is some  $C \in \Delta$  such that  $M \models C$ .

If  $M \not\models C$  for some  $C \in \Gamma$ , then  $C \in \Theta$  as well since  $\Theta = A, \Gamma$ , and so  $M \not\models C$  for some  $C \in \Theta$ . Similarly, if  $M \models C$  for some

$C \in \Delta$ , as  $C \in \Xi$ ,  $M \vDash C$  for some  $C \in \Xi$ . Consequently,  $\Theta \Rightarrow \Xi$  is valid.

2. The last inference is  $\neg$ L: Then the premise of the last inference is  $\Gamma \Rightarrow \Delta, A$  and the conclusion is  $\neg A, \Gamma \Rightarrow \Delta$ , i.e., the derivation ends in

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ \Gamma \Rightarrow \Delta, A \end{array}}{\neg A, \Gamma \Rightarrow \Delta} \neg\text{L}$$

and  $\Theta = \neg A, \Gamma$  while  $\Xi = \Delta$ .

The induction hypothesis tells us that  $\Gamma \Rightarrow \Delta, A$  is valid, i.e., for every  $M$ , either (a) for some  $C \in \Gamma$ ,  $M \vDash C$ , or (b) for some  $C \in \Delta$ ,  $M \vDash C$ , or (c)  $M \vDash A$ . We want to show that  $\Theta \Rightarrow \Xi$  is also valid. Let  $M$  be a structure. If (a) holds, then there is  $C \in \Gamma$  so that  $M \vDash C$ , but  $C \in \Theta$  as well. If (b) holds, there is  $C \in \Delta$  such that  $M \vDash C$ , but  $C \in \Xi$  as well. Finally, if  $M \vDash A$ , then  $M \vDash \neg A$ . Since  $\neg A \in \Theta$ , there is  $C \in \Theta$  such that  $M \vDash C$ . Consequently,  $\Theta \Rightarrow \Xi$  is valid.

3. The last inference is  $\neg$ R: Exercise.
4. The last inference is  $\wedge$ L: There are two variants:  $A \wedge B$  may be inferred on the left from  $A$  or from  $B$  on the left side of the premise. In the first case, the  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ A, \Gamma \Rightarrow \Delta \end{array}}{A \wedge B, \Gamma \Rightarrow \Delta} \wedge\text{L}$$

and  $\Theta = A \wedge B, \Gamma$  while  $\Xi = \Delta$ . Consider a structure  $M$ . Since by induction hypothesis,  $A, \Gamma \Rightarrow \Delta$  is valid, (a)  $M \vDash A$ , (b)  $M \vDash C$  for some  $C \in \Gamma$ , or (c)  $M \vDash C$  for some  $C \in \Delta$ . In



case (a),  $M \not\models A \wedge B$ , so there is  $C \in \Theta$  (namely,  $A \wedge B$ ) such that  $M \not\models C$ . In case (b), there is  $C \in \Gamma$  such that  $M \not\models C$ , and  $C \in \Theta$  as well. In case (c), there is  $C \in \Delta$  such that  $M \models C$ , and  $C \in \Xi$  as well since  $\Xi = \Delta$ . So in each case,  $M$  satisfies  $A \wedge B, \Gamma \Rightarrow \Delta$ . Since  $M$  was arbitrary,  $\Gamma \Rightarrow \Delta$  is valid. The case where  $A \wedge B$  is inferred from  $B$  is handled the same, changing  $A$  to  $B$ .

5. The last inference is  $\vee R$ : There are two variants:  $A \vee B$  may be inferred on the right from  $A$  or from  $B$  on the right side of the premise. In the first case,  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ \Gamma \Rightarrow \Delta, A \end{array}}{\Gamma \Rightarrow \Delta, A \vee B} \vee R$$

Now  $\Theta = \Gamma$  and  $\Xi = \Delta, A \vee B$ . Consider a structure  $M$ . Since  $\Gamma \Rightarrow \Delta, A$  is valid, (a)  $M \models A$ , (b)  $M \not\models C$  for some  $C \in \Gamma$ , or (c)  $M \models C$  for some  $C \in \Delta$ . In case (a),  $M \models A \vee B$ . In case (b), there is  $C \in \Gamma$  such that  $M \not\models C$ . In case (c), there is  $C \in \Delta$  such that  $M \models C$ . So in each case,  $M$  satisfies  $\Gamma \Rightarrow \Delta, A \vee B$ , i.e.,  $\Theta \Rightarrow \Xi$ . Since  $M$  was arbitrary,  $\Theta \Rightarrow \Xi$  is valid. The case where  $A \vee B$  is inferred from  $B$  is handled the same, changing  $A$  to  $B$ .

6. The last inference is  $\rightarrow R$ : Then  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ A, \Gamma \Rightarrow \Delta, B \end{array}}{\Gamma \Rightarrow \Delta, A \rightarrow B} \rightarrow R$$

Again, the induction hypothesis says that the premise is valid; we want to show that the conclusion is valid as well. Let  $M$  be arbitrary. Since  $A, \Gamma \Rightarrow \Delta, B$  is valid, at least one

of the following cases obtains: (a)  $M \not\models A$ , (b)  $M \models B$ , (c)  $M \not\models C$  for some  $C \in \Gamma$ , or (d)  $M \models C$  for some  $C \in \Delta$ . In cases (a) and (b),  $M \models A \rightarrow B$  and so there is a  $C \in \Delta, A \rightarrow B$  such that  $M \models C$ . In case (c), for some  $C \in \Gamma$ ,  $M \not\models C$ . In case (d), for some  $C \in \Delta$ ,  $M \models C$ . In each case,  $M$  satisfies  $\Gamma \Rightarrow \Delta, A \rightarrow B$ . Since  $M$  was arbitrary,  $\Gamma \Rightarrow \Delta, A \rightarrow B$  is valid.

7. The last inference is  $\forall\text{L}$ : Then there is a formula  $A(x)$  and a closed term  $t$  such that  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ A(t), \Gamma \Rightarrow \Delta \end{array}}{\forall x A(x), \Gamma \Rightarrow \Delta} \forall\text{L}$$

We want to show that the conclusion  $\forall x A(x), \Gamma \Rightarrow \Delta$  is valid. Consider a structure  $M$ . Since the premise  $A(t), \Gamma \Rightarrow \Delta$  is valid, (a)  $M \not\models A(t)$ , (b)  $M \not\models C$  for some  $C \in \Gamma$ , or (c)  $M \models C$  for some  $C \in \Delta$ . In case (a), by **Proposition 7.30**, if  $M \models \forall x A(x)$ , then  $M \models A(t)$ . Since  $M \not\models A(t)$ ,  $M \not\models \forall x A(x)$ . In case (b) and (c),  $M$  also satisfies  $\forall x A(x), \Gamma \Rightarrow \Delta$ . Since  $M$  was arbitrary,  $\forall x A(x), \Gamma \Rightarrow \Delta$  is valid.

8. The last inference is  $\exists\text{R}$ : Exercise.
9. The last inference is  $\forall\text{R}$ : Then there is a formula  $A(x)$  and a constant symbol  $a$  such that  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \Gamma \Rightarrow \Delta, A(a) \end{array}}{\Gamma \Rightarrow \Delta, \forall x A(x)} \forall\text{R}$$

where the eigenvariable condition is satisfied, i.e.,  $a$  does not occur in  $A(x)$ ,  $\Gamma$ , or  $\Delta$ . By induction hypothesis, the

premise of the last inference is valid. We have to show that the conclusion is valid as well, i.e., that for any structure  $M$ , (a)  $M \vDash \forall x A(x)$ , (b)  $M \not\vDash C$  for some  $C \in \Gamma$ , or (c)  $M \vDash C$  for some  $C \in \Delta$ .

Suppose  $M$  is an arbitrary structure. If (b) or (c) holds, we are done, so suppose neither holds: for all  $C \in \Gamma$ ,  $M \vDash C$ , and for all  $C \in \Delta$ ,  $M \not\vDash C$ . We have to show that (a) holds, i.e.,  $M \vDash \forall x A(x)$ . By **Proposition 7.18**, it suffices to show that  $M, s \vDash A(x)$  for all variable assignments  $s$ . So let  $s$  be an arbitrary variable assignment. Consider the structure  $M'$  which is just like  $M$  except  $a^{M'} = s(x)$ . By **Corollary 7.20**, for any  $C \in \Gamma$ ,  $M' \vDash C$  since  $a$  does not occur in  $\Gamma$ , and for any  $C \in \Delta$ ,  $M' \not\vDash C$ . But the premise is valid, so  $M' \vDash A(a)$ . By **Proposition 7.17**,  $M', s \vDash A(a)$ , since  $A(a)$  is a sentence. Now  $s \sim_x s$  with  $s(x) = \text{Val}_s^{M'}(a)$ , since we've defined  $M'$  in just this way. So **Proposition 7.22** applies, and we get  $M', s \vDash A(x)$ . Since  $a$  does not occur in  $A(x)$ , by **Proposition 7.19**,  $M, s \vDash A(x)$ . Since  $s$  was arbitrary, we've completed the proof that  $M, s \vDash A(x)$  for all variable assignments.

10. The last inference is  $\exists\text{L}$ : Exercise.

Now let's consider the possible inferences with two premises.

1. The last inference is a cut: then  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \Gamma \Rightarrow \Delta, A \end{array} \quad \begin{array}{c} \vdots \\ A, \Pi \Rightarrow \Lambda \end{array}}{\Gamma, \Pi \Rightarrow \Delta, \Lambda} \text{Cut}$$

Let  $M$  be a structure. By induction hypothesis, the premises are valid, so  $M$  satisfies both premises. We distinguish two cases: (a)  $M \not\vDash A$  and (b)  $M \vDash A$ . In case (a), in order for  $M$  to satisfy the left premise, it must satisfy  $\Gamma \Rightarrow \Delta$ . But then

it also satisfies the conclusion. In case (b), in order for  $M$  to satisfy the right premise, it must satisfy  $\Pi \setminus \Delta$ . Again,  $M$  satisfies the conclusion.

2. The last inference is  $\wedge R$ . Then  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \Gamma \Rightarrow \Delta, A \end{array} \quad \begin{array}{c} \vdots \\ \Gamma \Rightarrow \Delta, B \end{array}}{\Gamma \Rightarrow \Delta, A \wedge B} \wedge R$$

Consider a structure  $M$ . If  $M$  satisfies  $\Gamma \Rightarrow \Delta$ , we are done. So suppose it doesn't. Since  $\Gamma \Rightarrow \Delta, A$  is valid by induction hypothesis,  $M \vDash A$ . Similarly, since  $\Gamma \Rightarrow \Delta, B$  is valid,  $M \vDash B$ . But then  $M \vDash A \wedge B$ .

3. The last inference is  $\vee L$ : Exercise.
4. The last inference is  $\rightarrow L$ . Then  $\pi$  ends in

$$\frac{\begin{array}{c} \vdots \\ \Gamma \Rightarrow \Delta, A \end{array} \quad \begin{array}{c} \vdots \\ B, \Pi \Rightarrow \Lambda \end{array}}{A \rightarrow B, \Gamma, \Pi \Rightarrow \Delta, \Lambda} \rightarrow L$$

Again, consider a structure  $M$  and suppose  $M$  doesn't satisfy  $\Gamma, \Pi \Rightarrow \Delta, \Lambda$ . We have to show that  $M \not\vDash A \rightarrow B$ . If  $M$  doesn't satisfy  $\Gamma, \Pi \Rightarrow \Delta, \Lambda$ , it satisfies neither  $\Gamma \Rightarrow \Delta$  nor  $\Pi \Rightarrow \Lambda$ . Since,  $\Gamma \Rightarrow \Delta, A$  is valid, we have  $M \vDash A$ . Since  $B, \Pi \Rightarrow \Lambda$  is valid, we have  $M \not\vDash B$ . But then  $M \not\vDash A \rightarrow B$ , which is what we wanted to show.  $\square$

**Corollary 10.29.** *If  $\vdash A$  then  $A$  is valid.*

**Corollary 10.30.** *If  $\Gamma \vdash A$  then  $\Gamma \vDash A$ .*

*Proof.* If  $\Gamma \vdash A$  then for some finite subset  $\Gamma_0 \subseteq \Gamma$ , there is a derivation of  $\Gamma_0 \Rightarrow A$ . By **Theorem 10.28**, every structure  $M$  either makes some  $B \in \Gamma_0$  false or makes  $A$  true. Hence, if  $M \vDash \Gamma$  then also  $M \vDash A$ .  $\square$

**Corollary 10.31.** *If  $\Gamma$  is satisfiable, then it is consistent.*

*Proof.* We prove the contrapositive. Suppose that  $\Gamma$  is not consistent. Then there is a finite  $\Gamma_0 \subseteq \Gamma$  and a derivation of  $\Gamma_0 \Rightarrow \perp$ . By **Theorem 10.28**,  $\Gamma_0 \Rightarrow \perp$  is valid. In other words, for every structure  $M$ , there is  $C \in \Gamma_0$  so that  $M \not\vDash C$ , and since  $\Gamma_0 \subseteq \Gamma$ , that  $C$  is also in  $\Gamma$ . Thus, no  $M$  satisfies  $\Gamma$ , and  $\Gamma$  is not satisfiable.  $\square$

### 10.13 Derivations with Identity predicate

Derivations with identity predicate require additional initial sequents and inference rules.

**Definition 10.32 (Initial sequents for =).** If  $t$  is a closed term, then  $\Rightarrow t = t$  is an initial sequent.

The rules for = are ( $t_1$  and  $t_2$  are closed terms):

$$\frac{t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_1)}{t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_2)} = \qquad \frac{t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_2)}{t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_1)} =$$

**Example 10.33.** If  $s$  and  $t$  are closed terms, then  $s = t, A(s) \vdash A(t)$ :

$$\frac{\frac{A(s) \Rightarrow A(s)}{s = t, A(s) \Rightarrow A(s)} \text{WL}}{s = t, A(s) \Rightarrow A(t)} =$$

This may be familiar as the principle of substitutability of identicals, or Leibniz' Law.

**LK** proves that  $=$  is symmetric and transitive:

$$\frac{\Rightarrow t_1 = t_1}{\frac{t_1 = t_2 \Rightarrow t_1 = t_1}{t_1 = t_2 \Rightarrow t_2 = t_1}} \text{WL} \qquad \frac{\frac{t_1 = t_2 \Rightarrow t_1 = t_2}{t_2 = t_3, t_1 = t_2 \Rightarrow t_1 = t_2} \text{WL}}{\frac{t_2 = t_3, t_1 = t_2 \Rightarrow t_1 = t_3}{t_1 = t_2, t_2 = t_3 \Rightarrow t_1 = t_3}} \text{XL}$$

In the derivation on the left, the formula  $x = t_1$  is our  $A(x)$ . On the right, we take  $A(x)$  to be  $t_1 = x$ .

## 10.14 Soundness with Identity predicate

**Proposition 10.34.** *LK with initial sequents and rules for identity is sound.*

*Proof.* Initial sequents of the form  $\Rightarrow t = t$  are valid, since for every structure  $M$ ,  $M \vDash t = t$ . (Note that we assume the term  $t$  to be closed, i.e., it contains no variables, so variable assignments are irrelevant).

Suppose the last inference in a derivation is  $=$ . Then the premise is  $t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_1)$  and the conclusion is  $t_1 = t_2, \Gamma \Rightarrow \Delta, A(t_2)$ . Consider a structure  $M$ . We need to show that the conclusion is valid, i.e., if  $M \vDash t_1 = t_2$  and  $M \vDash \Gamma$ , then either  $M \vDash C$  for some  $C \in \Delta$  or  $M \vDash A(t_2)$ .

By induction hypothesis, the premise is valid. This means that if  $M \vDash t_1 = t_2$  and  $M \vDash \Gamma$  either (a) for some  $C \in \Delta$ ,  $M \vDash C$  or (b)  $M \vDash A(t_1)$ . In case (a) we are done. Consider case (b). Let  $s$  be a variable assignment with  $s(x) = \text{Val}^M(t_1)$ . By **Proposition 7.17**,  $M, s \vDash A(t_1)$ . Since  $s \sim_x s$ , by **Proposition 7.22**,  $M, s \vDash A(x)$ . since  $M \vDash t_1 = t_2$ , we have  $\text{Val}^M(t_1) = \text{Val}^M(t_2)$ , and hence  $s(x) = \text{Val}^M(t_2)$ . By applying **Proposition 7.22** again, we also have  $M, s \vDash A(t_2)$ . By **Proposition 7.17**,  $M \vDash A(t_2)$ .  $\square$

## Summary

**Proof systems** provide purely syntactic methods for characterizing consequence and compatibility between sentences. The **sequent calculus** is one such proof system. A **derivation** in it consists of a tree of sequents (a sequent  $\Gamma \Rightarrow \Delta$  consists of two sequences of formulas separated by  $\Rightarrow$ ). The topmost sequents in a derivation are **initial sequents** of the form  $A \Rightarrow A$ . All other sequents, for the derivation to be correct, must be correctly justified by one of a number of **inference rules**. These come in pairs; a rule for operating on the left and on the right side of a sequent for each connective and quantifier. For instance, if a sequent  $\Gamma \Rightarrow \Delta, A \rightarrow B$  is justified by the  $\rightarrow R$  rule, the preceding sequent (the **premise**) must be  $A, \Gamma \Rightarrow \Delta, B$ . Some rules also allow the order or number of sentences in a sequent to be manipulated, e.g., the XR rule allows two formulas on the right side of a sequent to be switched.

If there is a derivation of the sequent  $\Rightarrow A$ , we say  $A$  is a **theorem** and write  $\vdash A$ . If there is a derivation of  $\Gamma_0 \Rightarrow A$  where every  $B$  in  $\Gamma_0$  is in  $\Gamma$ , we say  $A$  is **derivable from  $\Gamma$**  and write  $\Gamma \vdash A$ . If there is a derivation of  $\Gamma_0 \Rightarrow$  where every  $B$  in  $\Gamma_0$  is in  $\Gamma$ , we say  $\Gamma$  is **inconsistent**, otherwise **consistent**. These notions are interrelated, e.g.,  $\Gamma \vdash A$  iff  $\Gamma \cup \{\neg A\}$  is inconsistent. They are also related to the corresponding semantic notions, e.g., if  $\Gamma \vdash A$  then  $\Gamma \vDash A$ . This property of proof systems—what can be derived from  $\Gamma$  is guaranteed to be entailed by  $\Gamma$ —is called **soundness**. The **soundness theorem** is proved by induction on the length of derivations, showing that each individual inference preserves validity of the conclusion sequent provided the premise sequents are valid.

## Problems

**Problem 10.1.** Give derivations of the following sequents:

1.  $A \wedge (B \wedge C) \Rightarrow (A \wedge B) \wedge C$ .

2.  $A \vee (B \vee C) \Rightarrow (A \vee B) \vee C.$
3.  $A \rightarrow (B \rightarrow C) \Rightarrow B \rightarrow (A \rightarrow C).$
4.  $A \Rightarrow \neg\neg A.$

**Problem 10.2.** Give derivations of the following sequents:

1.  $(A \vee B) \rightarrow C \Rightarrow A \rightarrow C.$
2.  $(A \rightarrow C) \wedge (B \rightarrow C) \Rightarrow (A \vee B) \rightarrow C.$
3.  $\Rightarrow \neg(A \wedge \neg A).$
4.  $B \rightarrow A \Rightarrow \neg A \rightarrow \neg B.$
5.  $\Rightarrow (A \rightarrow \neg A) \rightarrow \neg A.$
6.  $\Rightarrow \neg(A \rightarrow B) \rightarrow \neg B.$
7.  $A \rightarrow C \Rightarrow \neg(A \wedge \neg C).$
8.  $A \wedge \neg C \Rightarrow \neg(A \rightarrow C).$
9.  $A \vee B, \neg B \Rightarrow A.$
10.  $\neg A \vee \neg B \Rightarrow \neg(A \wedge B).$
11.  $\Rightarrow (\neg A \wedge \neg B) \rightarrow \neg(A \vee B).$
12.  $\Rightarrow \neg(A \vee B) \rightarrow (\neg A \wedge \neg B).$

**Problem 10.3.** Give derivations of the following sequents:

1.  $\neg(A \rightarrow B) \Rightarrow A.$
2.  $\neg(A \wedge B) \Rightarrow \neg A \vee \neg B.$
3.  $A \rightarrow B \Rightarrow \neg A \vee B.$
4.  $\Rightarrow \neg\neg A \rightarrow A.$
5.  $A \rightarrow B, \neg A \rightarrow B \Rightarrow B.$



$$6. (A \wedge B) \rightarrow C \Rightarrow (A \rightarrow C) \vee (B \rightarrow C).$$

$$7. (A \rightarrow B) \rightarrow A \Rightarrow A.$$

$$8. \Rightarrow (A \rightarrow B) \vee (B \rightarrow C).$$

(These all require the CR rule.)

**Problem 10.4.** Give derivations of the following sequents:

$$1. \Rightarrow (\forall x A(x) \wedge \forall y B(y)) \rightarrow \forall z (A(z) \wedge B(z)).$$

$$2. \Rightarrow (\exists x A(x) \vee \exists y B(y)) \rightarrow \exists z (A(z) \vee B(z)).$$

$$3. \forall x (A(x) \rightarrow B) \Rightarrow \exists y A(y) \rightarrow B.$$

$$4. \forall x \neg A(x) \Rightarrow \neg \exists x A(x).$$

$$5. \Rightarrow \neg \exists x A(x) \rightarrow \forall x \neg A(x).$$

$$6. \Rightarrow \neg \exists x \forall y ((A(x, y) \rightarrow \neg A(y, y)) \wedge (\neg A(y, y) \rightarrow A(x, y))).$$

**Problem 10.5.** Give derivations of the following sequents:

$$1. \Rightarrow \neg \forall x A(x) \rightarrow \exists x \neg A(x).$$

$$2. (\forall x A(x) \rightarrow B) \Rightarrow \exists y (A(y) \rightarrow B).$$

$$3. \Rightarrow \exists x (A(x) \rightarrow \forall y A(y)).$$

(These all require the CR rule.)

**Problem 10.6.** Prove **Proposition 10.16**

**Problem 10.7.** Prove that  $\Gamma \vdash \neg A$  iff  $\Gamma \cup \{A\}$  is inconsistent.

**Problem 10.8.** Complete the proof of **Theorem 10.28**.

**Problem 10.9.** Give derivations of the following sequents:

$$1. \Rightarrow \forall x \forall y ((x = y \wedge A(x)) \rightarrow A(y))$$

$$2. \exists x A(x) \wedge \forall y \forall z ((A(y) \wedge A(z)) \rightarrow y = z) \Rightarrow \exists x (A(x) \wedge \forall y (A(y) \rightarrow y = x))$$

## CHAPTER 11

# *Natural Deduction*

### 11.1 Rules and Derivations

Natural deduction systems are meant to closely parallel the informal reasoning used in mathematical proof (hence it is somewhat “natural”). Natural deduction proofs begin with assumptions. Inference rules are then applied. Assumptions are “discharged” by the  $\neg$ Intro,  $\rightarrow$ Intro,  $\forall$ Elim and  $\exists$ Elim inference rules, and the label of the discharged assumption is placed beside the inference for clarity.

**Definition 11.1 (Assumption).** An *assumption* is any sentence in the topmost position of any branch.

Derivations in natural deduction are certain trees of sentences, where the topmost sentences are assumptions, and if a sentence stands below one, two, or three other sequents, it must follow correctly by a rule of inference. The sentences at the top of the inference are called the *premises* and the sentence below the *conclusion* of the inference. The rules come in pairs, an introduction and an elimination rule for each logical operator. They introduce a logical operator in the conclusion or remove

a logical operator from a premise of the rule. Some of the rules allow an assumption of a certain type to be *discharged*. To indicate which assumption is discharged by which inference, we also assign labels to both the assumption and the inference. This is indicated by writing the assumption as “[ $A$ ] <sup>$n$</sup> .”

It is customary to consider rules for all the logical operators  $\wedge$ ,  $\vee$ ,  $\rightarrow$ ,  $\neg$ , and  $\perp$ , even if some of those are defined.

## 11.2 Propositional Rules

### Rules for $\wedge$

$$\frac{A \quad B}{A \wedge B} \wedge\text{Intro} \qquad \frac{A \wedge B}{A} \wedge\text{Elim}$$

$$\frac{A \wedge B}{B} \wedge\text{Elim}$$

### Rules for $\vee$

$$\frac{A}{A \vee B} \vee\text{Intro}$$

$$\frac{B}{A \vee B} \vee\text{Intro}$$

$$^n \frac{A \vee B \quad \begin{array}{c} [A]^n \\ \vdots \\ C \end{array} \quad \begin{array}{c} [B]^n \\ \vdots \\ C \end{array}}{C} \vee\text{Elim}$$

### Rules for $\rightarrow$

$$^n \frac{\begin{array}{c} [A]^n \\ \vdots \\ B \end{array}}{A \rightarrow B} \rightarrow\text{Intro}$$

$$\frac{A \rightarrow B \quad A}{B} \rightarrow\text{Elim}$$

### Rules for $\neg$

$$\begin{array}{c}
 [A]^n \\
 \vdots \\
 \perp \\
 n \frac{\perp}{\neg A} \neg\text{Intro}
 \end{array}
 \qquad
 \frac{\neg A \quad A}{\perp} \neg\text{Elim}$$

### Rules for $\perp$

$$\frac{\perp}{A} \perp_I
 \qquad
 \begin{array}{c}
 [\neg A]^n \\
 \vdots \\
 \perp \\
 n \frac{\perp}{A} \perp_C
 \end{array}$$

Note that  $\neg$ Intro and  $\perp_C$  are very similar: The difference is that  $\neg$ Intro derives a negated sentence  $\neg A$  but  $\perp_C$  a positive sentence  $A$ .

Whenever a rule indicates that some assumption may be discharged, we take this to be a permission, but not a requirement. E.g., in the  $\rightarrow$ Intro rule, we may discharge any number of assumptions of the form  $A$  in the derivation of the premise  $B$ , including zero.

## 11.3 Quantifier Rules

### Rules for $\forall$

$$\frac{A(a)}{\forall x A(x)} \forall\text{Intro}
 \qquad
 \frac{\forall x A(x)}{A(t)} \forall\text{Elim}$$

In the rules for  $\forall$ ,  $t$  is a closed term (a term that does not contain any variables), and  $a$  is a constant symbol which does not occur in the conclusion  $\forall x A(x)$ , or in any assumption which

is undischarged in the derivation ending with the premise  $A(a)$ . We call  $a$  the *eigenvariable* of the  $\forall$ Intro inference.<sup>1</sup>

### Rules for $\exists$

$$\frac{A(t)}{\exists x A(x)} \exists\text{Intro} \qquad \frac{\begin{array}{c} [A(a)]^n \\ \vdots \\ C \end{array}}{\exists x A(x)} \exists\text{Elim} \qquad \frac{n \quad C}{C}$$

Again,  $t$  is a closed term, and  $a$  is a constant which does not occur in the premise  $\exists x A(x)$ , in the conclusion  $C$ , or any assumption which is undischarged in the derivations ending with the two premises (other than the assumptions  $A(a)$ ). We call  $a$  the *eigenvariable* of the  $\exists$ Elim inference.

The condition that an eigenvariable neither occur in the premises nor in any assumption that is undischarged in the derivations leading to the premises for the  $\forall$ Intro or  $\exists$ Elim inference is called the *eigenvariable condition*.

Recall the convention that when  $A$  is a formula with the variable  $x$  free, we indicate this by writing  $A(x)$ . In the same context,  $A(t)$  then is short for  $A[t/x]$ . So we could also write the  $\exists$ Intro rule as:

$$\frac{A[t/x]}{\exists x A} \exists\text{Intro}$$

Note that  $t$  may already occur in  $A$ , e.g.,  $A$  might be  $P(t, x)$ . Thus, inferring  $\exists x P(t, x)$  from  $P(t, t)$  is a correct application of  $\exists$ Intro—you may “replace” one or more, and not necessarily all, occurrences of  $t$  in the premise by the bound variable  $x$ . However, the eigenvariable conditions in  $\forall$ Intro and  $\exists$ Elim require that the constant symbol  $a$  does not occur in  $A$ . So, you cannot correctly infer  $\forall x P(a, x)$  from  $P(a, a)$  using  $\forall$ Intro.

<sup>1</sup>We use the term “eigenvariable” even though  $a$  in the above rule is a constant. This has historical reasons.

In  $\exists$ Intro and  $\forall$ Elim there are no restrictions, and the term  $t$  can be anything, so we do not have to worry about any conditions. On the other hand, in the  $\exists$ Elim and  $\forall$ Intro rules, the eigenvariable condition requires that the constant symbol  $a$  does not occur anywhere in the conclusion or in an undischarged assumption. The condition is necessary to ensure that the system is sound, i.e., only derives sentences from undischarged assumptions from which they follow. Without this condition, the following would be allowed:

$$\frac{\exists x A(x) \quad \frac{[A(a)]^1}{\forall x A(x)} * \forall \text{Intro}}{\forall x A(x)} \exists \text{Elim}$$

However,  $\exists x A(x) \not\equiv \forall x A(x)$ .

As the elimination rules for quantifiers only allow substituting closed terms for variables, it follows that any formula that can be derived from a set of sentences is itself a sentence.

## 11.4 Derivations

We've said what an assumption is, and we've given the rules of inference. Derivations in natural deduction are inductively generated from these: each derivation either is an assumption on its own, or consists of one, two, or three derivations followed by a correct inference.

**Definition 11.2 (Derivation).** A *derivation* of a sentence  $A$  from assumptions  $\Gamma$  is a finite tree of sentences satisfying the following conditions:

1. The topmost sentences of the tree are either in  $\Gamma$  or are discharged by an inference in the tree.
2. The bottommost sentence of the tree is  $A$ .
3. Every sentence in the tree except the sentence  $A$  at the bot-

tom is a premise of a correct application of an inference rule whose conclusion stands directly below that sentence in the tree.

We then say that  $A$  is the *conclusion* of the derivation and  $\Gamma$  its undischarged assumptions.

If a derivation of  $A$  from  $\Gamma$  exists, we say that  $A$  is *derivable* from  $\Gamma$ , or in symbols:  $\Gamma \vdash A$ . If there is a derivation of  $A$  in which every assumption is discharged, we write  $\vdash A$ .

**Example 11.3.** Every assumption on its own is a derivation. So, e.g.,  $A$  by itself is a derivation, and so is  $B$  by itself. We can obtain a new derivation from these by applying, say, the  $\wedge$ Intro rule,

$$\frac{A \quad B}{A \wedge B} \wedge\text{Intro}$$

These rules are meant to be general: we can replace the  $A$  and  $B$  in it with any sentences, e.g., by  $C$  and  $D$ . Then the conclusion would be  $C \wedge D$ , and so

$$\frac{C \quad D}{C \wedge D} \wedge\text{Intro}$$

is a correct derivation. Of course, we can also switch the assumptions, so that  $D$  plays the role of  $A$  and  $C$  that of  $B$ . Thus,

$$\frac{D \quad C}{D \wedge C} \wedge\text{Intro}$$

is also a correct derivation.

We can now apply another rule, say,  $\rightarrow$ Intro, which allows us to conclude a conditional and allows us to discharge any assumption that is identical to the antecedent of that conditional. So both of the following would be correct derivations:

$$1 \frac{\frac{[C]^1 \quad D}{C \wedge D} \wedge\text{Intro}}{C \rightarrow (C \wedge D)} \rightarrow\text{Intro} \quad 1 \frac{\frac{C \quad [D]^1}{C \wedge D} \wedge\text{Intro}}{D \rightarrow (C \wedge D)} \rightarrow\text{Intro}$$

They show, respectively, that  $D \vdash C \rightarrow (C \wedge D)$  and  $C \vdash D \rightarrow (C \wedge D)$ .

Remember that discharging of assumptions is a permission, not a requirement: we don't have to discharge the assumptions. In particular, we can apply a rule even if the assumptions are not present in the derivation. For instance, the following is legal, even though there is no assumption  $A$  to be discharged:

$$1 \frac{B}{A \rightarrow B} \rightarrow\text{Intro}$$

## 11.5 Examples of Derivations

**Example 11.4.** Let's give a derivation of the sentence  $(A \wedge B) \rightarrow A$ .

We begin by writing the desired conclusion at the bottom of the derivation.

$$\frac{}{(A \wedge B) \rightarrow A}$$

Next, we need to figure out what kind of inference could result in a sentence of this form. The main operator of the conclusion is  $\rightarrow$ , so we'll try to arrive at the conclusion using the  $\rightarrow$ Intro rule. It is best to write down the assumptions involved and label the inference rules as you progress, so it is easy to see whether all assumptions have been discharged at the end of the proof.

$$1 \frac{\begin{array}{c} [A \wedge B]^1 \\ \vdots \\ \vdots \\ A \end{array}}{(A \wedge B) \rightarrow A} \rightarrow\text{Intro}$$

We now need to fill in the steps from the assumption  $A \wedge B$  to  $A$ . Since we only have one connective to deal with,  $\wedge$ , we must use the  $\wedge$  elim rule. This gives us the following proof:



$$1 \frac{\frac{[A \wedge B]^1}{A} \wedge\text{Elim}}{(A \wedge B) \rightarrow A} \rightarrow\text{Intro}$$

We now have a correct derivation of  $(A \wedge B) \rightarrow A$ .

**Example 11.5.** Now let's give a derivation of  $(\neg A \vee B) \rightarrow (A \rightarrow B)$ .

We begin by writing the desired conclusion at the bottom of the derivation.

$$\overline{(\neg A \vee B) \rightarrow (A \rightarrow B)}$$

To find a logical rule that could give us this conclusion, we look at the logical connectives in the conclusion:  $\neg$ ,  $\vee$ , and  $\rightarrow$ . We only care at the moment about the first occurrence of  $\rightarrow$  because it is the main operator of the sentence in the end-sequent, while  $\neg$ ,  $\vee$  and the second occurrence of  $\rightarrow$  are inside the scope of another connective, so we will take care of those later. We therefore start with the  $\rightarrow$ Intro rule. A correct application must look like this:

$$1 \frac{\begin{array}{c} [\neg A \vee B]^1 \\ \vdots \\ \vdots \\ A \rightarrow B \end{array}}{(\neg A \vee B) \rightarrow (A \rightarrow B)} \rightarrow\text{Intro}$$

This leaves us with two possibilities to continue. Either we can keep working from the bottom up and look for another application of the  $\rightarrow$ Intro rule, or we can work from the top down and apply a  $\vee$ Elim rule. Let us apply the latter. We will use the assumption  $\neg A \vee B$  as the leftmost premise of  $\vee$ Elim. For a valid application of  $\vee$ Elim, the other two premises must be identical to the conclusion  $A \rightarrow B$ , but each may be derived in turn from another assumption, namely one of the two disjuncts of  $\neg A \vee B$ . So our derivation will look like this:

$$\begin{array}{c}
 \begin{array}{c} [\neg A]^2 \\ \vdots \\ A \rightarrow B \end{array} \quad \begin{array}{c} [B]^2 \\ \vdots \\ A \rightarrow B \end{array} \\
 \hline
 2 \frac{[\neg A \vee B]^1 \quad A \rightarrow B \quad A \rightarrow B}{A \rightarrow B} \vee\text{Elim} \\
 \hline
 1 \frac{A \rightarrow B}{(\neg A \vee B) \rightarrow (A \rightarrow B)} \rightarrow\text{Intro}
 \end{array}$$

In each of the two branches on the right, we want to derive  $A \rightarrow B$ , which is best done using  $\rightarrow$ Intro.

$$\begin{array}{c}
 \begin{array}{c} [\neg A]^2, [A]^3 \\ \vdots \\ B \end{array} \quad \begin{array}{c} [B]^2, [A]^4 \\ \vdots \\ B \end{array} \\
 \hline
 2 \frac{[\neg A \vee B]^1 \quad 3 \frac{B}{A \rightarrow B} \rightarrow\text{Intro} \quad 4 \frac{B}{A \rightarrow B} \rightarrow\text{Intro}}{A \rightarrow B} \vee\text{Elim} \\
 \hline
 1 \frac{A \rightarrow B}{(\neg A \vee B) \rightarrow (A \rightarrow B)} \rightarrow\text{Intro}
 \end{array}$$

For the two missing parts of the derivation, we need derivations of  $B$  from  $\neg A$  and  $A$  in the middle, and from  $A$  and  $B$  on the left. Let's take the former first.  $\neg A$  and  $A$  are the two premises of  $\neg$ Elim:

$$\begin{array}{c}
 \frac{[\neg A]^2 \quad [A]^3}{\perp} \neg\text{Elim} \\
 \vdots \\
 B
 \end{array}$$

By using  $\perp_I$ , we can obtain  $B$  as a conclusion and complete the branch.

$$\begin{array}{c}
 \begin{array}{c} [\neg A]^2 \quad [A]^3 \\ \vdots \\ \perp \end{array} \quad \begin{array}{c} [B]^2, [A]^4 \\ \vdots \\ B \end{array} \\
 \hline
 2 \frac{[\neg A \vee B]^1 \quad 3 \frac{\perp}{A \rightarrow B} \perp_I \quad 4 \frac{B}{A \rightarrow B} \rightarrow\text{Intro}}{A \rightarrow B} \vee\text{Elim} \\
 \hline
 1 \frac{A \rightarrow B}{(\neg A \vee B) \rightarrow (A \rightarrow B)} \rightarrow\text{Intro}
 \end{array}$$

Let's now look at the rightmost branch. Here it's important to realize that the definition of derivation *allows assumptions to be discharged* but *does not require* them to be. In other words, if we can derive  $B$  from one of the assumptions  $A$  and  $B$  without using the other, that's ok. And to derive  $B$  from  $B$  is trivial:  $B$  by itself is such a derivation, and no inferences are needed. So we can simply delete the assumption  $A$ .

$$\begin{array}{c}
 \frac{[\neg A]^2 \quad [A]^3}{\perp} \neg\text{Elim} \\
 \frac{\perp}{B} \perp\text{I} \\
 \frac{B}{A \rightarrow B} \rightarrow\text{Intro} \\
 \frac{[B]^2}{A \rightarrow B} \rightarrow\text{Intro} \\
 \frac{[\neg A \vee B]^1}{A \rightarrow B} \vee\text{Elim} \\
 \frac{A \rightarrow B}{(\neg A \vee B) \rightarrow (A \rightarrow B)} \rightarrow\text{Intro}
 \end{array}$$

Note that in the finished derivation, the rightmost  $\rightarrow$ Intro inference does not actually discharge any assumptions.

**Example 11.6.** So far we have not needed the  $\perp_C$  rule. It is special in that it allows us to discharge an assumption that isn't a sub-formula of the conclusion of the rule. It is closely related to the  $\perp_I$  rule. In fact, the  $\perp_I$  rule is a special case of the  $\perp_C$  rule—there is a logic called “intuitionistic logic” in which only  $\perp_I$  is allowed. The  $\perp_C$  rule is a last resort when nothing else works. For instance, suppose we want to derive  $A \vee \neg A$ . Our usual strategy would be to attempt to derive  $A \vee \neg A$  using  $\vee$ Intro. But this would require us to derive either  $A$  or  $\neg A$  from no assumptions, and this can't be done.  $\perp_C$  to the rescue!

$$\begin{array}{c}
 [\neg(A \vee \neg A)]^1 \\
 \vdots \\
 \perp \\
 \frac{\perp}{A \vee \neg A} \perp_C
 \end{array}$$

Now we're looking for a derivation of  $\perp$  from  $\neg(A \vee \neg A)$ . Since  $\perp$  is the conclusion of  $\neg$ Elim we might try that:

$$\frac{\begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ \neg A \end{array} \quad \begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ A \end{array}}{1 \frac{\perp}{A \vee \neg A} \perp_C} \neg\text{Elim}$$

Our strategy for finding a derivation of  $\neg A$  calls for an application of  $\neg$ -Intro:

$$\frac{\begin{array}{c} [\neg(A \vee \neg A)]^1, [A]^2 \\ \vdots \\ 2 \frac{\perp}{\neg A} \neg\text{Intro} \end{array} \quad \begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ A \end{array}}{1 \frac{\perp}{A \vee \neg A} \perp_C} \neg\text{Elim}$$

Here, we can get  $\perp$  easily by applying  $\neg$ -Elim to the assumption  $\neg(A \vee \neg A)$  and  $A \vee \neg A$  which follows from our new assumption  $A$  by  $\vee$ -Intro:

$$\frac{\begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ 2 \frac{\perp}{\neg A} \neg\text{Intro} \end{array} \quad \frac{\begin{array}{c} [A]^2 \\ A \vee \neg A \end{array} \vee\text{Intro}}{\neg\text{Elim}} \quad \begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ A \end{array}}{1 \frac{\perp}{A \vee \neg A} \perp_C} \neg\text{Elim}$$

On the right side we use the same strategy, except we get  $A$  by  $\perp_C$ :

$$\frac{\begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ 2 \frac{\perp}{\neg A} \neg\text{Intro} \end{array} \quad \frac{\begin{array}{c} [A]^2 \\ A \vee \neg A \end{array} \vee\text{Intro}}{\neg\text{Elim}} \quad \frac{\begin{array}{c} [\neg(A \vee \neg A)]^1 \\ \vdots \\ 3 \frac{\perp}{A} \perp_C \end{array} \quad \frac{[\neg A]^3}{A \vee \neg A} \vee\text{Intro}}{\neg\text{Elim}}}{1 \frac{\perp}{A \vee \neg A} \perp_C}$$

## 11.6 Derivations with Quantifiers

**Example 11.7.** When dealing with quantifiers, we have to make sure not to violate the eigenvariable condition, and sometimes this requires us to play around with the order of carrying out certain inferences. In general, it helps to try and take care of rules subject to the eigenvariable condition first (they will be lower down in the finished proof).

Let's see how we'd give a derivation of the formula  $\exists x \neg A(x) \rightarrow \neg \forall x A(x)$ . Starting as usual, we write

$$\overline{\exists x \neg A(x) \rightarrow \neg \forall x A(x)}$$

We start by writing down what it would take to justify that last step using the  $\rightarrow$ Intro rule.

$$\begin{array}{c} [\exists x \neg A(x)]^1 \\ \vdots \\ \vdots \\ \neg \forall x A(x) \\ 1 \frac{\quad}{\exists x \neg A(x) \rightarrow \neg \forall x A(x)} \rightarrow \text{Intro} \end{array}$$

Since there is no obvious rule to apply to  $\neg \forall x A(x)$ , we will proceed by setting up the derivation so we can use the  $\exists$ Elim rule. Here we must pay attention to the eigenvariable condition, and choose a constant that does not appear in  $\exists x A(x)$  or any assumptions that it depends on. (Since no constant symbols appear, however, any choice will do fine.)

$$\begin{array}{c} [\neg A(a)]^2 \\ \vdots \\ \vdots \\ [\exists x \neg A(x)]^1 \quad \neg \forall x A(x) \\ 2 \frac{\quad}{\neg \forall x A(x)} \exists \text{Elim} \\ 1 \frac{\quad}{\exists x \neg A(x) \rightarrow \neg \forall x A(x)} \rightarrow \text{Intro} \end{array}$$

In order to derive  $\neg \forall x A(x)$ , we will attempt to use the  $\neg$ Intro rule: this requires that we derive a contradiction, possibly using

$\forall x A(x)$  as an additional assumption. Of course, this contradiction may involve the assumption  $\neg A(a)$  which will be discharged by the  $\exists$ Elim inference. We can set it up as follows:

$$\begin{array}{c}
 [\neg A(a)]^2, [\forall x A(x)]^3 \\
 \vdots \\
 \perp \\
 \frac{3 \quad \frac{\neg \forall x A(x)}{\neg \forall x A(x)} \neg\text{Intro}}{\exists x \neg A(x)} \exists\text{Elim} \\
 \frac{2 \quad \frac{[\exists x \neg A(x)]^1}{\neg \forall x A(x)}}{\exists x \neg A(x) \rightarrow \neg \forall x A(x)} \rightarrow\text{Intro}
 \end{array}$$

It looks like we are close to getting a contradiction. The easiest rule to apply is the  $\forall$ Elim, which has no eigenvariable conditions. Since we can use any term we want to replace the universally quantified  $x$ , it makes the most sense to continue using  $a$  so we can reach a contradiction.

$$\begin{array}{c}
 [\forall x A(x)]^3 \\
 \frac{[\neg A(a)]^2 \quad \frac{A(a)}{A(a)} \forall\text{Elim}}{\neg A(a)} \neg\text{Elim} \\
 \frac{3 \quad \frac{\perp}{\neg \forall x A(x)} \neg\text{Intro}}{\exists x \neg A(x)} \exists\text{Elim} \\
 \frac{2 \quad \frac{[\exists x \neg A(x)]^1}{\neg \forall x A(x)}}{\exists x \neg A(x) \rightarrow \neg \forall x A(x)} \rightarrow\text{Intro}
 \end{array}$$

It is important, especially when dealing with quantifiers, to double check at this point that the eigenvariable condition has not been violated. Since the only rule we applied that is subject to the eigenvariable condition was  $\exists$ Elim, and the eigenvariable  $a$  does not occur in any assumptions it depends on, this is a correct derivation.

**Example 11.8.** Sometimes we may derive a formula from other formulas. In these cases, we may have undischarged assumptions. It is important to keep track of our assumptions as well as the end goal.

Let's see how we'd give a derivation of the formula  $\exists x C(x, b)$  from the assumptions  $\exists x (A(x) \wedge B(x))$  and  $\forall x (B(x) \rightarrow C(x, b))$ . Starting as usual, we write the conclusion at the bottom.

$$\overline{\exists x C(x, b)}$$

We have two premises to work with. To use the first, i.e., try to find a derivation of  $\exists x C(x, b)$  from  $\exists x (A(x) \wedge B(x))$  we would use the  $\exists$ Elim rule. Since it has an eigenvariable condition, we will apply that rule first. We get the following:

$$\begin{array}{c} [A(a) \wedge B(a)]^1 \\ \vdots \\ \exists x C(x, b) \end{array} \quad \begin{array}{c} \exists x (A(x) \wedge B(x)) \\ \exists x C(x, b) \end{array} \quad \exists\text{Elim}$$

$$1 \frac{\quad}{\exists x C(x, b)}$$

The two assumptions we are working with share  $B$ . It may be useful at this point to apply  $\wedge$ Elim to separate out  $B(a)$ .

$$\begin{array}{c} [A(a) \wedge B(a)]^1 \\ \wedge\text{Elim} \\ B(a) \\ \vdots \\ \exists x C(x, b) \end{array} \quad \begin{array}{c} \exists x (A(x) \wedge B(x)) \\ \exists x C(x, b) \end{array} \quad \exists\text{Elim}$$

$$1 \frac{\quad}{\exists x C(x, b)}$$

The second assumption we have to work with is  $\forall x (B(x) \rightarrow C(x, b))$ . Since there is no eigenvariable condition we can instantiate  $x$  with the constant symbol  $a$  using  $\forall$ Elim to get  $B(a) \rightarrow C(a, b)$ . We now have both  $B(a) \rightarrow C(a, b)$  and  $B(a)$ . Our next move should be a straightforward application of the  $\rightarrow$ Elim rule.

$$\begin{array}{c}
 \frac{\frac{\forall x (B(x) \rightarrow C(x, b))}{B(a) \rightarrow C(a, b)} \forall\text{Elim} \quad \frac{[A(a) \wedge B(a)]^1}{B(a)} \wedge\text{Elim}}{C(a, b)} \rightarrow\text{Elim} \\
 \vdots \\
 1 \frac{\exists x (A(x) \wedge B(x)) \quad \exists x C(x, b)}{\exists x C(x, b)} \exists\text{Elim}
 \end{array}$$

We are so close! One application of  $\exists\text{Intro}$  and we have reached our goal.

$$\begin{array}{c}
 \frac{\frac{\forall x (B(x) \rightarrow C(x, b))}{B(a) \rightarrow C(a, b)} \forall\text{Elim} \quad \frac{[A(a) \wedge B(a)]^1}{B(a)} \wedge\text{Elim}}{C(a, b)} \rightarrow\text{Elim} \\
 \frac{\exists x (A(x) \wedge B(x)) \quad \frac{C(a, b)}{\exists x C(x, b)} \exists\text{Intro}}{\exists x C(x, b)} \exists\text{Elim} \\
 1 \frac{\exists x (A(x) \wedge B(x))}{\exists x C(x, b)} \exists\text{Elim}
 \end{array}$$

Since we ensured at each step that the eigenvariable conditions were not violated, we can be confident that this is a correct derivation.

**Example 11.9.** Give a derivation of the formula  $\neg\forall x A(x)$  from the assumptions  $\forall x A(x) \rightarrow \exists y B(y)$  and  $\neg\exists y B(y)$ . Starting as usual, we write the target formula at the bottom.

$$\overline{\neg\forall x A(x)}$$

The last line of the derivation is a negation, so let's try using  $\neg\text{Intro}$ . This will require that we figure out how to derive a contradiction.

$$\begin{array}{c}
 [\forall x A(x)]^1 \\
 \vdots \\
 \perp \\
 1 \frac{\perp}{\neg\forall x A(x)} \neg\text{Intro}
 \end{array}$$



So far so good. We can use  $\forall$ Elim but it's not obvious if that will help us get to our goal. Instead, let's use one of our assumptions.  $\forall x A(x) \rightarrow \exists y B(y)$  together with  $\forall x A(x)$  will allow us to use the  $\rightarrow$ Elim rule.

$$\frac{\forall x A(x) \rightarrow \exists y B(y) \quad [\forall x A(x)]^1}{\exists y B(y)} \rightarrow\text{Elim}$$

$$\vdots$$

$$1 \frac{\perp}{\neg \forall x A(x)} \neg\text{Intro}$$

We now have one final assumption to work with, and it looks like this will help us reach a contradiction by using  $\neg$ Elim.

$$\frac{\neg \exists y B(y) \quad \frac{\forall x A(x) \rightarrow \exists y B(y) \quad [\forall x A(x)]^1}{\exists y B(y)} \rightarrow\text{Elim}}{1 \frac{\perp}{\neg \forall x A(x)} \neg\text{Intro}} \neg\text{Elim}$$

## 11.7 Proof-Theoretic Notions

Just as we've defined a number of important semantic notions (validity, entailment, satisfiability), we now define corresponding *proof-theoretic notions*. These are not defined by appeal to satisfaction of sentences in structures, but by appeal to the derivability or non-derivability of certain sentences from others. It was an important discovery that these notions coincide. That they do is the content of the *soundness* and *completeness theorems*.

**Definition 11.10 (Theorems).** A sentence  $A$  is a *theorem* if there is a derivation of  $A$  in natural deduction in which all assumptions are discharged. We write  $\vdash A$  if  $A$  is a theorem and  $\not\vdash A$  if it is not.

**Definition 11.11 (Derivability).** A sentence  $A$  is *derivable* from a set of sentences  $\Gamma$ ,  $\Gamma \vdash A$ , if there is a derivation with conclusion  $A$  and in which every assumption is either discharged or is in  $\Gamma$ . If  $A$  is not derivable from  $\Gamma$  we write  $\Gamma \not\vdash A$ .

**Definition 11.12 (Consistency).** A set of sentences  $\Gamma$  is *inconsistent* iff  $\Gamma \vdash \perp$ . If  $\Gamma$  is not inconsistent, i.e., if  $\Gamma \not\vdash \perp$ , we say it is *consistent*.

**Proposition 11.13 (Reflexivity).** If  $A \in \Gamma$ , then  $\Gamma \vdash A$ .

*Proof.* The assumption  $A$  by itself is a derivation of  $A$  where every undischarged assumption (i.e.,  $A$ ) is in  $\Gamma$ .  $\square$

**Proposition 11.14 (Monotonicity).** If  $\Gamma \subseteq \Delta$  and  $\Gamma \vdash A$ , then  $\Delta \vdash A$ .

*Proof.* Any derivation of  $A$  from  $\Gamma$  is also a derivation of  $A$  from  $\Delta$ .  $\square$

**Proposition 11.15 (Transitivity).** If  $\Gamma \vdash A$  and  $\{A\} \cup \Delta \vdash B$ , then  $\Gamma \cup \Delta \vdash B$ .

*Proof.* If  $\Gamma \vdash A$ , there is a derivation  $\delta_0$  of  $A$  with all undischarged assumptions in  $\Gamma$ . If  $\{A\} \cup \Delta \vdash B$ , then there is a derivation  $\delta_1$  of  $B$  with all undischarged assumptions in  $\{A\} \cup \Delta$ . Now consider:

$$\begin{array}{c}
 \Delta, [A]^1 \\
 \vdots \\
 \vdots \delta_1 \\
 \vdots \\
 B \\
 \hline
 1 \frac{A \rightarrow B}{A \rightarrow B} \rightarrow \text{Intro}
 \end{array}
 \quad
 \begin{array}{c}
 \Gamma \\
 \vdots \\
 \vdots \delta_0 \\
 \vdots \\
 A \\
 \hline
 \rightarrow \text{Elim}
 \end{array}
 \quad
 \frac{\quad}{B}$$

The undischarged assumptions are now all among  $\Gamma \cup \Delta$ , so this shows  $\Gamma \cup \Delta \vdash B$ .  $\square$

When  $\Gamma = \{A_1, A_2, \dots, A_k\}$  is a finite set we may use the simplified notation  $A_1, A_2, \dots, A_k \vdash B$  for  $\Gamma \vdash B$ , in particular  $A \vdash B$  means that  $\{A\} \vdash B$ .

Note that if  $\Gamma \vdash A$  and  $A \vdash B$ , then  $\Gamma \vdash B$ . It follows also that if  $A_1, \dots, A_n \vdash B$  and  $\Gamma \vdash A_i$  for each  $i$ , then  $\Gamma \vdash B$ .

**Proposition 11.16.** *The following are equivalent.*

1.  $\Gamma$  is inconsistent.
2.  $\Gamma \vdash A$  for every sentence  $A$ .
3.  $\Gamma \vdash A$  and  $\Gamma \vdash \neg A$  for some sentence  $A$ .

*Proof.* Exercise.  $\square$

**Proposition 11.17 (Compactness).** 1. *If  $\Gamma \vdash A$  then there is a finite subset  $\Gamma_0 \subseteq \Gamma$  such that  $\Gamma_0 \vdash A$ .*

2. *If every finite subset of  $\Gamma$  is consistent, then  $\Gamma$  is consistent.*

*Proof.* 1. If  $\Gamma \vdash A$ , then there is a derivation  $\delta$  of  $A$  from  $\Gamma$ . Let  $\Gamma_0$  be the set of undischarged assumptions of  $\delta$ . Since any derivation is finite,  $\Gamma_0$  can only contain finitely many sentences. So,  $\delta$  is a derivation of  $A$  from a finite  $\Gamma_0 \subseteq \Gamma$ .

2. This is the contrapositive of (1) for the special case  $A \equiv \perp$ .

$\square$

## 11.8 Derivability and Consistency

We will now establish a number of properties of the derivability relation. They are independently interesting, but each will play a role in the proof of the completeness theorem.

**Proposition 11.18.** *If  $\Gamma \vdash A$  and  $\Gamma \cup \{A\}$  is inconsistent, then  $\Gamma$  is inconsistent.*

*Proof.* Let the derivation of  $A$  from  $\Gamma$  be  $\delta_1$  and the derivation of  $\perp$  from  $\Gamma \cup \{A\}$  be  $\delta_2$ . We can then derive:

$$\frac{\begin{array}{c} \Gamma, [A]^1 \\ \vdots \\ \delta_2 \\ \vdots \\ \perp \\ \hline 1 \frac{\perp}{\neg A} \neg\text{Intro} \end{array}}{\perp} \quad \frac{\begin{array}{c} \Gamma \\ \vdots \\ \delta_1 \\ \vdots \\ A \\ \hline \neg\text{Elim} \end{array}}{A} \neg\text{Elim}$$

In the new derivation, the assumption  $A$  is discharged, so it is a derivation from  $\Gamma$ .  $\square$

**Proposition 11.19.**  *$\Gamma \vdash A$  iff  $\Gamma \cup \{\neg A\}$  is inconsistent.*

*Proof.* First suppose  $\Gamma \vdash A$ , i.e., there is a derivation  $\delta_0$  of  $A$  from undischarged assumptions  $\Gamma$ . We obtain a derivation of  $\perp$  from  $\Gamma \cup \{\neg A\}$  as follows:

$$\frac{\begin{array}{c} \Gamma \\ \vdots \\ \delta_0 \\ \vdots \\ A \\ \hline \neg A \quad \neg\text{Elim} \end{array}}{\perp} \neg\text{Elim}$$

Now assume  $\Gamma \cup \{\neg A\}$  is inconsistent, and let  $\delta_1$  be the corresponding derivation of  $\perp$  from undischarged assumptions in  $\Gamma \cup \{\neg A\}$ . We obtain a derivation of  $A$  from  $\Gamma$  alone by using  $\perp_C$ :

$$\frac{\begin{array}{c} \Gamma, [\neg A]^1 \\ \vdots \\ \delta_1 \\ \vdots \\ \perp \\ \hline 1 \frac{\perp}{A} \perp_C \end{array}}{A} \perp_C$$

$\square$

**Proposition 11.20.** *If  $\Gamma \vdash A$  and  $\neg A \in \Gamma$ , then  $\Gamma$  is inconsistent.*

*Proof.* Suppose  $\Gamma \vdash A$  and  $\neg A \in \Gamma$ . Then there is a derivation  $\delta$  of  $A$  from  $\Gamma$ . Consider this simple application of the  $\neg$ -Elim rule:

$$\frac{\neg A \quad \begin{array}{c} \Gamma \\ \vdots \\ \delta \\ \vdots \\ A \end{array}}{\perp} \neg\text{Elim}$$

Since  $\neg A \in \Gamma$ , all undischarged assumptions are in  $\Gamma$ , this shows that  $\Gamma \vdash \perp$ .  $\square$

**Proposition 11.21.** *If  $\Gamma \cup \{A\}$  and  $\Gamma \cup \{\neg A\}$  are both inconsistent, then  $\Gamma$  is inconsistent.*

*Proof.* There are derivations  $\delta_1$  and  $\delta_2$  of  $\perp$  from  $\Gamma \cup \{A\}$  and  $\perp$  from  $\Gamma \cup \{\neg A\}$ , respectively. We can then derive

$$\frac{\begin{array}{c} \Gamma, [\neg A]^2 \\ \vdots \\ \delta_2 \\ \vdots \\ \perp \end{array} \quad \begin{array}{c} \Gamma, [A]^1 \\ \vdots \\ \delta_1 \\ \vdots \\ \perp \end{array}}{\perp} \frac{\begin{array}{c} 2 \frac{\perp}{\neg\neg A} \neg\text{Intro} \quad 1 \frac{\perp}{\neg A} \neg\text{Intro} \\ \neg\text{Elim} \end{array}}$$

Since the assumptions  $A$  and  $\neg A$  are discharged, this is a derivation of  $\perp$  from  $\Gamma$  alone. Hence  $\Gamma$  is inconsistent.  $\square$

## 11.9 Derivability and the Propositional Connectives

We establish that the derivability relation  $\vdash$  of natural deduction is strong enough to establish some basic facts involving the propositional connectives, such as that  $A \wedge B \vdash A$  and  $A, A \rightarrow B \vdash B$  (modus ponens). These facts are needed for the proof of the completeness theorem.

**Proposition 11.22.** 1. Both  $A \wedge B \vdash A$  and  $A \wedge B \vdash B$

2.  $A, B \vdash A \wedge B$ .

*Proof.* 1. We can derive both

$$\frac{A \wedge B}{A} \wedge\text{Elim} \qquad \frac{A \wedge B}{B} \wedge\text{Elim}$$

2. We can derive:

$$\frac{A \quad B}{A \wedge B} \wedge\text{Intro} \quad \square$$

**Proposition 11.23.** 1.  $A \vee B, \neg A, \neg B$  is inconsistent.

2. Both  $A \vdash A \vee B$  and  $B \vdash A \vee B$ .

*Proof.* 1. Consider the following derivation:

$$1 \frac{A \vee B \quad \frac{\frac{\neg A \quad [A]^1}{\perp} \neg\text{Elim} \quad \frac{\neg B \quad [B]^1}{\perp} \neg\text{Elim}}{\perp} \vee\text{Elim}}{\perp}$$

This is a derivation of  $\perp$  from undischarged assumptions  $A \vee B$ ,  $\neg A$ , and  $\neg B$ .

2. We can derive both

$$\frac{A}{A \vee B} \vee\text{Intro} \qquad \frac{B}{A \vee B} \vee\text{Intro} \quad \square$$

**Proposition 11.24.** 1.  $A, A \rightarrow B \vdash B$ .

2. Both  $\neg A \vdash A \rightarrow B$  and  $B \vdash A \rightarrow B$ .

*Proof.* 1. We can derive:

$$\frac{A \rightarrow B \quad A}{B} \rightarrow\text{Elim}$$

2. This is shown by the following two derivations:

$$\frac{\neg A \quad [A]^1}{\frac{\perp}{B} \perp_I} \neg\text{Elim} \qquad \frac{B}{A \rightarrow B} \rightarrow\text{Intro}$$

$$1 \frac{\perp}{A \rightarrow B} \rightarrow\text{Intro}$$

Note that  $\rightarrow\text{Intro}$  may, but does not have to, discharge the assumption  $A$ .  $\square$

## 11.10 Derivability and the Quantifiers

The completeness theorem also requires that the natural deduction rules yield the facts about  $\vdash$  established in this section.

**Theorem 11.25.** *If  $c$  is a constant not occurring in  $\Gamma$  or  $A(x)$  and  $\Gamma \vdash A(c)$ , then  $\Gamma \vdash \forall x A(x)$ .*

*Proof.* Let  $\delta$  be a derivation of  $A(c)$  from  $\Gamma$ . By adding a  $\forall\text{Intro}$  inference, we obtain a derivation of  $\forall x A(x)$ . Since  $c$  does not occur in  $\Gamma$  or  $A(x)$ , the eigenvariable condition is satisfied.  $\square$

**Proposition 11.26.** 1.  $A(t) \vdash \exists x A(x)$ .

2.  $\forall x A(x) \vdash A(t)$ .

*Proof.* 1. The following is a derivation of  $\exists x A(x)$  from  $A(t)$ :

$$\frac{A(t)}{\exists x A(x)} \exists\text{Intro}$$

2. The following is a derivation of  $A(t)$  from  $\forall x A(x)$ :

$$\frac{\forall x A(x)}{A(t)} \forall\text{Elim}$$

$\square$

## 11.11 Soundness

A derivation system, such as natural deduction, is *sound* if it cannot derive things that do not actually follow. Soundness is thus a kind of guaranteed safety property for derivation systems. Depending on which proof theoretic property is in question, we would like to know for instance, that

1. every derivable sentence is valid;
2. if a sentence is derivable from some others, it is also a consequence of them;
3. if a set of sentences is inconsistent, it is unsatisfiable.

These are important properties of a derivation system. If any of them do not hold, the derivation system is deficient—it would derive too much. Consequently, establishing the soundness of a derivation system is of the utmost importance.

**Theorem 11.27 (Soundness).** *If  $A$  is derivable from the undischarged assumptions  $\Gamma$ , then  $\Gamma \vDash A$ .*

*Proof.* Let  $\delta$  be a derivation of  $A$ . We proceed by induction on the number of inferences in  $\delta$ .

For the induction basis we show the claim if the number of inferences is 0. In this case,  $\delta$  consists only of a single sentence  $A$ , i.e., an assumption. That assumption is undischarged, since assumptions can only be discharged by inferences, and there are no inferences. So, any structure  $M$  that satisfies all of the undischarged assumptions of the proof also satisfies  $A$ .

Now for the inductive step. Suppose that  $\delta$  contains  $n$  inferences. The premise(s) of the lowermost inference are derived using sub-derivations, each of which contains fewer than  $n$  inferences. We assume the induction hypothesis: The premises of the lowermost inference follow from the undischarged assumptions of the sub-derivations ending in those premises. We have to show



that the conclusion  $A$  follows from the undischarged assumptions of the entire proof.

We distinguish cases according to the type of the lowermost inference. First, we consider the possible inferences with only one premise.

1. Suppose that the last inference is  $\neg$ -Intro: The derivation has the form

$$\begin{array}{c} \Gamma, [A]^n \\ \vdots \\ \delta_1 \\ \vdots \\ \perp \\ \hline \neg A \quad \neg\text{-Intro} \end{array}$$

By inductive hypothesis,  $\perp$  follows from the undischarged assumptions  $\Gamma \cup \{A\}$  of  $\delta_1$ . Consider a structure  $M$ . We need to show that, if  $M \models \Gamma$ , then  $M \models \neg A$ . Suppose for reductio that  $M \models \Gamma$ , but  $M \not\models \neg A$ , i.e.,  $M \models A$ . This would mean that  $M \models \Gamma \cup \{A\}$ . This is contrary to our inductive hypothesis. So,  $M \models \neg A$ .

2. The last inference is  $\wedge$ Elim: There are two variants:  $A$  or  $B$  may be inferred from the premise  $A \wedge B$ . Consider the first case. The derivation  $\delta$  looks like this:

$$\begin{array}{c} \Gamma \\ \vdots \\ \delta_1 \\ \vdots \\ A \wedge B \\ \hline A \quad \wedge\text{Elim} \end{array}$$

By inductive hypothesis,  $A \wedge B$  follows from the undischarged assumptions  $\Gamma$  of  $\delta_1$ . Consider a structure  $M$ . We need to show that, if  $M \models \Gamma$ , then  $M \models A$ . Suppose  $M \models \Gamma$ . By our inductive hypothesis ( $\Gamma \models A \wedge B$ ), we know that  $M \models A \wedge B$ . By definition,  $M \models A \wedge B$  iff  $M \models A$  and  $M \models B$ .

(The case where  $B$  is inferred from  $A \wedge B$  is handled similarly.)

3. The last inference is  $\vee$ Intro: There are two variants:  $A \vee B$  may be inferred from the premise  $A$  or the premise  $B$ . Consider the first case. The derivation has the form

$$\frac{\begin{array}{c} \Gamma \\ \vdots \\ \delta_1 \\ \vdots \\ A \end{array}}{A \vee B} \vee\text{Intro}$$

By inductive hypothesis,  $A$  follows from the undischarged assumptions  $\Gamma$  of  $\delta_1$ . Consider a structure  $M$ . We need to show that, if  $M \vDash \Gamma$ , then  $M \vDash A \vee B$ . Suppose  $M \vDash \Gamma$ ; then  $M \vDash A$  since  $\Gamma \vDash A$  (the inductive hypothesis). So it must also be the case that  $M \vDash A \vee B$ . (The case where  $A \vee B$  is inferred from  $B$  is handled similarly.)

4. The last inference is  $\rightarrow$ Intro:  $A \rightarrow B$  is inferred from a subproof with assumption  $A$  and conclusion  $B$ , i.e.,

$${}^n \frac{\begin{array}{c} \Gamma, [A]^n \\ \vdots \\ \delta_1 \\ \vdots \\ B \end{array}}{A \rightarrow B} \rightarrow\text{Intro}$$

By inductive hypothesis,  $B$  follows from the undischarged assumptions of  $\delta_1$ , i.e.,  $\Gamma \cup \{A\} \vDash B$ . Consider a structure  $M$ . The undischarged assumptions of  $\delta$  are just  $\Gamma$ , since  $A$  is discharged at the last inference. So we need to show that  $\Gamma \vDash A \rightarrow B$ . For reductio, suppose that for some structure  $M$ ,  $M \vDash \Gamma$  but  $M \not\vDash A \rightarrow B$ . So,  $M \vDash A$  and  $M \not\vDash B$ . But by hypothesis,  $B$  is a consequence of  $\Gamma \cup \{A\}$ , i.e.,  $M \vDash B$ , which is a contradiction. So,  $\Gamma \vDash A \rightarrow B$ .

5. The last inference is  $\perp_I$ : Here,  $\delta$  ends in

$$\frac{\begin{array}{c} \Gamma \\ \vdots \\ \delta_1 \\ \vdots \\ \perp \\ A \end{array}}{\perp_I}$$

By induction hypothesis,  $\Gamma \vDash \perp$ . We have to show that  $\Gamma \vDash A$ . Suppose not; then for some  $M$  we have  $M \vDash \Gamma$  and  $M \not\vDash A$ . But we always have  $M \not\vDash \perp$ , so this would mean that  $\Gamma \not\vDash \perp$ , contrary to the induction hypothesis.

6. The last inference is  $\perp_C$ : Exercise.
7. The last inference is  $\forall$ Intro: Then  $\delta$  has the form

$$\frac{\begin{array}{c} \Gamma \\ \vdots \\ \delta_1 \\ \vdots \\ A(a) \end{array}}{\forall x A(x)} \forall\text{Intro}$$

The premise  $A(a)$  is a consequence of the undischarged assumptions  $\Gamma$  by induction hypothesis. Consider some structure,  $M$ , such that  $M \vDash \Gamma$ . We need to show that  $M \vDash \forall x A(x)$ . Since  $\forall x A(x)$  is a sentence, this means we have to show that for every variable assignment  $s$ ,  $M, s \vDash A(x)$  (**Proposition 7.18**). Since  $\Gamma$  consists entirely of sentences,  $M, s \vDash B$  for all  $B \in \Gamma$  by **Definition 7.11**. Let  $M'$  be like  $M$  except that  $a^{M'} = s(x)$ . Since  $a$  does not occur in  $\Gamma$ ,  $M' \vDash \Gamma$  by **Corollary 7.20**. Since  $\Gamma \vDash A(a)$ ,  $M' \vDash A(a)$ . Since  $A(a)$  is a sentence,  $M', s \vDash A(a)$  by **Proposition 7.17**.  $M', s \vDash A(x)$  iff  $M' \vDash A(a)$  by **Proposition 7.22** (recall that  $A(a)$  is just  $A(x)[a/x]$ ). So,  $M', s \vDash A(x)$ . Since  $a$  does not occur in  $A(x)$ , by **Proposition 7.19**,  $M, s \vDash A(x)$ . But  $s$  was an arbitrary variable assignment, so  $M \vDash \forall x A(x)$ .

8. The last inference is  $\exists$ Intro: Exercise.
9. The last inference is  $\forall$ Elim: Exercise.

Now let's consider the possible inferences with several premises:  $\forall$ Elim,  $\wedge$ Intro,  $\rightarrow$ Elim, and  $\exists$ Elim.

1. The last inference is  $\wedge$ Intro.  $A \wedge B$  is inferred from the premises  $A$  and  $B$  and  $\delta$  has the form

$$\frac{\begin{array}{c} \Gamma_1 \\ \vdots \\ \delta_1 \\ \vdots \\ A \end{array} \quad \begin{array}{c} \Gamma_2 \\ \vdots \\ \delta_2 \\ \vdots \\ B \end{array}}{A \wedge B} \wedge\text{Intro}$$

By induction hypothesis,  $A$  follows from the undischarged assumptions  $\Gamma_1$  of  $\delta_1$  and  $B$  follows from the undischarged assumptions  $\Gamma_2$  of  $\delta_2$ . The undischarged assumptions of  $\delta$  are  $\Gamma_1 \cup \Gamma_2$ , so we have to show that  $\Gamma_1 \cup \Gamma_2 \vDash A \wedge B$ . Consider a structure  $M$  with  $M \vDash \Gamma_1 \cup \Gamma_2$ . Since  $M \vDash \Gamma_1$ , it must be the case that  $M \vDash A$  as  $\Gamma_1 \vDash A$ , and since  $M \vDash \Gamma_2$ ,  $M \vDash B$  since  $\Gamma_2 \vDash B$ . Together,  $M \vDash A \wedge B$ .

2. The last inference is  $\forall$ Elim: Exercise.
3. The last inference is  $\rightarrow$ Elim.  $B$  is inferred from the premises  $A \rightarrow B$  and  $A$ . The derivation  $\delta$  looks like this:

$$\frac{\begin{array}{c} \Gamma_1 \\ \vdots \\ \delta_1 \\ \vdots \\ A \rightarrow B \end{array} \quad \begin{array}{c} \Gamma_2 \\ \vdots \\ \delta_2 \\ \vdots \\ A \end{array}}{B} \rightarrow\text{Elim}$$

By induction hypothesis,  $A \rightarrow B$  follows from the undischarged assumptions  $\Gamma_1$  of  $\delta_1$  and  $A$  follows from the undischarged assumptions  $\Gamma_2$  of  $\delta_2$ . Consider a structure  $M$ . We

need to show that, if  $M \vDash \Gamma_1 \cup \Gamma_2$ , then  $M \vDash B$ . Suppose  $M \vDash \Gamma_1 \cup \Gamma_2$ . Since  $\Gamma_1 \vDash A \rightarrow B$ ,  $M \vDash A \rightarrow B$ . Since  $\Gamma_2 \vDash A$ , we have  $M \vDash A$ . This means that  $M \vDash B$  (For if  $M \not\vDash B$ , since  $M \vDash A$ , we'd have  $M \not\vDash A \rightarrow B$ , contradicting  $M \vDash A \rightarrow B$ ).

4. The last inference is  $\neg$ Elim: Exercise.
5. The last inference is  $\exists$ Elim: Exercise. □

**Corollary 11.28.** *If  $\vdash A$ , then  $A$  is valid.*

**Corollary 11.29.** *If  $\Gamma$  is satisfiable, then it is consistent.*

*Proof.* We prove the contrapositive. Suppose that  $\Gamma$  is not consistent. Then  $\Gamma \vdash \perp$ , i.e., there is a derivation of  $\perp$  from undischarged assumptions in  $\Gamma$ . By **Theorem 11.27**, any structure  $M$  that satisfies  $\Gamma$  must satisfy  $\perp$ . Since  $M \not\vDash \perp$  for every structure  $M$ , no  $M$  can satisfy  $\Gamma$ , i.e.,  $\Gamma$  is not satisfiable. □

## 11.12 Derivations with Identity predicate

Derivations with identity predicate require additional inference rules.

$$\frac{}{t = t} =\text{Intro} \qquad \frac{t_1 = t_2 \quad A(t_1)}{A(t_2)} =\text{Elim}$$

$$\frac{t_1 = t_2 \quad A(t_2)}{A(t_1)} =\text{Elim}$$

In the above rules,  $t$ ,  $t_1$ , and  $t_2$  are closed terms. The  $=$ Intro rule allows us to derive any identity statement of the form  $t = t$  outright, from no assumptions.

**Example 11.30.** If  $s$  and  $t$  are closed terms, then  $A(s), s = t \vdash A(t)$ :

$$\frac{s = t \quad A(s)}{A(t)} =\text{Elim}$$

This may be familiar as the “principle of substitutability of identicals,” or Leibniz’ Law.

**Example 11.31.** We derive the sentence

$$\forall x \forall y ((A(x) \wedge A(y)) \rightarrow x = y)$$

from the sentence

$$\exists x \forall y (A(y) \rightarrow y = x)$$

We develop the derivation backwards:

$$\begin{array}{c} \exists x \forall y (A(y) \rightarrow y = x) \quad [A(a) \wedge A(b)]^1 \\ \vdots \\ a = b \\ \frac{1 \quad \frac{\frac{((A(a) \wedge A(b)) \rightarrow a = b)}{\forall y ((A(a) \wedge A(y)) \rightarrow a = y)} \rightarrow\text{Intro}}{\forall x \forall y ((A(x) \wedge A(y)) \rightarrow x = y)} \forall\text{Intro}}{\forall y ((A(a) \wedge A(y)) \rightarrow a = y)} \forall\text{Intro} \end{array}$$

We’ll now have to use the main assumption: since it is an existential formula, we use  $\exists\text{Elim}$  to derive the intermediary conclusion  $a = b$ .

$$\begin{array}{c} [\forall y (A(y) \rightarrow y = c)]^2 \\ [A(a) \wedge A(b)]^1 \\ \vdots \\ a = b \\ \frac{2 \quad \frac{\frac{\frac{\frac{\frac{\exists x \forall y (A(y) \rightarrow y = x)}{\forall y ((A(a) \wedge A(y)) \rightarrow a = y)} \forall\text{Intro}}{\forall x \forall y ((A(x) \wedge A(y)) \rightarrow x = y)} \forall\text{Intro}}{\frac{1 \quad \frac{\frac{((A(a) \wedge A(b)) \rightarrow a = b)}{\forall y ((A(a) \wedge A(y)) \rightarrow a = y)} \rightarrow\text{Intro}}{\forall x \forall y ((A(x) \wedge A(y)) \rightarrow x = y)} \forall\text{Intro}}{\frac{a = b}{\exists\text{Elim}}}}}}}} \end{array}$$

The sub-derivation on the top right is completed by using its assumptions to show that  $a = c$  and  $b = c$ . This requires two separate derivations. The derivation for  $a = c$  is as follows:

$$\frac{\frac{[\forall y (A(y) \rightarrow y = c)]^2}{A(a) \rightarrow a = c} \forall\text{Elim} \quad \frac{[A(a) \wedge A(b)]^1}{A(a)} \wedge\text{Elim}}{a = c} \rightarrow\text{Elim}$$

From  $a = c$  and  $b = c$  we derive  $a = b$  by =Elim.

### 11.13 Soundness with Identity predicate

**Proposition 11.32.** *Natural deduction with rules for = is sound.*

*Proof.* Any formula of the form  $t = t$  is valid, since for every structure  $M$ ,  $M \models t = t$ . (Note that we assume the term  $t$  to be closed, i.e., it contains no variables, so variable assignments are irrelevant).

Suppose the last inference in a derivation is =Elim, i.e., the derivation has the following form:

$$\frac{\begin{array}{c} \Gamma_1 \\ \vdots \\ \delta_1 \\ \vdots \\ t_1 = t_2 \end{array} \quad \begin{array}{c} \Gamma_2 \\ \vdots \\ \delta_2 \\ \vdots \\ A(t_1) \end{array}}{A(t_2)} =\text{Elim}$$

The premises  $t_1 = t_2$  and  $A(t_1)$  are derived from undischarged assumptions  $\Gamma_1$  and  $\Gamma_2$ , respectively. We want to show that  $A(t_2)$  follows from  $\Gamma_1 \cup \Gamma_2$ . Consider a structure  $M$  with  $M \models \Gamma_1 \cup \Gamma_2$ . By induction hypothesis,  $M \models A(t_1)$  and  $M \models t_1 = t_2$ . Therefore,  $\text{Val}^M(t_1) = \text{Val}^M(t_2)$ . Let  $s$  be any variable assignment, and  $m = \text{Val}^M(t_1) = \text{Val}^M(t_2)$ . By **Proposition 7.22**,  $M, s \models A(t_1)$  iff  $M, s[m/x] \models A(x)$  iff  $M, s \models A(t_2)$ . Since  $M \models A(t_1)$ , we have  $M \models A(t_2)$ .  $\square$

## Summary

**Proof systems** provide purely syntactic methods for characterizing consequence and compatibility between sentences. **Natural deduction** is one such proof system. A **derivation** in it consists of a tree formulas. The topmost formulas in a derivation are **assumptions**. All other formulas, for the derivation to be correct, must be correctly justified by one of a number of **inference rules**. These come in pairs; an introduction and an elimination rule for each connective and quantifier. For instance, if a formula  $A$  is justified by a  $\rightarrow$ Elim rule, the preceding formulas (the **premises**) must be  $B \rightarrow A$  and  $B$  (for some  $B$ ). Some inference rules also allow assumptions to be **discharged**. For instance, if  $A \rightarrow B$  is inferred from  $B$  using  $\rightarrow$ Intro, any occurrences of  $A$  as assumptions in the derivation leading to the premise  $B$  may be discharged, and is given a label that is also recorded at the inference.

If there is a derivation with end formula  $A$  and all assumptions are discharged, we say  $A$  is a **theorem** and write  $\vdash A$ . If all undischarged assumptions are in some set  $\Gamma$ , we say  $A$  is **derivable from  $\Gamma$**  and write  $\Gamma \vdash A$ . If  $\Gamma \vdash \perp$  we say  $\Gamma$  is **inconsistent**, otherwise **consistent**. These notions are interrelated, e.g.,  $\Gamma \vdash A$  iff  $\Gamma \cup \{\neg A\}$  is inconsistent. They are also related to the corresponding semantic notions, e.g., if  $\Gamma \vdash A$  then  $\Gamma \models A$ . This property of proof systems—what can be derived from  $\Gamma$  is guaranteed to be entailed by  $\Gamma$ —is called **soundness**. The **soundness theorem** is proved by induction on the length of derivations, showing that each individual inference preserves entailment of its conclusion from open assumptions provided its premises are entailed by their undischarged assumptions.

## Problems

**Problem 11.1.** Give derivations that show the following:

1.  $A \wedge (B \wedge C) \vdash (A \wedge B) \wedge C$ .



2.  $A \vee (B \vee C) \vdash (A \vee B) \vee C.$
3.  $A \rightarrow (B \rightarrow C) \vdash B \rightarrow (A \rightarrow C).$
4.  $A \vdash \neg\neg A.$

**Problem 11.2.** Give derivations that show the following:

1.  $(A \vee B) \rightarrow C \vdash A \rightarrow C.$
2.  $(A \rightarrow C) \wedge (B \rightarrow C) \vdash (A \vee B) \rightarrow C.$
3.  $\vdash \neg(A \wedge \neg A).$
4.  $B \rightarrow A \vdash \neg A \rightarrow \neg B.$
5.  $\vdash (A \rightarrow \neg A) \rightarrow \neg A.$
6.  $\vdash \neg(A \rightarrow B) \rightarrow \neg B.$
7.  $A \rightarrow C \vdash \neg(A \wedge \neg C).$
8.  $A \wedge \neg C \vdash \neg(A \rightarrow C).$
9.  $A \vee B, \neg B \vdash A.$
10.  $\neg A \vee \neg B \vdash \neg(A \wedge B).$
11.  $\vdash (\neg A \wedge \neg B) \rightarrow \neg(A \vee B).$
12.  $\vdash \neg(A \vee B) \rightarrow (\neg A \wedge \neg B).$

**Problem 11.3.** Give derivations that show the following:

1.  $\neg(A \rightarrow B) \vdash A.$
2.  $\neg(A \wedge B) \vdash \neg A \vee \neg B.$
3.  $A \rightarrow B \vdash \neg A \vee B.$
4.  $\vdash \neg\neg A \rightarrow A.$
5.  $A \rightarrow B, \neg A \rightarrow B \vdash B.$

6.  $(A \wedge B) \rightarrow C \vdash (A \rightarrow C) \vee (B \rightarrow C).$

7.  $(A \rightarrow B) \rightarrow A \vdash A.$

8.  $\vdash (A \rightarrow B) \vee (B \rightarrow C).$

(These all require the  $\perp_C$  rule.)

**Problem 11.4.** Give derivations that show the following:

1.  $\vdash (\forall x A(x) \wedge \forall y B(y)) \rightarrow \forall z (A(z) \wedge B(z)).$

2.  $\vdash (\exists x A(x) \vee \exists y B(y)) \rightarrow \exists z (A(z) \vee B(z)).$

3.  $\forall x (A(x) \rightarrow B) \vdash \exists y A(y) \rightarrow B.$

4.  $\forall x \neg A(x) \vdash \neg \exists x A(x).$

5.  $\vdash \neg \exists x A(x) \rightarrow \forall x \neg A(x).$

6.  $\vdash \neg \exists x \forall y ((A(x, y) \rightarrow \neg A(y, y)) \wedge (\neg A(y, y) \rightarrow A(x, y))).$

**Problem 11.5.** Give derivations that show the following:

1.  $\vdash \neg \forall x A(x) \rightarrow \exists x \neg A(x).$

2.  $(\forall x A(x) \rightarrow B) \vdash \exists y (A(y) \rightarrow B).$

3.  $\vdash \exists x (A(x) \rightarrow \forall y A(y)).$

(These all require the  $\perp_C$  rule.)

**Problem 11.6.** Prove **Proposition 11.16**

**Problem 11.7.** Prove that  $\Gamma \vdash \neg A$  iff  $\Gamma \cup \{A\}$  is inconsistent.

**Problem 11.8.** Complete the proof of **Theorem 11.27**.

**Problem 11.9.** Prove that  $=$  is both symmetric and transitive, i.e., give derivations of  $\forall x \forall y (x = y \rightarrow y = x)$  and  $\forall x \forall y \forall z ((x = y \wedge y = z) \rightarrow x = z)$

**Problem 11.10.** Give derivations of the following formulas:

1.  $\forall x \forall y ((x = y \wedge A(x)) \rightarrow A(y))$
2.  $\exists x A(x) \wedge \forall y \forall z ((A(y) \wedge A(z)) \rightarrow y = z) \rightarrow \exists x (A(x) \wedge \forall y (A(y) \rightarrow y = x))$

## CHAPTER 12

# *The Completeness Theorem*

### 12.1 Introduction

The completeness theorem is one of the most fundamental results about logic. It comes in two formulations, the equivalence of which we'll prove. In its first formulation it says something fundamental about the relationship between semantic consequence and our derivation system: if a sentence  $A$  follows from some sentences  $\Gamma$ , then there is also a derivation that establishes  $\Gamma \vdash A$ . Thus, the derivation system is as strong as it can possibly be without proving things that don't actually follow.

In its second formulation, it can be stated as a model existence result: every consistent set of sentences is satisfiable. Consistency is a proof-theoretic notion: it says that our derivation system is unable to produce certain derivations. But who's to say that just because there are no derivations of a certain sort from  $\Gamma$ , it's guaranteed that there is a structure  $M$ ? Before the completeness theorem was first proved—in fact before we had the

derivation systems we now do—the great German mathematician David Hilbert held the view that consistency of mathematical theories guarantees the existence of the objects they are about. He put it as follows in a letter to Gottlob Frege:

If the arbitrarily given axioms do not contradict one another with all their consequences, then they are true and the things defined by the axioms exist. This is for me the criterion of truth and existence.

Frege vehemently disagreed. The second formulation of the completeness theorem shows that Hilbert was right in at least the sense that if the axioms are consistent, then *some* structure exists that makes them all true.

These aren't the only reasons the completeness theorem—or rather, its proof—is important. It has a number of important consequences, some of which we'll discuss separately. For instance, since any derivation that shows  $\Gamma \vdash A$  is finite and so can only use finitely many of the sentences in  $\Gamma$ , it follows by the completeness theorem that if  $A$  is a consequence of  $\Gamma$ , it is already a consequence of a finite subset of  $\Gamma$ . This is called *compactness*. Equivalently, if every finite subset of  $\Gamma$  is consistent, then  $\Gamma$  itself must be consistent.

Although the compactness theorem follows from the completeness theorem via the detour through derivations, it is also possible to use the *the proof of* the completeness theorem to establish it directly. For what the proof does is take a set of sentences with a certain property—consistency—and constructs a structure out of this set that has certain properties (in this case, that it satisfies the set). Almost the very same construction can be used to directly establish compactness, by starting from “finitely satisfiable” sets of sentences instead of consistent ones. The construction also yields other consequences, e.g., that any satisfiable set of sentences has a finite or countably infinite model. (This result is called the Löwenheim–Skolem theorem.) In general, the construction of structures from sets of sentences is used often in logic, and sometimes even in philosophy.

## 12.2 Outline of the Proof

The proof of the completeness theorem is a bit complex, and upon first reading it, it is easy to get lost. So let us outline the proof. The first step is a shift of perspective, that allows us to see a route to a proof. When completeness is thought of as “whenever  $\Gamma \vDash A$  then  $\Gamma \vdash A$ ,” it may be hard to even come up with an idea: for to show that  $\Gamma \vdash A$  we have to find a derivation, and it does not look like the hypothesis that  $\Gamma \vDash A$  helps us for this in any way. For some proof systems it is possible to directly construct a derivation, but we will take a slightly different approach. The shift in perspective required is this: completeness can also be formulated as: “if  $\Gamma$  is consistent, it is satisfiable.” Perhaps we can use the information in  $\Gamma$  together with the hypothesis that it is consistent to construct a structure that satisfies every sentence in  $\Gamma$ . After all, we know what kind of structure we are looking for: one that is as  $\Gamma$  describes it!

If  $\Gamma$  contains only atomic sentences, it is easy to construct a model for it. Suppose the atomic sentences are all of the form  $P(a_1, \dots, a_n)$  where the  $a_i$  are constant symbols. All we have to do is come up with a domain  $|M|$  and an assignment for  $P$  so that  $M \vDash P(a_1, \dots, a_n)$ . But that’s not very hard: put  $|M| = \mathbb{N}$ ,  $c_i^M = i$ , and for every  $P(a_1, \dots, a_n) \in \Gamma$ , put the tuple  $\langle k_1, \dots, k_n \rangle$  into  $P^M$ , where  $k_i$  is the index of the constant symbol  $a_i$  (i.e.,  $a_i \equiv c_{k_i}$ ).

Now suppose  $\Gamma$  contains some formula  $\neg B$ , with  $B$  atomic. We might worry that the construction of  $M$  interferes with the possibility of making  $\neg B$  true. But here’s where the consistency of  $\Gamma$  comes in: if  $\neg B \in \Gamma$ , then  $B \notin \Gamma$ , or else  $\Gamma$  would be inconsistent. And if  $B \notin \Gamma$ , then according to our construction of  $M$ ,  $M \not\vDash B$ , so  $M \vDash \neg B$ . So far so good.

What if  $\Gamma$  contains complex, non-atomic formulas? Say it contains  $A \wedge B$ . To make that true, we should proceed as if both  $A$  and  $B$  were in  $\Gamma$ . And if  $A \vee B \in \Gamma$ , then we will have to make at least one of them true, i.e., proceed as if one of them was in  $\Gamma$ .

This suggests the following idea: we add additional formulas

to  $\Gamma$  so as to (a) keep the resulting set consistent and (b) make sure that for every possible atomic sentence  $A$ , either  $A$  is in the resulting set, or  $\neg A$  is, and (c) such that, whenever  $A \wedge B$  is in the set, so are both  $A$  and  $B$ , if  $A \vee B$  is in the set, at least one of  $A$  or  $B$  is also, etc. We keep doing this (potentially forever). Call the set of all formulas so added  $\Gamma^*$ . Then our construction above would provide us with a structure  $M$  for which we could prove, by induction, that it satisfies all sentences in  $\Gamma^*$ , and hence also all sentence in  $\Gamma$  since  $\Gamma \subseteq \Gamma^*$ . It turns out that guaranteeing (a) and (b) is enough. A set of sentences for which (b) holds is called *complete*. So our task will be to extend the consistent set  $\Gamma$  to a consistent and complete set  $\Gamma^*$ .

There is one wrinkle in this plan: if  $\exists x A(x) \in \Gamma$  we would hope to be able to pick some constant symbol  $c$  and add  $A(c)$  in this process. But how do we know we can always do that? Perhaps we only have a few constant symbols in our language, and for each one of them we have  $\neg A(c) \in \Gamma$ . We can't also add  $A(c)$ , since this would make the set inconsistent, and we wouldn't know whether  $M$  has to make  $A(c)$  or  $\neg A(c)$  true. Moreover, it might happen that  $\Gamma$  contains only sentences in a language that has no constant symbols at all (e.g., the language of set theory).

The solution to this problem is to simply add infinitely many constants at the beginning, plus sentences that connect them with the quantifiers in the right way. (Of course, we have to verify that this cannot introduce an inconsistency.)

Our original construction works well if we only have constant symbols in the atomic sentences. But the language might also contain function symbols. In that case, it might be tricky to find the right functions on  $\mathbb{N}$  to assign to these function symbols to make everything work. So here's another trick: instead of using  $i$  to interpret  $c_i$ , just take the set of constant symbols itself as the domain. Then  $M$  can assign every constant symbol to itself:  $c_i^M = c_i$ . But why not go all the way: let  $|M|$  be all *terms* of the language! If we do this, there is an obvious assignment of functions (that take terms as arguments and have terms as values) to function symbols: we assign to the function symbol  $f_i^n$  the

function which, given  $n$  terms  $t_1, \dots, t_n$  as input, produces the term  $f_i^n(t_1, \dots, t_n)$  as value.

The last piece of the puzzle is what to do with  $=$ . The predicate symbol  $=$  has a fixed interpretation:  $M \models t = t'$  iff  $\text{Val}^M(t) = \text{Val}^M(t')$ . Now if we set things up so that the value of a term  $t$  is  $t$  itself, then this structure will make *no* sentence of the form  $t = t'$  true unless  $t$  and  $t'$  are one and the same term. And of course this is a problem, since basically every interesting theory in a language with function symbols will have as theorems sentences  $t = t'$  where  $t$  and  $t'$  are not the same term (e.g., in theories of arithmetic:  $(0 + 0) = 0$ ). To solve this problem, we change the domain of  $M$ : instead of using terms as the objects in  $|M|$ , we use sets of terms, and each set is so that it contains all those terms which the sentences in  $\Gamma$  require to be equal. So, e.g., if  $\Gamma$  is a theory of arithmetic, one of these sets will contain:  $0, (0 + 0), (0 \times 0)$ , etc. This will be the set we assign to  $0$ , and it will turn out that this set is also the value of all the terms in it, e.g., also of  $(0 + 0)$ . Therefore, the sentence  $(0 + 0) = 0$  will be true in this revised structure.

So here's what we'll do. First we investigate the properties of complete consistent sets, in particular we prove that a complete consistent set contains  $A \wedge B$  iff it contains both  $A$  and  $B$ ,  $A \vee B$  iff it contains at least one of them, etc. (**Proposition 12.2**). Then we define and investigate "saturated" sets of sentences. A saturated set is one which contains conditionals that link each quantified sentence to instances of it (**Definition 12.5**). We show that any consistent set  $\Gamma$  can always be extended to a saturated set  $\Gamma'$  (**Lemma 12.6**). If a set is consistent, saturated, and complete it also has the property that it contains  $\exists x A(x)$  iff it contains  $A(t)$  for some closed term  $t$  and  $\forall x A(x)$  iff it contains  $A(t)$  for all closed terms  $t$  (**Proposition 12.7**). We'll then take the saturated consistent set  $\Gamma'$  and show that it can be extended to a saturated, consistent, and complete set  $\Gamma^*$  (**Lemma 12.8**). This set  $\Gamma^*$  is what we'll use to define our term model  $M(\Gamma^*)$ . The term model has the set of closed terms as its domain, and the interpretation of its predicate symbols is given by the atomic sentences



in  $\Gamma^*$  (Definition 12.9). We'll use the properties of saturated, complete consistent sets to show that indeed  $M(\Gamma^*) \models A$  iff  $A \in \Gamma^*$  (Lemma 12.12), and thus in particular,  $M(\Gamma^*) \models \Gamma$ . Finally, we'll consider how to define a term model if  $\Gamma$  contains  $=$  as well (Definition 12.16) and show that it satisfies  $\Gamma^*$  (Lemma 12.19).

### 12.3 Complete Consistent Sets of Sentences

**Definition 12.1 (Complete set).** A set  $\Gamma$  of sentences is *complete* iff for any sentence  $A$ , either  $A \in \Gamma$  or  $\neg A \in \Gamma$ .

Complete sets of sentences leave no questions unanswered. For any sentence  $A$ ,  $\Gamma$  “says” if  $A$  is true or false. The importance of complete sets extends beyond the proof of the completeness theorem. A theory which is complete and axiomatizable, for instance, is always decidable.

Complete consistent sets are important in the completeness proof since we can guarantee that every consistent set of sentences  $\Gamma$  is contained in a complete consistent set  $\Gamma^*$ . A complete consistent set contains, for each sentence  $A$ , either  $A$  or its negation  $\neg A$ , but not both. This is true in particular for atomic sentences, so from a complete consistent set in a language suitably expanded by constant symbols, we can construct a structure where the interpretation of predicate symbols is defined according to which atomic sentences are in  $\Gamma^*$ . This structure can then be shown to make all sentences in  $\Gamma^*$  (and hence also all those in  $\Gamma$ ) true. The proof of this latter fact requires that  $\neg A \in \Gamma^*$  iff  $A \notin \Gamma^*$ ,  $(A \vee B) \in \Gamma^*$  iff  $A \in \Gamma^*$  or  $B \in \Gamma^*$ , etc.

In what follows, we will often tacitly use the properties of reflexivity, monotonicity, and transitivity of  $\vdash$  (see sections 10.8 and 11.7).

**Proposition 12.2.** *Suppose  $\Gamma$  is complete and consistent. Then:*

1. *If  $\Gamma \vdash A$ , then  $A \in \Gamma$ .*

2.  $A \wedge B \in \Gamma$  iff both  $A \in \Gamma$  and  $B \in \Gamma$ .
3.  $A \vee B \in \Gamma$  iff either  $A \in \Gamma$  or  $B \in \Gamma$ .
4.  $A \rightarrow B \in \Gamma$  iff either  $A \notin \Gamma$  or  $B \in \Gamma$ .

*Proof.* Let us suppose for all of the following that  $\Gamma$  is complete and consistent.

1. If  $\Gamma \vdash A$ , then  $A \in \Gamma$ .

Suppose that  $\Gamma \vdash A$ . Suppose to the contrary that  $A \notin \Gamma$ . Since  $\Gamma$  is complete,  $\neg A \in \Gamma$ . By **Propositions 10.20** and **11.20**,  $\Gamma$  is inconsistent. This contradicts the assumption that  $\Gamma$  is consistent. Hence, it cannot be the case that  $A \notin \Gamma$ , so  $A \in \Gamma$ .

2. Exercise.

3. First we show that if  $A \vee B \in \Gamma$ , then either  $A \in \Gamma$  or  $B \in \Gamma$ . Suppose  $A \vee B \in \Gamma$  but  $A \notin \Gamma$  and  $B \notin \Gamma$ . Since  $\Gamma$  is complete,  $\neg A \in \Gamma$  and  $\neg B \in \Gamma$ . By **Propositions 10.23** and **11.23**, item (1),  $\Gamma$  is inconsistent, a contradiction. Hence, either  $A \in \Gamma$  or  $B \in \Gamma$ .

For the reverse direction, suppose that  $A \in \Gamma$  or  $B \in \Gamma$ . By **Propositions 10.23** and **11.23**, item (2),  $\Gamma \vdash A \vee B$ . By (1),  $A \vee B \in \Gamma$ , as required.

4. Exercise. □

## 12.4 Henkin Expansion

Part of the challenge in proving the completeness theorem is that the model we construct from a complete consistent set  $\Gamma$  must make all the quantified formulas in  $\Gamma$  true. In order to guarantee this, we use a trick due to Leon Henkin. In essence, the

trick consists in expanding the language by infinitely many constant symbols and adding, for each formula with one free variable  $A(x)$  a formula of the form  $\exists x A(x) \rightarrow A(c)$ , where  $c$  is one of the new constant symbols. When we construct the structure satisfying  $\Gamma$ , this will guarantee that each true existential sentence has a witness among the new constants.

**Proposition 12.3.** *If  $\Gamma$  is consistent in  $\mathcal{L}$  and  $\mathcal{L}'$  is obtained from  $\mathcal{L}$  by adding a countably infinite set of new constant symbols  $c_0, c_1, \dots$ , then  $\Gamma$  is consistent in  $\mathcal{L}'$ .*

**Definition 12.4 (Saturated set).** A set  $\Gamma$  of formulas of a language  $\mathcal{L}$  is *saturated* iff for each formula  $A(x) \in \text{Frm}(\mathcal{L})$  with one free variable  $x$  there is a constant symbol  $c \in \mathcal{L}$  such that  $\exists x A(x) \rightarrow A(c) \in \Gamma$ .

The following definition will be used in the proof of the next theorem.

**Definition 12.5.** Let  $\mathcal{L}'$  be as in **Proposition 12.3**. Fix an enumeration  $A_0(x_0), A_1(x_1), \dots$  of all formulas  $A_i(x_i)$  of  $\mathcal{L}'$  in which one variable ( $x_i$ ) occurs free. We define the sentences  $D_n$  by induction on  $n$ .

Let  $c_0$  be the first constant symbol among the  $d_i$  we added to  $\mathcal{L}$  which does not occur in  $A_0(x_0)$ . Assuming that  $D_0, \dots, D_{n-1}$  have already been defined, let  $c_n$  be the first among the new constant symbols  $d_i$  that occurs neither in  $D_0, \dots, D_{n-1}$  nor in  $A_n(x_n)$ .

Now let  $D_n$  be the formula  $\exists x_n A_n(x_n) \rightarrow A_n(c_n)$ .

**Lemma 12.6.** *Every consistent set  $\Gamma$  can be extended to a saturated consistent set  $\Gamma'$ .*

*Proof.* Given a consistent set of sentences  $\Gamma$  in a language  $\mathcal{L}$ , expand the language by adding a countably infinite set of new constant symbols to form  $\mathcal{L}'$ . By **Proposition 12.3**,  $\Gamma$  is still consistent in the richer language. Further, let  $D_i$  be as in **Definition 12.5**. Let

$$\begin{aligned}\Gamma_0 &= \Gamma \\ \Gamma_{n+1} &= \Gamma_n \cup \{D_n\}\end{aligned}$$

i.e.,  $\Gamma_{n+1} = \Gamma \cup \{D_0, \dots, D_n\}$ , and let  $\Gamma' = \bigcup_n \Gamma_n$ .  $\Gamma'$  is clearly saturated.

If  $\Gamma'$  were inconsistent, then for some  $n$ ,  $\Gamma_n$  would be inconsistent (Exercise: explain why). So to show that  $\Gamma'$  is consistent it suffices to show, by induction on  $n$ , that each set  $\Gamma_n$  is consistent.

The induction basis is simply the claim that  $\Gamma_0 = \Gamma$  is consistent, which is the hypothesis of the theorem. For the induction step, suppose that  $\Gamma_n$  is consistent but  $\Gamma_{n+1} = \Gamma_n \cup \{D_n\}$  is inconsistent. Recall that  $D_n$  is  $\exists x_n A_n(x_n) \rightarrow A_n(c_n)$ , where  $A_n(x_n)$  is a formula of  $\mathcal{L}'$  with only the variable  $x_n$  free. By the way we've chosen the  $c_n$  (see **Definition 12.5**),  $c_n$  does not occur in  $A_n(x_n)$  nor in  $\Gamma_n$ .

If  $\Gamma_n \cup \{D_n\}$  is inconsistent, then  $\Gamma_n \vdash \neg D_n$ , and hence both of the following hold:

$$\Gamma_n \vdash \exists x_n A_n(x_n) \quad \Gamma_n \vdash \neg A_n(c_n)$$

Since  $c_n$  does not occur in  $\Gamma_n$  or in  $A_n(x_n)$ , **Theorems 10.25** and **11.25** applies. From  $\Gamma_n \vdash \neg A_n(c_n)$ , we obtain  $\Gamma_n \vdash \forall x_n \neg A_n(x_n)$ . Thus we have that both  $\Gamma_n \vdash \exists x_n A_n(x_n)$  and  $\Gamma_n \vdash \forall x_n \neg A_n(x_n)$ , so  $\Gamma_n$  itself is inconsistent. (Note that  $\forall x_n \neg A_n(x_n) \vdash \neg \exists x_n A_n(x_n)$ .) Contradiction:  $\Gamma_n$  was supposed to be consistent. Hence  $\Gamma_n \cup \{D_n\}$  is consistent.  $\square$

We'll now show that *complete*, consistent sets which are saturated have the property that it contains a universally quantified sentence iff it contains all its instances and it contains an existentially quantified sentence iff it contains at least one instance. We'll

use this to show that the structure we'll generate from a complete, consistent, saturated set makes all its quantified sentences true.

**Proposition 12.7.** *Suppose  $\Gamma$  is complete, consistent, and saturated.*

1.  $\exists x A(x) \in \Gamma$  iff  $A(t) \in \Gamma$  for at least one closed term  $t$ .
2.  $\forall x A(x) \in \Gamma$  iff  $A(t) \in \Gamma$  for all closed terms  $t$ .

*Proof.* 1. First suppose that  $\exists x A(x) \in \Gamma$ . Because  $\Gamma$  is saturated,  $(\exists x A(x) \rightarrow A(c)) \in \Gamma$  for some constant symbol  $c$ . By Propositions 10.24 and 11.24, item (1), and Proposition 12.2(1),  $A(c) \in \Gamma$ .

For the other direction, saturation is not necessary: Suppose  $A(t) \in \Gamma$ . Then  $\Gamma \vdash \exists x A(x)$  by Propositions 10.26 and 11.26, item (1). By Proposition 12.2(1),  $\exists x A(x) \in \Gamma$ .

2. Exercise. □

## 12.5 Lindenbaum's Lemma

We now prove a lemma that shows that any consistent set of sentences is contained in some set of sentences which is not just consistent, but also complete. The proof works by adding one sentence at a time, guaranteeing at each step that the set remains consistent. We do this so that for every  $A$ , either  $A$  or  $\neg A$  gets added at some stage. The union of all stages in that construction then contains either  $A$  or its negation  $\neg A$  and is thus complete. It is also consistent, since we make sure at each stage not to introduce an inconsistency.

**Lemma 12.8 (Lindenbaum's Lemma).** *Every consistent set  $\Gamma$  in a language  $\mathcal{L}$  can be extended to a complete and consistent set  $\Gamma^*$ .*

*Proof.* Let  $\Gamma$  be consistent. Let  $A_0, A_1, \dots$  be an enumeration of all the sentences of  $\mathcal{L}$ . Define  $\Gamma_0 = \Gamma$ , and

$$\Gamma_{n+1} = \begin{cases} \Gamma_n \cup \{A_n\} & \text{if } \Gamma_n \cup \{A_n\} \text{ is consistent;} \\ \Gamma_n \cup \{\neg A_n\} & \text{otherwise.} \end{cases}$$

Let  $\Gamma^* = \bigcup_{n \geq 0} \Gamma_n$ .

Each  $\Gamma_n$  is consistent:  $\Gamma_0$  is consistent by definition. If  $\Gamma_{n+1} = \Gamma_n \cup \{A_n\}$ , this is because the latter is consistent. If it isn't,  $\Gamma_{n+1} = \Gamma_n \cup \{\neg A_n\}$ . We have to verify that  $\Gamma_n \cup \{\neg A_n\}$  is consistent. Suppose it's not. Then *both*  $\Gamma_n \cup \{A_n\}$  and  $\Gamma_n \cup \{\neg A_n\}$  are inconsistent. This means that  $\Gamma_n$  would be inconsistent by **Propositions 10.21** and **11.21**, contrary to the induction hypothesis.

For every  $n$  and every  $i < n$ ,  $\Gamma_i \subseteq \Gamma_n$ . This follows by a simple induction on  $n$ . For  $n = 0$ , there are no  $i < 0$ , so the claim holds automatically. For the inductive step, suppose it is true for  $n$ . We have  $\Gamma_{n+1} = \Gamma_n \cup \{A_n\}$  or  $= \Gamma_n \cup \{\neg A_n\}$  by construction. So  $\Gamma_n \subseteq \Gamma_{n+1}$ . If  $i < n$ , then  $\Gamma_i \subseteq \Gamma_n$  by inductive hypothesis, and so  $\Gamma_i \subseteq \Gamma_{n+1}$  by transitivity of  $\subseteq$ .

From this it follows that every finite subset of  $\Gamma^*$  is a subset of  $\Gamma_n$  for some  $n$ , since each  $B \in \Gamma^*$  not already in  $\Gamma_0$  is added at some stage  $i$ . If  $n$  is the last one of these, then all  $B$  in the finite subset are in  $\Gamma_n$ . So, every finite subset of  $\Gamma^*$  is consistent. By **Propositions 10.17** and **11.17**,  $\Gamma^*$  is consistent.

Every sentence of  $\text{Frm}(\mathcal{L})$  appears on the list used to define  $\Gamma^*$ . If  $A_n \notin \Gamma^*$ , then that is because  $\Gamma_n \cup \{A_n\}$  was inconsistent. But then  $\neg A_n \in \Gamma^*$ , so  $\Gamma^*$  is complete.  $\square$

## 12.6 Construction of a Model

Right now we are not concerned about  $=$ , i.e., we only want to show that a consistent set  $\Gamma$  of sentences not containing  $=$  is satisfiable. We first extend  $\Gamma$  to a consistent, complete, and saturated set  $\Gamma^*$ . In this case, the definition of a model  $M(\Gamma^*)$  is simple: We take the set of closed terms of  $\mathcal{L}'$  as the domain. We assign every

constant symbol to itself, and make sure that more generally, for every closed term  $t$ ,  $\text{Val}^{M(\Gamma^*)}(t) = t$ . The predicate symbols are assigned extensions in such a way that an atomic sentence is true in  $M(\Gamma^*)$  iff it is in  $\Gamma^*$ . This will obviously make all the atomic sentences in  $\Gamma^*$  true in  $M(\Gamma^*)$ . The rest are true provided the  $\Gamma^*$  we start with is consistent, complete, and saturated.

**Definition 12.9 (Term model).** Let  $\Gamma^*$  be a complete and consistent, saturated set of sentences in a language  $\mathcal{L}$ . The *term model*  $M(\Gamma^*)$  of  $\Gamma^*$  is the structure defined as follows:

1. The domain  $|M(\Gamma^*)|$  is the set of all closed terms of  $\mathcal{L}$ .
2. The interpretation of a constant symbol  $c$  is  $c$  itself:  
 $c^{M(\Gamma^*)} = c$ .
3. The function symbol  $f$  is assigned the function which, given as arguments the closed terms  $t_1, \dots, t_n$ , has as value the closed term  $f(t_1, \dots, t_n)$ :

$$f^{M(\Gamma^*)}(t_1, \dots, t_n) = f(t_1, \dots, t_n)$$

4. If  $R$  is an  $n$ -place predicate symbol, then

$$\langle t_1, \dots, t_n \rangle \in R^{M(\Gamma^*)} \text{ iff } R(t_1, \dots, t_n) \in \Gamma^*.$$

We will now check that we indeed have  $\text{Val}^{M(\Gamma^*)}(t) = t$ .

**Lemma 12.10.** *Let  $M(\Gamma^*)$  be the term model of Definition 12.9, then  $\text{Val}^{M(\Gamma^*)}(t) = t$ .*

*Proof.* The proof is by induction on  $t$ , where the base case, when  $t$  is a constant symbol, follows directly from the definition of the term model. For the induction step assume  $t_1, \dots, t_n$  are closed terms such that  $\text{Val}^{M(\Gamma^*)}(t_i) = t_i$  and that  $f$  is an  $n$ -ary function symbol. Then

$$\text{Val}^{M(\Gamma^*)}(f(t_1, \dots, t_n)) = f^{M(\Gamma^*)}(\text{Val}^{M(\Gamma^*)}(t_1), \dots, \text{Val}^{M(\Gamma^*)}(t_n))$$

$$\begin{aligned}
 &= f^{M(\Gamma^*)}(t_1, \dots, t_n) \\
 &= f(t_1, \dots, t_n),
 \end{aligned}$$

and so by induction this holds for every closed term  $t$ .  $\square$

A structure  $M$  may make an existentially quantified sentence  $\exists x A(x)$  true without there being an instance  $A(t)$  that it makes true. A structure  $M$  may make all instances  $A(t)$  of a universally quantified sentence  $\forall x A(x)$  true, without making  $\forall x A(x)$  true. This is because in general not every element of  $|M|$  is the value of a closed term ( $M$  may not be covered). This is the reason the satisfaction relation is defined via variable assignments. However, for our term model  $M(\Gamma^*)$  this wouldn't be necessary—because it is covered. This is the content of the next result.

**Proposition 12.11.** *Let  $M(\Gamma^*)$  be the term model of [Definition 12.9](#).*

1.  $M(\Gamma^*) \models \exists x A(x)$  iff  $M(\Gamma^*) \models A(t)$  for at least one closed term  $t$ .
2.  $M(\Gamma^*) \models \forall x A(x)$  iff  $M(\Gamma^*) \models A(t)$  for all closed terms  $t$ .

*Proof.* 1. By [Proposition 7.18](#),  $M(\Gamma^*) \models \exists x A(x)$  iff for at least one variable assignment  $s$ ,  $M(\Gamma^*), s \models A(x)$ . As  $|M(\Gamma^*)|$  consists of the closed terms of  $\mathcal{L}$ , this is the case iff there is at least one closed term  $t$  such that  $s(x) = t$  and  $M(\Gamma^*), s \models A(x)$ . By [Proposition 7.22](#),  $M(\Gamma^*), s \models A(x)$  iff  $M(\Gamma^*), s \models A(t)$ , where  $s(x) = t$ . By [Proposition 7.17](#),  $M(\Gamma^*), s \models A(t)$  iff  $M(\Gamma^*) \models A(t)$ , since  $A(t)$  is a sentence.

2. Exercise.  $\square$

**Lemma 12.12 (Truth Lemma).** *Suppose  $A$  does not contain  $=$ . Then  $M(\Gamma^*) \models A$  iff  $A \in \Gamma^*$ .*

*Proof.* We prove both directions simultaneously, and by induction on  $A$ .



1.  $A \equiv \perp$ :  $M(\Gamma^*) \not\models \perp$  by definition of satisfaction. On the other hand,  $\perp \notin \Gamma^*$  since  $\Gamma^*$  is consistent.
  2.  $A \equiv R(t_1, \dots, t_n)$ :  $M(\Gamma^*) \models R(t_1, \dots, t_n)$  iff  $\langle t_1, \dots, t_n \rangle \in R^{M(\Gamma^*)}$  (by the definition of satisfaction) iff  $R(t_1, \dots, t_n) \in \Gamma^*$  (by the construction of  $M(\Gamma^*)$ ).
  3.  $A \equiv \neg B$ :  $M(\Gamma^*) \models A$  iff  $M(\Gamma^*) \not\models B$  (by definition of satisfaction). By induction hypothesis,  $M(\Gamma^*) \not\models B$  iff  $B \notin \Gamma^*$ . Since  $\Gamma^*$  is consistent and complete,  $B \notin \Gamma^*$  iff  $\neg B \in \Gamma^*$ .
  4.  $A \equiv B \wedge C$ : exercise.
  5.  $A \equiv B \vee C$ :  $M(\Gamma^*) \models A$  iff  $M(\Gamma^*) \models B$  or  $M(\Gamma^*) \models C$  (by definition of satisfaction) iff  $B \in \Gamma^*$  or  $C \in \Gamma^*$  (by induction hypothesis). This is the case iff  $(B \vee C) \in \Gamma^*$  (by **Proposition 12.2(3)**).
  6.  $A \equiv B \rightarrow C$ : exercise.
  7.  $A \equiv \forall x B(x)$ : exercise.
  8.  $A \equiv \exists x B(x)$ :  $M(\Gamma^*) \models A$  iff  $M(\Gamma^*) \models B(t)$  for at least one term  $t$  (**Proposition 12.11**). By induction hypothesis, this is the case iff  $B(t) \in \Gamma^*$  for at least one term  $t$ . By **Proposition 12.7**, this in turn is the case iff  $\exists x B(x) \in \Gamma^*$ .
- 

## 12.7 Identity

The construction of the term model given in the preceding section is enough to establish completeness for first-order logic for sets  $\Gamma$  that do not contain  $=$ . The term model satisfies every  $A \in \Gamma^*$  which does not contain  $=$  (and hence all  $A \in \Gamma$ ). It does not work, however, if  $=$  is present. The reason is that  $\Gamma^*$  then may contain a sentence  $t = t'$ , but in the term model the value of any term is that term itself. Hence, if  $t$  and  $t'$  are different terms,

their values in the term model—i.e.,  $t$  and  $t'$ , respectively—are different, and so  $t = t'$  is false. We can fix this, however, using a construction known as “factoring.”

**Definition 12.13.** Let  $\Gamma^*$  be a consistent and complete set of sentences in  $\mathcal{L}$ . We define the relation  $\approx$  on the set of closed terms of  $\mathcal{L}$  by

$$t \approx t' \quad \text{iff} \quad t = t' \in \Gamma^*$$

**Proposition 12.14.** *The relation  $\approx$  has the following properties:*

1.  $\approx$  is reflexive.
2.  $\approx$  is symmetric.
3.  $\approx$  is transitive.
4. If  $t \approx t'$ ,  $f$  is a function symbol, and  $t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n$  are closed terms, then

$$f(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \approx f(t_1, \dots, t_{i-1}, t', t_{i+1}, \dots, t_n).$$

5. If  $t \approx t'$ ,  $R$  is a predicate symbol, and  $t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n$  are closed terms, then

$$R(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) \in \Gamma^* \quad \text{iff} \\ R(t_1, \dots, t_{i-1}, t', t_{i+1}, \dots, t_n) \in \Gamma^*.$$

*Proof.* Since  $\Gamma^*$  is consistent and complete,  $t = t' \in \Gamma^*$  iff  $\Gamma^* \vdash t = t'$ . Thus it is enough to show the following:

1.  $\Gamma^* \vdash t = t$  for all closed terms  $t$ .
2. If  $\Gamma^* \vdash t = t'$  then  $\Gamma^* \vdash t' = t$ .
3. If  $\Gamma^* \vdash t = t'$  and  $\Gamma^* \vdash t' = t''$ , then  $\Gamma^* \vdash t = t''$ .

4. If  $\Gamma^* \vdash t = t'$ , then

$$\Gamma^* \vdash f(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n) = f(t_1, \dots, t_{i-1}, t', t_{i+1}, \dots, t_n)$$

for every  $n$ -place function symbol  $f$  and closed terms  $t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n$ .

5. If  $\Gamma^* \vdash t = t'$  and  $\Gamma^* \vdash R(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$ , then  $\Gamma^* \vdash R(t_1, \dots, t_{i-1}, t', t_{i+1}, \dots, t_n)$  for every  $n$ -place predicate symbol  $R$  and closed terms  $t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_n$ .  $\square$

**Definition 12.15.** Suppose  $\Gamma^*$  is a consistent and complete set in a language  $\mathcal{L}$ ,  $t$  is a closed term, and  $\approx$  as in the previous definition. Then:

$$[t]_{\approx} = \{t' : t' \in \text{Trm}(\mathcal{L}), t \approx t'\}$$

and  $\text{Trm}(\mathcal{L})/\approx = \{[t]_{\approx} : t \in \text{Trm}(\mathcal{L})\}$ .

**Definition 12.16.** Let  $M = M(\Gamma^*)$  be the term model for  $\Gamma^*$  from [Definition 12.9](#). Then  $M/\approx$  is the following structure:

1.  $|M/\approx| = \text{Trm}(\mathcal{L})/\approx$ .
2.  $c^{M/\approx} = [c]_{\approx}$
3.  $f^{M/\approx}([t_1]_{\approx}, \dots, [t_n]_{\approx}) = [f(t_1, \dots, t_n)]_{\approx}$
4.  $\langle [t_1]_{\approx}, \dots, [t_n]_{\approx} \rangle \in R^{M/\approx}$  iff  $M \models R(t_1, \dots, t_n)$ , i.e., iff  $R(t_1, \dots, t_n) \in \Gamma^*$ .

Note that we have defined  $f^{M/\approx}$  and  $R^{M/\approx}$  for elements of  $\text{Trm}(\mathcal{L})/\approx$  by referring to them as  $[t]_{\approx}$ , i.e., via *representatives*  $t \in [t]_{\approx}$ . We have to make sure that these definitions do not depend on the choice of these representatives, i.e., that for some other choices  $t'$  which determine the same equivalence classes ( $[t]_{\approx} =$

$[t']_{\approx}$ ), the definitions yield the same result. For instance, if  $R$  is a one-place predicate symbol, the last clause of the definition says that  $[t]_{\approx} \in R^{M/\approx}$  iff  $M \vDash R(t)$ . If for some other term  $t'$  with  $t \approx t'$ ,  $M \not\vDash R(t)$ , then the definition would require  $[t']_{\approx} \notin R^{M/\approx}$ . If  $t \approx t'$ , then  $[t]_{\approx} = [t']_{\approx}$ , but we can't have both  $[t]_{\approx} \in R^{M/\approx}$  and  $[t]_{\approx} \notin R^{M/\approx}$ . However, **Proposition 12.14** guarantees that this cannot happen.

**Proposition 12.17.**  *$M/\approx$  is well defined, i.e., if  $t_1, \dots, t_n, t'_1, \dots, t'_n$  are closed terms, and  $t_i \approx t'_i$  then*

$$1. [f(t_1, \dots, t_n)]_{\approx} = [f(t'_1, \dots, t'_n)]_{\approx}, \text{ i.e.,}$$

$$f(t_1, \dots, t_n) \approx f(t'_1, \dots, t'_n)$$

and

$$2. M \vDash R(t_1, \dots, t_n) \text{ iff } M \vDash R(t'_1, \dots, t'_n), \text{ i.e.,}$$

$$R(t_1, \dots, t_n) \in \Gamma^* \text{ iff } R(t'_1, \dots, t'_n) \in \Gamma^*.$$

*Proof.* Follows from **Proposition 12.14** by induction on  $n$ . □

As in the case of the term model, before proving the truth lemma we need the following lemma.

**Lemma 12.18.** *Let  $M = M(\Gamma^*)$ , then  $\text{Val}^{M/\approx}(t) = [t]_{\approx}$ .*

*Proof.* The proof is similar to that of **Lemma 12.10**. □

**Lemma 12.19.**  *$M/\approx \vDash A$  iff  $A \in \Gamma^*$  for all sentences  $A$ .*

*Proof.* By induction on  $A$ , just as in the proof of **Lemma 12.12**. The only case that needs additional attention is when  $A \equiv t = t'$ .

$$M/\approx \vDash t = t' \text{ iff } [t]_{\approx} = [t']_{\approx} \text{ (by definition of } M/\approx)$$

$$\text{iff } t \approx t' \text{ (by definition of } [t]_{\approx})$$

$$\text{iff } t = t' \in \Gamma^* \text{ (by definition of } \approx). \quad \square$$

Note that while  $M(\Gamma^*)$  is always countable and infinite,  $M/\approx$  may be finite, since it may turn out that there are only finitely many classes  $[t]_{\approx}$ . This is to be expected, since  $\Gamma$  may contain sentences which require any structure in which they are true to be finite. For instance,  $\forall x \forall y x = y$  is a consistent sentence, but is satisfied only in structures with a domain that contains exactly one element.

## 12.8 The Completeness Theorem

Let's combine our results: we arrive at the completeness theorem.

**Theorem 12.20 (Completeness Theorem).** *Let  $\Gamma$  be a set of sentences. If  $\Gamma$  is consistent, it is satisfiable.*

*Proof.* Suppose  $\Gamma$  is consistent. By [Lemma 12.6](#), there is a saturated consistent set  $\Gamma' \supseteq \Gamma$ . By [Lemma 12.8](#), there is a  $\Gamma^* \supseteq \Gamma'$  which is consistent and complete. Since  $\Gamma' \subseteq \Gamma^*$ , for each formula  $A(x)$ ,  $\Gamma^*$  contains a sentence of the form  $\exists x A(x) \rightarrow A(c)$  and so  $\Gamma^*$  is saturated. If  $\Gamma$  does not contain  $=$ , then by [Lemma 12.12](#),  $M(\Gamma^*) \models A$  iff  $A \in \Gamma^*$ . From this it follows in particular that for all  $A \in \Gamma$ ,  $M(\Gamma^*) \models A$ , so  $\Gamma$  is satisfiable. If  $\Gamma$  does contain  $=$ , then by [Lemma 12.19](#), for all sentences  $A$ ,  $M/\approx \models A$  iff  $A \in \Gamma^*$ . In particular,  $M/\approx \models A$  for all  $A \in \Gamma$ , so  $\Gamma$  is satisfiable.  $\square$

**Corollary 12.21 (Completeness Theorem, Second Version).** *For all  $\Gamma$  and sentences  $A$ : if  $\Gamma \models A$  then  $\Gamma \vdash A$ .*

*Proof.* Note that the  $\Gamma$ 's in [Corollary 12.21](#) and [Theorem 12.20](#) are universally quantified. To make sure we do not confuse ourselves, let us restate [Theorem 12.20](#) using a different variable: for any set of sentences  $\Delta$ , if  $\Delta$  is consistent, it is satisfiable. By contraposition, if  $\Delta$  is not satisfiable, then  $\Delta$  is inconsistent. We will use this to prove the corollary.

Suppose that  $\Gamma \vDash A$ . Then  $\Gamma \cup \{\neg A\}$  is unsatisfiable by [Proposition 7.27](#). Taking  $\Gamma \cup \{\neg A\}$  as our  $\Delta$ , the previous version of [Theorem 12.20](#) gives us that  $\Gamma \cup \{\neg A\}$  is inconsistent. By [Propositions 10.19](#) and [11.19](#),  $\Gamma \vdash A$ .  $\square$

## 12.9 The Compactness Theorem

One important consequence of the completeness theorem is the compactness theorem. The compactness theorem states that if each *finite* subset of a set of sentences is satisfiable, the entire set is satisfiable—even if the set itself is infinite. This is far from obvious. There is nothing that seems to rule out, at first glance at least, the possibility of there being infinite sets of sentences which are contradictory, but the contradiction only arises, so to speak, from the infinite number. The compactness theorem says that such a scenario can be ruled out: there are no unsatisfiable infinite sets of sentences each finite subset of which is satisfiable. Like the completeness theorem, it has a version related to entailment: if an infinite set of sentences entails something, already a finite subset does.

**Definition 12.22.** A set  $\Gamma$  of formulas is *finitely satisfiable* iff every finite  $\Gamma_0 \subseteq \Gamma$  is satisfiable.

**Theorem 12.23 (Compactness Theorem).** *The following hold for any sentences  $\Gamma$  and  $A$ :*

1.  $\Gamma \vDash A$  iff there is a finite  $\Gamma_0 \subseteq \Gamma$  such that  $\Gamma_0 \vDash A$ .
2.  $\Gamma$  is satisfiable iff it is finitely satisfiable.

*Proof.* We prove (2). If  $\Gamma$  is satisfiable, then there is a structure  $M$  such that  $M \vDash A$  for all  $A \in \Gamma$ . Of course, this  $M$  also satisfies every finite subset of  $\Gamma$ , so  $\Gamma$  is finitely satisfiable.

Now suppose that  $\Gamma$  is finitely satisfiable. Then every finite subset  $\Gamma_0 \subseteq \Gamma$  is satisfiable. By soundness ([Corollaries 10.31](#)

and 11.29), every finite subset is consistent. Then  $\Gamma$  itself must be consistent by Propositions 10.17 and 11.17. By completeness (Theorem 12.20), since  $\Gamma$  is consistent, it is satisfiable.  $\square$

**Example 12.24.** In every model  $M$  of a theory  $\Gamma$ , each term  $t$  of course picks out an element of  $|M|$ . Can we guarantee that it is also true that every element of  $|M|$  is picked out by some term or other? In other words, are there theories  $\Gamma$  all models of which are covered? The compactness theorem shows that this is not the case if  $\Gamma$  has infinite models. Here's how to see this: Let  $M$  be an infinite model of  $\Gamma$ , and let  $c$  be a constant symbol not in the language of  $\Gamma$ . Let  $\Delta$  be the set of all sentences  $c \neq t$  for  $t$  a term in the language  $\mathcal{L}$  of  $\Gamma$ , i.e.,

$$\Delta = \{c \neq t : t \in \text{Trm}(\mathcal{L})\}.$$

A finite subset of  $\Gamma \cup \Delta$  can be written as  $\Gamma' \cup \Delta'$ , with  $\Gamma' \subseteq \Gamma$  and  $\Delta' \subseteq \Delta$ . Since  $\Delta'$  is finite, it can contain only finitely many terms. Let  $a \in |M|$  be an element of  $|M|$  not picked out by any of them, and let  $M'$  be the structure that is just like  $M$ , but also  $c^{M'} = a$ . Since  $a \neq \text{Val}^M(t)$  for all  $t$  occurring in  $\Delta'$ ,  $M' \models \Delta'$ . Since  $M \models \Gamma$ ,  $\Gamma' \subseteq \Gamma$ , and  $c$  does not occur in  $\Gamma$ , also  $M' \models \Gamma'$ . Together,  $M' \models \Gamma' \cup \Delta'$  for every finite subset  $\Gamma' \cup \Delta'$  of  $\Gamma \cup \Delta$ . So every finite subset of  $\Gamma \cup \Delta$  is satisfiable. By compactness,  $\Gamma \cup \Delta$  itself is satisfiable. So there are models  $M \models \Gamma \cup \Delta$ . Every such  $M$  is a model of  $\Gamma$ , but is not covered, since  $\text{Val}^M(c) \neq \text{Val}^M(t)$  for all terms  $t$  of  $\mathcal{L}$ .

**Example 12.25.** Consider a language  $\mathcal{L}$  containing the predicate symbol  $<$ , constant symbols  $0$ ,  $1$ , and function symbols  $+$ ,  $\times$ , and  $-$ . Let  $\Gamma$  be the set of all sentences in this language true in the structure  $\mathbb{Q}$  with domain  $\mathbb{Q}$  and the obvious interpretations.  $\Gamma$  is the set of all sentences of  $\mathcal{L}$  true about the rational numbers. Of course, in  $\mathbb{Q}$  (and even in  $\mathbb{R}$ ), there are no numbers  $r$  which are greater than  $0$  but less than  $1/k$  for all  $k \in \mathbb{Z}^+$ . Such a number, if it existed, would be an *infinitesimal*: non-zero, but infinitely small. The compactness theorem can be used to show

that there are models of  $\Gamma$  in which infinitesimals exist. We do not have a function symbol for division in our language (division by zero is undefined, and function symbols have to be interpreted by total functions). However, we can still express that  $r < 1/k$ , since this is the case iff  $r \cdot k < 1$ . Now let  $c$  be a new constant symbol and let  $\Delta$  be

$$\{0 < c\} \cup \{c \times \bar{k} < 1 : k \in \mathbb{Z}^+\}$$

(where  $\bar{k} = (1 + (1 + \dots + (1 + 1) \dots))$  with  $k$  1's). For any finite subset  $\Delta_0$  of  $\Delta$  there is a  $K$  such that for all the sentences  $c \times \bar{k} < 1$  in  $\Delta_0$  have  $k < K$ . If we expand  $\mathcal{Q}$  to  $\mathcal{Q}'$  with  $c^{\mathcal{Q}'} = 1/K$  we have that  $\mathcal{Q}' \models \Gamma \cup \Delta_0$ , and so  $\Gamma \cup \Delta$  is finitely satisfiable (Exercise: prove this in detail). By compactness,  $\Gamma \cup \Delta$  is satisfiable. Any model  $\mathcal{S}$  of  $\Gamma \cup \Delta$  contains an infinitesimal, namely  $c^{\mathcal{S}}$ .

**Example 12.26.** We know that first-order logic with identity predicate can express that the size of the domain must have some minimal size: The sentence  $A_{\geq n}$  (which says “there are at least  $n$  distinct objects”) is true only in structures where  $|M|$  has at least  $n$  objects. So if we take

$$\Delta = \{A_{\geq n} : n \geq 1\}$$

then any model of  $\Delta$  must be infinite. Thus, we can guarantee that a theory only has infinite models by adding  $\Delta$  to it: the models of  $\Gamma \cup \Delta$  are all and only the infinite models of  $\Gamma$ .

So first-order logic can express infinitude. The compactness theorem shows that it cannot express finitude, however. For suppose some set of sentences  $\Lambda$  were satisfied in all and only finite structures. Then  $\Delta \cup \Lambda$  is finitely satisfiable. Why? Suppose  $\Delta' \cup \Lambda' \subseteq \Delta \cup \Lambda$  is finite with  $\Delta' \subseteq \Delta$  and  $\Lambda' \subseteq \Lambda$ . Let  $n$  be the largest number such that  $A_{\geq n} \in \Delta'$ .  $\Lambda$ , being satisfied in all finite structures, has a model  $M$  with finitely many but  $\geq n$  elements. But then  $M \models \Delta' \cup \Lambda'$ . By compactness,  $\Delta \cup \Lambda$  has an infinite model, contradicting the assumption that  $\Lambda$  is satisfied only in finite structures.



## 12.10 A Direct Proof of the Compactness Theorem

We can prove the Compactness Theorem directly, without appealing to the Completeness Theorem, using the same ideas as in the proof of the completeness theorem. In the proof of the Completeness Theorem we started with a consistent set  $\Gamma$  of sentences, expanded it to a consistent, saturated, and complete set  $\Gamma^*$  of sentences, and then showed that in the term model  $M(\Gamma^*)$  constructed from  $\Gamma^*$ , all sentences of  $\Gamma$  are true, so  $\Gamma$  is satisfiable.

We can use the same method to show that a finitely satisfiable set of sentences is satisfiable. We just have to prove the corresponding versions of the results leading to the truth lemma where we replace “consistent” with “finitely satisfiable.”

**Proposition 12.27.** *Suppose  $\Gamma$  is complete and finitely satisfiable. Then:*

1.  $(A \wedge B) \in \Gamma$  iff both  $A \in \Gamma$  and  $B \in \Gamma$ .
2.  $(A \vee B) \in \Gamma$  iff either  $A \in \Gamma$  or  $B \in \Gamma$ .
3.  $(A \rightarrow B) \in \Gamma$  iff either  $A \notin \Gamma$  or  $B \in \Gamma$ .

**Lemma 12.28.** *Every finitely satisfiable set  $\Gamma$  can be extended to a saturated finitely satisfiable set  $\Gamma'$ .*

**Proposition 12.29.** *Suppose  $\Gamma$  is complete, finitely satisfiable, and saturated.*

1.  $\exists x A(x) \in \Gamma$  iff  $A(t) \in \Gamma$  for at least one closed term  $t$ .
2.  $\forall x A(x) \in \Gamma$  iff  $A(t) \in \Gamma$  for all closed terms  $t$ .

**Lemma 12.30.** *Every finitely satisfiable set  $\Gamma$  can be extended to a complete and finitely satisfiable set  $\Gamma^*$ .*

**Theorem 12.31 (Compactness).**  *$\Gamma$  is satisfiable if and only if it is finitely satisfiable.*

*Proof.* If  $\Gamma$  is satisfiable, then there is a structure  $M$  such that  $M \models A$  for all  $A \in \Gamma$ . Of course, this  $M$  also satisfies every finite subset of  $\Gamma$ , so  $\Gamma$  is finitely satisfiable.

Now suppose that  $\Gamma$  is finitely satisfiable. By [Lemma 12.28](#), there is a finitely satisfiable, saturated set  $\Gamma' \supseteq \Gamma$ . By [Lemma 12.30](#),  $\Gamma'$  can be extended to a complete and finitely satisfiable set  $\Gamma^*$ , and  $\Gamma^*$  is still saturated. Construct the term model  $M(\Gamma^*)$  as in [Definition 12.9](#). Note that [Proposition 12.11](#) did not rely on the fact that  $\Gamma^*$  is consistent (or complete or saturated, for that matter), but just on the fact that  $M(\Gamma^*)$  is covered. The proof of the Truth Lemma ([Lemma 12.12](#)) goes through if we replace references to [Proposition 12.2](#) and [Proposition 12.7](#) by references to [Proposition 12.27](#) and [Proposition 12.29](#)  $\square$

## 12.11 The Löwenheim–Skolem Theorem

The Löwenheim–Skolem Theorem says that if a theory has an infinite model, then it also has a model that is at most countably infinite. An immediate consequence of this fact is that first-order logic cannot express that the size of a structure is uncountable: any sentence or set of sentences satisfied in all uncountable structures is also satisfied in some countable structure.

**Theorem 12.32.** *If  $\Gamma$  is consistent then it has a countable model, i.e., it is satisfiable in a structure whose domain is either finite or countably infinite.*

*Proof.* If  $\Gamma$  is consistent, the structure  $M$  delivered by the proof of the completeness theorem has a domain  $|M|$  that is no larger

than the set of the terms of the language  $\mathcal{L}$ . So  $M$  is at most countably infinite.  $\square$

**Theorem 12.33.** *If  $\Gamma$  is a consistent set of sentences in the language of first-order logic without identity, then it has a countably infinite model, i.e., it is satisfiable in a structure whose domain is infinite and countable.*

*Proof.* If  $\Gamma$  is consistent and contains no sentences in which identity appears, then the structure  $M$  delivered by the proof of the completeness theorem has a domain  $|M|$  identical to the set of terms of the language  $\mathcal{L}'$ . So  $M$  is countably infinite, since  $\text{Trm}(\mathcal{L}')$  is.  $\square$

**Example 12.34 (Skolem's Paradox).** Zermelo–Fraenkel set theory **ZFC** is a very powerful framework in which practically all mathematical statements can be expressed, including facts about the sizes of sets. So for instance, **ZFC** can prove that the set  $\mathbb{R}$  of real numbers is uncountable, it can prove Cantor's Theorem that the power set of any set is larger than the set itself, etc. If **ZFC** is consistent, its models are all infinite, and moreover, they all contain elements about which the theory says that they are uncountable, such as the element that makes true the theorem of **ZFC** that the power set of the natural numbers exists. By the Löwenheim–Skolem Theorem, **ZFC** also has countable models—models that contain “uncountable” sets but which themselves are countable.

## Summary

The **completeness theorem** is the converse of the **soundness theorem**. In one form it states that if  $\Gamma \models A$  then  $\Gamma \vdash A$ , in another that if  $\Gamma$  is consistent then it is satisfiable. We proved the second form (and derived the first from the second). The proof is involved and requires a number of steps. We start with a consistent set  $\Gamma$ . First we add infinitely many new constant symbols  $c_i$

as well as formulas of the form  $\exists x A(x) \rightarrow A(c)$  where each formula  $A(x)$  with a free variable in the expanded language is paired with one of the new constants. This results in a **saturated** consistent set of sentences containing  $\Gamma$ . It is still consistent. Now we take that set and extend it to a **complete consistent set**. A complete consistent set has the nice property that for any sentence  $A$ , either  $A$  or  $\neg A$  is in the set (but never both). Since we started from a saturated set, we now have a saturated, complete, consistent set of sentences  $\Gamma^*$  that includes  $\Gamma$ . From this set it is now possible to define a structure  $M$  such that  $M(\Gamma^*) \models A$  iff  $A \in \Gamma^*$ . In particular,  $M(\Gamma^*) \models \Gamma$ , i.e.,  $\Gamma$  is satisfiable. If  $\models$  is present, the construction is slightly more complex.

Two important corollaries follow from the completeness theorem. The **compactness theorem** states that  $\Gamma \models A$  iff  $\Gamma_0 \models A$  for some finite  $\Gamma_0 \subseteq \Gamma$ . An equivalent formulation is that  $\Gamma$  is satisfiable iff every finite  $\Gamma_0 \subseteq \Gamma$  is satisfiable. The compactness theorem is useful to prove the existence of structures with certain properties. For instance, we can use it to show that there are infinite models for every theory which has arbitrarily large finite models. This means in particular that finitude cannot be expressed in first-order logic. The second corollary, the **Löwenheim-Skolem Theorem**, states that every satisfiable  $\Gamma$  has a countable model. It in turn shows that uncountability cannot be expressed in first-order logic.

## Problems

- Problem 12.1.** Complete the proof of [Proposition 12.2](#).
- Problem 12.2.** Complete the proof of [Proposition 12.11](#).
- Problem 12.3.** Complete the proof of [Lemma 12.12](#).
- Problem 12.4.** Complete the proof of [Proposition 12.14](#).
- Problem 12.5.** Complete the proof of [Lemma 12.18](#).

**Problem 12.6.** Use **Corollary 12.21** to prove **Theorem 12.20**, thus showing that the two formulations of the completeness theorem are equivalent.

**Problem 12.7.** In order for a derivation system to be complete, its rules must be strong enough to prove every unsatisfiable set inconsistent. Which of the rules of derivation were necessary to prove completeness? Are any of these rules not used anywhere in the proof? In order to answer these questions, make a list or diagram that shows which of the rules of derivation were used in which results that lead up to the proof of **Theorem 12.20**. Be sure to note any tacit uses of rules in these proofs.

**Problem 12.8.** Prove (1) of **Theorem 12.23**.

**Problem 12.9.** In the standard model of arithmetic  $N$ , there is no element  $k \in |N|$  which satisfies every formula  $\bar{n} < x$  (where  $\bar{n}$  is  $0'\dots'$  with  $n$  's). Use the compactness theorem to show that the set of sentences in the language of arithmetic which are true in the standard model of arithmetic  $N$  are also true in a structure  $N'$  that contains an element which *does* satisfy every formula  $\bar{n} < x$ .

**Problem 12.10.** Prove **Proposition 12.27**. Avoid the use of  $\vdash$ .

**Problem 12.11.** Prove **Lemma 12.28**. (Hint: The crucial step is to show that if  $\Gamma_n$  is finitely satisfiable, so is  $\Gamma_n \cup \{D_n\}$ , without any appeal to derivations or consistency.)

**Problem 12.12.** Prove **Proposition 12.29**.

**Problem 12.13.** Prove **Lemma 12.30**. (Hint: the crucial step is to show that if  $\Gamma_n$  is finitely satisfiable, then either  $\Gamma_n \cup \{A_n\}$  or  $\Gamma_n \cup \{\neg A_n\}$  is finitely satisfiable.)

**Problem 12.14.** Write out the complete proof of the Truth Lemma (**Lemma 12.12**) in the version required for the proof of **Theorem 12.31**.

## CHAPTER 13

# *Beyond First-order Logic*

### 13.1 Overview

First-order logic is not the only system of logic of interest: there are many extensions and variations of first-order logic. A logic typically consists of the formal specification of a language, usually, but not always, a deductive system, and usually, but not always, an intended semantics. But the technical use of the term raises an obvious question: what do logics that are not first-order logic have to do with the word “logic,” used in the intuitive or philosophical sense? All of the systems described below are designed to model reasoning of some form or another; can we say what makes them logical?

No easy answers are forthcoming. The word “logic” is used in different ways and in different contexts, and the notion, like that of “truth,” has been analyzed from numerous philosophical stances. For example, one might take the goal of logical reasoning to be the determination of which statements are necessarily

true, true a priori, true independent of the interpretation of the nonlogical terms, true by virtue of their form, or true by linguistic convention; and each of these conceptions requires a good deal of clarification. Even if one restricts one's attention to the kind of logic used in mathematics, there is little agreement as to its scope. For example, in the *Principia Mathematica*, Russell and Whitehead tried to develop mathematics on the basis of logic, in the *logicist* tradition begun by Frege. Their system of logic was a form of higher-type logic similar to the one described below. In the end they were forced to introduce axioms which, by most standards, do not seem purely logical (notably, the axiom of infinity, and the axiom of reducibility), but one might nonetheless hold that some forms of higher-order reasoning should be accepted as logical. In contrast, Quine, whose ontology does not admit "propositions" as legitimate objects of discourse, argues that second-order and higher-order logic are really manifestations of set theory in sheep's clothing; in other words, systems involving quantification over predicates are not purely logical.

For now, it is best to leave such philosophical issues for a rainy day, and simply think of the systems below as formal idealizations of various kinds of reasoning, logical or otherwise.

## 13.2 Many-Sorted Logic

In first-order logic, variables and quantifiers range over a single domain. But it is often useful to have multiple (disjoint) domains: for example, you might want to have a domain of numbers, a domain of geometric objects, a domain of functions from numbers to numbers, a domain of abelian groups, and so on.

Many-sorted logic provides this kind of framework. One starts with a list of "sorts"—the "sort" of an object indicates the "domain" it is supposed to inhabit. One then has variables and quantifiers for each sort, and (usually) an identity predicate for each sort. Functions and relations are also "typed" by the sorts of objects they can take as arguments. Otherwise, one keeps the

usual rules of first-order logic, with versions of the quantifier-rules repeated for each sort.

For example, to study international relations we might choose a language with two sorts of objects, French citizens and German citizens. We might have a unary relation, “drinks wine,” for objects of the first sort; another unary relation, “eats wurst,” for objects of the second sort; and a binary relation, “forms a multinational married couple,” which takes two arguments, where the first argument is of the first sort and the second argument is of the second sort. If we use variables  $a, b, c$  to range over French citizens and  $x, y, z$  to range over German citizens, then

$$\forall a \forall x [(MarriedTo(a, x) \rightarrow (DrinksWine(a) \vee \neg EatsWurst(x)))]$$

asserts that if any French person is married to a German, either the French person drinks wine or the German doesn’t eat wurst.

Many-sorted logic can be embedded in first-order logic in a natural way, by lumping all the objects of the many-sorted domains together into one first-order domain, using unary predicate symbols to keep track of the sorts, and relativizing quantifiers. For example, the first-order language corresponding to the example above would have unary predicate symbols “*German*” and “*French*,” in addition to the other relations described, with the sort requirements erased. A sorted quantifier  $\forall x A$ , where  $x$  is a variable of the German sort, translates to

$$\forall x (German(x) \rightarrow A).$$

We need to add axioms that insure that the sorts are separate—e.g.,  $\forall x \neg (German(x) \wedge French(x))$ —as well as axioms that guarantee that “drinks wine” only holds of objects satisfying the predicate  $French(x)$ , etc. With these conventions and axioms, it is not difficult to show that many-sorted sentences translate to first-order sentences, and many-sorted derivations translate to first-order derivations. Also, many-sorted structures “translate” to corresponding first-order structures and vice-versa, so we also have a completeness theorem for many-sorted logic.



### 13.3 Second-Order logic

The language of second-order logic allows one to quantify not just over a domain of individuals, but over relations on that domain as well. Given a first-order language  $\mathcal{L}$ , for each  $k$  one adds variables  $R$  which range over  $k$ -ary relations, and allows quantification over those variables. If  $R$  is a variable for a  $k$ -ary relation, and  $t_1, \dots, t_k$  are ordinary (first-order) terms,  $R(t_1, \dots, t_k)$  is an atomic formula. Otherwise, the set of formulas is defined just as in the case of first-order logic, with additional clauses for second-order quantification. Note that we only have the identity predicate for first-order terms: if  $R$  and  $S$  are relation variables of the same arity  $k$ , we can define  $R = S$  to be an abbreviation for

$$\forall x_1 \dots \forall x_k (R(x_1, \dots, x_k) \leftrightarrow S(x_1, \dots, x_k)).$$

The rules for second-order logic simply extend the quantifier rules to the new second order variables. Here, however, one has to be a little bit careful to explain how these variables interact with the predicate symbols of  $\mathcal{L}$ , and with formulas of  $\mathcal{L}$  more generally. At the bare minimum, relation variables count as terms, so one has inferences of the form

$$A(R) \vdash \exists R A(R)$$

But if  $\mathcal{L}$  is the language of arithmetic with a constant relation symbol  $<$ , one would also expect the following inference to be valid:

$$x < y \vdash \exists R R(x, y)$$

or for a given formula  $A$ ,

$$A(x_1, \dots, x_k) \vdash \exists R R(x_1, \dots, x_k)$$

More generally, we might want to allow inferences of the form

$$A[\lambda \vec{x}. B(\vec{x})/R] \vdash \exists R A$$

where  $A[\lambda \vec{x}. B(\vec{x})/R]$  denotes the result of replacing every atomic formula of the form  $Rt_1, \dots, t_k$  in  $A$  by  $B(t_1, \dots, t_k)$ . This last rule

is equivalent to having a *comprehension schema*, i.e., an axiom of the form

$$\exists R \forall x_1, \dots, x_k (A(x_1, \dots, x_k) \leftrightarrow R(x_1, \dots, x_k)),$$

one for each formula  $A$  in the second-order language, in which  $R$  is not a free variable. (Exercise: show that if  $R$  is allowed to occur in  $A$ , this schema is inconsistent!)

When logicians refer to the “axioms of second-order logic” they usually mean the minimal extension of first-order logic by second-order quantifier rules together with the comprehension schema. But it is often interesting to study weaker subsystems of these axioms and rules. For example, note that in its full generality the axiom schema of comprehension is *impredicative*: it allows one to assert the existence of a relation  $R(x_1, \dots, x_k)$  that is “defined” by a formula with second-order quantifiers; and these quantifiers range over the set of all such relations—a set which includes  $R$  itself! Around the turn of the twentieth century, a common reaction to Russell’s paradox was to lay the blame on such definitions, and to avoid them in developing the foundations of mathematics. If one prohibits the use of second-order quantifiers in the formula  $A$ , one has a *predicative* form of comprehension, which is somewhat weaker.

From the semantic point of view, one can think of a second-order structure as consisting of a first-order structure for the language, coupled with a set of relations on the domain over which the second-order quantifiers range (more precisely, for each  $k$  there is a set of relations of arity  $k$ ). Of course, if comprehension is included in the derivation system, then we have the added requirement that there are enough relations in the “second-order part” to satisfy the comprehension axioms—otherwise the derivation system is not sound! One easy way to insure that there are enough relations around is to take the second-order part to consist of *all* the relations on the first-order part. Such a structure is called *full*, and, in a sense, is really the “intended structure” for the language. If we restrict our attention to full structures we have

what is known as the *full* second-order semantics. In that case, specifying a structure boils down to specifying the first-order part, since the contents of the second-order part follow from that implicitly.

To summarize, there is some ambiguity when talking about second-order logic. In terms of the derivation system, one might have in mind either

1. A “minimal” second-order derivation system, together with some comprehension axioms.
2. The “standard” second-order derivation system, with full comprehension.

In terms of the semantics, one might be interested in either

1. The “weak” semantics, where a structure consists of a first-order part, together with a second-order part big enough to satisfy the comprehension axioms.
2. The “standard” second-order semantics, in which one considers full structures only.

When logicians do not specify the derivation system or the semantics they have in mind, they are usually referring to the second item on each list. The advantage to using this semantics is that, as we will see, it gives us categorical descriptions of many natural mathematical structures; at the same time, the derivation system is quite strong, and sound for this semantics. The drawback is that the derivation system is *not* complete for the semantics; in fact, *no* effectively given derivation system is complete for the full second-order semantics. On the other hand, we will see that the derivation system *is* complete for the weakened semantics; this implies that if a sentence is not provable, then there is *some* structure, not necessarily the full one, in which it is false.

The language of second-order logic is quite rich. One can identify unary relations with subsets of the domain, and so in

particular you can quantify over these sets; for example, one can express induction for the natural numbers with a single axiom

$$\forall R ((R(0) \wedge \forall x (R(x) \rightarrow R(x'))) \rightarrow \forall x R(x)).$$

If one takes the language of arithmetic to have symbols  $0, \prime, +, \times$  and  $<$ , one can add the following axioms to describe their behavior:

1.  $\forall x \neg x' = 0$
2.  $\forall x \forall y (s(x) = s(y) \rightarrow x = y)$
3.  $\forall x (x + 0) = x$
4.  $\forall x \forall y (x + y') = (x + y)'$
5.  $\forall x (x \times 0) = 0$
6.  $\forall x \forall y (x \times y') = ((x \times y) + x)$
7.  $\forall x \forall y (x < y \leftrightarrow \exists z y = (x + z'))$

It is not difficult to show that these axioms, together with the axiom of induction above, provide a categorical description of the structure  $\mathbb{N}$ , the standard model of arithmetic, provided we are using the full second-order semantics. Given any structure  $M$  in which these axioms are true, define a function  $f$  from  $\mathbb{N}$  to the domain of  $M$  using ordinary recursion on  $\mathbb{N}$ , so that  $f(0) = 0^M$  and  $f(x+1) = \prime^M(f(x))$ . Using ordinary induction on  $\mathbb{N}$  and the fact that axioms (1) and (2) hold in  $M$ , we see that  $f$  is injective. To see that  $f$  is surjective, let  $P$  be the set of elements of  $|M|$  that are in the range of  $f$ . Since  $M$  is full,  $P$  is in the second-order domain. By the construction of  $f$ , we know that  $0^M$  is in  $P$ , and that  $P$  is closed under  $\prime^M$ . The fact that the induction axiom holds in  $M$  (in particular, for  $P$ ) guarantees that  $P$  is equal to the entire first-order domain of  $M$ . This shows that  $f$  is a bijection. Showing that  $f$  is a homomorphism is no more difficult, using ordinary induction on  $\mathbb{N}$  repeatedly.

In set-theoretic terms, a function is just a special kind of relation; for example, a unary function  $f$  can be identified with a binary relation  $R$  satisfying  $\forall x \exists! y R(x, y)$ . As a result, one can quantify over functions too. Using the full semantics, one can then define the class of infinite structures to be the class of structures  $M$  for which there is an injective function from the domain of  $M$  to a proper subset of itself:

$$\exists f (\forall x \forall y (f(x) = f(y) \rightarrow x = y) \wedge \exists y \forall x f(x) \neq y).$$

The negation of this sentence then defines the class of finite structures.

In addition, one can define the class of well-orderings, by adding the following to the definition of a linear ordering:

$$\forall P (\exists x P(x) \rightarrow \exists x (P(x) \wedge \forall y (y < x \rightarrow \neg P(y)))).$$

This asserts that every non-empty set has a least element, modulo the identification of “set” with “one-place relation”. For another example, one can express the notion of connectedness for graphs, by saying that there is no nontrivial separation of the vertices into disconnected parts:

$$\neg \exists A (\exists x A(x) \wedge \exists y \neg A(y) \wedge \forall w \forall z ((A(w) \wedge \neg A(z)) \rightarrow \neg R(w, z))).$$

For yet another example, you might try as an exercise to define the class of finite structures whose domain has even size. More strikingly, one can provide a categorical description of the real numbers as a complete ordered field containing the rationals.

In short, second-order logic is much more expressive than first-order logic. That’s the good news; now for the bad. We have already mentioned that there is no effective derivation system that is complete for the full second-order semantics. For better or for worse, many of the properties of first-order logic are absent, including compactness and the Löwenheim–Skolem theorems.

On the other hand, if one is willing to give up the full second-order semantics in terms of the weaker one, then the minimal

second-order derivation system is complete for this semantics. In other words, if we read  $\vdash$  as “proves in the minimal system” and  $\models$  as “logically implies in the weaker semantics”, we can show that whenever  $\Gamma \models A$  then  $\Gamma \vdash A$ . If one wants to include specific comprehension axioms in the derivation system, one has to restrict the semantics to second-order structures that satisfy these axioms: for example, if  $\Delta$  consists of a set of comprehension axioms (possibly all of them), we have that if  $\Gamma \cup \Delta \models A$ , then  $\Gamma \cup \Delta \vdash A$ . In particular, if  $A$  is not provable using the comprehension axioms we are considering, then there is a model of  $\neg A$  in which these comprehension axioms nonetheless hold.

The easiest way to see that the completeness theorem holds for the weaker semantics is to think of second-order logic as a many-sorted logic, as follows. One sort is interpreted as the ordinary “first-order” domain, and then for each  $k$  we have a domain of “relations of arity  $k$ .” We take the language to have built-in relation symbols “ $true_k(R, x_1, \dots, x_k)$ ” which is meant to assert that  $R$  holds of  $x_1, \dots, x_k$ , where  $R$  is a variable of the sort “ $k$ -ary relation” and  $x_1, \dots, x_k$  are objects of the first-order sort.

With this identification, the weak second-order semantics is essentially the usual semantics for many-sorted logic; and we have already observed that many-sorted logic can be embedded in first-order logic. Modulo the translations back and forth, then, the weaker conception of second-order logic is really a form of first-order logic in disguise, where the domain contains both “objects” and “relations” governed by the appropriate axioms.

## 13.4 Higher-Order logic

Passing from first-order logic to second-order logic enabled us to talk about sets of objects in the first-order domain, within the formal language. Why stop there? For example, third-order logic should enable us to deal with sets of sets of objects, or perhaps even sets which contain both objects and sets of objects. And fourth-order logic will let us talk about sets of objects of that kind.

As you may have guessed, one can iterate this idea arbitrarily.

In practice, higher-order logic is often formulated in terms of functions instead of relations. (Modulo the natural identifications, this difference is inessential.) Given some basic “sorts”  $A$ ,  $B$ ,  $C$ , ... (which we will now call “types”), we can create new ones by stipulating

If  $\sigma$  and  $\tau$  are finite types then so is  $\sigma \rightarrow \tau$ .

Think of types as syntactic “labels,” which classify the objects we want in our domain;  $\sigma \rightarrow \tau$  describes those objects that are functions which take objects of type  $\sigma$  to objects of type  $\tau$ . For example, we might want to have a type  $\Omega$  of truth values, “true” and “false,” and a type  $\mathbb{N}$  of natural numbers. In that case, you can think of objects of type  $\mathbb{N} \rightarrow \Omega$  as unary relations, or subsets of  $\mathbb{N}$ ; objects of type  $\mathbb{N} \rightarrow \mathbb{N}$  are functions from natural numbers to natural numbers; and objects of type  $(\mathbb{N} \rightarrow \mathbb{N}) \rightarrow \mathbb{N}$  are “functionals,” that is, higher-type functions that take functions to numbers.

As in the case of second-order logic, one can think of higher-order logic as a kind of many-sorted logic, where there is a sort for each type of object we want to consider. But it is usually clearer just to define the syntax of higher-type logic from the ground up. For example, we can define a set of finite types inductively, as follows:

1.  $\mathbb{N}$  is a finite type.
2. If  $\sigma$  and  $\tau$  are finite types, then so is  $\sigma \rightarrow \tau$ .
3. If  $\sigma$  and  $\tau$  are finite types, so is  $\sigma \times \tau$ .

Intuitively,  $\mathbb{N}$  denotes the type of the natural numbers,  $\sigma \rightarrow \tau$  denotes the type of functions from  $\sigma$  to  $\tau$ , and  $\sigma \times \tau$  denotes the type of pairs of objects, one from  $\sigma$  and one from  $\tau$ . We can then define a set of terms inductively, as follows:

1. For each type  $\sigma$ , there is a stock of variables  $x, y, z, \dots$  of type  $\sigma$

2.  $0$  is a term of type  $\mathbb{N}$
3.  $S$  (successor) is a term of type  $\mathbb{N} \rightarrow \mathbb{N}$
4. If  $s$  is a term of type  $\sigma$ , and  $t$  is a term of type  $\mathbb{N} \rightarrow (\sigma \rightarrow \sigma)$ , then  $R_{st}$  is a term of type  $\mathbb{N} \rightarrow \sigma$
5. If  $s$  is a term of type  $\tau \rightarrow \sigma$  and  $t$  is a term of type  $\tau$ , then  $s(t)$  is a term of type  $\sigma$
6. If  $s$  is a term of type  $\sigma$  and  $x$  is a variable of type  $\tau$ , then  $\lambda x. s$  is a term of type  $\tau \rightarrow \sigma$ .
7. If  $s$  is a term of type  $\sigma$  and  $t$  is a term of type  $\tau$ , then  $\langle s, t \rangle$  is a term of type  $\sigma \times \tau$ .
8. If  $s$  is a term of type  $\sigma \times \tau$  then  $p_1(s)$  is a term of type  $\sigma$  and  $p_2(s)$  is a term of type  $\tau$ .

Intuitively,  $R_{st}$  denotes the function defined recursively by

$$\begin{aligned} R_{st}(0) &= s \\ R_{st}(x+1) &= t(x, R_{st}(x)), \end{aligned}$$

$\langle s, t \rangle$  denotes the pair whose first component is  $s$  and whose second component is  $t$ , and  $p_1(s)$  and  $p_2(s)$  denote the first and second elements (“projections”) of  $s$ . Finally,  $\lambda x. s$  denotes the function  $f$  defined by

$$f(x) = s$$

for any  $x$  of type  $\sigma$ ; so item (6) gives us a form of comprehension, enabling us to define functions using terms. Formulas are built up from identity predicate statements  $s = t$  between terms of the same type, the usual propositional connectives, and higher-type quantification. One can then take the axioms of the system to be the basic equations governing the terms defined above, together with the usual rules of logic with quantifiers and identity predicate.



If one augments the finite type system with a type  $\Omega$  of truth values, one has to include axioms which govern its use as well. In fact, if one is clever, one can get rid of complex formulas entirely, replacing them with terms of type  $\Omega$ ! The proof system can then be modified accordingly. The result is essentially the *simple theory of types* set forth by Alonzo Church in the 1930s.

As in the case of second-order logic, there are different versions of higher-type semantics that one might want to use. In the full version, variables of type  $\sigma \rightarrow \tau$  range over the set of *all* functions from the objects of type  $\sigma$  to objects of type  $\tau$ . As you might expect, this semantics is too strong to admit a complete, effective derivation system. But one can consider a weaker semantics, in which a structure consists of sets of elements  $T_\tau$  for each type  $\tau$ , together with appropriate operations for application, projection, etc. If the details are carried out correctly, one can obtain completeness theorems for the kinds of derivation systems described above.

Higher-type logic is attractive because it provides a framework in which we can embed a good deal of mathematics in a natural way: starting with  $\mathbb{N}$ , one can define real numbers, continuous functions, and so on. It is also particularly attractive in the context of intuitionistic logic, since the types have clear “constructive” interpretations. In fact, one can develop constructive versions of higher-type semantics (based on intuitionistic, rather than classical logic) that clarify these constructive interpretations quite nicely, and are, in many ways, more interesting than the classical counterparts.

## 13.5 Intuitionistic Logic

In contrast to second-order and higher-order logic, intuitionistic first-order logic represents a restriction of the classical version, intended to model a more “constructive” kind of reasoning. The following examples may serve to illustrate some of the underlying motivations.

Suppose someone came up to you one day and announced that they had determined a natural number  $x$ , with the property that if  $x$  is prime, the Riemann hypothesis is true, and if  $x$  is composite, the Riemann hypothesis is false. Great news! Whether the Riemann hypothesis is true or not is one of the big open questions of mathematics, and here they seem to have reduced the problem to one of calculation, that is, to the determination of whether a specific number is prime or not.

What is the magic value of  $x$ ? They describe it as follows:  $x$  is the natural number that is equal to 7 if the Riemann hypothesis is true, and 9 otherwise.

Angrily, you demand your money back. From a classical point of view, the description above does in fact determine a unique value of  $x$ ; but what you really want is a value of  $x$  that is given *explicitly*.

To take another, perhaps less contrived example, consider the following question. We know that it is possible to raise an irrational number to a rational power, and get a rational result. For example,  $\sqrt{2}^2 = 2$ . What is less clear is whether or not it is possible to raise an irrational number to an *irrational* power, and get a rational result. The following theorem answers this in the affirmative:

**Theorem 13.1.** *There are irrational numbers  $a$  and  $b$  such that  $a^b$  is rational.*

*Proof.* Consider  $\sqrt{2}^{\sqrt{2}}$ . If this is rational, we are done: we can let  $a = b = \sqrt{2}$ . Otherwise, it is irrational. Then we have

$$(\sqrt{2}^{\sqrt{2}})^{\sqrt{2}} = \sqrt{2}^{\sqrt{2} \cdot \sqrt{2}} = \sqrt{2}^2 = 2,$$

which is certainly rational. So, in this case, let  $a$  be  $\sqrt{2}^{\sqrt{2}}$ , and let  $b$  be  $\sqrt{2}$ . □

Does this constitute a valid proof? Most mathematicians feel that it does. But again, there is something a little bit unsatisfying

here: we have proved the existence of a pair of real numbers with a certain property, without being able to say *which* pair of numbers it is. It is possible to prove the same result, but in such a way that the pair  $a, b$  is given in the proof: take  $a = \sqrt{3}$  and  $b = \log_3 4$ . Then

$$a^b = \sqrt{3}^{\log_3 4} = 3^{1/2 \cdot \log_3 4} = (3^{\log_3 4})^{1/2} = 4^{1/2} = 2,$$

since  $3^{\log_3 x} = x$ .

Intuitionistic logic is designed to model a kind of reasoning where moves like the one in the first proof are disallowed. Proving the existence of an  $x$  satisfying  $A(x)$  means that you have to give a specific  $x$ , and a proof that it satisfies  $A$ , like in the second proof. Proving that  $A$  or  $B$  holds requires that you can prove one or the other.

Formally speaking, intuitionistic first-order logic is what you get if you restrict a derivation system for first-order logic in a certain way. Similarly, there are intuitionistic versions of second-order or higher-order logic. From the mathematical point of view, these are just formal deductive systems, but, as already noted, they are intended to model a kind of mathematical reasoning. One can take this to be the kind of reasoning that is justified on a certain philosophical view of mathematics (such as Brouwer's intuitionism); one can take it to be a kind of mathematical reasoning which is more "concrete" and satisfying (along the lines of Bishop's constructivism); and one can argue about whether or not the formal description captures the informal motivation. But whatever philosophical positions we may hold, we can study intuitionistic logic as a formally presented logic; and for whatever reasons, many mathematical logicians find it interesting to do so.

There is an informal constructive interpretation of the intuitionist connectives, usually known as the BHK interpretation (named after Brouwer, Heyting, and Kolmogorov). It runs as follows: a proof of  $A \wedge B$  consists of a proof of  $A$  paired with a proof of  $B$ ; a proof of  $A \vee B$  consists of either a proof of  $A$ , or a proof of  $B$ , where we have explicit information as to which is the

case; a proof of  $A \rightarrow B$  consists of a procedure, which transforms a proof of  $A$  to a proof of  $B$ ; a proof of  $\forall x A(x)$  consists of a procedure which returns a proof of  $A(x)$  for any value of  $x$ ; and a proof of  $\exists x A(x)$  consists of a value of  $x$ , together with a proof that this value satisfies  $A$ . One can describe the interpretation in computational terms known as the “Curry–Howard isomorphism” or the “formulas-as-types paradigm”: think of a formula as specifying a certain kind of data type, and proofs as computational objects of these data types that enable us to see that the corresponding formula is true.

Intuitionistic logic is often thought of as being classical logic “minus” the law of the excluded middle. This following theorem makes this more precise.

**Theorem 13.2.** *Intuitionistically, the following axiom schemata are equivalent:*

1.  $(\neg A \rightarrow \perp) \rightarrow A$ .
2.  $A \vee \neg A$
3.  $\neg\neg A \rightarrow A$

Obtaining instances of one schema from either of the others is a good exercise in intuitionistic logic.

The first deductive systems for intuitionistic propositional logic, put forth as formalizations of Brouwer’s intuitionism, are due, independently, to Kolmogorov, Glivenko, and Heyting. The first formalization of intuitionistic first-order logic (and parts of intuitionist mathematics) is due to Heyting. Though a number of classically valid schemata are not intuitionistically valid, many are.

The *double-negation translation* describes an important relationship between classical and intuitionist logic. It is defined inductively follows (think of  $A^N$  as the “intuitionist” translation of the classical formula  $A$ ):

$$A^N \equiv \neg\neg A \quad \text{for atomic formulas } A$$

$$\begin{aligned} (A \wedge B)^N &\equiv (A^N \wedge B^N) \\ (A \vee B)^N &\equiv \neg\neg(A^N \vee B^N) \\ (A \rightarrow B)^N &\equiv (A^N \rightarrow B^N) \\ (\forall x A)^N &\equiv \forall x A^N \\ (\exists x A)^N &\equiv \neg\neg\exists x A^N \end{aligned}$$

Kolmogorov and Glivenko had versions of this translation for propositional logic; for predicate logic, it is due to Gödel and Gentzen, independently. We have

**Theorem 13.3.**    1.  $A \leftrightarrow A^N$  is provable classically  
 2. If  $A$  is provable classically, then  $A^N$  is provable intuitionistically.

We can now envision the following dialogue. Classical mathematician: “I’ve proved  $A$ !” Intuitionist mathematician: “Your proof isn’t valid. What you’ve really proved is  $A^N$ .” Classical mathematician: “Fine by me!” As far as the classical mathematician is concerned, the intuitionist is just splitting hairs, since the two are equivalent. But the intuitionist insists there is a difference.

Note that the above translation concerns pure logic only; it does not address the question as to what the appropriate *nonlogical* axioms are for classical and intuitionistic mathematics, or what the relationship is between them. But the following slight extension of the theorem above provides some useful information:

**Theorem 13.4.** If  $\Gamma$  proves  $A$  classically,  $\Gamma^N$  proves  $A^N$  intuitionistically.

In other words, if  $A$  is provable from some hypotheses classically, then  $A^N$  is provable from their double-negation translations.

To show that a sentence or propositional formula is intuitionistically valid, all you have to do is provide a proof. But how can

you show that it is not valid? For that purpose, we need a semantics that is sound, and preferably complete. A semantics due to Kripke nicely fits the bill.

We can play the same game we did for classical logic: define the semantics, and prove soundness and completeness. It is worthwhile, however, to note the following distinction. In the case of classical logic, the semantics was the “obvious” one, in a sense implicit in the meaning of the connectives. Though one can provide some intuitive motivation for Kripke semantics, the latter does not offer the same feeling of inevitability. In addition, the notion of a classical structure is a natural mathematical one, so we can either take the notion of a structure to be a tool for studying classical first-order logic, or take classical first-order logic to be a tool for studying mathematical structures. In contrast, Kripke structures can only be viewed as a logical construct; they don’t seem to have independent mathematical interest.

A Kripke structure  $\mathfrak{M} = \langle W, R, V \rangle$  for a propositional language consists of a set  $W$ , partial order  $R$  on  $W$  with a least element, and an “monotone” assignment of propositional variables to the elements of  $W$ . The intuition is that the elements of  $W$  represent “worlds,” or “states of knowledge”; an element  $v \geq u$  represents a “possible future state” of  $u$ ; and the propositional variables assigned to  $u$  are the propositions that are known to be true in state  $u$ . The forcing relation  $\mathfrak{M}, w \Vdash A$  then extends this relationship to arbitrary formulas in the language; read  $\mathfrak{M}, w \Vdash A$  as “ $A$  is true in state  $w$ .” The relationship is defined inductively, as follows:

1.  $\mathfrak{M}, w \Vdash p_i$  iff  $p_i$  is one of the propositional variables assigned to  $w$ .
2.  $\mathfrak{M}, w \not\Vdash \perp$ .
3.  $\mathfrak{M}, w \Vdash (A \wedge B)$  iff  $\mathfrak{M}, w \Vdash A$  and  $\mathfrak{M}, w \Vdash B$ .
4.  $\mathfrak{M}, w \Vdash (A \vee B)$  iff  $\mathfrak{M}, w \Vdash A$  or  $\mathfrak{M}, w \Vdash B$ .

5.  $\mathfrak{M}, w \Vdash (A \rightarrow B)$  iff, whenever  $w' \geq w$  and  $\mathfrak{M}, w' \Vdash A$ , then  $\mathfrak{M}, w' \Vdash B$ .

It is a good exercise to try to show that  $\neg(p \wedge q) \rightarrow (\neg p \vee \neg q)$  is not intuitionistically valid, by cooking up a Kripke structure that provides a counterexample.

## 13.6 Modal Logics

Consider the following example of a conditional sentence:

If Jeremy is alone in that room, then he is drunk and naked and dancing on the chairs.

This is an example of a conditional assertion that may be materially true but nonetheless misleading, since it seems to suggest that there is a stronger link between the antecedent and conclusion other than simply that either the antecedent is false or the consequent true. That is, the wording suggests that the claim is not only true in this particular world (where it may be trivially true, because Jeremy is not alone in the room), but that, moreover, the conclusion *would have* been true *had* the antecedent been true. In other words, one can take the assertion to mean that the claim is true not just in this world, but in any “possible” world; or that it is *necessarily* true, as opposed to just true in this particular world.

Modal logic was designed to make sense of this kind of necessity. One obtains modal propositional logic from ordinary propositional logic by adding a box operator; which is to say, if  $A$  is a formula, so is  $\Box A$ . Intuitively,  $\Box A$  asserts that  $A$  is *necessarily* true, or true in any possible world.  $\Diamond A$  is usually taken to be an abbreviation for  $\neg \Box \neg A$ , and can be read as asserting that  $A$  is *possibly* true. Of course, modality can be added to predicate logic as well.

Kripke structures can be used to provide a semantics for modal logic; in fact, Kripke first designed this semantics with

modal logic in mind. Rather than restricting to partial orders, more generally one has a set of “possible worlds,”  $P$ , and a binary “accessibility” relation  $R(x, y)$  between worlds. Intuitively,  $R(p, q)$  asserts that the world  $q$  is compatible with  $p$ ; i.e., if we are “in” world  $p$ , we have to entertain the possibility that the world could have been like  $q$ .

Modal logic is sometimes called an “intensional” logic, as opposed to an “extensional” one. The intended semantics for an extensional logic, like classical logic, will only refer to a single world, the “actual” one; while the semantics for an “intensional” logic relies on a more elaborate ontology. In addition to structuring necessity, one can use modality to structure other linguistic constructions, reinterpreting  $\Box$  and  $\Diamond$  according to the application. For example:

1. In provability logic,  $\Box A$  is read “ $A$  is provable” and  $\Diamond A$  is read “ $A$  is consistent.”
2. In epistemic logic, one might read  $\Box A$  as “I know  $A$ ” or “I believe  $A$ .”
3. In temporal logic, one can read  $\Box A$  as “ $A$  is always true” and  $\Diamond A$  as “ $A$  is sometimes true.”

One would like to augment logic with rules and axioms dealing with modality. For example, the system **S4** consists of the ordinary axioms and rules of propositional logic, together with the following axioms:

$$\Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$$

$$\Box A \rightarrow A$$

$$\Box A \rightarrow \Box \Box A$$

as well as a rule, “from  $A$  conclude  $\Box A$ .” **S5** adds the following axiom:

$$\Diamond A \rightarrow \Box \Diamond A$$



Variations of these axioms may be suitable for different applications; for example,  $S_5$  is usually taken to characterize the notion of logical necessity. And the nice thing is that one can usually find a semantics for which the derivation system is sound and complete by restricting the accessibility relation in the Kripke structures in natural ways. For example, **S4** corresponds to the class of Kripke structures in which the accessibility relation is reflexive and transitive. **S5** corresponds to the class of Kripke structures in which the accessibility relation is *universal*, which is to say that every world is accessible from every other; so  $\Box A$  holds if and only if  $A$  holds in every world.

## 13.7 Other Logics

As you may have gathered by now, it is not hard to design a new logic. You too can create your own a syntax, make up a deductive system, and fashion a semantics to go with it. You might have to be a bit clever if you want the derivation system to be complete for the semantics, and it might take some effort to convince the world at large that your logic is truly interesting. But, in return, you can enjoy hours of good, clean fun, exploring your logic's mathematical and computational properties.

Recent decades have witnessed a veritable explosion of formal logics. Fuzzy logic is designed to model reasoning about vague properties. Probabilistic logic is designed to model reasoning about uncertainty. Default logics and nonmonotonic logics are designed to model defeasible forms of reasoning, which is to say, "reasonable" inferences that can later be overturned in the face of new information. There are epistemic logics, designed to model reasoning about knowledge; causal logics, designed to model reasoning about causal relationships; and even "deontic" logics, which are designed to model reasoning about moral and ethical obligations. Depending on whether the primary motivation for introducing these systems is philosophical, mathematical, or computational, you may find such creatures studied under the

rubric of mathematical logic, philosophical logic, artificial intelligence, cognitive science, or elsewhere.

The list goes on and on, and the possibilities seem endless. We may never attain Leibniz' dream of reducing all of human reason to calculation—but that can't stop us from trying.

**PART III**

*Incompleteness*

## CHAPTER 14

# *Introduction to Incompleteness*

### 14.1 Historical Background

In this section, we will briefly discuss historical developments that will help put the incompleteness theorems in context. In particular, we will give a very sketchy overview of the history of mathematical logic; and then say a few words about the history of the foundations of mathematics.

The phrase “mathematical logic” is ambiguous. One can interpret the word “mathematical” as describing the subject matter, as in, “the logic of mathematics,” denoting the principles of mathematical reasoning; or as describing the methods, as in “the mathematics of logic,” denoting a mathematical study of the principles of reasoning. The account that follows involves mathematical logic in both senses, often at the same time.

The study of logic began, essentially, with Aristotle, who lived approximately 384–322 BCE. His *Categories*, *Prior analytics*, and *Posterior analytics* include systematic studies of the principles of scientific reasoning, including a thorough and systematic study of the syllogism.

Aristotle’s logic dominated scholastic philosophy through the middle ages; indeed, as late as the eighteenth century, Kant main-

tained that Aristotle's logic was perfect and in no need of revision. But the theory of the syllogism is far too limited to model anything but the most superficial aspects of mathematical reasoning. A century earlier, Leibniz, a contemporary of Newton's, imagined a complete "calculus" for logical reasoning, and made some rudimentary steps towards designing such a calculus, essentially describing a version of propositional logic.

The nineteenth century was a watershed for logic. In 1854 George Boole wrote *The Laws of Thought*, with a thorough algebraic study of propositional logic that is not far from modern presentations. In 1879 Gottlob Frege published his *Begriffsschrift* (Concept writing) which extends propositional logic with quantifiers and relations, and thus includes first-order logic. In fact, Frege's logical systems included higher-order logic as well, and more. In his *Basic Laws of Arithmetic*, Frege set out to show that all of arithmetic could be derived in his *Begriffsschrift* from purely logical assumption. Unfortunately, these assumptions turned out to be inconsistent, as Russell showed in 1902. But setting aside the inconsistent axiom, Frege more or less invented modern logic singlehandedly, a startling achievement. Quantificational logic was also developed independently by algebraically-minded thinkers after Boole, including Peirce and Schröder.

Let us now turn to developments in the foundations of mathematics. Of course, since logic plays an important role in mathematics, there is a good deal of interaction with the developments just described. For example, Frege developed his logic with the explicit purpose of showing that all of mathematics could be based solely on his logical framework; in particular, he wished to show that mathematics consists of a priori *analytic* truths instead of, as Kant had maintained, a priori *synthetic* ones.

Many take the birth of mathematics proper to have occurred with the Greeks. Euclid's *Elements*, written around 300 B.C., is already a mature representative of Greek mathematics, with its emphasis on rigor and precision. The definitions and proofs in Euclid's *Elements* survive more or less intact in high school geometry textbooks today (to the extent that geometry is still taught in

high schools). This model of mathematical reasoning has been held to be a paradigm for rigorous argumentation not only in mathematics but in branches of philosophy as well. (Spinoza even presented moral and religious arguments in the Euclidean style, which is strange to see!)

Calculus was invented by Newton and Leibniz in the seventeenth century. (A fierce priority dispute raged for centuries, but most scholars today hold that the two developments were for the most part independent.) Calculus involves reasoning about, for example, infinite sums of infinitely small quantities; these features fueled criticism by Bishop Berkeley, who argued that belief in God was no less rational than the mathematics of his time. The methods of calculus were widely used in the eighteenth century, for example by Leonhard Euler, who used calculations involving infinite sums with dramatic results.

In the nineteenth century, mathematicians tried to address Berkeley's criticisms by putting calculus on a firmer foundation. Efforts by Cauchy, Weierstrass, Bolzano, and others led to our contemporary definitions of limits, continuity, differentiation, and integration in terms of "epsilon and deltas," in other words, devoid of any reference to infinitesimals. Later in the century, mathematicians tried to push further, and explain all aspects of calculus, including the real numbers themselves, in terms of the natural numbers. (Kronecker: "God created the whole numbers, all else is the work of man.") In 1872, Dedekind wrote "Continuity and the irrational numbers," where he showed how to "construct" the real numbers as sets of rational numbers (which, as you know, can be viewed as pairs of natural numbers); in 1888 he wrote "Was sind und was sollen die Zahlen" (roughly, "What are the natural numbers, and what should they be?") which aimed to explain the natural numbers in purely "logical" terms. In 1887 Kronecker wrote "Über den Zahlbegriff" ("On the concept of number") where he spoke of representing all mathematical object in terms of the integers; in 1889 Giuseppe Peano gave formal, symbolic axioms for the natural numbers.

The end of the nineteenth century also brought a new bold-

ness in dealing with the infinite. Before then, infinitary objects and structures (like the set of natural numbers) were treated gingerly; “infinitely many” was understood as “as many as you want,” and “approaches in the limit” was understood as “gets as close as you want.” But Georg Cantor showed that it was possible to take the infinite at face value. Work by Cantor, Dedekind, and others help to introduce the general set-theoretic understanding of mathematics that is now widely accepted.

This brings us to twentieth century developments in logic and foundations. In 1902 Russell discovered the paradox in Frege’s logical system. In 1904 Zermelo proved Cantor’s well-ordering principle, using the so-called “axiom of choice”; the legitimacy of this axiom prompted a good deal of debate. Between 1910 and 1913 the three volumes of Russell and Whitehead’s *Principia Mathematica* appeared, extending the Fregean program of establishing mathematics on logical grounds. Unfortunately, Russell and Whitehead were forced to adopt two principles that seemed hard to justify as purely logical: an axiom of infinity and an axiom of “reducibility.” In the 1900’s Poincaré criticized the use of “impredicative definitions” in mathematics, and in the 1910’s Brouwer began proposing to refound all of mathematics in an “intuitionistic” basis, which avoided the use of the law of the excluded middle ( $A \vee \neg A$ ).

Strange days indeed! The program of reducing all of mathematics to logic is now referred to as “logicism,” and is commonly viewed as having failed, due to the difficulties mentioned above. The program of developing mathematics in terms of intuitionistic mental constructions is called “intuitionism,” and is viewed as posing overly severe restrictions on everyday mathematics. Around the turn of the century, David Hilbert, one of the most influential mathematicians of all time, was a strong supporter of the new, abstract methods introduced by Cantor and Dedekind: “no one will drive us from the paradise that Cantor has created for us.” At the same time, he was sensitive to foundational criticisms of these new methods (oddly enough, now called “classical”). He proposed a way of having one’s cake and eating

it too:

1. Represent classical methods with formal axioms and rules; represent mathematical questions as formulas in an axiomatic system.
2. Use safe, “finitary” methods to prove that these formal deductive systems are consistent.

Hilbert’s work went a long way toward accomplishing the first goal. In 1899, he had done this for geometry in his celebrated book *Foundations of geometry*. In subsequent years, he and a number of his students and collaborators worked on other areas of mathematics to do what Hilbert had done for geometry. Hilbert himself gave axiom systems for arithmetic and analysis. Zermelo gave an axiomatization of set theory, which was expanded on by Fraenkel, Skolem, von Neumann, and others. By the mid-1920s, there were two approaches that laid claim to the title of an axiomatization of “all” of mathematics, the *Principia mathematica* of Russell and Whitehead, and what came to be known as Zermelo–Fraenkel set theory.

In 1921, Hilbert set out on a research project to establish the goal of proving these systems to be consistent. He was aided in this project by several of his students, in particular Bernays, Ackermann, and later Gentzen. The basic idea for accomplishing this goal was to cast the question of the possibility of a derivation of an inconsistency in mathematics as a combinatorial problem about possible sequences of symbols, namely possible sequences of sentences which meet the criterion of being a correct derivation of, say,  $A \wedge \neg A$  from the axioms of an axiom system for arithmetic, analysis, or set theory. A proof of the impossibility of such a sequence of symbols would—since it is itself a mathematical proof—be formalizable in these axiomatic systems. In other words, there would be some sentence  $\text{Con}$  which states that, say, arithmetic is consistent. Moreover, this sentence should be provable in the systems in question, especially if its proof requires only very restricted, “finitary” means.



The second aim, that the axiom systems developed would settle every mathematical question, can be made precise in two ways. In one way, we can formulate it as follows: For any sentence  $A$  in the language of an axiom system for mathematics, either  $A$  or  $\neg A$  is provable from the axioms. If this were true, then there would be no sentences which can neither be proved nor refuted on the basis of the axioms, no questions which the axioms do not settle. An axiom system with this property is called *complete*. Of course, for any given sentence it might still be a difficult task to determine which of the two alternatives holds. But in principle there should be a method to do so. In fact, for the axiom and derivation systems considered by Hilbert, completeness would imply that such a method exists—although Hilbert did not realize this. The second way to interpret the question would be this stronger requirement: that there be a mechanical, computational method which would determine, for a given sentence  $A$ , whether it is derivable from the axioms or not.

In 1931, Gödel proved the two “incompleteness theorems,” which showed that this program could not succeed. There is no axiom system for mathematics which is complete, specifically, the sentence that expresses the consistency of the axioms is a sentence which can neither be proved nor refuted.

This struck a lethal blow to Hilbert’s original program. However, as is so often the case in mathematics, it also opened up exciting new avenues for research. If there is no one, all-encompassing formal system of mathematics, it makes sense to develop more circumscribed systems and investigate what can be proved in them. It also makes sense to develop less restricted methods of proof for establishing the consistency of these systems, and to find ways to measure how hard it is to prove their consistency. Since Gödel showed that (almost) every formal system has questions it cannot settle, it makes sense to look for “interesting” questions a given formal system cannot settle, and to figure out how strong a formal system has to be to settle them. To the present day, logicians have been pursuing these questions in a new mathematical discipline, the theory of proofs.

## 14.2 Definitions

In order to carry out Hilbert's project of formalizing mathematics and showing that such a formalization is consistent and complete, the first order of business would be that of picking a language, logical framework, and a system of axioms. For our purposes, let us suppose that mathematics can be formalized in a first-order language, i.e., that there is some set of constant symbols, function symbols, and predicate symbols which, together with the connectives and quantifiers of first-order logic, allow us to express the claims of mathematics. Most people agree that such a language exists: the language of set theory, in which  $\in$  is the only non-logical symbol. That such a simple language is so expressive is of course a very implausible claim at first sight, and it took a lot of work to establish that practically of all mathematics can be expressed in this very austere vocabulary. To keep things simple, for now, let's restrict our discussion to arithmetic, so the part of mathematics that just deals with the natural numbers  $\mathbb{N}$ . The natural language in which to express facts of arithmetic is  $\mathcal{L}_A$ .  $\mathcal{L}_A$  contains a single two-place predicate symbol  $<$ , a single constant symbol  $0$ , one one-place function symbol  $\iota$ , and two two-place function symbols  $+$  and  $\times$ .

**Definition 14.1.** A set of sentences  $\Gamma$  is a *theory* if it is closed under entailment, i.e., if  $\Gamma = \{A : \Gamma \vDash A\}$ .

There are two easy ways to specify theories. One is as the set of sentences true in some structure. For instance, consider the structure for  $\mathcal{L}_A$  in which the domain is  $\mathbb{N}$  and all non-logical symbols are interpreted as you would expect.

**Definition 14.2.** The *standard model of arithmetic* is the structure  $N$  defined as follows:

1.  $|N| = \mathbb{N}$

2.  $0^N = 0$
3.  $r^N(n) = n + 1$  for all  $n \in \mathbb{N}$
4.  $+^N(n, m) = n + m$  for all  $n, m \in \mathbb{N}$
5.  $\times^N(n, m) = n \cdot m$  for all  $n, m \in \mathbb{N}$
6.  $<^N = \{\langle n, m \rangle : n \in \mathbb{N}, m \in \mathbb{N}, n < m\}$

Note the difference between  $\times$  and  $\cdot$ :  $\times$  is a symbol in the language of arithmetic. Of course, we've chosen it to remind us of multiplication, but  $\times$  is not the multiplication operation but a two-place function symbol (officially,  $f_1^2$ ). By contrast,  $\cdot$  is the ordinary multiplication function. When you see something like  $n \cdot m$ , we mean the product of the numbers  $n$  and  $m$ ; when you see something like  $x \times y$  we are talking about a term in the language of arithmetic. In the standard model, the function symbol times is interpreted as the function  $\cdot$  on the natural numbers. For addition, we use  $+$  as both the function symbol of the language of arithmetic, and the addition function on the natural numbers. Here you have to use the context to determine what is meant.

**Definition 14.3.** The theory of *true arithmetic* is the set of sentences satisfied in the standard model of arithmetic, i.e.,

$$\mathbf{TA} = \{A : \mathbb{N} \models A\}.$$

$\mathbf{TA}$  is a theory, for whenever  $\mathbf{TA} \models A$ ,  $A$  is satisfied in every structure which satisfies  $\mathbf{TA}$ . Since  $M \models \mathbf{TA}$ ,  $M \models A$ , and so  $A \in \mathbf{TA}$ .

The other way to specify a theory  $\Gamma$  is as the set of sentences entailed by some set of sentences  $\Gamma_0$ . In that case,  $\Gamma$  is the “closure” of  $\Gamma_0$  under entailment. Specifying a theory this way is only interesting if  $\Gamma_0$  is explicitly specified, e.g., if the elements of  $\Gamma_0$  are listed. At the very least,  $\Gamma_0$  has to be decidable, i.e., there has to be a computable test for when a sentence counts as an

element of  $\Gamma_0$  or not. We call the sentences in  $\Gamma_0$  *axioms* for  $\Gamma$ , and  $\Gamma$  *axiomatized* by  $\Gamma_0$ .

**Definition 14.4.** A theory  $\Gamma$  is *axiomatized* by  $\Gamma_0$  iff

$$\Gamma = \{A : \Gamma_0 \vDash A\}$$

**Definition 14.5.** The theory  $\mathbf{Q}$  axiomatized by the following sentences is known as “Robinson’s  $\mathbf{Q}$ ” and is a very simple theory of arithmetic.

$$\forall x \forall y (x' = y' \rightarrow x = y) \quad (Q_1)$$

$$\forall x 0 \neq x' \quad (Q_2)$$

$$\forall x (x = 0 \vee \exists y x = y') \quad (Q_3)$$

$$\forall x (x + 0) = x \quad (Q_4)$$

$$\forall x \forall y (x + y') = (x + y)' \quad (Q_5)$$

$$\forall x (x \times 0) = 0 \quad (Q_6)$$

$$\forall x \forall y (x \times y') = ((x \times y) + x) \quad (Q_7)$$

$$\forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y) \quad (Q_8)$$

The set of sentences  $\{Q_1, \dots, Q_8\}$  are the axioms of  $\mathbf{Q}$ , so  $\mathbf{Q}$  consists of all sentences entailed by them:

$$\mathbf{Q} = \{A : \{Q_1, \dots, Q_8\} \vDash A\}.$$

**Definition 14.6.** Suppose  $A(x)$  is a formula in  $\mathcal{L}_A$  with free variables  $x$  and  $y_1, \dots, y_n$ . Then any sentence of the form

$$\forall y_1 \dots \forall y_n ((A(0) \wedge \forall x (A(x) \rightarrow A(x'))) \rightarrow \forall x A(x))$$

is an instance of the *induction schema*.

*Peano arithmetic*  $\mathbf{PA}$  is the theory axiomatized by the axioms

of  $\mathbf{Q}$  together with all instances of the induction schema.

Every instance of the induction schema is true in  $N$ . This is easiest to see if the formula  $A$  only has one free variable  $x$ . Then  $A(x)$  defines a subset  $X_A$  of  $\mathbb{N}$  in  $N$ .  $X_A$  is the set of all  $n \in \mathbb{N}$  such that  $N, s \models A(x)$  when  $s(x) = n$ . The corresponding instance of the induction schema is

$$((A(0) \wedge \forall x (A(x) \rightarrow A(x')))) \rightarrow \forall x A(x)).$$

If its antecedent is true in  $N$ , then  $0 \in X_A$  and, whenever  $n \in X_A$ , so is  $n + 1$ . Since  $0 \in X_A$ , we get  $1 \in X_A$ . With  $1 \in X_A$  we get  $2 \in X_A$ . And so on. So for every  $n \in \mathbb{N}$ ,  $n \in X_A$ . But this means that  $\forall x A(x)$  is satisfied in  $N$ .

Both  $\mathbf{Q}$  and  $\mathbf{PA}$  are axiomatized theories. The big question is, how strong are they? For instance, can  $\mathbf{PA}$  prove all the truths about  $\mathbb{N}$  that can be expressed in  $\mathcal{L}_A$ ? Specifically, do the axioms of  $\mathbf{PA}$  settle all the questions that can be formulated in  $\mathcal{L}_A$ ?

Another way to put this is to ask: Is  $\mathbf{PA} = \mathbf{TA}$ ?  $\mathbf{TA}$  obviously does prove (i.e., it includes) all the truths about  $\mathbb{N}$ , and it settles all the questions that can be formulated in  $\mathcal{L}_A$ , since if  $A$  is a sentence in  $\mathcal{L}_A$ , then either  $N \models A$  or  $N \models \neg A$ , and so either  $\mathbf{TA} \models A$  or  $\mathbf{TA} \models \neg A$ . Call such a theory *complete*.

**Definition 14.7.** A theory  $\Gamma$  is *complete* iff for every sentence  $A$  in its language, either  $\Gamma \models A$  or  $\Gamma \models \neg A$ .

By the Completeness Theorem,  $\Gamma \models A$  iff  $\Gamma \vdash A$ , so  $\Gamma$  is complete iff for every sentence  $A$  in its language, either  $\Gamma \vdash A$  or  $\Gamma \vdash \neg A$ .

Another question we are led to ask is this: Is there a computational procedure we can use to test if a sentence is in  $\mathbf{TA}$ , in  $\mathbf{PA}$ , or even just in  $\mathbf{Q}$ ? We can make this more precise by defining when a set (e.g., a set of sentences) is decidable.

**Definition 14.8.** A set  $X$  is *decidable* iff there is a computational procedure which on input  $x$  returns 1 if  $x \in X$  and 0 otherwise.

So our question becomes: Is **TA** (**PA**, **Q**) decidable?

The answer to all these questions will be: no. None of these theories are decidable. However, this phenomenon is not specific to these particular theories. In fact, *any* theory that satisfies certain conditions is subject to the same results. One of these conditions, which **Q** and **PA** satisfy, is that they are axiomatized by a decidable set of axioms.

**Definition 14.9.** A theory is *axiomatizable* if it is axiomatized by a decidable set of axioms.

**Example 14.10.** Any theory axiomatized by a finite set of sentences is axiomatizable, since any finite set is decidable. Thus, **Q**, for instance, is axiomatizable.

Schematically axiomatized theories like **PA** are also axiomatizable. For to test if  $B$  is among the axioms of **PA**, i.e., to compute the function  $\chi_X$  where  $\chi_X(B) = 1$  if  $B$  is an axiom of **PA** and  $= 0$  otherwise, we can do the following: First, check if  $B$  is one of the axioms of **Q**. If it is, the answer is “yes” and the value of  $\chi_X(B) = 1$ . If not, test if it is an instance of the induction schema. This can be done systematically; in this case, perhaps it’s easiest to see that it can be done as follows: Any instance of the induction schema begins with a number of universal quantifiers, and then a sub-formula that is a conditional. The consequent of that conditional is  $\forall x A(x, y_1, \dots, y_n)$  where  $x$  and  $y_1, \dots, y_n$  are all the free variables of  $A$  and the initial quantifiers of  $B$  bind the variables  $y_1, \dots, y_n$ . Once we have extracted this  $A$  and checked that its free variables match the variables bound by the universal quantifiers at the front and  $\forall x$ , we go on to check that the antecedent of the conditional matches

$$A(0, y_1, \dots, y_n) \wedge \forall x (A(x, y_1, \dots, y_n) \rightarrow A(x', y_1, \dots, y_n))$$

Again, if it does,  $B$  is an instance of the induction schema, and if it doesn't,  $B$  isn't.

In answering this question—and the more general question of which theories are complete or decidable—it will be useful to consider also the following definition. Recall that a set  $X$  is countable iff it is empty or if there is a surjective function  $f: \mathbb{N} \rightarrow X$ . Such a function is called an enumeration of  $X$ .

**Definition 14.11.** A set  $X$  is called *computably enumerable* (c.e. for short) iff it is empty or it has a computable enumeration.

In addition to axiomatizability, another condition on theories to which the incompleteness theorems apply will be that they are strong enough to prove basic facts about computable functions and decidable relations. By “basic facts,” we mean sentences which express what the values of computable functions are for each of their arguments. And by “strong enough” we mean that the theories in question count these sentences among its theorems. For instance, consider a prototypical computable function: addition. The value of  $+$  for arguments 2 and 3 is 5, i.e.,  $2+3=5$ . A sentence in the language of arithmetic that expresses that the value of  $+$  for arguments 2 and 3 is 5 is:  $(\bar{2} + \bar{3}) = \bar{5}$ . And, e.g.,  $\mathbf{Q}$  proves this sentence. More generally, we would like there to be, for each computable function  $f(x_1, x_2)$  a formula  $A_f(x_1, x_2, y)$  in  $\mathcal{L}_A$  such that  $\mathbf{Q} \vdash A_f(\bar{n}_1, \bar{n}_2, \bar{m})$  whenever  $f(n_1, n_2) = m$ . In this way,  $\mathbf{Q}$  proves that the value of  $f$  for arguments  $n_1, n_2$  is  $m$ . In fact, we require that it proves a bit more, namely that no other number is the value of  $f$  for arguments  $n_1, n_2$ . And the same goes for decidable relations. This is made precise in the following two definitions.

**Definition 14.12.** A formula  $A(x_1, \dots, x_k, y)$  represents the function  $f: \mathbb{N}^k \rightarrow \mathbb{N}$  in  $\Gamma$  iff whenever  $f(n_1, \dots, n_k) = m$ , then

1.  $\Gamma \vdash A(\bar{n}_1, \dots, \bar{n}_k, \bar{m})$ , and

$$2. \Gamma \vdash \forall y (A(\overline{n_1}, \dots, \overline{n_k}, y) \rightarrow y = \overline{m}).$$

**Definition 14.13.** A formula  $A(x_1, \dots, x_k)$  represents the relation  $R \subseteq \mathbb{N}^k$  iff,

1. whenever  $R(n_1, \dots, n_k)$ ,  $\Gamma \vdash A(\overline{n_1}, \dots, \overline{n_k})$ , and
2. whenever not  $R(n_1, \dots, n_k)$ ,  $\Gamma \vdash \neg A(\overline{n_1}, \dots, \overline{n_k})$ .

A theory is “strong enough” for the incompleteness theorems to apply if it represents all computable functions and all decidable relations.  $\mathbf{Q}$  and its extensions satisfy this condition, but it will take us a while to establish this—it’s a non-trivial fact about the kinds of things  $\mathbf{Q}$  can prove, and it’s hard to show because  $\mathbf{Q}$  has only a few axioms from which we’ll have to prove all these facts. However,  $\mathbf{Q}$  is a very weak theory. So although it’s hard to prove that  $\mathbf{Q}$  represents all computable functions, most interesting theories are stronger than  $\mathbf{Q}$ , i.e., prove more than  $\mathbf{Q}$  does. And if  $\mathbf{Q}$  proves something, any stronger theory does; since  $\mathbf{Q}$  represents all computable functions, every stronger theory does. This means that many interesting theories meet this condition of the incompleteness theorems. So our hard work will pay off, since it shows that the incompleteness theorems apply to a wide range of theories. Certainly, any theory aiming to formalize “all of mathematics” must prove everything that  $\mathbf{Q}$  proves, since it should at the very least be able to capture the results of elementary computations. So any theory that is a candidate for a theory of “all of mathematics” will be one to which the incompleteness theorems apply.

### 14.3 Overview of Incompleteness Results

Hilbert expected that mathematics could be formalized in an axiomatizable theory which it would be possible to prove complete and decidable. Moreover, he aimed to prove the consistency of



this theory with very weak, “finitary,” means, which would defend classical mathematics against the challenges of intuitionism. Gödel’s incompleteness theorems showed that these goals cannot be achieved.

Gödel’s first incompleteness theorem showed that a version of Russell and Whitehead’s *Principia Mathematica* is not complete. But the proof was actually very general and applies to a wide variety of theories. This means that it wasn’t just that *Principia Mathematica* did not manage to completely capture mathematics, but that *no* acceptable theory does. It took a while to isolate the features of theories that suffice for the incompleteness theorems to apply, and to generalize Gödel’s proof to apply make it depend only on these features. But we are now in a position to state a very general version of the first incompleteness theorem for theories in the language  $\mathcal{L}_A$  of arithmetic.

**Theorem 14.14.** *If  $\Gamma$  is a consistent and axiomatizable theory in  $\mathcal{L}_A$  which represents all computable functions and decidable relations, then  $\Gamma$  is not complete.*

To say that  $\Gamma$  is not complete is to say that for at least one sentence  $A$ ,  $\Gamma \not\vdash A$  and  $\Gamma \not\vdash \neg A$ . Such a sentence is called *independent* (of  $\Gamma$ ). We can in fact relatively quickly prove that there must be independent sentences. But the power of Gödel’s proof of the theorem lies in the fact that it exhibits a *specific example* of such an independent sentence. The intriguing construction produces a sentence  $G_\Gamma$ , called a *Gödel sentence* for  $\Gamma$ , which is unprovable because in  $\Gamma$ ,  $G_\Gamma$  is equivalent to the claim that  $G_\Gamma$  is unprovable in  $\Gamma$ . It does so *constructively*, i.e., given an axiomatization of  $\Gamma$  and a description of the derivation system, the proof gives a method for actually writing down  $G_\Gamma$ .

The construction in Gödel’s proof requires that we find a way to express in  $\mathcal{L}_A$  the properties of and operations on terms and formulas of  $\mathcal{L}_A$  itself. These include properties such as “ $A$  is a sentence,” “ $\delta$  is a derivation of  $A$ ,” and operations such as  $A[t/x]$ . This way must (a) express these properties and relations

via a “coding” of symbols and sequences thereof (which is what terms, formulas, derivations, etc. are) as natural numbers (which is what  $\mathcal{L}_A$  can talk about). It must (b) do this in such a way that  $\Gamma$  will prove the relevant facts, so we must show that these properties are coded by decidable properties of natural numbers and the operations correspond to computable functions on natural numbers. This is called “arithmetization of syntax.”

Before we investigate how syntax can be arithmetized, however, we will consider the condition that  $\Gamma$  is “strong enough,” i.e., represents all computable functions and decidable relations. This requires that we give a precise definition of “computable.” This can be done in a number of ways, e.g., via the model of Turing machines, or as those functions computable by programs in some general-purpose programming language. Since our aim is to represent these functions and relations in a theory in the language  $\mathcal{L}_A$ , however, it is best to pick a simple definition of computability of just numerical functions. This is the notion of *recursive function*. So we will first discuss the recursive functions. We will then show that  $\mathbf{Q}$  already represents all recursive functions and relations. This will allow us to apply the incompleteness theorem to specific theories such as  $\mathbf{Q}$  and  $\mathbf{PA}$ , since we will have established that these are examples of theories that are “strong enough.”

The end result of the arithmetization of syntax is a formula  $\text{Prov}_\Gamma(x)$  which, via the coding of formulas as numbers, expresses provability from the axioms of  $\Gamma$ . Specifically, if  $A$  is coded by the number  $n$ , and  $\Gamma \vdash A$ , then  $\Gamma \vdash \text{Prov}_\Gamma(\bar{n})$ . This “provability predicate” for  $\Gamma$  allows us also to express, in a certain sense, the consistency of  $\Gamma$  as a sentence of  $\mathcal{L}_A$ : let the “consistency statement” for  $\Gamma$  be the sentence  $\neg\text{Prov}_\Gamma(\bar{n})$ , where we take  $n$  to be the code of a contradiction, e.g., of  $\perp$ . The second incompleteness theorem states that consistent axiomatizable theories also do not prove their own consistency statements. The conditions required for this theorem to apply are a bit more stringent than just that the theory represents all computable functions and decidable relations, but we will show that  $\mathbf{PA}$  satisfies them.

## 14.4 Undecidability and Incompleteness

Gödel's proof of the incompleteness theorems require arithmetization of syntax. But even without that we can obtain some nice results just on the assumption that a theory represents all decidable relations. The proof is a diagonal argument similar to the proof of the undecidability of the halting problem.

**Theorem 14.15.** *If  $\Gamma$  is a consistent theory that represents every decidable relation, then  $\Gamma$  is not decidable.*

*Proof.* Suppose  $\Gamma$  were decidable. We show that if  $\Gamma$  represents every decidable relation, it must be inconsistent.

Decidable properties (one-place relations) are represented by formulas with one free variable. Let  $A_0(x), A_1(x), \dots$ , be a computable enumeration of all such formulas. Now consider the following set  $D \subseteq \mathbb{N}$ :

$$D = \{n : \Gamma \vdash \neg A_n(\bar{n})\}$$

The set  $D$  is decidable, since we can test if  $n \in D$  by first computing  $A_n(x)$ , and from this  $\neg A_n(\bar{n})$ . Obviously, substituting the term  $\bar{n}$  for every free occurrence of  $x$  in  $A_n(x)$  and prefixing  $A(\bar{n})$  by  $\neg$  is a mechanical matter. By assumption,  $\Gamma$  is decidable, so we can test if  $\neg A(\bar{n}) \in \Gamma$ . If it is,  $n \in D$ , and if it isn't,  $n \notin D$ . So  $D$  is likewise decidable.

Since  $\Gamma$  represents all decidable properties, it represents  $D$ . And the formulas which represent  $D$  in  $\Gamma$  are all among  $A_0(x), A_1(x), \dots$ . So let  $d$  be a number such that  $A_d(x)$  represents  $D$  in  $\Gamma$ . If  $d \notin D$ , then, since  $A_d(x)$  represents  $D$ ,  $\Gamma \vdash \neg A_d(\bar{d})$ . But that means that  $d$  meets the defining condition of  $D$ , and so  $d \in D$ . This contradicts  $d \notin D$ . So by indirect proof,  $d \in D$ .

Since  $d \in D$ , by the definition of  $D$ ,  $\Gamma \vdash \neg A_d(\bar{d})$ . On the other hand, since  $A_d(x)$  represents  $D$  in  $\Gamma$ ,  $\Gamma \vdash A_d(\bar{d})$ . Hence,  $\Gamma$  is inconsistent.  $\square$

The preceding theorem shows that no consistent theory that represents all decidable relations can be decidable. We will show

that  $\mathbf{Q}$  does represent all decidable relations; this means that all theories that include  $\mathbf{Q}$ , such as  $\mathbf{PA}$  and  $\mathbf{TA}$ , also do, and hence also are not decidable. (Since all these theories are true in the standard model, they are all consistent.)

We can also use this result to obtain a weak version of the first incompleteness theorem. Any theory that is axiomatizable and complete is decidable. Consistent theories that are axiomatizable and represent all decidable properties then cannot be complete.

**Theorem 14.16.** *If  $\Gamma$  is axiomatizable and complete it is decidable.*

*Proof.* Any inconsistent theory is decidable, since inconsistent theories contain all sentences, so the answer to the question “is  $A \in \Gamma$ ” is always “yes,” i.e., can be decided.

So suppose  $\Gamma$  is consistent, and furthermore is axiomatizable, and complete. Since  $\Gamma$  is axiomatizable, it is computably enumerable. For we can enumerate all the correct derivations from the axioms of  $\Gamma$  by a computable function. From a correct derivation we can compute the sentence it derives, and so together there is a computable function that enumerates all theorems of  $\Gamma$ . A sentence is a theorem of  $\Gamma$  iff  $\neg A$  is not a theorem, since  $\Gamma$  is consistent and complete. We can therefore decide if  $A \in \Gamma$  as follows. Enumerate all theorems of  $\Gamma$ . When  $A$  appears on this list, we know that  $\Gamma \vdash A$ . When  $\neg A$  appears on this list, we know that  $\Gamma \not\vdash A$ . Since  $\Gamma$  is complete, one of these cases eventually obtains, so the procedure eventually produces an answer.  $\square$

**Corollary 14.17.** *If  $\Gamma$  is consistent, axiomatizable, and represents every decidable property, it is not complete.*

*Proof.* If  $\Gamma$  were complete, it would be decidable by the previous theorem (since it is axiomatizable and consistent). But since  $\Gamma$  represents every decidable property, it is not decidable, by the first theorem.  $\square$

Once we have established that, e.g.,  $\mathbf{Q}$ , represents all decidable properties, the corollary tells us that  $\mathbf{Q}$  must be incomplete. However, its proof does not provide an example of an independent sentence; it merely shows that such a sentence must exist. For this, we have to arithmetize syntax and follow Gödel's original proof idea. And of course, we still have to show the first claim, namely that  $\mathbf{Q}$  does, in fact, represent all decidable properties.

It should be noted that not every *interesting* theory is incomplete or undecidable. There are many theories that are sufficiently strong to describe interesting mathematical facts that do not satisfy the conditions of Gödel's result. For instance,  $\mathbf{Pres} = \{A \in \mathcal{L}_{A^+} : \mathbf{N} \models A\}$ , the set of sentences of the language of arithmetic without  $\times$  true in the standard model, is both complete and decidable. This theory is called Presburger arithmetic, and proves all the truths about natural numbers that can be formulated just with  $0$ ,  $\prime$ , and  $+$ .

## Summary

Hilbert's program aimed to show that all of mathematics could be formalized in an axiomatized theory in a formal language, such as the language of arithmetic or of set theory. He believed that such a theory would be **complete**. That is, for every sentence  $A$ , either  $\mathbf{T} \vdash A$  or  $\mathbf{T} \vdash \neg A$ . In this sense then,  $\mathbf{T}$  would have settled every mathematical question: it would either prove that it's true or that it's false. If Hilbert had been right, it would also have turned out that mathematics is **decidable**. That's because any axiomatizable theory is **computably enumerable**, i.e., there is a computable function that lists all its theorems. We can test if a sentence  $A$  is a theorem by listing all of them until we find  $A$  (in which it is a theorem) or  $\neg A$  (in which case it isn't). Alas, Hilbert was wrong. Gödel proved that no axiomatizable, consistent theory that is "strong enough" is complete. That's the **first incompleteness theorem**. The requirement that the theory be "strong enough" amounts to it representing all computable func-

tions and relations. Specifically, the very weak theory  $\mathbf{Q}$  satisfies this property, and any theory that is at least as strong as  $\mathbf{Q}$  also does. He also showed—that is the **second incompleteness theorem**—that the sentence that expresses the consistency of the theory is itself undecidable in it, i.e., the theory proves neither it nor its negation. So Hilbert’s further aim of finding “finitary” consistency proof of all of mathematics cannot be realized. For any finitary consistency proof would, presumably, be formalizable in a theory that captures all of mathematics. Finally, we established that theories that represent all computable functions and relations are not **decidable**. Note that although axiomatizability and completeness implies decidability, incompleteness does not imply undecidability. So this result shows that the second of Hilbert’s goals, namely that there be a procedure that decides if  $\mathbf{T} \vdash A$  or not, can also not be achieved, at least not for theories at least as strong as  $\mathbf{Q}$ .

## Problems

**Problem 14.1.** Show that  $\mathbf{TA} = \{A : N \vDash A\}$  is not axiomatizable. You may assume that  $\mathbf{TA}$  represents all decidable properties.

## CHAPTER 15

# *Recursive Functions*

### 15.1 Introduction

In order to develop a mathematical theory of computability, one has to, first of all, develop a *model* of computability. We now think of computability as the kind of thing that computers do, and computers work with symbols. But at the beginning of the development of theories of computability, the paradigmatic example of computation was *numerical* computation. Mathematicians were always interested in number-theoretic functions, i.e., functions  $f: \mathbb{N}^n \rightarrow \mathbb{N}$  that can be computed. So it is not surprising that at the beginning of the theory of computability, it was such functions that were studied. The most familiar examples of computable numerical functions, such as addition, multiplication, exponentiation (of natural numbers) share an interesting feature: they can be defined *recursively*. It is thus quite natural to attempt a general definition of *computable function* on the basis of recursive definitions. Among the many possible ways to define number-theoretic functions recursively, one particularly simple pattern of definition here becomes central: so-called *primitive recursion*.

In addition to computable functions, we might be interested

in computable sets and relations. A set is computable if we can compute the answer to whether or not a given number is an element of the set, and a relation is computable iff we can compute whether or not a tuple  $\langle n_1, \dots, n_k \rangle$  is an element of the relation. By considering the *characteristic function* of a set or relation, discussion of computable sets and relations can be subsumed under that of computable functions. Thus we can define primitive recursive relations as well, e.g., the relation “ $n$  evenly divides  $m$ ” is a primitive recursive relation.

Primitive recursive functions—those that can be defined using just primitive recursion—are not, however, the only computable number-theoretic functions. Many generalizations of primitive recursion have been considered, but the most powerful and widely-accepted additional way of computing functions is by unbounded search. This leads to the definition of *partial recursive functions*, and a related definition to *general recursive functions*. General recursive functions are computable and total, and the definition characterizes exactly the partial recursive functions that happen to be total. Recursive functions can simulate every other model of computation (Turing machines, lambda calculus, etc.) and so represent one of the many accepted models of computation.

## 15.2 Primitive Recursion

A characteristic of the natural numbers is that every natural number can be reached from 0 by applying the successor operation  $+1$  finitely many times—any natural number is either 0 or the successor of ... the successor of 0. One way to specify a function  $h: \mathbb{N} \rightarrow \mathbb{N}$  that makes use of this fact is this: (a) specify what the value of  $h$  is for argument 0, and (b) also specify how to, given the value of  $h(x)$ , compute the value of  $h(x + 1)$ . For (a) tells us directly what  $h(0)$  is, so  $h$  is defined for 0. Now, using the instruction given by (b) for  $x = 0$ , we can compute  $h(1) = h(0 + 1)$  from  $h(0)$ . Using the same instructions for  $x = 1$ , we compute  $h(2) = h(1 + 1)$  from  $h(1)$ , and so on. For every natural num-



ber  $x$ , we'll eventually reach the step where we define  $h(x)$  from  $h(x+1)$ , and so  $h(x)$  is defined for all  $x \in \mathbb{N}$ .

For instance, suppose we specify  $h: \mathbb{N} \rightarrow \mathbb{N}$  by the following two equations:

$$\begin{aligned}h(0) &= 1 \\h(x+1) &= 2 \cdot h(x)\end{aligned}$$

If we already know how to multiply, then these equations give us the information required for (a) and (b) above. By successively applying the second equation, we get that

$$\begin{aligned}h(1) &= 2 \cdot h(0) = 2, \\h(2) &= 2 \cdot h(1) = 2 \cdot 2, \\h(3) &= 2 \cdot h(2) = 2 \cdot 2 \cdot 2, \\&\vdots\end{aligned}$$

We see that the function  $h$  we have specified is  $h(x) = 2^x$ .

The characteristic feature of the natural numbers guarantees that there is only one function  $h$  that meets these two criteria. A pair of equations like these is called a *definition by primitive recursion* of the function  $h$ . It is so-called because we define  $h$  “recursively,” i.e., the definition, specifically the second equation, involves  $h$  itself on the right-hand-side. It is “primitive” because in defining  $h(x+1)$  we only use the value  $h(x)$ , i.e., the immediately preceding value. This is the simplest way of defining a function on  $\mathbb{N}$  recursively.

We can define even more fundamental functions like addition and multiplication by primitive recursion. In these cases, however, the functions in question are 2-place. We fix one of the argument places, and use the other for the recursion. E.g, to define  $\text{add}(x, y)$  we can fix  $x$  and define the value first for  $y = 0$  and then for  $y + 1$  in terms of  $y$ . Since  $x$  is fixed, it will appear on the left and on the right side of the defining equations.

$$\text{add}(x, 0) = x$$

$$\text{add}(x, y + 1) = \text{add}(x, y) + 1$$

These equations specify the value of  $\text{add}$  for all  $x$  and  $y$ . To find  $\text{add}(2, 3)$ , for instance, we apply the defining equations for  $x = 2$ , using the first to find  $\text{add}(2, 0) = 2$ , then using the second to successively find  $\text{add}(2, 1) = 2 + 1 = 3$ ,  $\text{add}(2, 2) = 3 + 1 = 4$ ,  $\text{add}(2, 3) = 4 + 1 = 5$ .

In the definition of  $\text{add}$  we used  $+$  on the right-hand-side of the second equation, but only to add 1. In other words, we used the successor function  $\text{succ}(z) = z + 1$  and applied it to the previous value  $\text{add}(x, y)$  to define  $\text{add}(x, y + 1)$ . So we can think of the recursive definition as given in terms of a single function which we apply to the previous value. However, it doesn't hurt—and sometimes is necessary—to allow the function to depend not just on the previous value but also on  $x$  and  $y$ . Consider:

$$\begin{aligned}\text{mult}(x, 0) &= 0 \\ \text{mult}(x, y + 1) &= \text{add}(\text{mult}(x, y), x)\end{aligned}$$

This is a primitive recursive definition of a function  $\text{mult}$  by applying the function  $\text{add}$  to both the preceding value  $\text{mult}(x, y)$  and the first argument  $x$ . It also defines the function  $\text{mult}(x, y)$  for all arguments  $x$  and  $y$ . For instance,  $\text{mult}(2, 3)$  is determined by successively computing  $\text{mult}(2, 0)$ ,  $\text{mult}(2, 1)$ ,  $\text{mult}(2, 2)$ , and  $\text{mult}(2, 3)$ :

$$\begin{aligned}\text{mult}(2, 0) &= 0 \\ \text{mult}(2, 1) &= \text{mult}(2, 0 + 1) = \text{add}(\text{mult}(2, 0), 2) = \text{add}(0, 2) = 2 \\ \text{mult}(2, 2) &= \text{mult}(2, 1 + 1) = \text{add}(\text{mult}(2, 1), 2) = \text{add}(2, 2) = 4 \\ \text{mult}(2, 3) &= \text{mult}(2, 2 + 1) = \text{add}(\text{mult}(2, 2), 2) = \text{add}(4, 2) = 6\end{aligned}$$

The general pattern then is this: to give a primitive recursive definition of a function  $h(x_0, \dots, x_{k-1}, y)$ , we provide two equations. The first defines the value of  $h(x_0, \dots, x_{k-1}, 0)$  without reference to  $h$ . The second defines the value of  $h(x_0, \dots, x_{k-1}, y + 1)$  in terms of  $h(x_0, \dots, x_{k-1}, y)$ , the other arguments  $x_0, \dots, x_{k-1}$ ,

and  $y$ . Only the immediately preceding value of  $h$  may be used in that second equation. If we think of the operations given by the right-hand-sides of these two equations as themselves being functions  $f$  and  $g$ , then the general pattern to define a new function  $h$  by primitive recursion is this:

$$\begin{aligned}h(x_0, \dots, x_{k-1}, 0) &= f(x_0, \dots, x_{k-1}) \\h(x_0, \dots, x_{k-1}, y + 1) &= g(x_0, \dots, x_{k-1}, y, h(x_0, \dots, x_{k-1}, y))\end{aligned}$$

In the case of add, we have  $k = 1$  and  $f(x_0) = x_0$  (the identity function), and  $g(x_0, y, z) = z + 1$  (the 3-place function that returns the successor of its third argument):

$$\begin{aligned}\text{add}(x_0, 0) &= f(x_0) = x_0 \\ \text{add}(x_0, y + 1) &= g(x_0, y, \text{add}(x_0, y)) = \text{succ}(\text{add}(x_0, y))\end{aligned}$$

In the case of mult, we have  $f(x_0) = 0$  (the constant function always returning 0) and  $g(x_0, y, z) = \text{add}(z, x_0)$  (the 3-place function that returns the sum of its last and first argument):

$$\begin{aligned}\text{mult}(x_0, 0) &= f(x_0) = 0 \\ \text{mult}(x_0, y + 1) &= g(x_0, y, \text{mult}(x_0, y)) = \text{add}(\text{mult}(x_0, y), x_0)\end{aligned}$$

### 15.3 Composition

If  $f$  and  $g$  are two one-place functions of natural numbers, we can compose them:  $h(x) = g(f(x))$ . The new function  $h(x)$  is then defined by *composition* from the functions  $f$  and  $g$ . We'd like to generalize this to functions of more than one argument.

Here's one way of doing this: suppose  $f$  is a  $k$ -place function, and  $g_0, \dots, g_{k-1}$  are  $k$  functions which are all  $n$ -place. Then we can define a new  $n$ -place function  $h$  as follows:

$$h(x_0, \dots, x_{n-1}) = f(g_0(x_0, \dots, x_{n-1}), \dots, g_{k-1}(x_0, \dots, x_{n-1}))$$

If  $f$  and all  $g_i$  are computable, so is  $h$ : To compute  $h(x_0, \dots, x_{n-1})$ , first compute the values  $y_i = g_i(x_0, \dots, x_{n-1})$  for

each  $i = 0, \dots, k - 1$ . Then feed these values into  $f$  to compute  $h(x_0, \dots, x_{k-1}) = f(y_0, \dots, y_{k-1})$ .

This may seem like an overly restrictive characterization of what happens when we compute a new function using some existing ones. For one thing, sometimes we do not use all the arguments of a function, as when we defined  $g(x, y, z) = \text{succ}(z)$  for use in the primitive recursive definition of add. Suppose we are allowed use of the following functions:

$$P_i^n(x_0, \dots, x_{n-1}) = x_i$$

The functions  $P_i^k$  are called *projection* functions:  $P_i^n$  is an  $n$ -place function. Then  $g$  can be defined by

$$g(x, y, z) = \text{succ}(P_2^3(x, y, z)).$$

Here the role of  $f$  is played by the 1-place function  $\text{succ}$ , so  $k = 1$ . And we have one 3-place function  $P_2^3$  which plays the role of  $g_0$ . The result is a 3-place function that returns the successor of the third argument.

The projection functions also allow us to define new functions by reordering or identifying arguments. For instance, the function  $h(x) = \text{add}(x, x)$  can be defined by

$$h(x_0) = \text{add}(P_0^1(x_0), P_0^1(x_0)).$$

Here  $k = 2$ ,  $n = 1$ , the role of  $f(y_0, y_1)$  is played by  $\text{add}$ , and the roles of  $g_0(x_0)$  and  $g_1(x_0)$  are both played by  $P_0^1(x_0)$ , the one-place projection function (aka the identity function).

If  $f(y_0, y_1)$  is a function we already have, we can define the function  $h(x_0, x_1) = f(x_1, x_0)$  by

$$h(x_0, x_1) = f(P_1^2(x_0, x_1), P_0^2(x_0, x_1)).$$

Here  $k = 2$ ,  $n = 2$ , and the roles of  $g_0$  and  $g_1$  are played by  $P_1^2$  and  $P_0^2$ , respectively.

You may also worry that  $g_0, \dots, g_{k-1}$  are all required to have the same arity  $n$ . (Remember that the *arity* of a function is the

number of arguments; an  $n$ -place function has arity  $n$ .) But adding the projection functions provides the desired flexibility. For example, suppose  $f$  and  $g$  are 3-place functions and  $h$  is the 2-place function defined by

$$h(x, y) = f(x, g(x, x, y), y).$$

The definition of  $h$  can be rewritten with the projection functions, as

$$h(x, y) = f(P_0^2(x, y), g(P_0^2(x, y), P_0^2(x, y), P_1^2(x, y)), P_1^2(x, y)).$$

Then  $h$  is the composition of  $f$  with  $P_0^2$ ,  $l$ , and  $P_1^2$ , where

$$l(x, y) = g(P_0^2(x, y), P_0^2(x, y), P_1^2(x, y)),$$

i.e.,  $l$  is the composition of  $g$  with  $P_0^2$ ,  $P_0^2$ , and  $P_1^2$ .

## 15.4 Primitive Recursion Functions

Let us record again how we can define new functions from existing ones using primitive recursion and composition.

**Definition 15.1.** Suppose  $f$  is a  $k$ -place function ( $k \geq 1$ ) and  $g$  is a  $(k + 2)$ -place function. The function defined by *primitive recursion from  $f$  and  $g$*  is the  $(k + 1)$ -place function  $h$  defined by the equations

$$\begin{aligned} h(x_0, \dots, x_{k-1}, 0) &= f(x_0, \dots, x_{k-1}) \\ h(x_0, \dots, x_{k-1}, y + 1) &= g(x_0, \dots, x_{k-1}, y, h(x_0, \dots, x_{k-1}, y)) \end{aligned}$$

**Definition 15.2.** Suppose  $f$  is a  $k$ -place function, and  $g_0, \dots, g_{k-1}$  are  $k$  functions which are all  $n$ -place. The function defined by *composition from  $f$  and  $g_0, \dots, g_{k-1}$*  is the  $n$ -place function  $h$

defined by

$$h(x_0, \dots, x_{n-1}) = f(g_0(x_0, \dots, x_{n-1}), \dots, g_{k-1}(x_0, \dots, x_{n-1})).$$

In addition to succ and the projection functions

$$P_i^n(x_0, \dots, x_{n-1}) = x_i,$$

for each natural number  $n$  and  $i < n$ , we will include among the primitive recursive functions the function  $\text{zero}(x) = 0$ .

**Definition 15.3.** The set of primitive recursive functions is the set of functions from  $\mathbb{N}^n$  to  $\mathbb{N}$ , defined inductively by the following clauses:

1. zero is primitive recursive.
2. succ is primitive recursive.
3. Each projection function  $P_i^n$  is primitive recursive.
4. If  $f$  is a  $k$ -place primitive recursive function and  $g_0, \dots, g_{k-1}$  are  $n$ -place primitive recursive functions, then the composition of  $f$  with  $g_0, \dots, g_{k-1}$  is primitive recursive.
5. If  $f$  is a  $k$ -place primitive recursive function and  $g$  is a  $k + 2$ -place primitive recursive function, then the function defined by primitive recursion from  $f$  and  $g$  is primitive recursive.

Put more concisely, the set of primitive recursive functions is the smallest set containing zero, succ, and the projection functions  $P_j^n$ , and which is closed under composition and primitive recursion.

Another way of describing the set of primitive recursive functions is by defining it in terms of “stages.” Let  $S_0$  denote the set of starting functions: zero, succ, and the projections. These are the primitive recursive functions of stage 0. Once a stage  $S_i$  has been

defined, let  $S_{i+1}$  be the set of all functions you get by applying a single instance of composition or primitive recursion to functions already in  $S_i$ . Then

$$S = \bigcup_{i \in \mathbb{N}} S_i$$

is the set of all primitive recursive functions

Let us verify that `add` is a primitive recursive function.

**Proposition 15.4.** *The addition function  $\text{add}(x, y) = x + y$  is primitive recursive.*

*Proof.* We already have a primitive recursive definition of `add` in terms of two functions  $f$  and  $g$  which matches the format of

**Definition 15.1:**

$$\begin{aligned} \text{add}(x_0, 0) &= f(x_0) = x_0 \\ \text{add}(x_0, y + 1) &= g(x_0, y, \text{add}(x_0, y)) = \text{succ}(\text{add}(x_0, y)) \end{aligned}$$

So `add` is primitive recursive provided  $f$  and  $g$  are as well.  $f(x_0) = x_0 = P_0^1(x_0)$ , and the projection functions count as primitive recursive, so  $f$  is primitive recursive. The function  $g$  is the three-place function  $g(x_0, y, z)$  defined by

$$g(x_0, y, z) = \text{succ}(z).$$

This does not yet tell us that  $g$  is primitive recursive, since  $g$  and `succ` are not quite the same function: `succ` is one-place, and  $g$  has to be three-place. But we can define  $g$  “officially” by composition as

$$g(x_0, y, z) = \text{succ}(P_2^3(x_0, y, z))$$

Since `succ` and  $P_2^3$  count as primitive recursive functions,  $g$  does as well, since it can be defined by composition from primitive recursive functions. □

**Proposition 15.5.** *The multiplication function  $\text{mult}(x, y) = x \cdot y$  is primitive recursive.*

*Proof.* Exercise. □

**Example 15.6.** Here's our very first example of a primitive recursive definition:

$$\begin{aligned}h(0) &= 1 \\h(y + 1) &= 2 \cdot h(y).\end{aligned}$$

This function cannot fit into the form required by **Definition 15.1**, since  $k = 0$ . The definition also involves the constants 1 and 2. To get around the first problem, let's introduce a dummy argument and define the function  $h'$ :

$$\begin{aligned}h'(x_0, 0) &= f(x_0) = 1 \\h'(x_0, y + 1) &= g(x_0, y, h'(x_0, y)) = 2 \cdot h'(x_0, y).\end{aligned}$$

The function  $f(x_0) = 1$  can be defined from `succ` and `zero` by composition:  $f(x_0) = \text{succ}(\text{zero}(x_0))$ . The function  $g$  can be defined by composition from  $g'(z) = 2 \cdot z$  and projections:

$$g(x_0, y, z) = g'(P_2^3(x_0, y, z))$$

and  $g'$  in turn can be defined by composition as

$$g'(z) = \text{mult}(g''(z), P_0^1(z))$$

and

$$g''(z) = \text{succ}(f(z)),$$

where  $f$  is as above:  $f(z) = \text{succ}(\text{zero}(z))$ . Now that we have  $h'$ , we can use composition again to let  $h(y) = h'(P_0^1(y), P_0^1(y))$ . This shows that  $h$  can be defined from the basic functions using a sequence of compositions and primitive recursions, so  $h$  is primitive recursive.



## 15.5 Primitive Recursion Notations

One advantage to having the precise inductive description of the primitive recursive functions is that we can be systematic in describing them. For example, we can assign a “notation” to each such function, as follows. Use symbols  $\text{zero}$ ,  $\text{succ}$ , and  $P_i^n$  for zero, successor, and the projections. Now suppose  $h$  is defined by composition from a  $k$ -place function  $f$  and  $n$ -place functions  $g_0, \dots, g_{k-1}$ , and we have assigned notations  $F, G_0, \dots, G_{k-1}$  to the latter functions. Then, using a new symbol  $\text{Comp}_{k,n}$ , we can denote the function  $h$  by  $\text{Comp}_{k,n}[F, G_0, \dots, G_{k-1}]$ .

For functions defined by primitive recursion, we can use analogous notations. Suppose the  $(k+1)$ -ary function  $h$  is defined by primitive recursion from the  $k$ -ary function  $f$  and the  $(k+2)$ -ary function  $g$ , and the notations assigned to  $f$  and  $g$  are  $F$  and  $G$ , respectively. Then the notation assigned to  $h$  is  $\text{Rec}_k[F, G]$ .

Recall that the addition function is defined by primitive recursion as

$$\text{add}(x_0, 0) = P_0^1(x_0) = x_0$$

$$\text{add}(x_0, y + 1) = \text{succ}(P_2^3(x_0, y, \text{add}(x_0, y))) = \text{add}(x_0, y) + 1$$

Here the role of  $f$  is played by  $P_0^1$ , and the role of  $g$  is played by  $\text{succ}(P_2^3(x_0, y, z))$ , which is assigned the notation  $\text{Comp}_{1,3}[\text{succ}, P_2^3]$  as it is the result of defining a function by composition from the 1-ary function  $\text{succ}$  and the 3-ary function  $P_2^3$ . With this setup, we can denote the addition function by

$$\text{Rec}_1[P_0^1, \text{Comp}_{1,3}[\text{succ}, P_2^3]].$$

Having these notations sometimes proves useful, e.g., when enumerating primitive recursive functions.

## 15.6 Primitive Recursive Functions are Computable

Suppose a function  $h$  is defined by primitive recursion

$$\begin{aligned}h(\vec{x}, 0) &= f(\vec{x}) \\h(\vec{x}, y + 1) &= g(\vec{x}, y, h(\vec{x}, y))\end{aligned}$$

and suppose the functions  $f$  and  $g$  are computable. (We use  $\vec{x}$  to abbreviate  $x_0, \dots, x_{k-1}$ .) Then  $h(\vec{x}, 0)$  can obviously be computed, since it is just  $f(\vec{x})$  which we assume is computable.  $h(\vec{x}, 1)$  can then also be computed, since  $1 = 0 + 1$  and so  $h(\vec{x}, 1)$  is just

$$h(\vec{x}, 1) = g(\vec{x}, 0, h(\vec{x}, 0)) = g(\vec{x}, 0, f(\vec{x})).$$

We can go on in this way and compute

$$\begin{aligned}h(\vec{x}, 2) &= g(\vec{x}, 1, h(\vec{x}, 1)) = g(\vec{x}, 1, g(\vec{x}, 0, f(\vec{x}))) \\h(\vec{x}, 3) &= g(\vec{x}, 2, h(\vec{x}, 2)) = g(\vec{x}, 2, g(\vec{x}, 1, g(\vec{x}, 0, f(\vec{x})))) \\h(\vec{x}, 4) &= g(\vec{x}, 3, h(\vec{x}, 3)) = g(\vec{x}, 3, g(\vec{x}, 2, g(\vec{x}, 1, g(\vec{x}, 0, f(\vec{x})))))) \\&\vdots\end{aligned}$$

Thus, to compute  $h(\vec{x}, y)$  in general, successively compute  $h(\vec{x}, 0)$ ,  $h(\vec{x}, 1)$ ,  $\dots$ , until we reach  $h(\vec{x}, y)$ .

Thus, a primitive recursive definition yields a new computable function if the functions  $f$  and  $g$  are computable. Composition of functions also results in a computable function if the functions  $f$  and  $g_i$  are computable.

Since the basic functions zero, succ, and  $P_i^n$  are computable, and composition and primitive recursion yield computable functions from computable functions, this means that every primitive recursive function is computable.

## 15.7 Examples of Primitive Recursive Functions

We already have some examples of primitive recursive functions: the addition and multiplication functions `add` and `mult`. The identity function  $\text{id}(x) = x$  is primitive recursive, since it is just  $P_0^1$ . The constant functions  $\text{const}_n(x) = n$  are primitive recursive since they can be defined from `zero` and `succ` by successive composition. This is useful when we want to use constants in primitive recursive definitions, e.g., if we want to define the function  $f(x) = 2 \cdot x$  can obtain it by composition from  $\text{const}_2(x)$  and multiplication as  $f(x) = \text{mult}(\text{const}_2(x), P_0^1(x))$ . We'll make use of this trick from now on.

**Proposition 15.7.** *The exponentiation function  $\text{exp}(x, y) = x^y$  is primitive recursive.*

*Proof.* We can define `exp` primitive recursively as

$$\begin{aligned}\text{exp}(x, 0) &= 1 \\ \text{exp}(x, y + 1) &= \text{mult}(x, \text{exp}(x, y)).\end{aligned}$$

Strictly speaking, this is not a recursive definition from primitive recursive functions. Officially, though, we have:

$$\begin{aligned}\text{exp}(x, 0) &= f(x) \\ \text{exp}(x, y + 1) &= g(x, y, \text{exp}(x, y)).\end{aligned}$$

where

$$\begin{aligned}f(x) &= \text{succ}(\text{zero}(x)) = 1 \\ g(x, y, z) &= \text{mult}(P_0^3(x, y, z), P_2^3(x, y, z)) = x \cdot z\end{aligned}$$

and so  $f$  and  $g$  are defined from primitive recursive functions by composition.  $\square$

**Proposition 15.8.** *The predecessor function  $\text{pred}(y)$  defined by*

$$\text{pred}(y) = \begin{cases} 0 & \text{if } y = 0 \\ y - 1 & \text{otherwise} \end{cases}$$

*is primitive recursive.*

*Proof.* Note that

$$\begin{aligned} \text{pred}(0) &= 0 \text{ and} \\ \text{pred}(y + 1) &= y. \end{aligned}$$

This is almost a primitive recursive definition. It does not, strictly speaking, fit into the pattern of definition by primitive recursion, since that pattern requires at least one extra argument  $x$ . It is also odd in that it does not actually use  $\text{pred}(y)$  in the definition of  $\text{pred}(y + 1)$ . But we can first define  $\text{pred}'(x, y)$  by

$$\begin{aligned} \text{pred}'(x, 0) &= \text{zero}(x) = 0, \\ \text{pred}'(x, y + 1) &= P_1^3(x, y, \text{pred}'(x, y)) = y. \end{aligned}$$

and then define  $\text{pred}$  from it by composition, e.g., as  $\text{pred}(x) = \text{pred}'(\text{zero}(x), P_0^1(x))$ .  $\square$

**Proposition 15.9.** *The factorial function  $\text{fac}(x) = x! = 1 \cdot 2 \cdot 3 \cdots x$  is primitive recursive.*

*Proof.* The obvious primitive recursive definition is

$$\begin{aligned} \text{fac}(0) &= 1 \\ \text{fac}(y + 1) &= \text{fac}(y) \cdot (y + 1). \end{aligned}$$

Officially, we have to first define a two-place function  $h$

$$\begin{aligned} h(x, 0) &= \text{const}_1(x) \\ h(x, y + 1) &= g(x, y, h(x, y)) \end{aligned}$$

where  $g(x, y, z) = \text{mult}(P_2^3(x, y, z), \text{succ}(P_1^3(x, y, z)))$  and then let

$$\text{fac}(y) = h(P_0^1(y), P_0^1(y)) = h(y, y).$$

From now on we'll be a bit more *laissez-faire* and not give the official definitions by composition and primitive recursion.  $\square$

**Proposition 15.10.** *Truncated subtraction,  $x \dot{-} y$ , defined by*

$$x \dot{-} y = \begin{cases} 0 & \text{if } x < y \\ x - y & \text{otherwise} \end{cases}$$

*is primitive recursive.*

*Proof.* We have:

$$\begin{aligned} x \dot{-} 0 &= x \\ x \dot{-} (y + 1) &= \text{pred}(x \dot{-} y) \end{aligned} \quad \square$$

**Proposition 15.11.** *The distance between  $x$  and  $y$ ,  $|x - y|$ , is primitive recursive.*

*Proof.* We have  $|x - y| = (x \dot{-} y) + (y \dot{-} x)$ , so the distance can be defined by composition from  $+$  and  $\dot{-}$ , which are primitive recursive.  $\square$

**Proposition 15.12.** *The maximum of  $x$  and  $y$ ,  $\max(x, y)$ , is primitive recursive.*

*Proof.* We can define  $\max(x, y)$  by composition from  $+$  and  $\dot{-}$  by

$$\max(x, y) = x + (y \dot{-} x).$$

If  $x$  is the maximum, i.e.,  $x \geq y$ , then  $y \dot{-} x = 0$ , so  $x + (y \dot{-} x) = x + 0 = x$ . If  $y$  is the maximum, then  $y \dot{-} x = y - x$ , and so  $x + (y \dot{-} x) = x + (y - x) = y$ .  $\square$

**Proposition 15.13.** *The minimum of  $x$  and  $y$ ,  $\min(x, y)$ , is primitive recursive.*

*Proof.* Exercise. □

**Proposition 15.14.** *The set of primitive recursive functions is closed under the following two operations:*

1. *Finite sums: if  $f(\vec{x}, z)$  is primitive recursive, then so is the function*

$$g(\vec{x}, y) = \sum_{z=0}^y f(\vec{x}, z).$$

2. *Finite products: if  $f(\vec{x}, z)$  is primitive recursive, then so is the function*

$$h(\vec{x}, y) = \prod_{z=0}^y f(\vec{x}, z).$$

*Proof.* For example, finite sums are defined recursively by the equations

$$\begin{aligned} g(\vec{x}, 0) &= f(\vec{x}, 0) \\ g(\vec{x}, y + 1) &= g(\vec{x}, y) + f(\vec{x}, y + 1). \end{aligned} \quad \square$$

## 15.8 Primitive Recursive Relations

**Definition 15.15.** A relation  $R(\vec{x})$  is said to be primitive recursive if its characteristic function,

$$\chi_R(\vec{x}) = \begin{cases} 1 & \text{if } R(\vec{x}) \\ 0 & \text{otherwise} \end{cases}$$

is primitive recursive.

In other words, when one speaks of a primitive recursive relation  $R(\vec{x})$ , one is referring to a relation of the form  $\chi_R(\vec{x}) = 1$ , where  $\chi_R$  is a primitive recursive function which, on any input, returns either 1 or 0. For example, the relation  $\text{IsZero}(x)$ , which holds if and only if  $x = 0$ , corresponds to the function  $\chi_{\text{IsZero}}$ , defined using primitive recursion by

$$\begin{aligned}\chi_{\text{IsZero}}(0) &= 1, \\ \chi_{\text{IsZero}}(x + 1) &= 0.\end{aligned}$$

It should be clear that one can compose relations with other primitive recursive functions. So the following are also primitive recursive:

1. The equality relation,  $x = y$ , defined by  $\text{IsZero}(|x - y|)$
2. The less-than relation,  $x \leq y$ , defined by  $\text{IsZero}(x \dot{-} y)$

**Proposition 15.16.** *The set of primitive recursive relations is closed under Boolean operations, that is, if  $P(\vec{x})$  and  $Q(\vec{x})$  are primitive recursive, so are*

1.  $\neg P(\vec{x})$
2.  $P(\vec{x}) \wedge Q(\vec{x})$
3.  $P(\vec{x}) \vee Q(\vec{x})$
4.  $P(\vec{x}) \rightarrow Q(\vec{x})$

*Proof.* Suppose  $P(\vec{x})$  and  $Q(\vec{x})$  are primitive recursive, i.e., their characteristic functions  $\chi_P$  and  $\chi_Q$  are. We have to show that the characteristic functions of  $\neg P(\vec{x})$ , etc., are also primitive recursive.

$$\chi_{\neg P}(\vec{x}) = \begin{cases} 0 & \text{if } \chi_P(\vec{x}) = 1 \\ 1 & \text{otherwise} \end{cases}$$

We can define  $\chi_{\neg P}(\vec{x})$  as  $1 \dot{-} \chi_P(\vec{x})$ .

$$\chi_{P \wedge Q}(\vec{x}) = \begin{cases} 1 & \text{if } \chi_P(\vec{x}) = \chi_Q(\vec{x}) = 1 \\ 0 & \text{otherwise} \end{cases}$$

We can define  $\chi_{P \wedge Q}(\vec{x})$  as  $\chi_P(\vec{x}) \cdot \chi_Q(\vec{x})$  or as  $\min(\chi_P(\vec{x}), \chi_Q(\vec{x}))$ . Similarly,

$$\begin{aligned} \chi_{P \vee Q}(\vec{x}) &= \max(\chi_P(\vec{x}), \chi_Q(\vec{x})) \text{ and} \\ \chi_{P \rightarrow Q}(\vec{x}) &= \max(1 \dot{-} \chi_P(\vec{x}), \chi_Q(\vec{x})). \quad \square \end{aligned}$$

**Proposition 15.17.** *The set of primitive recursive relations is closed under bounded quantification, i.e., if  $R(\vec{x}, z)$  is a primitive recursive relation, then so are the relations*

$$\begin{aligned} (\forall z < y) R(\vec{x}, z) \text{ and} \\ (\exists z < y) R(\vec{x}, z). \end{aligned}$$

$(\forall z < y) R(\vec{x}, z)$  holds of  $\vec{x}$  and  $y$  if and only if  $R(\vec{x}, z)$  holds for every  $z$  less than  $y$ , and similarly for  $(\exists z < y) R(\vec{x}, z)$ .

*Proof.* By convention, we take  $(\forall z < 0) R(\vec{x}, z)$  to be true (for the trivial reason that there are no  $z$  less than 0) and  $(\exists z < 0) R(\vec{x}, z)$  to be false. A bounded universal quantifier functions just like a finite product or iterated minimum, i.e., if  $P(\vec{x}, y) \Leftrightarrow (\forall z < y) R(\vec{x}, z)$  then  $\chi_P(\vec{x}, y)$  can be defined by

$$\begin{aligned} \chi_P(\vec{x}, 0) &= 1 \\ \chi_P(\vec{x}, y + 1) &= \min(\chi_P(\vec{x}, y), \chi_R(\vec{x}, y)). \end{aligned}$$

Bounded existential quantification can similarly be defined using max. Alternatively, it can be defined from bounded universal quantification, using the equivalence  $(\exists z < y) R(\vec{x}, z) \Leftrightarrow \neg(\forall z < y) \neg R(\vec{x}, z)$ . Note that, for example, a bounded quantifier of the form  $(\exists x \leq y) \dots x \dots$  is equivalent to  $(\exists x < y + 1) \dots x \dots$ .  $\square$



Another useful primitive recursive function is the conditional function,  $\text{cond}(x, y, z)$ , defined by

$$\text{cond}(x, y, z) = \begin{cases} y & \text{if } x = 0 \\ z & \text{otherwise.} \end{cases}$$

This is defined recursively by

$$\begin{aligned} \text{cond}(0, y, z) &= y, \\ \text{cond}(x + 1, y, z) &= z. \end{aligned}$$

One can use this to justify definitions of primitive recursive functions by cases from primitive recursive relations:

**Proposition 15.18.** *If  $g_0(\vec{x}), \dots, g_m(\vec{x})$  are primitive recursive functions, and  $R_0(\vec{x}), \dots, R_{m-1}(\vec{x})$  are primitive recursive relations, then the function  $f$  defined by*

$$f(\vec{x}) = \begin{cases} g_0(\vec{x}) & \text{if } R_0(\vec{x}) \\ g_1(\vec{x}) & \text{if } R_1(\vec{x}) \text{ and not } R_0(\vec{x}) \\ \vdots \\ g_{m-1}(\vec{x}) & \text{if } R_{m-1}(\vec{x}) \text{ and none of the previous hold} \\ g_m(\vec{x}) & \text{otherwise} \end{cases}$$

*is also primitive recursive.*

*Proof.* When  $m = 1$ , this is just the function defined by

$$f(\vec{x}) = \text{cond}(\chi_{-R_0}(\vec{x}), g_0(\vec{x}), g_1(\vec{x})).$$

For  $m$  greater than 1, one can just compose definitions of this form.  $\square$

## 15.9 Bounded Minimization

It is often useful to define a function as the least number satisfying some property or relation  $P$ . If  $P$  is decidable, we can

compute this function simply by trying out all the possible numbers,  $0, 1, 2, \dots$ , until we find the least one satisfying  $P$ . This kind of unbounded search takes us out of the realm of primitive recursive functions. However, if we're only interested in the least number *less than some independently given bound*, we stay primitive recursive. In other words, and a bit more generally, suppose we have a primitive recursive relation  $R(x, z)$ . Consider the function that maps  $x$  and  $y$  to the least  $z < y$  such that  $R(x, z)$ . It, too, can be computed, by testing whether  $R(x, 0), R(x, 1), \dots, R(x, y - 1)$ . But why is it primitive recursive?

**Proposition 15.19.** *If  $R(\vec{x}, z)$  is primitive recursive, so is the function  $m_R(\vec{x}, y)$  which returns the least  $z$  less than  $y$  such that  $R(\vec{x}, z)$  holds, if there is one, and  $y$  otherwise. We will write the function  $m_R$  as*

$$(\min z < y) R(\vec{x}, z),$$

*Proof.* Note that there can be no  $z < 0$  such that  $R(\vec{x}, z)$  since there is no  $z < 0$  at all. So  $m_R(\vec{x}, 0) = 0$ .

In case the bound is of the form  $y + 1$  we have three cases:

1. There is a  $z < y$  such that  $R(\vec{x}, z)$ , in which case  $m_R(\vec{x}, y + 1) = m_R(\vec{x}, y)$ .
2. There is no such  $z < y$  but  $R(\vec{x}, y)$  holds, then  $m_R(\vec{x}, y + 1) = y$ .
3. There is no  $z < y + 1$  such that  $R(\vec{x}, z)$ , then  $m_R(\vec{x}, y + 1) = y + 1$ .

So we can define  $m_R(\vec{x}, 0)$  by primitive recursion as follows:

$$m_R(\vec{x}, 0) = 0$$

$$m_R(\vec{x}, y + 1) = \begin{cases} m_R(\vec{x}, y) & \text{if } m_R(\vec{x}, y) \neq y \\ y & \text{if } m_R(\vec{x}, y) = y \text{ and } R(\vec{x}, y) \\ y + 1 & \text{otherwise.} \end{cases}$$

Note that there is a  $z < y$  such that  $R(\vec{x}, z)$  iff  $m_R(\vec{x}, y) \neq y$ .  $\square$

## 15.10 Primes

Bounded quantification and bounded minimization provide us with a good deal of machinery to show that natural functions and relations are primitive recursive. For example, consider the relation “ $x$  divides  $y$ ”, written  $x \mid y$ . The relation  $x \mid y$  holds if division of  $y$  by  $x$  is possible without remainder, i.e., if  $y$  is an integer multiple of  $x$ . (If it doesn’t hold, i.e., the remainder when dividing  $x$  by  $y$  is  $> 0$ , we write  $x \nmid y$ .) In other words,  $x \mid y$  iff for some  $z$ ,  $x \cdot z = y$ . Obviously, any such  $z$ , if it exists, must be  $\leq y$ . So, we have that  $x \mid y$  iff for some  $z \leq y$ ,  $x \cdot z = y$ . We can define the relation  $x \mid y$  by bounded existential quantification from  $=$  and multiplication by

$$x \mid y \Leftrightarrow (\exists z \leq y) (x \cdot z) = y.$$

We’ve thus shown that  $x \mid y$  is primitive recursive.

A natural number  $x$  is *prime* if it is neither 0 nor 1 and is only divisible by 1 and itself. In other words, prime numbers are such that, whenever  $y \mid x$ , either  $y = 1$  or  $y = x$ . To test if  $x$  is prime, we only have to check if  $y \mid x$  for all  $y \leq x$ , since if  $y > x$ , then automatically  $y \nmid x$ . So, the relation  $\text{Prime}(x)$ , which holds iff  $x$  is prime, can be defined by

$$\text{Prime}(x) \Leftrightarrow x \geq 2 \wedge (\forall y \leq x) (y \mid x \rightarrow y = 1 \vee y = x)$$

and is thus primitive recursive.

The primes are 2, 3, 5, 7, 11, etc. Consider the function  $p(x)$  which returns the  $x$ th prime in that sequence, i.e.,  $p(0) = 2$ ,  $p(1) = 3$ ,  $p(2) = 5$ , etc. (For convenience we will often write  $p(x)$  as  $p_x$  ( $p_0 = 2$ ,  $p_1 = 3$ , etc.))

If we had a function  $\text{nextPrime}(x)$ , which returns the first prime number larger than  $x$ ,  $p$  can be easily defined using primitive recursion:

$$\begin{aligned} p(0) &= 2 \\ p(x+1) &= \text{nextPrime}(p(x)) \end{aligned}$$

Since  $\text{nextPrime}(x)$  is the least  $y$  such that  $y > x$  and  $y$  is prime, it can be easily computed by unbounded search. But it can also be defined by bounded minimization, thanks to a result due to Euclid: there is always a prime number between  $x$  and  $x! + 1$ .

$$\text{nextPrime}(x) = (\min y \leq x! + 1) (y > x \wedge \text{Prime}(y)).$$

This shows, that  $\text{nextPrime}(x)$  and hence  $p(x)$  are (not just computable but) primitive recursive.

(If you're curious, here's a quick proof of Euclid's theorem. Suppose  $p_n$  is the largest prime  $\leq x$  and consider the product  $p = p_0 \cdot p_1 \cdot \dots \cdot p_n$  of all primes  $\leq x$ . Either  $p + 1$  is prime or there is a prime between  $x$  and  $p + 1$ . Why? Suppose  $p + 1$  is not prime. Then some prime number  $q \mid p + 1$  where  $q < p + 1$ . None of the primes  $\leq x$  divide  $p + 1$ . (By definition of  $p$ , each of the primes  $p_i \leq x$  divides  $p$ , i.e., with remainder 0. So, each of the primes  $p_i \leq x$  divides  $p + 1$  with remainder 1, and so  $p_i \nmid p + 1$ .) Hence,  $q$  is a prime  $> x$  and  $< p + 1$ . And  $p \leq x!$ , so there is a prime  $> x$  and  $\leq x! + 1$ .)

## 15.11 Sequences

The set of primitive recursive functions is remarkably robust. But we will be able to do even more once we have developed a adequate means of handling *sequences*. We will identify finite sequences of natural numbers with natural numbers in the following way: the sequence  $\langle a_0, a_1, a_2, \dots, a_k \rangle$  corresponds to the number

$$p_0^{a_0+1} \cdot p_1^{a_1+1} \cdot p_2^{a_2+1} \cdot \dots \cdot p_k^{a_k+1}.$$

We add one to the exponents to guarantee that, for example, the sequences  $\langle 2, 7, 3 \rangle$  and  $\langle 2, 7, 3, 0, 0 \rangle$  have distinct numeric codes. We can take both 0 and 1 to code the empty sequence; for concreteness, let  $\Lambda$  denote 0.

The reason that this coding of sequences works is the so-called Fundamental Theorem of Arithmetic: every natural number  $n \geq$

2 can be written in one and only one way in the form

$$n = p_0^{a_0} \cdot p_1^{a_1} \cdot \dots \cdot p_k^{a_k}$$

with  $a_k \geq 1$ . This guarantees that the mapping  $\langle \rangle(a_0, \dots, a_k) = \langle a_0, \dots, a_k \rangle$  is injective: different sequences are mapped to different numbers; to each number only at most one sequence corresponds.

We'll now show that the operations of determining the length of a sequence, determining its  $i$ th element, appending an element to a sequence, and concatenating two sequences, are all primitive recursive.

**Proposition 15.20.** *The function  $\text{len}(s)$ , which returns the length of the sequence  $s$ , is primitive recursive.*

*Proof.* Let  $R(i, s)$  be the relation defined by

$$R(i, s) \text{ iff } p_i \mid s \wedge p_{i+1} \nmid s.$$

$R$  is clearly primitive recursive. Whenever  $s$  is the code of a non-empty sequence, i.e.,

$$s = p_0^{a_0+1} \cdot \dots \cdot p_k^{a_k+1},$$

$R(i, s)$  holds if  $p_i$  is the largest prime such that  $p_i \mid s$ , i.e.,  $i = k$ . The length of  $s$  thus is  $i+1$  iff  $p_i$  is the largest prime that divides  $s$ , so we can let

$$\text{len}(s) = \begin{cases} 0 & \text{if } s = 0 \text{ or } s = 1 \\ 1 + (\min i < s) R(i, s) & \text{otherwise} \end{cases}$$

We can use bounded minimization, since there is only one  $i$  that satisfies  $R(s, i)$  when  $s$  is a code of a sequence, and if  $i$  exists it is less than  $s$  itself.  $\square$

**Proposition 15.21.** *The function  $\text{append}(s, a)$ , which returns the result of appending  $a$  to the sequence  $s$ , is primitive recursive.*

*Proof.*  $\text{append}$  can be defined by:

$$\text{append}(s, a) = \begin{cases} 2^{a+1} & \text{if } s = 0 \text{ or } s = 1 \\ s \cdot p_{\text{len}(s)}^{a+1} & \text{otherwise.} \end{cases} \quad \square$$

**Proposition 15.22.** *The function  $\text{element}(s, i)$ , which returns the  $i$ th element of  $s$  (where the initial element is called the 0th), or 0 if  $i$  is greater than or equal to the length of  $s$ , is primitive recursive.*

*Proof.* Note that  $a$  is the  $i$ th element of  $s$  iff  $p_i^{a+1}$  is the largest power of  $p_i$  that divides  $s$ , i.e.,  $p_i^{a+1} \mid s$  but  $p_i^{a+2} \nmid s$ . So:

$$\text{element}(s, i) = \begin{cases} 0 & \text{if } i \geq \text{len}(s) \\ (\min a < s) (p_i^{a+2} \nmid s) & \text{otherwise.} \end{cases} \quad \square$$

Instead of using the official names for the functions defined above, we introduce a more compact notation. We will use  $(s)_i$  instead of  $\text{element}(s, i)$ , and  $\langle s_0, \dots, s_k \rangle$  to abbreviate

$$\text{append}(\text{append}(\dots \text{append}(\Lambda, s_0) \dots), s_k).$$

Note that if  $s$  has length  $k$ , the elements of  $s$  are  $(s)_0, \dots, (s)_{k-1}$ .

**Proposition 15.23.** *The function  $\text{concat}(s, t)$ , which concatenates two sequences, is primitive recursive.*

*Proof.* We want a function  $\text{concat}$  with the property that

$$\text{concat}(\langle a_0, \dots, a_k \rangle, \langle b_0, \dots, b_l \rangle) = \langle a_0, \dots, a_k, b_0, \dots, b_l \rangle.$$

We'll use a "helper" function  $\text{hconcat}(s, t, n)$  which concatenates the first  $n$  symbols of  $t$  to  $s$ . This function can be defined by primitive recursion as follows:

$$\text{hconcat}(s, t, 0) = s$$

$$\text{hconcat}(s, t, n + 1) = \text{append}(\text{hconcat}(s, t, n), (t)_n)$$

Then we can define `concat` by

$$\text{concat}(s, t) = \text{hconcat}(s, t, \text{len}(t)). \quad \square$$

We will write  $s \frown t$  instead of `concat`( $s, t$ ).

It will be useful for us to be able to bound the numeric code of a sequence in terms of its length and its largest element. Suppose  $s$  is a sequence of length  $k$ , each element of which is less than or equal to some number  $x$ . Then  $s$  has at most  $k$  prime factors, each at most  $p_{k-1}$ , and each raised to at most  $x + 1$  in the prime factorization of  $s$ . In other words, if we define

$$\text{sequenceBound}(x, k) = p_{k-1}^{k \cdot (x+1)},$$

then the numeric code of the sequence  $s$  described above is at most `sequenceBound`( $x, k$ ).

Having such a bound on sequences gives us a way of defining new functions using bounded search. For example, we can define `concat` using bounded search. All we need to do is write down a primitive recursive *specification* of the object (number of the concatenated sequence) we are looking for, and a bound on how far to look. The following works:

$$\begin{aligned} \text{concat}(s, t) = & (\min v < \text{sequenceBound}(s + t, \text{len}(s) + \text{len}(t))) \\ & (\text{len}(v) = \text{len}(s) + \text{len}(t) \wedge \\ & (\forall i < \text{len}(s)) ((v)_i = (s)_i) \wedge \\ & (\forall j < \text{len}(t)) ((v)_{\text{len}(s)+j} = (t)_j)) \end{aligned}$$

**Proposition 15.24.** *The function `subseq`( $s, i, n$ ) which returns the subsequence of  $s$  of length  $n$  beginning at the  $i$ th element, is primitive recursive.*

*Proof.* Exercise. □

## 15.12 Trees

Sometimes it is useful to represent trees as natural numbers, just like we can represent sequences by numbers and properties of and operations on them by primitive recursive relations and functions on their codes. We'll use sequences and their codes to do this. A tree can be either a single node (possibly with a label) or else a node (possibly with a label) connected to a number of subtrees. The node is called the *root* of the tree, and the subtrees it is connected to its *immediate subtrees*.

We code trees recursively as a sequence  $\langle k, d_1, \dots, d_k \rangle$ , where  $k$  is the number of immediate subtrees and  $d_1, \dots, d_k$  the codes of the immediate subtrees. If the nodes have labels, they can be included after the immediate subtrees. So a tree consisting just of a single node with label  $l$  would be coded by  $\langle 0, l \rangle$ , and a tree consisting of a root (labelled  $l_1$ ) connected to two single nodes (labelled  $l_2, l_3$ ) would be coded by  $\langle 2, \langle 0, l_2 \rangle, \langle 0, l_3 \rangle, l_1 \rangle$ .

**Proposition 15.25.** *The function  $\text{SubtreeSeq}(t)$ , which returns the code of a sequence the elements of which are the codes of all subtrees of the tree with code  $t$ , is primitive recursive.*

*Proof.* First note that  $\text{ISubtrees}(t) = \text{subseq}(t, 1, (t)_0)$  is primitive recursive and returns the codes of the immediate subtrees of a tree  $t$ . Now we can define a helper function  $\text{hSubtreeSeq}(t, n)$  which computes the sequence of all subtrees which are  $n$  nodes removed from the root. The sequence of subtrees of  $t$  which is 0 nodes removed from the root—in other words, begins at the root of  $t$ —is the sequence consisting just of  $t$ . To obtain a sequence of all level  $n+1$  subtrees of  $t$ , we concatenate the level  $n$  subtrees with a sequence consisting of all immediate subtrees of the level  $n$  subtrees. To get a list of all these, note that if  $f(x)$  is a primitive recursive function returning codes of sequences, then  $g_f(s, k) = f((s)_0) \frown \dots \frown f((s)_k)$  is also primitive recursive:

$$g(s, 0) = f((s)_0)$$



$$g(s, k + 1) = g(s, k) \frown f((s)_{k+1})$$

For instance, if  $s$  is a sequence of trees, then  $h(s) = g_{\text{ISubtrees}}(s, \text{len}(s))$  gives the sequence of the immediate subtrees of the elements of  $s$ . We can use it to define  $\text{hSubtreeSeq}$  by

$$\text{hSubtreeSeq}(t, 0) = \langle t \rangle$$

$$\text{hSubtreeSeq}(t, n + 1) = \text{hSubtreeSeq}(t, n) \frown h(\text{hSubtreeSeq}(t, n)).$$

The maximum level of subtrees in a tree coded by  $t$ , i.e., the maximum distance between the root and a leaf node, is bounded by the code  $t$ . So a sequence of codes of all subtrees of the tree coded by  $t$  is given by  $\text{hSubtreeSeq}(t, t)$ .  $\square$

### 15.13 Other Recursions

Using pairing and sequencing, we can justify more exotic (and useful) forms of primitive recursion. For example, it is often useful to define two functions simultaneously, such as in the following definition:

$$\begin{aligned} h_0(\vec{x}, 0) &= f_0(\vec{x}) \\ h_1(\vec{x}, 0) &= f_1(\vec{x}) \\ h_0(\vec{x}, y + 1) &= g_0(\vec{x}, y, h_0(\vec{x}, y), h_1(\vec{x}, y)) \\ h_1(\vec{x}, y + 1) &= g_1(\vec{x}, y, h_0(\vec{x}, y), h_1(\vec{x}, y)) \end{aligned}$$

This is an instance of *simultaneous recursion*. Another useful way of defining functions is to give the value of  $h(\vec{x}, y + 1)$  in terms of *all* the values  $h(\vec{x}, 0), \dots, h(\vec{x}, y)$ , as in the following definition:

$$\begin{aligned} h(\vec{x}, 0) &= f(\vec{x}) \\ h(\vec{x}, y + 1) &= g(\vec{x}, y, \langle h(\vec{x}, 0), \dots, h(\vec{x}, y) \rangle). \end{aligned}$$

The following schema captures this idea more succinctly:

$$h(\vec{x}, y) = g(\vec{x}, y, \langle h(\vec{x}, 0), \dots, h(\vec{x}, y - 1) \rangle)$$

with the understanding that the last argument to  $g$  is just the empty sequence when  $y$  is 0. In either formulation, the idea is that in computing the “successor step,” the function  $h$  can make use of the entire sequence of values computed so far. This is known as a *course-of-values* recursion. For a particular example, it can be used to justify the following type of definition:

$$h(\vec{x}, y) = \begin{cases} g(\vec{x}, y, h(\vec{x}, k(\vec{x}, y))) & \text{if } k(\vec{x}, y) < y \\ f(\vec{x}) & \text{otherwise} \end{cases}$$

In other words, the value of  $h$  at  $y$  can be computed in terms of the value of  $h$  at *any* previous value, given by  $k$ .

You should think about how to obtain these functions using ordinary primitive recursion. One final version of primitive recursion is more flexible in that one is allowed to change the *parameters* (side values) along the way:

$$\begin{aligned} h(\vec{x}, 0) &= f(\vec{x}) \\ h(\vec{x}, y + 1) &= g(\vec{x}, y, h(k(\vec{x}, y), y)) \end{aligned}$$

This, too, can be simulated with ordinary primitive recursion. (Doing so is tricky. For a hint, try unwinding the computation by hand.)

## 15.14 Non-Primitive Recursive Functions

The primitive recursive functions do not exhaust the intuitively computable functions. It should be intuitively clear that we can make a list of all the unary primitive recursive functions,  $f_0, f_1, f_2, \dots$  such that we can effectively compute the value of  $f_x$  on input  $y$ ; in other words, the function  $g(x, y)$ , defined by

$$g(x, y) = f_x(y)$$

is computable. But then so is the function

$$h(x) = g(x, x) + 1$$

$$= f_x(x) + 1.$$

For each primitive recursive function  $f_i$ , the value of  $h$  and  $f_i$  differ at  $i$ . So  $h$  is computable, but not primitive recursive; and one can say the same about  $g$ . This is an “effective” version of Cantor’s diagonalization argument.

One can provide more explicit examples of computable functions that are not primitive recursive. For example, let the notation  $g^n(x)$  denote  $g(g(\dots g(x)))$ , with  $n$   $g$ ’s in all; and define a sequence  $g_0, g_1, \dots$  of functions by

$$\begin{aligned} g_0(x) &= x + 1 \\ g_{n+1}(x) &= g_n^x(x) \end{aligned}$$

You can confirm that each function  $g_n$  is primitive recursive. Each successive function grows much faster than the one before;  $g_1(x)$  is equal to  $2x$ ,  $g_2(x)$  is equal to  $2^x \cdot x$ , and  $g_3(x)$  grows roughly like an exponential stack of  $x$  2’s. The Ackermann–Péter function is essentially the function  $G(x) = g_x(x)$ , and one can show that this grows faster than any primitive recursive function.

Let us return to the issue of enumerating the primitive recursive functions. Remember that we have assigned symbolic notations to each primitive recursive function; so it suffices to enumerate notations. We can assign a natural number  $\#(F)$  to each notation  $F$ , recursively, as follows:

$$\begin{aligned} \#(0) &= \langle 0 \rangle \\ \#(S) &= \langle 1 \rangle \\ \#(P_i^n) &= \langle 2, n, i \rangle \\ \#(\text{Comp}_{k,l}[H, G_0, \dots, G_{k-1}]) &= \langle 3, k, l, \#(H), \#(G_0), \dots, \#(G_{k-1}) \rangle \\ \#(\text{Rec}_l[G, H]) &= \langle 4, l, \#(G), \#(H) \rangle \end{aligned}$$

Here we are using the fact that every sequence of numbers can be viewed as a natural number, using the codes from the last section. The upshot is that every code is assigned a natural number. Of course, some sequences (and hence some numbers) do not

correspond to notations; but we can let  $f_i$  be the unary primitive recursive function with notation coded as  $i$ , if  $i$  codes such a notation; and the constant 0 function otherwise. The net result is that we have an explicit way of enumerating the unary primitive recursive functions.

(In fact, some functions, like the constant zero function, will appear more than once on the list. This is not just an artifact of our coding, but also a result of the fact that the constant zero function has more than one notation. We will later see that one can not computably avoid these repetitions; for example, there is no computable function that decides whether or not a given notation represents the constant zero function.)

We can now take the function  $g(x, y)$  to be given by  $f_x(y)$ , where  $f_x$  refers to the enumeration we have just described. How do we know that  $g(x, y)$  is computable? Intuitively, this is clear: to compute  $g(x, y)$ , first “unpack”  $x$ , and see if it is a notation for a unary function. If it is, compute the value of that function on input  $y$ .

You may already be convinced that (with some work!) one can write a program (say, in Java or C++) that does this; and now we can appeal to the Church–Turing thesis, which says that anything that, intuitively, is computable can be computed by a Turing machine.

Of course, a more direct way to show that  $g(x, y)$  is computable is to describe a Turing machine that computes it, explicitly. This would, in particular, avoid the Church–Turing thesis and appeals to intuition. Soon we will have built up enough machinery to show that  $g(x, y)$  is computable, appealing to a model of computation that can be *simulated* on a Turing machine: namely, the recursive functions.

## 15.15 Partial Recursive Functions

To motivate the definition of the recursive functions, note that our proof that there are computable functions that are not primi-

tive recursive actually establishes much more. The argument was simple: all we used was the fact that it is possible to enumerate functions  $f_0, f_1, \dots$  such that, as a function of  $x$  and  $y$ ,  $f_x(y)$  is computable. So the argument applies to *any class of functions that can be enumerated in such a way*. This puts us in a bind: we would like to describe the computable functions explicitly; but any explicit description of a collection of computable functions cannot be exhaustive!

The way out is to allow *partial* functions to come into play. We will see that it *is* possible to enumerate the partial computable functions. In fact, we already pretty much know that this is the case, since it is possible to enumerate Turing machines in a systematic way. We will come back to our diagonal argument later, and explore why it does not go through when partial functions are included.

The question is now this: what do we need to add to the primitive recursive functions to obtain all the partial recursive functions? We need to do two things:

1. Modify our definition of the primitive recursive functions to allow for partial functions as well.
2. *Add* something to the definition, so that some new partial functions are included.

The first is easy. As before, we will start with zero, successor, and projections, and close under composition and primitive recursion. The only difference is that we have to modify the definitions of composition and primitive recursion to allow for the possibility that some of the terms in the definition are not defined. If  $f$  and  $g$  are partial functions, we will write  $f(x) \downarrow$  to mean that  $f$  is defined at  $x$ , i.e.,  $x$  is in the domain of  $f$ ; and  $f(x) \uparrow$  to mean the opposite, i.e., that  $f$  is not defined at  $x$ . We will use  $f(x) \simeq g(x)$  to mean that either  $f(x)$  and  $g(x)$  are both undefined, or they are both defined and equal. We will use these notations for more complicated terms as well. We will adopt the

convention that if  $h$  and  $g_0, \dots, g_k$  all are partial functions, then

$$h(g_0(\vec{x}), \dots, g_k(\vec{x}))$$

is defined if and only if each  $g_i$  is defined at  $\vec{x}$ , and  $h$  is defined at  $g_0(\vec{x}), \dots, g_k(\vec{x})$ . With this understanding, the definitions of composition and primitive recursion for partial functions is just as above, except that we have to replace “=” by “ $\simeq$ ”.

What we will add to the definition of the primitive recursive functions to obtain partial functions is the *unbounded search operator*. If  $f(x, \vec{z})$  is any partial function on the natural numbers, define  $\mu x f(x, \vec{z})$  to be

the least  $x$  such that  $f(0, \vec{z}), f(1, \vec{z}), \dots, f(x, \vec{z})$  are all defined, and  $f(x, \vec{z}) = 0$ , if such an  $x$  exists

with the understanding that  $\mu x f(x, \vec{z})$  is undefined otherwise. This defines  $\mu x f(x, \vec{z})$  uniquely.

Note that our definition makes no reference to Turing machines, or algorithms, or any specific computational model. But like composition and primitive recursion, there is an operational, computational intuition behind unbounded search. When it comes to the computability of a partial function, arguments where the function is undefined correspond to inputs for which the computation does not halt. The procedure for computing  $\mu x f(x, \vec{z})$  will amount to this: compute  $f(0, \vec{z}), f(1, \vec{z}), f(2, \vec{z})$  until a value of 0 is returned. If any of the intermediate computations do not halt, however, neither does the computation of  $\mu x f(x, \vec{z})$ .

If  $R(x, \vec{z})$  is any relation,  $\mu x R(x, \vec{z})$  is defined to be  $\mu x (1 \dot{-} \chi_R(x, \vec{z}))$ . In other words,  $\mu x R(x, \vec{z})$  returns the least value of  $x$  such that  $R(x, \vec{z})$  holds. So, if  $f(x, \vec{z})$  is a total function,  $\mu x f(x, \vec{z})$  is the same as  $\mu x (f(x, \vec{z}) = 0)$ . But note that our original definition is more general, since it allows for the possibility that  $f(x, \vec{z})$  is not everywhere defined (whereas, in contrast, the characteristic function of a relation is always total).

**Definition 15.26.** The set of *partial recursive functions* is the smallest set of partial functions from the natural numbers to the natural numbers (of various arities) containing zero, successor, and projections, and closed under composition, primitive recursion, and unbounded search.

Of course, some of the partial recursive functions will happen to be total, i.e., defined for every argument.

**Definition 15.27.** The set of *recursive functions* is the set of partial recursive functions that are total.

A recursive function is sometimes called “total recursive” to emphasize that it is defined everywhere.

## 15.16 The Normal Form Theorem

**Theorem 15.28 (Kleene’s Normal Form Theorem).** *There is a primitive recursive relation  $T(e, x, s)$  and a primitive recursive function  $U(s)$ , with the following property: if  $f$  is any partial recursive function, then for some  $e$ ,*

$$f(x) \simeq U(\mu s T(e, x, s))$$

*for every  $x$ .*

The proof of the normal form theorem is involved, but the basic idea is simple. Every partial recursive function has an *index*  $e$ , intuitively, a number coding its program or definition. If  $f(x) \downarrow$ , the computation can be recorded systematically and coded by some number  $s$ , and the fact that  $s$  codes the computation of  $f$  on input  $x$  can be checked primitive recursively using only  $x$  and the definition  $e$ . Consequently, the relation  $T$ , “the function with index  $e$  has a computation for input  $x$ , and  $s$  codes this computation,” is primitive recursive. Given the full record of the com-

putation  $s$ , the “upshot” of  $s$  is the value of  $f(x)$ , and it can be obtained from  $s$  primitive recursively as well.

The normal form theorem shows that only a single unbounded search is required for the definition of any partial recursive function. Basically, we can search through all numbers until we find one that codes a computation of the function with index  $e$  for input  $x$ . We can use the numbers  $e$  as “names” of partial recursive functions, and write  $\varphi_e$  for the function  $f$  defined by the equation in the theorem. Note that any partial recursive function can have more than one index—in fact, every partial recursive function has infinitely many indices.

### 15.17 The Halting Problem

The *halting problem* in general is the problem of deciding, given the specification  $e$  (e.g., program) of a computable function and a number  $n$ , whether the computation of the function on input  $n$  halts, i.e., produces a result. Famously, Alan Turing proved that this problem itself cannot be solved by a computable function, i.e., the function

$$h(e, n) = \begin{cases} 1 & \text{if computation } e \text{ halts on input } n \\ 0 & \text{otherwise,} \end{cases}$$

is not computable.

In the context of partial recursive functions, the role of the specification of a program may be played by the index  $e$  given in Kleene’s normal form theorem. If  $f$  is a partial recursive function, any  $e$  for which the equation in the normal form theorem holds, is an index of  $f$ . Given a number  $e$ , the normal form theorem states that

$$\varphi_e(x) \simeq U(\mu s T(e, x, s))$$

is partial recursive, and for every partial recursive  $f: \mathbb{N} \rightarrow \mathbb{N}$ , there is an  $e \in \mathbb{N}$  such that  $\varphi_e(x) \simeq f(x)$  for all  $x \in \mathbb{N}$ . In fact,



for each such  $f$  there is not just one, but infinitely many such  $e$ . The *halting function*  $h$  is defined by

$$h(e, x) = \begin{cases} 1 & \text{if } \varphi_e(x) \downarrow \\ 0 & \text{otherwise.} \end{cases}$$

Note that  $h(e, x) = 0$  if  $\varphi_e(x) \uparrow$ , but also when  $e$  is not the index of a partial recursive function at all.

**Theorem 15.29.** *The halting function  $h$  is not partial recursive.*

*Proof.* If  $h$  were partial recursive, we could define

$$d(y) = \begin{cases} 1 & \text{if } h(y, y) = 0 \\ \mu x \ x \neq x & \text{otherwise.} \end{cases}$$

Since no number  $x$  satisfies  $x \neq x$ , there is no  $\mu x \ x \neq x$ , and so  $d(y) \uparrow$  iff  $h(y, y) \neq 0$ . From this definition it follows that

1.  $d(y) \downarrow$  iff  $\varphi_y(y) \uparrow$  or  $y$  is not the index of a partial recursive function.
2.  $d(y) \uparrow$  iff  $\varphi_y(y) \downarrow$ .

If  $h$  were partial recursive, then  $d$  would be partial recursive as well. Thus, by the Kleene normal form theorem, it has an index  $e_d$ . Consider the value of  $h(e_d, e_d)$ . There are two possible cases, 0 and 1.

1. If  $h(e_d, e_d) = 1$  then  $\varphi_{e_d}(e_d) \downarrow$ . But  $\varphi_{e_d} \simeq d$ , and  $d(e_d)$  is defined iff  $h(e_d, e_d) = 0$ . So  $h(e_d, e_d) \neq 1$ .
2. If  $h(e_d, e_d) = 0$  then either  $e_d$  is not the index of a partial recursive function, or it is and  $\varphi_{e_d}(e_d) \uparrow$ . But again,  $\varphi_{e_d} \simeq d$ , and  $d(e_d)$  is undefined iff  $\varphi_{e_d}(e_d) \downarrow$ .

The upshot is that  $e_d$  cannot, after all, be the index of a partial recursive function. But if  $h$  were partial recursive,  $d$  would be too, and so our definition of  $e_d$  as an index of it would be admissible. We must conclude that  $h$  cannot be partial recursive.  $\square$

## 15.18 General Recursive Functions

There is another way to obtain a set of total functions. Say a total function  $f(x, \vec{z})$  is *regular* if for every sequence of natural numbers  $\vec{z}$ , there is an  $x$  such that  $f(x, \vec{z}) = 0$ . In other words, the regular functions are exactly those functions to which one can apply unbounded search, and end up with a total function. One can, conservatively, restrict unbounded search to regular functions:

**Definition 15.30.** The set of *general recursive functions* is the smallest set of functions from the natural numbers to the natural numbers (of various arities) containing zero, successor, and projections, and closed under composition, primitive recursion, and unbounded search applied to *regular* functions.

Clearly every general recursive function is total. The difference between **Definition 15.30** and **Definition 15.27** is that in the latter one is allowed to use partial recursive functions along the way; the only requirement is that the function you end up with at the end is total. So the word “general,” a historic relic, is a misnomer; on the surface, **Definition 15.30** is *less* general than **Definition 15.27**. But, fortunately, the difference is illusory; though the definitions are different, the set of general recursive functions and the set of recursive functions are one and the same.

### Summary

In order to show that  $\mathbf{Q}$  represents all computable functions, we need a precise model of computability that we can take as the basis for a proof. There are, of course, many models of computability, such as Turing machines. One model that plays a significant role historically—it’s one of the first models proposed, and is also the one used by Gödel himself—is that of the **recursive functions**. The recursive functions are a class of arithmetical functions—that is, their domain and range are the natural numbers—that can be defined from a few basic functions using

a few operations. The basic functions are the constant zero function zero ( $\text{zero}(x) = 0$ ), the immediate successor function succ ( $\text{succ}(x) = x + 1$ ), and the projection functions. The operations are **composition**, **primitive recursion**, and **regular minimization**. Composition is simply a general version of “chaining together” functions: first apply one, then apply the other to the result. Primitive recursion defines a new function  $f$  from two functions  $g, h$  already defined, by stipulating that the value of  $f$  for 0 is given by  $g$ , and the value for any number  $n + 1$  is given by  $h$  applied to  $f(n)$ . Functions that can be defined using just these two principles are called **primitive recursive**. A relation is primitive recursive iff its characteristic function is. It turns out that a whole list of interesting functions and relations are primitive recursive (such as addition, multiplication, exponentiation, divisibility), and that we can define new primitive recursive functions and relations from old ones using principles such as bounded quantification and bounded minimization. In particular, this allowed us to show that we can deal with **sequences** of numbers in primitive recursive ways. That is, there is a way to “code” sequences of numbers as single numbers in such a way that we can compute the  $i$ -th element, the length, the concatenation of two sequences, etc., all using primitive recursive functions operating on these codes. To obtain all the computable functions, we finally added definition by **regular minimization** to composition and primitive recursion. A function  $g(x, y)$  is **regular** iff, for every  $y$  it takes the value 0 for at least one  $x$ . If  $f$  is regular, the least  $x$  such that  $g(x, y) = 0$  always exists, and can be found simply by computing all the values of  $g(0, y), g(1, y)$ , etc., until one of them is = 0. The resulting function  $f(y) = \mu x g(x, y) = 0$  is the function defined by regular minimization from  $g$ . It is always total and computable. The resulting set of functions are called **general recursive**. One version of the Church-Turing Thesis says that the computable arithmetical functions are exactly the general recursive ones.

## Problems

**Problem 15.1.** Prove [Proposition 15.5](#) by showing that the primitive recursive definition of mult can be put into the form required by [Definition 15.1](#) and showing that the corresponding functions  $f$  and  $g$  are primitive recursive.

**Problem 15.2.** Give the complete primitive recursive notation for mult.

**Problem 15.3.** Prove [Proposition 15.13](#).

**Problem 15.4.** Show that

$$f(x, y) = 2^{(2^{\dots^{2^x}})} \} y \text{ 2's}$$

is primitive recursive.

**Problem 15.5.** Show that integer division  $d(x, y) = \lfloor x/y \rfloor$  (i.e., division, where you disregard everything after the decimal point) is primitive recursive. When  $y = 0$ , we stipulate  $d(x, y) = 0$ . Give an explicit definition of  $d$  using primitive recursion and composition.

**Problem 15.6.** Show that the three place relation  $x \equiv y \pmod n$  (congruence modulo  $n$ ) is primitive recursive.

**Problem 15.7.** Suppose  $R(\vec{x}, z)$  is primitive recursive. Define the function  $m'_R(\vec{x}, y)$  which returns the least  $z$  less than  $y$  such that  $R(\vec{x}, z)$  holds, if there is one, and 0 otherwise, by primitive recursion from  $\chi_R$ .

**Problem 15.8.** Define integer division  $d(x, y)$  using bounded minimization.

**Problem 15.9.** Show that there is a primitive recursive function  $\text{sconcat}(s)$  with the property that

$$\text{sconcat}(\langle s_0, \dots, s_k \rangle) = s_0 \frown \dots \frown s_k.$$

**Problem 15.10.** Show that there is a primitive recursive function  $\text{tail}(s)$  with the property that

$$\begin{aligned}\text{tail}(\Lambda) &= 0 \text{ and} \\ \text{tail}(\langle s_0, \dots, s_k \rangle) &= \langle s_1, \dots, s_k \rangle.\end{aligned}$$

**Problem 15.11.** Prove **Proposition 15.24**.

**Problem 15.12.** The definition of  $\text{hSubtreeSeq}$  in the proof of **Proposition 15.25** in general includes repetitions. Give an alternative definition which guarantees that the code of a subtree occurs only once in the resulting list.

**Problem 15.13.** Define the remainder function  $r(x, y)$  by course-of-values recursion. (If  $x, y$  are natural numbers and  $y > 0$ ,  $r(x, y)$  is the number less than  $y$  such that  $x = z \times y + r(x, y)$  for some  $z$ . For definiteness, let's say that if  $y = 0$ ,  $r(x, 0) = 0$ .)

## CHAPTER 16

# *Arithmetization of Syntax*

### 16.1 Introduction

In order to connect computability and logic, we need a way to talk about the objects of logic (symbols, terms, formulas, derivations), operations on them, and their properties and relations, in a way amenable to computational treatment. We can do this directly, by considering computable functions and relations on symbols, sequences of symbols, and other objects built from them. Since the objects of logical syntax are all finite and built from a countable sets of symbols, this is possible for some models of computation. But other models of computation—such as the recursive functions—are restricted to numbers, their relations and functions. Moreover, ultimately we also want to be able to deal with syntax within certain theories, specifically, in theories formulated in the language of arithmetic. In these cases it is necessary to *arithmetize* syntax, i.e., to represent syntactic objects, operations on them, and their relations, as numbers, arithmetical functions, and arithmetical relations, respectively. The idea, which goes back to Leibniz, is to assign numbers to syntactic objects.

It is relatively straightforward to assign numbers to symbols as their “codes.” Some symbols pose a bit of a challenge, since,

e.g., there are infinitely many variables, and even infinitely many function symbols of each arity  $n$ . But of course it's possible to assign numbers to symbols systematically in such a way that, say,  $v_2$  and  $v_3$  are assigned different codes. Sequences of symbols (such as terms and formulas) are a bigger challenge. But if we can deal with sequences of numbers purely arithmetically (e.g., by the powers-of-primes coding of sequences), we can extend the coding of individual symbols to coding of sequences of symbols, and then further to sequences or other arrangements of formulas, such as derivations. This extended coding is called "Gödel numbering." Every term, formula, and derivation is assigned a Gödel number.

By coding sequences of symbols as sequences of their codes, and by choosing a system of coding sequences that can be dealt with using computable functions, we can then also deal with Gödel numbers using computable functions. In practice, all the relevant functions will be primitive recursive. For instance, computing the length of a sequence and computing the  $i$ -th element of a sequence from the code of the sequence are both primitive recursive. If the number coding the sequence is, e.g., the Gödel number of a formula  $A$ , we immediately see that the length of a formula and the (code of the)  $i$ -th symbol in a formula can also be computed from the Gödel number of  $A$ . It is a bit harder to prove that, e.g., the property of being the Gödel number of a correctly formed term or of a correct derivation is primitive recursive. It is nevertheless possible, because the sequences of interest (terms, formulas, derivations) are inductively defined.

As an example, consider the operation of substitution. If  $A$  is a formula,  $x$  a variable, and  $t$  a term, then  $A[t/x]$  is the result of replacing every free occurrence of  $x$  in  $A$  by  $t$ . Now suppose we have assigned Gödel numbers to  $A$ ,  $x$ ,  $t$ —say,  $k$ ,  $l$ , and  $m$ , respectively. The same scheme assigns a Gödel number to  $A[t/x]$ , say,  $n$ . This mapping—of  $k$ ,  $l$ , and  $m$  to  $n$ —is the arithmetical analog of the substitution operation. When the substitution operation maps  $A$ ,  $x$ ,  $t$  to  $A[t/x]$ , the arithmetized substitution function maps the Gödel numbers  $k$ ,  $l$ ,  $m$  to the Gödel number  $n$ . We will see that this function is primitive recursive.

Arithmetization of syntax is not just of abstract interest, although it was originally a non-trivial insight that languages like the language of arithmetic, which do not come with mechanisms for “talking about” languages can, after all, formalize complex properties of expressions. It is then just a small step to ask what a theory in this language, such as Peano arithmetic, can *prove* about its own language (including, e.g., whether sentences are provable or true). This leads us to the famous limitative theorems of Gödel (about unprovability) and Tarski (the undefinability of truth). But the trick of arithmetizing syntax is also important in order to prove some important results in computability theory, e.g., about the computational power of theories or the relationship between different models of computability. The arithmetization of syntax serves as a model for arithmetizing other objects and properties. For instance, it is similarly possible to arithmetize configurations and computations (say, of Turing machines). This makes it possible to simulate computations in one model (e.g., Turing machines) in another (e.g., recursive functions).

## 16.2 Coding Symbols

The basic language  $\mathcal{L}$  of first order logic makes use of the symbols

$$\perp \quad \neg \quad \vee \quad \wedge \quad \rightarrow \quad \forall \quad \exists \quad = \quad ( \quad ) \quad ,$$

together with countable sets of variables and constant symbols, and countable sets of function symbols and predicate symbols of arbitrary arity. We can assign *codes* to each of these symbols in such a way that every symbol is assigned a unique number as its code, and no two different symbols are assigned the same number. We know that this is possible since the set of all symbols is countable and so there is a bijection between it and the set of natural numbers. But we want to make sure that we can recover the symbol (as well as some information about it, e.g., the arity of a function symbol) from its code in a computable way. There are many possible ways of doing this, of course. Here is one such way,



which uses primitive recursive functions. (Recall that  $\langle n_0, \dots, n_k \rangle$  is the number coding the sequence of numbers  $n_0, \dots, n_k$ .)

**Definition 16.1.** If  $s$  is a symbol of  $\mathcal{L}$ , let the *symbol code*  $c_s$  be defined as follows:

1. If  $s$  is among the logical symbols,  $c_s$  is given by the following table:

$\perp$	$\neg$	$\vee$	$\wedge$	$\rightarrow$	$\forall$
$\langle 0, 0 \rangle$	$\langle 0, 1 \rangle$	$\langle 0, 2 \rangle$	$\langle 0, 3 \rangle$	$\langle 0, 4 \rangle$	$\langle 0, 5 \rangle$
$\exists$	$=$	$($	$)$	$,$	
$\langle 0, 6 \rangle$	$\langle 0, 7 \rangle$	$\langle 0, 8 \rangle$	$\langle 0, 9 \rangle$	$\langle 0, 10 \rangle$	

2. If  $s$  is the  $i$ -th variable  $v_i$ , then  $c_s = \langle 1, i \rangle$ .
3. If  $s$  is the  $i$ -th constant symbol  $c_i$ , then  $c_s = \langle 2, i \rangle$ .
4. If  $s$  is the  $i$ -th  $n$ -ary function symbol  $f_i^n$ , then  $c_s = \langle 3, n, i \rangle$ .
5. If  $s$  is the  $i$ -th  $n$ -ary predicate symbol  $P_i^n$ , then  $c_s = \langle 4, n, i \rangle$ .

**Proposition 16.2.** *The following relations are primitive recursive:*

1.  $\text{Fn}(x, n)$  iff  $x$  is the code of  $f_i^n$  for some  $i$ , i.e.,  $x$  is the code of an  $n$ -ary function symbol.
2.  $\text{Pred}(x, n)$  iff  $x$  is the code of  $P_i^n$  for some  $i$  or  $x$  is the code of  $=$  and  $n = 2$ , i.e.,  $x$  is the code of an  $n$ -ary predicate symbol.

**Definition 16.3.** If  $s_0, \dots, s_{n-1}$  is a sequence of symbols, its *Gödel number* is  $\langle c_{s_0}, \dots, c_{s_{n-1}} \rangle$ .

Note that *codes* and *Gödel numbers* are different things. For instance, the variable  $v_5$  has a code  $c_{v_5} = \langle 1, 5 \rangle = 2^2 \cdot 3^6$ . But the variable  $v_5$  considered as a term is also a sequence of symbols (of

length 1). The Gödel number  $\#v_5\#$  of the term  $v_5$  is  $\langle c_{v_5} \rangle = 2^{c_{v_5}+1} = 2^{2^2 \cdot 3^6 + 1}$ .

**Example 16.4.** Recall that if  $k_0, \dots, k_{n-1}$  is a sequence of numbers, then the code of the sequence  $\langle k_0, \dots, k_{n-1} \rangle$  in the power-of-primes coding is

$$2^{k_0+1} \cdot 3^{k_1+1} \cdot \dots \cdot p_{n-1}^{k_{n-1}},$$

where  $p_i$  is the  $i$ -th prime (starting with  $p_0 = 2$ ). So for instance, the formula  $v_0 = 0$ , or, more explicitly,  $(v_0, c_0)$ , has the Gödel number

$$\langle c_=(, c_0, c_{v_0}, c_0, c_0, c_0) \rangle.$$

Here,  $c_=($  is  $\langle 0, 7 \rangle = 2^{0+1} \cdot 3^{7+1}$ ,  $c_{v_0}$  is  $\langle 1, 0 \rangle = 2^{1+1} \cdot 3^{0+1}$ , etc. So  $\#(v_0, c_0)\#$  is

$$\begin{aligned} 2^{c_+=1} \cdot 3^{c_+=1} \cdot 5^{c_{v_0}+1} \cdot 7^{c_+=1} \cdot 11^{c_{c_0}+1} \cdot 13^{c_+=1} = \\ 2^{2^1 \cdot 3^8 + 1} \cdot 3^{2^1 \cdot 3^9 + 1} \cdot 5^{2^2 \cdot 3^1 + 1} \cdot 7^{2^1 \cdot 3^{11} + 1} \cdot 11^{2^3 \cdot 3^1 + 1} \cdot 13^{2^1 \cdot 3^{10} + 1} = \\ 2^{13 \cdot 123} \cdot 3^{39 \cdot 367} \cdot 5^{13} \cdot 7^{354 \cdot 295} \cdot 11^{25} \cdot 13^{118 \cdot 099}. \end{aligned}$$

## 16.3 Coding Terms

A term is simply a certain kind of sequence of symbols: it is built up inductively from constants and variables according to the formation rules for terms. Since sequences of symbols can be coded as numbers—using a coding scheme for the symbols plus a way to code sequences of numbers—assigning Gödel numbers to terms is not difficult. The challenge is rather to show that the property a number has if it is the Gödel number of a correctly formed term is computable, or in fact primitive recursive.

Variables and constant symbols are the simplest terms, and testing whether  $x$  is the Gödel number of such a term is easy:  $\text{Var}(x)$  holds if  $x$  is  $\#v_i\#$  for some  $i$ . In other words,  $x$  is a sequence of length 1 and its single element  $(x)_0$  is the code of some variable  $v_i$ , i.e.,  $x$  is  $\langle \langle 1, i \rangle \rangle$  for some  $i$ . Similarly,  $\text{Const}(x)$  holds

if  $x$  is  $\#C_i^\#$  for some  $i$ . Both of these relations are primitive recursive, since if such an  $i$  exists, it must be  $< x$ :

$$\begin{aligned}\text{Var}(x) &\Leftrightarrow (\exists i < x) x = \langle\langle 1, i \rangle\rangle \\ \text{Const}(x) &\Leftrightarrow (\exists i < x) x = \langle\langle 2, i \rangle\rangle\end{aligned}$$

**Proposition 16.5.** *The relations  $\text{Term}(x)$  and  $\text{CTerm}(x)$  which hold iff  $x$  is the Gödel number of a term or a closed term, respectively, are primitive recursive.*

*Proof.* A sequence of symbols  $s$  is a term iff there is a sequence  $s_0, \dots, s_{k-1} = s$  of terms which records how the term  $s$  was formed from constant symbols and variables according to the formation rules for terms. To express that such a putative formation sequence follows the formation rules it has to be the case that, for each  $i < k$ , either

1.  $s_i$  is a variable  $v_j$ , or
2.  $s_i$  is a constant symbol  $c_j$ , or
3.  $s_i$  is built from  $n$  terms  $t_1, \dots, t_n$  occurring prior to place  $i$  using an  $n$ -place function symbol  $f_j^n$ .

To show that the corresponding relation on Gödel numbers is primitive recursive, we have to express this condition primitive recursively, i.e., using primitive recursive functions, relations, and bounded quantification.

Suppose  $y$  is the number that codes the sequence  $s_0, \dots, s_{k-1}$ , i.e.,  $y = \langle\#s_0^\#, \dots, \#s_{k-1}^\#\rangle$ . It codes a formation sequence for the term with Gödel number  $x$  iff for all  $i < k$ :

1.  $\text{Var}((y)_i)$ , or
2.  $\text{Const}((y)_i)$ , or

3. there is an  $n$  and a number  $z = \langle z_1, \dots, z_n \rangle$  such that each  $z_l$  is equal to some  $(y)_{i'}$  for  $i' < i$  and

$$(y)_i = \#f_j^n(\# \frown \text{flatten}(z) \frown \#)^{\#},$$

and moreover  $(y)_{k-1} = x$ . (The function  $\text{flatten}(z)$  turns the sequence  $\langle \#t_1^{\#}, \dots, \#t_n^{\#} \rangle$  into  $\#t_1, \dots, t_n^{\#}$  and is primitive recursive.)

The indices  $j, n$ , the Gödel numbers  $z_l$  of the terms  $t_l$ , and the code  $z$  of the sequence  $\langle z_1, \dots, z_n \rangle$ , in (3) are all less than  $y$ . We can replace  $k$  above with  $\text{len}(y)$ . Hence we can express “ $y$  is the code of a formation sequence of the term with Gödel number  $x$ ” in a way that shows that this relation is primitive recursive.

We now just have to convince ourselves that there is a primitive recursive bound on  $y$ . But if  $x$  is the Gödel number of a term, it must have a formation sequence with at most  $\text{len}(x)$  terms (since every term in the formation sequence of  $s$  must start at some place in  $s$ , and no two subterms can start at the same place). The Gödel number of each subterm of  $s$  is of course  $\leq x$ . Hence, there always is a formation sequence with code  $\leq p_{k-1}^{k(x+1)}$ , where  $k = \text{len}(x)$ .

For  $\text{ClTerm}$ , simply leave out the clause for variables.  $\square$

**Proposition 16.6.** *The function  $\text{num}(n) = \#\bar{n}^{\#}$  is primitive recursive.*

*Proof.* We define  $\text{num}(n)$  by primitive recursion:

$$\begin{aligned} \text{num}(0) &= \#0^{\#} \\ \text{num}(n+1) &= \# \prime (\# \frown \text{num}(n) \frown \#)^{\#}. \end{aligned} \quad \square$$

## 16.4 Coding Formulas

**Proposition 16.7.** *The relation  $\text{Atom}(x)$  which holds iff  $x$  is the Gödel number of an atomic formula, is primitive recursive.*

*Proof.* The number  $x$  is the Gödel number of an atomic formula iff one of the following holds:

1. There are  $n, j < x$ , and  $z < x$  such that for each  $i < n$ ,  $\text{Term}((z)_i)$  and  $x =$

$$\#P_j^n(\# \frown \text{flatten}(z) \frown \#)^{\#}.$$

2. There are  $z_1, z_2 < x$  such that  $\text{Term}(z_1)$ ,  $\text{Term}(z_2)$ , and  $x =$

$$\#=(\# \frown z_1 \frown \#, \# \frown z_2 \frown \#)^{\#}.$$

3.  $x = \# \perp^{\#}$ . □

**Proposition 16.8.** *The relation  $\text{Frm}(x)$  which holds iff  $x$  is the Gödel number of a formula is primitive recursive.*

*Proof.* A sequence of symbols  $s$  is a formula iff there is formation sequence  $s_0, \dots, s_{k-1} = s$  of formula which records how  $s$  was formed from atomic formulas according to the formation rules. The code for each  $s_i$  (and indeed of the code of the sequence  $\langle s_0, \dots, s_{k-1} \rangle$ ) is less than the code  $x$  of  $s$ . □

**Proposition 16.9.** *The relation  $\text{FreeOcc}(x, z, i)$ , which holds iff the  $i$ -th symbol of the formula with Gödel number  $x$  is a free occurrence of the variable with Gödel number  $z$ , is primitive recursive.*

*Proof.* Exercise. □

**Proposition 16.10.** *The property  $\text{Sent}(x)$  which holds iff  $x$  is the Gödel number of a sentence is primitive recursive.*

*Proof.* A sentence is a formula without free occurrences of variables. So  $\text{Sent}(x)$  holds iff

$$(\forall i < \text{len}(x)) (\forall z < x) \\ (\exists j < z) z = \#v_j \rightarrow \neg \text{FreeOcc}(x, z, i). \quad \square$$

## 16.5 Substitution

Recall that substitution is the operation of replacing all free occurrences of a variable  $u$  in a formula  $A$  by a term  $t$ , written  $A[t/u]$ . This operation, when carried out on Gödel numbers of variables, formulas, and terms, is primitive recursive.

**Proposition 16.11.** *There is a primitive recursive function  $\text{Subst}(x, y, z)$  with the property that*

$$\text{Subst}(\#A\#, \#t\#, \#u\#) = \#A[t/u]\#.$$

*Proof.* We can then define a function  $\text{hSubst}$  by primitive recursion as follows:

$$\begin{aligned} \text{hSubst}(x, y, z, 0) &= \Lambda \\ \text{hSubst}(x, y, z, i + 1) &= \\ &\begin{cases} \text{hSubst}(x, y, z, i) \frown y & \text{if } \text{FreeOcc}(x, z, i) \\ \text{append}(\text{hSubst}(x, y, z, i), (x)_i) & \text{otherwise.} \end{cases} \end{aligned}$$

$\text{Subst}(x, y, z)$  can now be defined as  $\text{hSubst}(x, y, z, \text{len}(x))$ .  $\square$

**Proposition 16.12.** *The relation  $\text{FreeFor}(x, y, z)$ , which holds iff the term with Gödel number  $y$  is free for the variable with Gödel number  $z$  in the formula with Gödel number  $x$ , is primitive recursive.*

*Proof.* Exercise.  $\square$

## 16.6 Derivations in LK

In order to arithmetize derivations, we must represent derivations as numbers. Since derivations are trees of sequents where each inference carries also a label, a recursive representation is the most obvious approach: we represent a derivation as a tuple, the components of which are the end-sequent, the label, and the representations of the sub-derivations leading to the premises of the last inference.

**Definition 16.13.** If  $\Gamma$  is a finite sequence of sentences,  $\Gamma = \langle A_1, \dots, A_n \rangle$ , then  $\# \Gamma^\# = \langle \# A_1^\#, \dots, \# A_n^\# \rangle$ .

If  $\Gamma \Rightarrow \Delta$  is a sequent, then a Gödel number of  $\Gamma \Rightarrow \Delta$  is

$$\# \Gamma \Rightarrow \Delta^\# = \langle \# \Gamma^\#, \# \Delta^\# \rangle$$

If  $\pi$  is a derivation in **LK**, then  $\#\pi^\#$  is defined as follows:

1. If  $\pi$  consists only of the initial sequent  $\Gamma \Rightarrow \Delta$ , then  $\#\pi^\#$  is

$$\langle 0, \# \Gamma \Rightarrow \Delta^\# \rangle.$$

2. If  $\pi$  ends in an inference with one or two premises, has  $\Gamma \Rightarrow \Delta$  as its conclusion, and  $\pi_1$  and  $\pi_2$  are the immediate subproof ending in the premise of the last inference, then  $\#\pi^\#$  is

$$\langle 1, \# \pi_1^\#, \# \Gamma \Rightarrow \Delta^\#, k \rangle \text{ or}$$

$$\langle 2, \# \pi_1^\#, \# \pi_2^\#, \# \Gamma \Rightarrow \Delta^\#, k \rangle,$$

respectively, where  $k$  is given by the following table according to which rule was used in the last inference:

Rule:	WL	WR	CL	CR	XL	XR
$k$ :	1	2	3	4	5	6
Rule:	$\neg$ L	$\neg$ R	$\wedge$ L	$\wedge$ R	$\vee$ L	$\vee$ R
$k$ :	7	8	9	10	11	12
Rule:	$\rightarrow$ L	$\rightarrow$ R	$\forall$ L	$\forall$ R	$\exists$ L	$\exists$ R
$k$ :	13	14	15	16	17	18
Rule:	Cut	=				
$k$ :	19	20				

**Example 16.14.** Consider the very simple derivation

$$\frac{\frac{A \Rightarrow A}{A \wedge B \Rightarrow A} \wedge L}{\Rightarrow (A \wedge B) \rightarrow A} \rightarrow R$$

The Gödel number of the initial sequent would be  $p_0 = \langle 0, \#A \Rightarrow A^\# \rangle$ . The Gödel number of the derivation ending in the conclusion of  $\wedge L$  would be  $p_1 = \langle 1, p_0, \#A \wedge B \Rightarrow A^\#, 9 \rangle$  (1 since  $\wedge L$  has one premise, the Gödel number of the conclusion  $A \wedge B \Rightarrow A$ , and 9 is the number coding  $\wedge L$ ). The Gödel number of the entire derivation then is  $\langle 1, p_1, \#\Rightarrow (A \wedge B) \rightarrow A^\#, 14 \rangle$ , i.e.,

$$\langle 1, \langle 1, \langle 0, \#A \Rightarrow A^\# \rangle, \#A \wedge B \Rightarrow A^\#, 9 \rangle, \#\Rightarrow (A \wedge B) \rightarrow A^\#, 14 \rangle.$$

Having settled on a representation of derivations, we must also show that we can manipulate such derivations primitive recursively, and express their essential properties and relations so. Some operations are simple: e.g., given a Gödel number  $p$  of a derivation,  $\text{EndSeq}(p) = (p)_{(p)_0+1}$  gives us the Gödel number of its end-sequent and  $\text{LastRule}(p) = (p)_{(p)_0+2}$  the code of its last rule. The property  $\text{Sequent}(s)$  defined by

$$\text{len}(s) = 2 \wedge (\forall i < \text{len}((s)_0) + \text{len}((s)_1)) \text{Sent}(((s)_0 \frown (s)_1)_i)$$



holds of  $s$  iff  $s$  is the Gödel number of a sequent consisting of sentences. Some are much harder. We'll at least sketch how to do this. The goal is to show that the relation " $\pi$  is a derivation of  $A$  from  $\Gamma$ " is a primitive recursive relation of the Gödel numbers of  $\pi$  and  $A$ .

**Proposition 16.15.** *The property  $\text{Correct}(p)$  which holds iff the last inference in the derivation  $\pi$  with Gödel number  $p$  is correct, is primitive recursive.*

*Proof.*  $\Gamma \Rightarrow \Delta$  is an initial sequent if either there is a sentence  $A$  such that  $\Gamma \Rightarrow \Delta$  is  $A \Rightarrow A$ , or there is a term  $t$  such that  $\Gamma \Rightarrow \Delta$  is  $\emptyset \Rightarrow t = t$ . In terms of Gödel numbers,  $\text{InitSeq}(s)$  holds iff

$$(\exists x < s) (\text{Sent}(x) \wedge s = \langle \langle x \rangle, \langle x \rangle \rangle) \vee$$

$$(\exists t < s) (\text{Term}(t) \wedge s = \langle 0, \langle \# = (\# \frown t \frown \#, \# \frown t \frown \#) \rangle \rangle).$$

We also have to show that for each rule of inference  $R$  the relation  $\text{FollowsBy}_R(p)$  is primitive recursive, where  $\text{FollowsBy}_R(p)$  holds iff  $p$  is the Gödel number of derivation  $\pi$ , and the end-sequent of  $\pi$  follows by a correct application of  $R$  from the immediate sub-derivations of  $\pi$ .

A simple case is that of the  $\wedge R$  rule. If  $\pi$  ends in a correct  $\wedge R$  inference, it looks like this:

$$\frac{\begin{array}{c} \vdots \\ \pi_1 \\ \vdots \\ \Gamma \Rightarrow \Delta, A \end{array} \quad \begin{array}{c} \vdots \\ \pi_2 \\ \vdots \\ \Gamma \Rightarrow \Delta, B \end{array}}{\Gamma \Rightarrow \Delta, A \wedge B} \wedge R$$

So, the last inference in the derivation  $\pi$  is a correct application of  $\wedge R$  iff there are sequences of sentences  $\Gamma$  and  $\Delta$  as well as two sentences  $A$  and  $B$  such that the end-sequent of  $\pi_1$  is  $\Gamma \Rightarrow \Delta, A$ , the end-sequent of  $\pi_2$  is  $\Gamma \Rightarrow \Delta, B$ , and the end-sequent of  $\pi$  is  $\Gamma \Rightarrow \Delta, A \wedge B$ . We just have to translate this into Gödel numbers. If  $s = \# \Gamma \Rightarrow \Delta \#$  then  $(s)_0 = \# \Gamma \#$  and  $(s)_1 = \# \Delta \#$ . So,  $\text{FollowsBy}_{\wedge R}(p)$  holds iff

$$(\exists g < p) (\exists d < p) (\exists a < p) (\exists b < p)$$

$$\begin{aligned}
\text{EndSequent}(\mathcal{p}) &= \langle g, d \frown \langle \#(\# \frown a \frown \# \wedge \# \frown b \frown \#) \# \rangle \rangle \wedge \\
\text{EndSequent}((\mathcal{p})_1) &= \langle g, d \frown \langle a \rangle \rangle \wedge \\
\text{EndSequent}((\mathcal{p})_2) &= \langle g, d \frown \langle b \rangle \rangle \wedge \\
(\mathcal{p})_0 &= 2 \wedge \text{LastRule}(\mathcal{p}) = 10.
\end{aligned}$$

The individual lines express, respectively, “there is a sequence  $\Gamma$  with Gödel number  $g$ , there is a sequence  $\Delta$  with Gödel number  $d$ , a formula  $A$  with Gödel number  $a$ , and a formula  $B$  with Gödel number  $b$ ,” such that “the end-sequent of  $\pi$  is  $\Gamma \Rightarrow \Delta, A \wedge B$ ,” “the end-sequent of  $\pi_1$  is  $\Gamma \Rightarrow \Delta, A$ ,” “the end-sequent of  $\pi_2$  is  $\Gamma \Rightarrow \Delta, B$ ,” and “ $\pi$  has two immediate subderivations and the last inference rule is  $\wedge R$  (with number 10).”

The last inference in  $\pi$  is a correct application of  $\exists R$  iff there are sequences  $\Gamma$  and  $\Delta$ , a formula  $A$ , a variable  $x$ , and a term  $t$ , such that the end-sequent of  $\pi$  is  $\Gamma \Rightarrow \Delta, \exists x A$  and the end-sequent of  $\pi_1$  is  $\Gamma \Rightarrow \Delta, A[t/x]$ . So in terms of Gödel numbers, we have  $\text{FollowsBy}_{\exists R}(\mathcal{p})$  iff

$$\begin{aligned}
&(\exists g < \mathcal{p}) (\exists d < \mathcal{p}) (\exists a < \mathcal{p}) (\exists x < \mathcal{p}) (\exists t < \mathcal{p}) \\
&\text{EndSequent}(\mathcal{p}) = \langle g, d \frown \langle \# \exists \# \frown x \frown a \rangle \rangle \wedge \\
&\text{EndSequent}((\mathcal{p})_1) = \langle g, d \frown \langle \text{Subst}(a, t, x) \rangle \rangle \wedge \\
&(\mathcal{p})_0 = 1 \wedge \text{LastRule}(\mathcal{p}) = 18.
\end{aligned}$$

We then define  $\text{Correct}(\mathcal{p})$  as

$$\begin{aligned}
&\text{Sequent}(\text{EndSequent}(\mathcal{p})) \wedge \\
&[(\text{LastRule}(\mathcal{p}) = 1 \wedge \text{FollowsBy}_{\text{WL}}(\mathcal{p})) \vee \cdots \vee \\
&\quad (\text{LastRule}(\mathcal{p}) = 20 \wedge \text{FollowsBy}_{\_}(\mathcal{p})) \vee \\
&\quad (\mathcal{p})_0 = 0 \wedge \text{InitialSeq}(\text{EndSequent}(\mathcal{p}))]
\end{aligned}$$

The first line ensures that the end-sequent of  $d$  is actually a sequent consisting of sentences. The last line covers the case where  $\mathcal{p}$  is just an initial sequent.  $\square$

**Proposition 16.16.** *The relation  $\text{Deriv}(p)$  which holds if  $p$  is the Gödel number of a correct derivation  $\pi$ , is primitive recursive.*

*Proof.* A derivation  $\pi$  is correct if every one of its inferences is a correct application of a rule, i.e., if every one of its sub-derivations ends in a correct inference. So,  $\text{Deriv}(d)$  iff

$$(\forall i < \text{len}(\text{SubtreeSeq}(p))) \text{Correct}((\text{SubtreeSeq}(p))_i). \quad \square$$

**Proposition 16.17.** *Suppose  $\Gamma$  is a primitive recursive set of sentences. Then the relation  $\text{Prf}_\Gamma(x, y)$  expressing “ $x$  is the code of a derivation  $\pi$  of  $\Gamma_0 \Rightarrow A$  for some finite  $\Gamma_0 \subseteq \Gamma$  and  $y$  is the Gödel number of  $A$ ” is primitive recursive.*

*Proof.* Suppose “ $y \in \Gamma$ ” is given by the primitive recursive predicate  $R_\Gamma(y)$ . We have to show that  $\text{Prf}_\Gamma(x, y)$  which holds iff  $y$  is the Gödel number of a sentence  $A$  and  $x$  is the code of an **LK**-derivation with end-sequent  $\Gamma_0 \Rightarrow A$  is primitive recursive.

By the previous proposition, the property  $\text{Deriv}(x)$  which holds iff  $x$  is the code of a correct derivation  $\pi$  in **LK** is primitive recursive. If  $x$  is such a code, then  $\text{EndSequent}(x)$  is the code of the end-sequent of  $\pi$ , and so  $(\text{EndSequent}(x))_0$  is the code of the left side of the end sequent and  $(\text{EndSequent}(x))_1$  the right side. So we can express “the right side of the end-sequent of  $\pi$  is  $A$ ” as  $\text{len}((\text{EndSequent}(x))_1) = 1 \wedge ((\text{EndSequent}(x))_1)_0 = x$ . The left side of the end-sequent of  $\pi$  is of course automatically finite, we just have to express that every sentence in it is in  $\Gamma$ . Thus we can define  $\text{Prf}_\Gamma(x, y)$  by

$$\begin{aligned} \text{Prf}_\Gamma(x, y) \Leftrightarrow & \text{Deriv}(x) \wedge \\ & (\forall i < \text{len}((\text{EndSequent}(x))_0)) R_\Gamma(((\text{EndSequent}(x))_0)_i) \wedge \\ & \text{len}((\text{EndSequent}(x))_1) = 1 \wedge ((\text{EndSequent}(x))_1)_0 = y. \end{aligned}$$

## 16.7 Derivations in Natural Deduction

In order to arithmetize derivations, we must represent derivations as numbers. Since derivations are trees of formulas where each inference carries one or two labels, a recursive representation is the most obvious approach: we represent a derivation as a tuple, the components of which are the number of immediate sub-derivations leading to the premises of the last inference, the representations of these sub-derivations, and the end-formula, the discharge label of the last inference, and a number indicating the type of the last inference.

**Definition 16.18.** If  $\delta$  is a derivation in natural deduction, then  $\# \delta^\#$  is defined inductively as follows:

1. If  $\delta$  consists only of the assumption  $A$ , then  $\# \delta^\#$  is  $\langle 0, \# A^\#, n \rangle$ . The number  $n$  is 0 if it is an undischarged assumption, and the numerical label otherwise.
2. If  $\delta$  ends in an inference with one, two, or three premises, then  $\# \delta^\#$  is

$$\begin{aligned} &\langle 1, \# \delta_1^\#, \# A^\#, n, k \rangle, \\ &\langle 2, \# \delta_1^\#, \# \delta_2^\#, \# A^\#, n, k \rangle, \text{ or} \\ &\langle 3, \# \delta_1^\#, \# \delta_2^\#, \# \delta_3^\#, \# A^\#, n, k \rangle, \end{aligned}$$

respectively. Here  $\delta_1, \delta_2, \delta_3$  are the sub-derivations ending in the premise(s) of the last inference in  $\delta$ ,  $A$  is the conclusion of the last inference in  $\delta$ ,  $n$  is the discharge label of the last inference (0 if the inference does not discharge any assumptions), and  $k$  is given by the following table according to which rule was used in the last inference.

Rule:	$\wedge$ Intro	$\wedge$ Elim	$\vee$ Intro	$\vee$ Elim
$k$ :	1	2	3	4
Rule:	$\rightarrow$ Intro	$\rightarrow$ Elim	$\neg$ Intro	$\neg$ Elim
$k$ :	5	6	7	8
Rule:	$\perp_I$	$\perp_C$	$\forall$ Intro	$\forall$ Elim
$k$ :	9	10	11	12
Rule:	$\exists$ Intro	$\exists$ Elim	$=$ Intro	$=$ Elim
$k$ :	13	14	15	16

**Example 16.19.** Consider the very simple derivation

$$\begin{array}{c}
 [A \wedge B]^1 \\
 \hline
 A \\
 \hline
 (A \wedge B) \rightarrow A \quad \rightarrow\text{Intro} \\
 \wedge\text{Elim}
 \end{array}$$

The Gödel number of the assumption would be  $d_0 = \langle 0, \#A \wedge B\#, 1 \rangle$ . The Gödel number of the derivation ending in the conclusion of  $\wedge$ Elim would be  $d_1 = \langle 1, d_0, \#A\#, 0, 2 \rangle$  (1 since  $\wedge$ Elim has one premise, the Gödel number of conclusion  $A$ , 0 because no assumption is discharged, and 2 is the number coding  $\wedge$ Elim). The Gödel number of the entire derivation then is  $\langle 1, d_1, \#((A \wedge B) \rightarrow A)\#, 1, 5 \rangle$ , i.e.,

$$\langle 1, \langle 1, \langle 0, \#(A \wedge B)\#, 1 \rangle, \#A\#, 0, 2 \rangle, \#((A \wedge B) \rightarrow A)\#, 1, 5 \rangle.$$

Having settled on a representation of derivations, we must also show that we can manipulate Gödel numbers of such derivations primitive recursively, and express their essential properties and relations. Some operations are simple: e.g., given a Gödel number  $d$  of a derivation,  $\text{EndFmla}(d) = (d)_{(d)_0+1}$  gives us the Gödel number of its end-formula,  $\text{DischargeLabel}(d) = (d)_{(d)_0+2}$  gives us the discharge label and  $\text{LastRule}(d) = (d)_{(d)_0+3}$  the number indicating the type of the last inference. Some are much harder. We'll at least sketch how to do this. The goal is to show that the relation " $\delta$  is a derivation of  $A$  from  $\Gamma$ " is a primitive recursive relation of the Gödel numbers of  $\delta$  and  $A$ .

**Proposition 16.20.** *The following relations are primitive recursive:*

1. *A occurs as an assumption in  $\delta$  with label  $n$ .*
2. *All assumptions in  $\delta$  with label  $n$  are of the form  $A$  (i.e., we can discharge the assumption  $A$  using label  $n$  in  $\delta$ ).*

*Proof.* We have to show that the corresponding relations between Gödel numbers of formulas and Gödel numbers of derivations are primitive recursive.

1. We want to show that  $\text{Assum}(x, d, n)$ , which holds if  $x$  is the Gödel number of an assumption of the derivation with Gödel number  $d$  labelled  $n$ , is primitive recursive. This is the case if the derivation with Gödel number  $\langle 0, x, n \rangle$  is a sub-derivation of  $d$ . Note that the way we code derivations is a special case of the coding of trees introduced in [section 15.12](#), so the primitive recursive function  $\text{SubtreeSeq}(d)$  gives a sequence of Gödel numbers of all sub-derivations of  $d$  (of length at most  $d$ ). So we can define

$$\text{Assum}(x, d, n) \Leftrightarrow (\exists i < d) (\text{SubtreeSeq}(d))_i = \langle 0, x, n \rangle.$$

2. We want to show that  $\text{Discharge}(x, d, n)$ , which holds if all assumptions with label  $n$  in the derivation with Gödel number  $d$  all are the formula with Gödel number  $x$ . But this relation holds iff  $(\forall y < d) (\text{Assum}(y, d, n) \rightarrow y = x)$ .  $\square$

**Proposition 16.21.** *The property  $\text{Correct}(d)$  which holds iff the last inference in the derivation  $\delta$  with Gödel number  $d$  is correct, is primitive recursive.*

*Proof.* Here we have to show that for each rule of inference  $R$  the relation  $\text{FollowsBy}_R(d)$  is primitive recursive, where  $\text{FollowsBy}_R(d)$  holds iff  $d$  is the Gödel number of derivation  $\delta$ , and the end-formula of  $\delta$  follows by a correct application of  $R$  from the immediate sub-derivations of  $\delta$ .

A simple case is that of the  $\wedge$ Intro rule. If  $\delta$  ends in a correct  $\wedge$ Intro inference, it looks like this:

$$\frac{\begin{array}{c} \vdots \\ \delta_1 \\ \vdots \\ A \end{array} \quad \begin{array}{c} \vdots \\ \delta_2 \\ \vdots \\ B \end{array}}{A \wedge B} \wedge\text{Intro}$$

Then the Gödel number  $d$  of  $\delta$  is  $\langle 2, d_1, d_2, \#(A \wedge B)\#, 0, k \rangle$  where  $\text{EndFmla}(d_1) = \#A\#$ ,  $\text{EndFmla}(d_2) = \#B\#$ ,  $n = 0$ , and  $k = 1$ . So we can define  $\text{FollowsBy}_{\wedge\text{Intro}}(d)$  as

$$(d)_0 = 2 \wedge \text{DischargeLabel}(d) = 0 \wedge \text{LastRule}(d) = 1 \wedge \\ \text{EndFmla}(d) = \#(\# \frown \text{EndFmla}((d)_1) \frown \#\wedge\# \frown \text{EndFmla}((d)_2) \frown \#)\#.$$

Another simple example is the  $=$ Intro rule. Here the premise is an empty derivation, i.e.,  $(d)_1 = 0$ , and no discharge label, i.e.,  $n = 0$ . However,  $A$  must be of the form  $t = t$ , for a closed term  $t$ . Here, a primitive recursive definition is

$$(d)_0 = 1 \wedge (d)_1 = 0 \wedge \text{DischargeLabel}(d) = 0 \wedge \\ (\exists t < d) (\text{ClTerm}(t) \wedge \text{EndFmla}(d) = \#(\# \frown t \frown \# \frown t \frown \#)\#)$$

For a more complicated example,  $\text{FollowsBy}_{\rightarrow\text{Intro}}(d)$  holds iff the end-formula of  $\delta$  is of the form  $(A \rightarrow B)$ , where the end-formula of  $\delta_1$  is  $B$ , and any assumption in  $\delta$  labelled  $n$  is of the form  $A$ . We can express this primitive recursively by

$$(d)_0 = 1 \wedge \\ (\exists a < d) (\text{Discharge}(a, (d)_1, \text{DischargeLabel}(d)) \wedge \\ \text{EndFmla}(d) = \#(\# \frown a \frown \#\rightarrow\# \frown \text{EndFmla}((d)_1) \frown \#)\#)$$

(Think of  $a$  as the Gödel number of  $A$ ).

For another example, consider  $\exists$ Intro. Here, the last inference in  $\delta$  is correct iff there is a formula  $A$ , a closed term  $t$  and

a variable  $x$  such that  $A[t/x]$  is the end-formula of the derivation  $\delta_1$  and  $\exists x A$  is the conclusion of the last inference. So,  $\text{FollowsBy}_{\exists\text{Intro}}(d)$  holds iff

$$(d)_0 = 1 \wedge \text{DischargeLabel}(d) = 0 \wedge \\ (\exists a < d) (\exists x < d) (\exists t < d) (\text{CI}(\text{Term}(t)) \wedge \text{Var}(x) \wedge \\ \text{Subst}(a, t, x) = \text{EndFmla}((d)_1) \wedge \text{EndFmla}(d) = (\exists x \# \wedge x \wedge a)).$$

We then define  $\text{Correct}(d)$  as

$$\text{Sent}(\text{EndFmla}(d)) \wedge \\ (\text{LastRule}(d) = 1 \wedge \text{FollowsBy}_{\wedge\text{Intro}}(d)) \vee \dots \vee \\ (\text{LastRule}(d) = 16 \wedge \text{FollowsBy}_{=\text{Elim}}(d)) \vee \\ (\exists n < d) (\exists x < d) (d = \langle 0, x, n \rangle).$$

The first line ensures that the end-formula of  $d$  is a sentence. The last line covers the case where  $d$  is just an assumption.  $\square$

**Proposition 16.22.** *The relation  $\text{Deriv}(d)$  which holds if  $d$  is the Gödel number of a correct derivation  $\delta$ , is primitive recursive.*

*Proof.* A derivation  $\delta$  is correct if every one of its inferences is a correct application of a rule, i.e., if every one of its sub-derivations ends in a correct inference. So,  $\text{Deriv}(d)$  iff

$$(\forall i < \text{len}(\text{SubtreeSeq}(d))) \text{Correct}((\text{SubtreeSeq}(d))_i) \quad \square$$

**Proposition 16.23.** *The relation  $\text{OpenAssum}(z, d)$  that holds if  $z$  is the Gödel number of an undischarged assumption  $A$  of the derivation  $\delta$  with Gödel number  $d$ , is primitive recursive.*

*Proof.* An occurrence of an assumption is discharged if it occurs with label  $n$  in a sub-derivation of  $\delta$  that ends in a rule with discharge label  $n$ . So  $A$  is an undischarged assumption of  $\delta$  if at



least one of its occurrences is not discharged in  $\delta$ . We must be careful:  $\delta$  may contain both discharged and undischarged occurrences of  $A$ .

Consider a sequence  $\delta_0, \dots, \delta_k$  where  $\delta_0 = \delta$ ,  $\delta_k$  is the assumption  $[A]^n$  (for some  $n$ ), and  $\delta_{i+1}$  is an immediate sub-derivation of  $\delta_i$ . If such a sequence exists in which no  $\delta_i$  ends in an inference with discharge label  $n$ , then  $A$  is an undischarged assumption of  $\delta$ .

The primitive recursive function  $\text{SubtreeSeq}(d)$  provides us with a sequence of Gödel numbers of all sub-derivations of  $\delta$ . Any sequence of Gödel numbers of sub-derivations of  $\delta$  is a subsequence of it. Being a subsequence of is a primitive recursive relation:  $\text{Subseq}(s, s')$  holds iff  $(\forall i < \text{len}(s)) \exists j < \text{len}(s') (s)_i = (s')_j$ . Being an immediate sub-derivation is as well:  $\text{Subderiv}(d, d')$  iff  $(\exists j < (d')_0) d = (d')_j$ . So we can define  $\text{OpenAssum}(z, d)$  by

$$\begin{aligned} (\exists s < \text{SubtreeSeq}(d)) & (\text{Subseq}(s, \text{SubtreeSeq}(d)) \wedge (s)_0 = d \wedge \\ & (\exists n < d) ((s)_{\text{len}(s)-1} = \langle 0, z, n \rangle \wedge \\ & (\forall i < (\text{len}(s) - 1)) (\text{Subderiv}((s)_{i+1}, (s)_i)) \wedge \\ & \text{DischargeLabel}((s)_{i+1}) \neq n)). \quad \square \end{aligned}$$

**Proposition 16.24.** *Suppose  $\Gamma$  is a primitive recursive set of sentences. Then the relation  $\text{Prf}_\Gamma(x, y)$  expressing “ $x$  is the code of a derivation  $\delta$  of  $A$  from undischarged assumptions in  $\Gamma$  and  $y$  is the Gödel number of  $A$ ” is primitive recursive.*

*Proof.* Suppose “ $y \in \Gamma$ ” is given by the primitive recursive predicate  $R_\Gamma(y)$ . We have to show that  $\text{Prf}_\Gamma(x, y)$  which holds iff  $y$  is the Gödel number of a sentence  $A$  and  $x$  is the code of a natural deduction derivation with end formula  $A$  and all undischarged assumptions in  $\Gamma$  is primitive recursive.

By **Proposition 16.22**, the property  $\text{Deriv}(x)$  which holds iff  $x$  is the Gödel number of a correct derivation  $\delta$  in natural deduction is primitive recursive. Thus we can define  $\text{Prf}_\Gamma(x, y)$  by

$$\text{Prf}_\Gamma(x, y) \Leftrightarrow \text{Deriv}(x) \wedge \text{EndFmla}(x) = y \wedge$$

$$(\forall z < x) (\text{OpenAssum}(z, x) \rightarrow R_I(z)). \quad \square$$

## Summary

The proof of the incompleteness theorems requires that we have a way to talk about provability in a theory (such as **PA**) in the language of the theory itself, i.e., in the language of arithmetic. But the language of arithmetic only deals with numbers, not with formulas or derivations. The solution to this problem is to define a systematic mapping from formulas and derivations to numbers. The number associated with a formula or a derivation is called its **Gödel number**. If  $A$  is a formula,  $\#A^\#$  is its Gödel number. We showed that important operations on formulas turn into primitive recursive functions on the respective Gödel numbers. For instance,  $A[t/x]$ , the operation of substituting a term  $t$  for every free occurrence of  $x$  in  $A$ , corresponds to an arithmetical function  $\text{subst}(n, m, k)$  which, if applied to the Gödel numbers of  $A$ ,  $t$ , and  $x$ , yields the Gödel number of  $A[t/x]$ . In other words,  $\text{subst}(\#A^\#, \#t^\#, \#x^\#) = \#A[t/x]^\#$ . Likewise, properties of derivations turn into primitive recursive relations on the respective Gödel numbers. In particular, the property  $\text{Deriv}(n)$  that holds of  $n$  if it is the Gödel number of a correct derivation in natural deduction, is primitive recursive. Showing that these are primitive recursive required a fair amount of work, and at times some ingenuity, and depended essentially on the fact that operating with sequences is primitive recursive. If a theory **T** is decidable, then we can use  $\text{Deriv}$  to define a decidable relation  $\text{Prf}_T(n, m)$  which holds if  $n$  is the Gödel number of a derivation of the sentence with Gödel number  $m$  from **T**. This relation is primitive recursive if the set of axioms of **T** is, and merely general recursive if the axioms of **T** are decidable but not primitive recursive.

## Problems

**Problem 16.1.** Show that the function  $\text{flatten}(z)$ , which turns the sequence  $\langle \#t_1\#, \dots, \#t_n\# \rangle$  into  $\#t_1, \dots, t_n\#$ , is primitive recursive.

**Problem 16.2.** Give a detailed proof of [Proposition 16.8](#) along the lines of the first proof of [Proposition 16.5](#).

**Problem 16.3.** Prove [Proposition 16.9](#). You may make use of the fact that any substring of a formula which is a formula is a sub-formula of it.

**Problem 16.4.** Prove [Proposition 16.12](#)

**Problem 16.5.** Define the following properties as in [Proposition 16.15](#):

1.  $\text{FollowsBy}_{\text{Cut}}(p)$ ,
2.  $\text{FollowsBy}_{\rightarrow\text{L}}(p)$ ,
3.  $\text{FollowsBy}_{=} (p)$ ,
4.  $\text{FollowsBy}_{\forall\text{R}}(p)$ .

For the last one, you will have to also show that you can test primitive recursively if the last inference of the derivation with Gödel number  $p$  satisfies the eigenvariable condition, i.e., the eigenvariable  $a$  of the  $\forall\text{R}$  does not occur in the end-sequent.

**Problem 16.6.** Define the following properties as in [Proposition 16.21](#):

1.  $\text{FollowsBy}_{\rightarrow\text{Elim}}(d)$ ,
2.  $\text{FollowsBy}_{=\text{Elim}}(d)$ ,
3.  $\text{FollowsBy}_{\forall\text{Elim}}(d)$ ,
4.  $\text{FollowsBy}_{\forall\text{Intro}}(d)$ .

For the last one, you will have to also show that you can test primitive recursively if the last inference of the derivation with Gödel number  $d$  satisfies the eigenvariable condition, i.e., the eigenvariable  $a$  of the  $\forall$ Intro inference occurs neither in the end-formula of  $d$  nor in an open assumption of  $d$ . You may use the primitive recursive predicate `OpenAssum` from [Proposition 16.23](#) for this.

## CHAPTER 17

# *Representability in $\mathbf{Q}$*

### 17.1 Introduction

The incompleteness theorems apply to theories in which basic facts about computable functions can be expressed and proved. We will describe a very minimal such theory called “ $\mathbf{Q}$ ” (or, sometimes, “Robinson’s  $Q$ ,” after Raphael Robinson). We will say what it means for a function to be *representable* in  $\mathbf{Q}$ , and then we will prove the following:

A function is representable in  $\mathbf{Q}$  if and only if it is computable.

For one thing, this provides us with another model of computability. But we will also use it to show that the set  $\{A : \mathbf{Q} \vdash A\}$  is not decidable, by reducing the halting problem to it. By the time we are done, we will have proved much stronger things than this.

The language of  $\mathbf{Q}$  is the language of arithmetic;  $\mathbf{Q}$  consists of the following axioms (to be used in conjunction with the other axioms and rules of first-order logic with identity predicate):

$$\forall x \forall y (x' = y' \rightarrow x = y) \quad (Q_1)$$

$$\forall x 0 \neq x' \quad (Q_2)$$

$$\forall x (x = 0 \vee \exists y x = y') \quad (\mathbf{Q}_3)$$

$$\forall x (x + 0) = x \quad (\mathbf{Q}_4)$$

$$\forall x \forall y (x + y') = (x + y)' \quad (\mathbf{Q}_5)$$

$$\forall x (x \times 0) = 0 \quad (\mathbf{Q}_6)$$

$$\forall x \forall y (x \times y') = ((x \times y) + x) \quad (\mathbf{Q}_7)$$

$$\forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y) \quad (\mathbf{Q}_8)$$

For each natural number  $n$ , define the numeral  $\bar{n}$  to be the term  $0''\dots'$  where there are  $n$  tick marks in all. So,  $\bar{0}$  is the constant symbol  $0$  by itself,  $\bar{1}$  is  $0'$ ,  $\bar{2}$  is  $0''$ , etc.

As a theory of arithmetic,  $\mathbf{Q}$  is *extremely* weak; for example, you can't even prove very simple facts like  $\forall x x \neq x'$  or  $\forall x \forall y (x + y) = (y + x)$ . But we will see that much of the reason that  $\mathbf{Q}$  is so interesting is *because* it is so weak. In fact, it is just barely strong enough for the incompleteness theorem to hold. Another reason  $\mathbf{Q}$  is interesting is because it has a *finite* set of axioms.

A stronger theory than  $\mathbf{Q}$  (called *Peano arithmetic*  $\mathbf{PA}$ ) is obtained by adding a schema of induction to  $\mathbf{Q}$ :

$$(A(0) \wedge \forall x (A(x) \rightarrow A(x'))) \rightarrow \forall x A(x)$$

where  $A(x)$  is any formula. If  $A(x)$  contains free variables other than  $x$ , we add universal quantifiers to the front to bind all of them (so that the corresponding instance of the induction schema is a sentence). For instance, if  $A(x, y)$  also contains the variable  $y$  free, the corresponding instance is

$$\forall y ((A(0) \wedge \forall x (A(x) \rightarrow A(x'))) \rightarrow \forall x A(x))$$

Using instances of the induction schema, one can prove much more from the axioms of  $\mathbf{PA}$  than from those of  $\mathbf{Q}$ . In fact, it takes a good deal of work to find “natural” statements about the natural numbers that can't be proved in Peano arithmetic!

**Definition 17.1.** A function  $f(x_0, \dots, x_k)$  from the natural numbers to the natural numbers is said to be *representable in  $\mathbf{Q}$*  if there is a formula  $A_f(x_0, \dots, x_k, y)$  such that whenever  $f(n_0, \dots, n_k) = m$ ,  $\mathbf{Q}$  proves

1.  $A_f(\overline{n_0}, \dots, \overline{n_k}, \overline{m})$
2.  $\forall y (A_f(\overline{n_0}, \dots, \overline{n_k}, y) \rightarrow \overline{m} = y)$ .

There are other ways of stating the definition; for example, we could equivalently require that  $\mathbf{Q}$  proves  $\forall y (A_f(\overline{n_0}, \dots, \overline{n_k}, y) \leftrightarrow y = \overline{m})$ .

**Theorem 17.2.** *A function is representable in  $\mathbf{Q}$  if and only if it is computable.*

There are two directions to proving the theorem. The left-to-right direction is fairly straightforward once arithmetization of syntax is in place. The other direction requires more work. Here is the basic idea: we pick “general recursive” as a way of making “computable” precise, and show that every general recursive function is representable in  $\mathbf{Q}$ . Recall that a function is general recursive if it can be defined from zero, the successor function  $\text{succ}$ , and the projection functions  $P_i^n$ , using composition, primitive recursion, and regular minimization. So one way of showing that every general recursive function is representable in  $\mathbf{Q}$  is to show that the basic functions are representable, and whenever some functions are representable, then so are the functions defined from them using composition, primitive recursion, and regular minimization. In other words, we might show that the basic functions are representable, and that the representable functions are “closed under” composition, primitive recursion, and regular minimization. This guarantees that every general recursive function is representable.

It turns out that the step where we would show that representable functions are closed under primitive recursion is hard.

In order to avoid this step, we show first that in fact we can do without primitive recursion. That is, we show that every general recursive function can be defined from basic functions using composition and regular minimization alone. To do this, we show that primitive recursion can actually be done by a specific regular minimization. However, for this to work, we have to add some additional basic functions: addition, multiplication, and the characteristic function of the identity relation  $\chi_{=}$ . Then, we can prove the theorem by showing that all of *these* basic functions are representable in  $\mathbf{Q}$ , and the representable functions are closed under composition and regular minimization.

## 17.2 Functions Representable in $\mathbf{Q}$ are Computable

We'll prove that every function that is representable in  $\mathbf{Q}$  is computable. We first have to establish a lemma about functions representable in  $\mathbf{Q}$ .

**Lemma 17.3.** *If  $f(x_0, \dots, x_k)$  is representable in  $\mathbf{Q}$ , there is a formula  $A(x_0, \dots, x_k, y)$  such that*

$$\mathbf{Q} \vdash A_f(\bar{n}_0, \dots, \bar{n}_k, \bar{m}) \quad \text{iff} \quad m = f(n_0, \dots, n_k).$$

*Proof.* The “if” part is **Definition 17.1(1)**. The “only if” part is seen as follows: Suppose  $\mathbf{Q} \vdash A_f(\bar{n}_0, \dots, \bar{n}_k, \bar{m})$  but  $m \neq f(n_0, \dots, n_k)$ . Let  $l = f(n_0, \dots, n_k)$ . By **Definition 17.1(1)**,  $\mathbf{Q} \vdash A_f(\bar{n}_0, \dots, \bar{n}_k, \bar{l})$ . By **Definition 17.1(2)**,  $\forall y (A_f(\bar{n}_0, \dots, \bar{n}_k, y) \rightarrow \bar{l} = y)$ . Using logic and the assumption that  $\mathbf{Q} \vdash A_f(\bar{n}_0, \dots, \bar{n}_k, \bar{m})$ , we get that  $\mathbf{Q} \vdash \bar{l} = \bar{m}$ . On the other hand, by **Lemma 17.14**,  $\mathbf{Q} \vdash \bar{l} \neq \bar{m}$ . So  $\mathbf{Q}$  is inconsistent. But that is impossible, since  $\mathbf{Q}$  is satisfied by the standard model (see **Definition 14.2**),  $N \models \mathbf{Q}$ , and satisfiable theories are always consistent by the Soundness Theorem (**Corollaries 10.31** and **11.29**).  $\square$



**Lemma 17.4.** *Every function that is representable in  $\mathbf{Q}$  is computable.*

*Proof.* Let's first give the intuitive idea for why this is true. To compute  $f$ , we do the following. List all the possible derivations  $\delta$  in the language of arithmetic. This is possible to do mechanically. For each one, check if it is a derivation of a formula of the form  $A_f(\overline{n_0}, \dots, \overline{n_k}, \overline{m})$  (the formula representing  $f$  in  $\mathbf{Q}$  from Lemma 17.3). If it is,  $m = f(n_0, \dots, n_k)$  by Lemma 17.3, and we've found the value of  $f$ . The search terminates because  $\mathbf{Q} \vdash A_f(\overline{n_0}, \dots, \overline{n_k}, \overline{f(n_0, \dots, n_k)})$ , so eventually we find a  $\delta$  of the right sort.

This is not quite precise because our procedure operates on derivations and formulas instead of just on numbers, and we haven't explained exactly why "listing all possible derivations" is mechanically possible. But as we've seen, it is possible to code terms, formulas, and derivations by Gödel numbers. We've also introduced a precise model of computation, the general recursive functions. And we've seen that the relation  $\text{Prf}_{\mathbf{Q}}(d, y)$ , which holds iff  $d$  is the Gödel number of a derivation of the formula with Gödel number  $y$  from the axioms of  $\mathbf{Q}$ , is (primitive) recursive. Other primitive recursive functions we'll need are  $\text{num}$  (Proposition 16.6) and  $\text{Subst}$  (Proposition 16.11). From these, it is possible to define  $f$  by minimization; thus,  $f$  is recursive.

First, define

$$A(n_0, \dots, n_k, m) = \text{Subst}(\text{Subst}(\dots \text{Subst}({}^{\#}A_f^{\#}, \text{num}(n_0), {}^{\#}x_0^{\#}), \dots), \text{num}(n_k), {}^{\#}x_k^{\#}), \text{num}(m), {}^{\#}y^{\#})$$

This looks complicated, but it's just the function  $A(n_0, \dots, n_k, m) = {}^{\#}A_f(\overline{n_0}, \dots, \overline{n_k}, \overline{m})^{\#}$ .

Now, consider the relation  $R(n_0, \dots, n_k, s)$  which holds if  $(s)_0$  is the Gödel number of a derivation from  $\mathbf{Q}$  of  $A_f(\overline{n_0}, \dots, \overline{n_k}, \overline{(s)_1})$ :

$$R(n_0, \dots, n_k, s) \quad \text{iff} \quad \text{Prf}_{\mathbf{Q}}((s)_0, A(n_0, \dots, n_k, (s)_1))$$

If we can find an  $s$  such that  $R(n_0, \dots, n_k, s)$  hold, we have found a pair of numbers— $(s)_0$  and  $(s)_1$ —such that  $(s)_0$  is the Gödel number of a derivation of  $A_f(\overline{n}_0, \dots, \overline{n}_k, (s)_1)$ . So looking for  $s$  is like looking for the pair  $d$  and  $m$  in the informal proof. And a computable function that “looks for” such an  $s$  can be defined by regular minimization. Note that  $R$  is regular: for every  $n_0, \dots, n_k$ , there is a derivation  $\delta$  of  $\mathbf{Q} \vdash A_f(\overline{n}_0, \dots, \overline{n}_k, \overline{f(n_0, \dots, n_k)})$ , so  $R(n_0, \dots, n_k, s)$  holds for  $s = \langle \# \delta^\#, f(n_0, \dots, n_k) \rangle$ . So, we can write  $f$  as

$$f(n_0, \dots, n_k) = (\mu s R(n_0, \dots, n_k, s))_1. \quad \square$$

### 17.3 The Beta Function Lemma

In order to show that we can carry out primitive recursion if addition, multiplication, and  $\chi_ =$  are available, we need to develop functions that handle sequences. (If we had exponentiation as well, our task would be easier.) When we had primitive recursion, we could define things like the “ $n$ -th prime,” and pick a fairly straightforward coding. But here we do not have primitive recursion—in fact we want to show that we can do primitive recursion using minimization—so we need to be more clever.

**Lemma 17.5.** *There is a function  $\beta(d, i)$  such that for every sequence  $a_0, \dots, a_n$  there is a number  $d$ , such that for every  $i \leq n$ ,  $\beta(d, i) = a_i$ . Moreover,  $\beta$  can be defined from the basic functions using just composition and regular minimization.*

Think of  $d$  as coding the sequence  $\langle a_0, \dots, a_n \rangle$ , and  $\beta(d, i)$  returning the  $i$ -th element. (Note that this “coding” does *not* use the power-of-primes coding we’re already familiar with!). The lemma is fairly minimal; it doesn’t say we can concatenate sequences or append elements, or even that we can *compute*  $d$  from  $a_0, \dots, a_n$  using functions definable by composition and regular minimization. All it says is that there is a “decoding” function such that every sequence is “coded.”

The use of the notation  $\beta$  is Gödel's. To repeat, the hard part of proving the lemma is defining a suitable  $\beta$  using the seemingly restricted resources, i.e., using just composition and minimization—however, we're allowed to use addition, multiplication, and  $\chi_=\$ . There are various ways to prove this lemma, but one of the cleanest is still Gödel's original method, which used a number-theoretic fact called Sunzi's Theorem (traditionally, the "Chinese Remainder Theorem").

**Definition 17.6.** Two natural numbers  $a$  and  $b$  are *relatively prime* iff their greatest common divisor is 1; in other words, they have no other divisors in common.

**Definition 17.7.** Natural numbers  $a$  and  $b$  are *congruent modulo  $c$* ,  $a \equiv b \pmod{c}$ , iff  $c \mid (a - b)$ , i.e.,  $a$  and  $b$  have the same remainder when divided by  $c$ .

Here is Sunzi's Theorem:

**Theorem 17.8.** *Suppose  $x_0, \dots, x_n$  are (pairwise) relatively prime. Let  $y_0, \dots, y_n$  be any numbers. Then there is a number  $z$  such that*

$$\begin{aligned} z &\equiv y_0 \pmod{x_0} \\ z &\equiv y_1 \pmod{x_1} \\ &\vdots \\ z &\equiv y_n \pmod{x_n}. \end{aligned}$$

Here is how we will use Sunzi's Theorem: if  $x_0, \dots, x_n$  are bigger than  $y_0, \dots, y_n$  respectively, then we can take  $z$  to code the sequence  $\langle y_0, \dots, y_n \rangle$ . To recover  $y_i$ , we need only divide  $z$  by  $x_i$  and take the remainder. To use this coding, we will need to find suitable values for  $x_0, \dots, x_n$ .

A couple of observations will help us in this regard. Given  $y_0, \dots, y_n$ , let

$$j = \max(n, y_0, \dots, y_n) + 1,$$

and let

$$\begin{aligned}x_0 &= 1 + j! \\x_1 &= 1 + 2 \cdot j! \\x_2 &= 1 + 3 \cdot j! \\&\vdots \\x_n &= 1 + (n + 1) \cdot j!\end{aligned}$$

Then two things are true:

1.  $x_0, \dots, x_n$  are relatively prime.
2. For each  $i$ ,  $y_i < x_i$ .

To see that (1) is true, note that if  $p$  is a prime number and  $p \mid x_i$  and  $p \mid x_k$ , then  $p \mid 1 + (i + 1)j!$  and  $p \mid 1 + (k + 1)j!$ . But then  $p$  divides their difference,

$$(1 + (i + 1)j!) - (1 + (k + 1)j!) = (i - k)j!.$$

Since  $p$  divides  $1 + (i + 1)j!$ , it can't divide  $j!$  as well (otherwise, the first division would leave a remainder of 1). So  $p$  divides  $i - k$ , since  $p$  divides  $(i - k)j!$ . But  $|i - k|$  is at most  $n$ , and we have chosen  $j > n$ , so this implies that  $p \mid j!$ , again a contradiction. So there is no prime number dividing both  $x_i$  and  $x_k$ . Clause (2) is easy: we have  $y_i < j < j! < x_i$ .

Now let us prove the  $\beta$  function lemma. Remember that we can use 0, successor, plus, times,  $\chi_{=}$ , projections, and any function defined from them using composition and minimization applied to regular functions. We can also use a relation if its characteristic function is so definable. As before we can show that these relations are closed under Boolean combinations and bounded quantification; for example:

$$\begin{aligned}\text{not}(x) &= \chi_{\neq}(x, 0) \\(\min x \leq z) R(x, y) &= \mu x (R(x, y) \vee x = z)\end{aligned}$$

$$(\exists x \leq z) R(x, y) \Leftrightarrow R((\min x \leq z) R(x, y), y)$$

We can then show that all of the following are also definable without primitive recursion:

1. The pairing function,  $J(x, y) = \frac{1}{2}[(x + y)(x + y + 1)] + x$ ;
2. the projection functions

$$K(z) = (\min x \leq z) (\exists y \leq z) z = J(x, y),$$

$$L(z) = (\min y \leq z) (\exists x \leq z) z = J(x, y);$$

3. the less-than relation  $x < y$ ;
4. the divisibility relation  $x \mid y$ ;
5. the function  $\text{rem}(x, y)$  which returns the remainder when  $y$  is divided by  $x$ .

Now define

$$\beta^*(d_0, d_1, i) = \text{rem}(1 + (i + 1)d_1, d_0) \text{ and}$$

$$\beta(d, i) = \beta^*(K(d), L(d), i).$$

This is the function we want. Given  $a_0, \dots, a_n$  as above, let

$$j = \max(n, a_0, \dots, a_n) + 1,$$

and let  $d_1 = j!$ . By (1) above, we know that  $1 + d_1, 1 + 2d_1, \dots, 1 + (n + 1)d_1$  are relatively prime, and by (2) that all are greater than  $a_0, \dots, a_n$ . By Sunzi's Theorem there is a value  $d_0$  such that for each  $i$ ,

$$d_0 \equiv a_i \pmod{(1 + (i + 1)d_1)}$$

and so (because  $d_1$  is greater than  $a_i$ ),

$$a_i = \text{rem}(1 + (i + 1)d_1, d_0).$$

Let  $d = J(d_0, d_1)$ . Then for each  $i \leq n$ , we have

$$\beta(d, i) = \beta^*(d_0, d_1, i)$$

$$\begin{aligned}
 &= \text{rem}(1 + (i + 1)d_1, d_0) \\
 &= a_i
 \end{aligned}$$

which is what we need. This completes the proof of the  $\beta$ -function lemma.

## 17.4 Simulating Primitive Recursion

Now we can show that definition by primitive recursion can be “simulated” by regular minimization using the beta function. Suppose we have  $f(\vec{x})$  and  $g(\vec{x}, y, z)$ . Then the function  $h(x, \vec{z})$  defined from  $f$  and  $g$  by primitive recursion is

$$\begin{aligned}
 h(\vec{x}, 0) &= f(\vec{x}) \\
 h(\vec{x}, y + 1) &= g(\vec{x}, y, h(\vec{x}, y)).
 \end{aligned}$$

We need to show that  $h$  can be defined from  $f$  and  $g$  using just composition and regular minimization, using the basic functions and functions defined from them using composition and regular minimization (such as  $\beta$ ).

**Lemma 17.9.** *If  $h$  can be defined from  $f$  and  $g$  using primitive recursion, it can be defined from  $f$ ,  $g$ , the functions zero, succ,  $P_i^n$ , add, mult,  $\chi_=\$ , using composition and regular minimization.*

*Proof.* First, define an auxiliary function  $\hat{h}(\vec{x}, y)$  which returns the least number  $d$  such that  $d$  codes a sequence which satisfies

1.  $(d)_0 = f(\vec{x})$ , and
2. for each  $i < y$ ,  $(d)_{i+1} = g(\vec{x}, i, (d)_i)$ ,

where now  $(d)_i$  is short for  $\beta(d, i)$ . In other words,  $\hat{h}$  returns the sequence  $\langle h(\vec{x}, 0), h(\vec{x}, 1), \dots, h(\vec{x}, y) \rangle$ . We can write  $\hat{h}$  as

$$\hat{h}(\vec{x}, y) = \mu d (\beta(d, 0) = f(\vec{x}) \wedge (\forall i < y) \beta(d, i+1) = g(\vec{x}, i, \beta(d, i))).$$

Note: no primitive recursion is needed here, just minimization. The function we minimize is regular because of the beta function lemma [Lemma 17.5](#).

But now we have

$$h(\vec{x}, y) = \beta(\hat{h}(\vec{x}, y), y),$$

so  $h$  can be defined from the basic functions using just composition and regular minimization.  $\square$

## 17.5 Basic Functions are Representable in $\mathbf{Q}$

First we have to show that all the basic functions are representable in  $\mathbf{Q}$ . In the end, we need to show how to assign to each  $k$ -ary basic function  $f(x_0, \dots, x_{k-1})$  a formula  $A_f(x_0, \dots, x_{k-1}, y)$  that represents it.

We will be able to represent zero, successor, plus, times, the characteristic function for equality, and projections. In each case, the appropriate representing function is entirely straightforward; for example, zero is represented by the formula  $y = 0$ , successor is represented by the formula  $x'_0 = y$ , and addition is represented by the formula  $(x_0 + x_1) = y$ . The work involves showing that  $\mathbf{Q}$  can prove the relevant sentences; for example, saying that addition is represented by the formula above involves showing that for every pair of natural numbers  $m$  and  $n$ ,  $\mathbf{Q}$  proves

$$\begin{aligned} \bar{n} + \bar{m} &= \overline{n + m} \text{ and} \\ \forall y ((\bar{n} + \bar{m}) = y &\rightarrow y = \overline{n + m}). \end{aligned}$$

**Proposition 17.10.** *The zero function  $\text{zero}(x) = 0$  is represented in  $\mathbf{Q}$  by  $A_{\text{zero}}(x, y) \equiv y = 0$ .*

**Proposition 17.11.** *The successor function  $\text{succ}(x) = x + 1$  is represented in  $\mathbf{Q}$  by  $A_{\text{succ}}(x, y) \equiv y = x'$ .*

**Proposition 17.12.** *The projection function  $P_i^n(x_0, \dots, x_{n-1}) = x_i$  is represented in  $\mathbf{Q}$  by*

$$A_{P_i^n}(x_0, \dots, x_{n-1}, y) \equiv y = x_i.$$

**Proposition 17.13.** *The characteristic function of  $=$ ,*

$$\chi_{=(x_0, x_1)} = \begin{cases} 1 & \text{if } x_0 = x_1 \\ 0 & \text{otherwise} \end{cases}$$

*is represented in  $\mathbf{Q}$  by*

$$A_{\chi_{=}}(x_0, x_1, y) \equiv (x_0 = x_1 \wedge y = \bar{1}) \vee (x_0 \neq x_1 \wedge y = \bar{0}).$$

The proof requires the following lemma.

**Lemma 17.14.** *Given natural numbers  $n$  and  $m$ , if  $n \neq m$ , then  $\mathbf{Q} \vdash \bar{n} \neq \bar{m}$ .*

*Proof.* Use induction on  $n$  to show that for every  $m$ , if  $n \neq m$ , then  $\mathbf{Q} \vdash \bar{n} \neq \bar{m}$ .

In the base case,  $n = 0$ . If  $m$  is not equal to 0, then  $m = k + 1$  for some natural number  $k$ . We have an axiom that says  $\forall x \, 0 \neq x'$ . By a quantifier axiom, replacing  $x$  by  $\bar{k}$ , we can conclude  $0 \neq \bar{k}'$ . But  $\bar{k}'$  is just  $\bar{m}$ .

In the induction step, we can assume the claim is true for  $n$ , and consider  $n + 1$ . Let  $m$  be any natural number. There are two possibilities: either  $m = 0$  or for some  $k$  we have  $m = k + 1$ . The first case is handled as above. In the second case, suppose  $n + 1 \neq k + 1$ . Then  $n \neq k$ . By the induction hypothesis for  $n$  we have  $\mathbf{Q} \vdash \bar{n} \neq \bar{k}$ . We have an axiom that says  $\forall x \forall y \, x' = y' \rightarrow x =$



$y$ . Using a quantifier axiom, we have  $\bar{n}' = \bar{k}' \rightarrow \bar{n} = \bar{k}$ . Using propositional logic, we can conclude, in  $\mathbf{Q}$ ,  $\bar{n} \neq \bar{k} \rightarrow \bar{n}' \neq \bar{k}'$ . Using modus ponens, we can conclude  $\bar{n}' \neq \bar{k}'$ , which is what we want, since  $\bar{k}'$  is  $\bar{m}$ .  $\square$

Note that the lemma does not say much: in essence it says that  $\mathbf{Q}$  can prove that different numerals denote different objects. For example,  $\mathbf{Q}$  proves  $0'' \neq 0'''$ . But showing that this holds in general requires some care. Note also that although we are using induction, it is induction *outside* of  $\mathbf{Q}$ .

*Proof of Proposition 17.13.* If  $n = m$ , then  $\bar{n}$  and  $\bar{m}$  are the same term, and  $\chi_{=}(n, m) = 1$ . But  $\mathbf{Q} \vdash (\bar{n} = \bar{m} \wedge \bar{1} = \bar{1})$ , so it proves  $A_{=}(n, m, \bar{1})$ . If  $n \neq m$ , then  $\chi_{=}(n, m) = 0$ . By Lemma 17.14,  $\mathbf{Q} \vdash \bar{n} \neq \bar{m}$  and so also  $(\bar{n} \neq \bar{m} \wedge 0 = 0)$ . Thus  $\mathbf{Q} \vdash A_{=}(n, m, \bar{0})$ .

For the second part, we also have two cases. If  $n = m$ , we have to show that  $\mathbf{Q} \vdash \forall y (A_{=}(n, m, y) \rightarrow y = \bar{1})$ . Arguing informally, suppose  $A_{=}(n, m, y)$ , i.e.,

$$(\bar{n} = \bar{n} \wedge y = \bar{1}) \vee (\bar{n} \neq \bar{n} \wedge y = \bar{0})$$

The left disjunct implies  $y = \bar{1}$  by logic; the right contradicts  $\bar{n} = \bar{n}$  which is provable by logic.

Suppose, on the other hand, that  $n \neq m$ . Then  $A_{=}(n, m, y)$  is

$$(\bar{n} = \bar{m} \wedge y = \bar{1}) \vee (\bar{n} \neq \bar{m} \wedge y = \bar{0})$$

Here, the left disjunct contradicts  $\bar{n} \neq \bar{m}$ , which is provable in  $\mathbf{Q}$  by Lemma 17.14; the right disjunct entails  $y = \bar{0}$ .  $\square$

**Proposition 17.15.** *The addition function  $\text{add}(x_0, x_1) = x_0 + x_1$  is represented in  $\mathbf{Q}$  by*

$$A_{\text{add}}(x_0, x_1, y) \equiv y = (x_0 + x_1).$$

**Lemma 17.16.**  $\mathbf{Q} \vdash (\bar{n} + \bar{m}) = \overline{n + m}$

*Proof.* We prove this by induction on  $m$ . If  $m = 0$ , the claim is that  $\mathbf{Q} \vdash (\bar{n} + 0) = \bar{n}$ . This follows by axiom  $Q_4$ . Now suppose the claim for  $m$ ; let's prove the claim for  $m + 1$ , i.e., prove that  $\mathbf{Q} \vdash (\bar{n} + \overline{m + 1}) = \overline{n + m + 1}$ . Note that  $\overline{m + 1}$  is just  $\bar{m}'$ , and  $\overline{n + m + 1}$  is just  $\overline{n + m}'$ . By axiom  $Q_5$ ,  $\mathbf{Q} \vdash (\bar{n} + \bar{m}') = (\bar{n} + \bar{m})'$ . By induction hypothesis,  $\mathbf{Q} \vdash (\bar{n} + \bar{m}) = \overline{n + m}$ . So  $\mathbf{Q} \vdash (\bar{n} + \bar{m}') = \overline{n + m}'$ .  $\square$

*Proof of Proposition 17.15.* The formula  $A_{\text{add}}(x_0, x_1, y)$  representing add is  $y = (x_0 + x_1)$ . First we show that if  $\text{add}(n, m) = k$ , then  $\mathbf{Q} \vdash A_{\text{add}}(\bar{n}, \bar{m}, \bar{k})$ , i.e.,  $\mathbf{Q} \vdash \bar{k} = (\bar{n} + \bar{m})$ . But since  $k = n + m$ ,  $\bar{k}$  just is  $\overline{n + m}$ , and we've shown in Lemma 17.16 that  $\mathbf{Q} \vdash (\bar{n} + \bar{m}) = \overline{n + m}$ .

We also have to show that if  $\text{add}(n, m) = k$ , then

$$\mathbf{Q} \vdash \forall y (A_{\text{add}}(\bar{n}, \bar{m}, y) \rightarrow y = \bar{k}).$$

Suppose we have  $(\bar{n} + \bar{m}) = y$ . Since

$$\mathbf{Q} \vdash (\bar{n} + \bar{m}) = \overline{n + m},$$

we can replace the left side with  $\overline{n + m}$  and get  $\overline{n + m} = y$ , for arbitrary  $y$ .  $\square$

**Proposition 17.17.** *The multiplication function  $\text{mult}(x_0, x_1) = x_0 \cdot x_1$  is represented in  $\mathbf{Q}$  by*

$$A_{\text{mult}}(x_0, x_1, y) \equiv y = (x_0 \times x_1).$$

*Proof.* Exercise.  $\square$

**Lemma 17.18.**  $\mathbf{Q} \vdash (\overline{n} \times \overline{m}) = \overline{n \cdot m}$

*Proof.* Exercise. □

Recall that we use  $\times$  for the function symbol of the language of arithmetic, and  $\cdot$  for the ordinary multiplication operation on numbers. So  $\cdot$  can appear between expressions for numbers (such as in  $m \cdot n$ ) while  $\times$  appears only between terms of the language of arithmetic (such as in  $(\overline{m} \times \overline{n})$ ). Even more confusingly,  $+$  is used for both the function symbol and the addition operation. When it appears between terms—e.g., in  $(\overline{n} + \overline{m})$ —it is the 2-place function symbol of the language of arithmetic, and when it appears between numbers—e.g., in  $n + m$ —it is the addition operation. This includes the case  $\overline{n + m}$ : this is the standard numeral corresponding to the number  $n + m$ .

## 17.6 Composition is Representable in $\mathbf{Q}$

Suppose  $h$  is defined by

$$h(x_0, \dots, x_{l-1}) = f(g_0(x_0, \dots, x_{l-1}), \dots, g_{k-1}(x_0, \dots, x_{l-1})).$$

where we have already found formulas  $A_f, A_{g_0}, \dots, A_{g_{k-1}}$  representing the functions  $f$ , and  $g_0, \dots, g_{k-1}$ , respectively. We have to find a formula  $A_h$  representing  $h$ .

Let's start with a simple case, where all functions are 1-place, i.e., consider  $h(x) = f(g(x))$ . If  $A_f(y, z)$  represents  $f$ , and  $A_g(x, y)$  represents  $g$ , we need a formula  $A_h(x, z)$  that represents  $h$ . Note that  $h(x) = z$  iff there is a  $y$  such that both  $z = f(y)$  and  $y = g(x)$ . (If  $h(x) = z$ , then  $g(x)$  is such a  $y$ ; if such a  $y$  exists, then since  $y = g(x)$  and  $z = f(y)$ ,  $z = f(g(x))$ .) This suggests that  $\exists y (A_g(x, y) \wedge A_f(y, z))$  is a good candidate for  $A_h(x, z)$ . We just have to verify that  $\mathbf{Q}$  proves the relevant formulas.

**Proposition 17.19.** *If  $h(n) = m$ , then  $\mathbf{Q} \vdash A_h(\bar{n}, \bar{m})$ .*

*Proof.* Suppose  $h(n) = m$ , i.e.,  $f(g(n)) = m$ . Let  $k = g(n)$ . Then

$$\mathbf{Q} \vdash A_g(\bar{n}, \bar{k})$$

since  $A_g$  represents  $g$ , and

$$\mathbf{Q} \vdash A_f(\bar{k}, \bar{m})$$

since  $A_f$  represents  $f$ . Thus,

$$\mathbf{Q} \vdash A_g(\bar{n}, \bar{k}) \wedge A_f(\bar{k}, \bar{m})$$

and consequently also

$$\mathbf{Q} \vdash \exists y (A_g(\bar{n}, y) \wedge A_f(y, \bar{m})),$$

i.e.,  $\mathbf{Q} \vdash A_h(\bar{n}, \bar{m})$ . □

**Proposition 17.20.** *If  $h(n) = m$ , then  $\mathbf{Q} \vdash \forall z (A_h(\bar{n}, z) \rightarrow z = \bar{m})$ .*

*Proof.* Suppose  $h(n) = m$ , i.e.,  $f(g(n)) = m$ . Let  $k = g(n)$ . Then

$$\mathbf{Q} \vdash \forall y (A_g(\bar{n}, y) \rightarrow y = \bar{k})$$

since  $A_g$  represents  $g$ , and

$$\mathbf{Q} \vdash \forall z (A_f(\bar{k}, z) \rightarrow z = \bar{m})$$

since  $A_f$  represents  $f$ . Using just a little bit of logic, we can show that also

$$\mathbf{Q} \vdash \forall z (\exists y (A_g(\bar{n}, y) \wedge A_f(y, z)) \rightarrow z = \bar{m}).$$

i.e.,  $\mathbf{Q} \vdash \forall y (A_h(\bar{n}, y) \rightarrow y = \bar{m})$ . □

The same idea works in the more complex case where  $f$  and  $g_i$  have arity greater than 1.

**Proposition 17.21.** *If  $A_f(y_0, \dots, y_{k-1}, z)$  represents  $f(y_0, \dots, y_{k-1})$  in  $\mathbf{Q}$ , and  $A_{g_i}(x_0, \dots, x_{l-1}, y)$  represents  $g_i(x_0, \dots, x_{l-1})$  in  $\mathbf{Q}$ , then*

$$\exists y_0 \dots \exists y_{k-1} (A_{g_0}(x_0, \dots, x_{l-1}, y_0) \wedge \dots \wedge A_{g_{k-1}}(x_0, \dots, x_{l-1}, y_{k-1}) \wedge A_f(y_0, \dots, y_{k-1}, z))$$

*represents*

$$h(x_0, \dots, x_{l-1}) = f(g_0(x_0, \dots, x_{l-1}), \dots, g_{k-1}(x_0, \dots, x_{l-1})).$$

*Proof.* Exercise. □

## 17.7 Regular Minimization is Representable in $\mathbf{Q}$

Let's consider unbounded search. Suppose  $g(x, z)$  is regular and representable in  $\mathbf{Q}$ , say by the formula  $A_g(x, z, y)$ . Let  $f$  be defined by  $f(z) = \mu x [g(x, z) = 0]$ . We would like to find a formula  $A_f(z, y)$  representing  $f$ . The value of  $f(z)$  is that number  $x$  which (a) satisfies  $g(x, z) = 0$  and (b) is the least such, i.e., for any  $w < x$ ,  $g(w, z) \neq 0$ . So the following is a natural choice:

$$A_f(z, y) \equiv A_g(y, z, 0) \wedge \forall w (w < y \rightarrow \neg A_g(w, z, 0)).$$

In the general case, of course, we would have to replace  $z$  with  $z_0, \dots, z_k$ .

The proof, again, will involve some lemmas about things  $\mathbf{Q}$  is strong enough to prove.

**Lemma 17.22.** *For every constant symbol  $a$  and every natural number  $n$ ,*

$$\mathbf{Q} \vdash (a' + \bar{n}) = (a + \bar{n})'.$$

*Proof.* The proof is, as usual, by induction on  $n$ . In the base case,  $n = 0$ , we need to show that  $\mathbf{Q}$  proves  $(a' + 0) = (a + 0)'$ . But we

have:

$$\mathbf{Q} \vdash (a' + 0) = a' \quad \text{by axiom } Q_4 \quad (17.1)$$

$$\mathbf{Q} \vdash (a + 0) = a \quad \text{by axiom } Q_4 \quad (17.2)$$

$$\mathbf{Q} \vdash (a + 0)' = a' \quad \text{by eq. (17.2)} \quad (17.3)$$

$$\mathbf{Q} \vdash (a' + 0) = (a + 0)' \quad \text{by eq. (17.1) and eq. (17.3)}$$

In the induction step, we can assume that we have shown that  $\mathbf{Q} \vdash (a' + \bar{n}) = (a + \bar{n})'$ . Since  $\overline{n+1}$  is  $\bar{n}'$ , we need to show that  $\mathbf{Q}$  proves  $(a' + \bar{n}') = (a + \bar{n}')'$ . We have:

$$\mathbf{Q} \vdash (a' + \bar{n}') = (a' + \bar{n})' \quad \text{by axiom } Q_5 \quad (17.4)$$

$$\mathbf{Q} \vdash (a' + \bar{n}') = (a + \bar{n}')' \quad \text{inductive hypothesis} \quad (17.5)$$

$$\mathbf{Q} \vdash (a' + \bar{n})' = (a + \bar{n}')' \quad \text{by eq. (17.4) and eq. (17.5)}. \quad \square$$

It is again worth mentioning that this is weaker than saying that  $\mathbf{Q}$  proves  $\forall x \forall y (x' + y) = (x + y)'$ . Although this sentence is true in  $N$ ,  $\mathbf{Q}$  does not prove it.

**Lemma 17.23.**  $\mathbf{Q} \vdash \forall x \neg x < 0$ .

*Proof.* We give the proof informally (i.e., only giving hints as to how to construct the formal derivation).

We have to prove  $\neg a < 0$  for an arbitrary  $a$ . By the definition of  $<$ , we need to prove  $\neg \exists y (y' + a) = 0$  in  $\mathbf{Q}$ . We'll assume  $\exists y (y' + a) = 0$  and prove a contradiction. Suppose  $(b' + a) = 0$ . Using  $Q_3$ , we have that  $a = 0 \vee \exists y a = y'$ . We distinguish cases.

Case 1:  $a = 0$  holds. From  $(b' + a) = 0$ , we have  $(b' + 0) = 0$ . By axiom  $Q_4$  of  $\mathbf{Q}$ , we have  $(b' + 0) = b'$ , and hence  $b' = 0$ . But by axiom  $Q_2$  we also have  $b' \neq 0$ , a contradiction.

Case 2: For some  $c$ ,  $a = c'$ . But then we have  $(b' + c') = 0$ . By axiom  $Q_5$ , we have  $(b' + c)' = 0$ , again contradicting axiom  $Q_2$ .  $\square$

**Lemma 17.24.** *For every natural number  $n$ ,*

$$\mathbf{Q} \vdash \forall x (x < \overline{n+1} \rightarrow (x = 0 \vee \dots \vee x = \overline{n})).$$

*Proof.* We use induction on  $n$ . Let us consider the base case, when  $n = 0$ . In that case, we need to show  $a < \overline{1} \rightarrow a = 0$ , for arbitrary  $a$ . Suppose  $a < \overline{1}$ . Then by the defining axiom for  $<$ , we have  $\exists y (y' + a) = 0'$  (since  $\overline{1} \equiv 0'$ ).

Suppose  $b$  has that property, i.e., we have  $(b' + a) = 0'$ . We need to show  $a = 0$ . By axiom  $Q_3$ , we have either  $a = 0$  or that there is a  $c$  such that  $a = c'$ . In the former case, there is nothing to show. So suppose  $a = c'$ . Then we have  $(b' + c') = 0'$ . By axiom  $Q_5$  of  $\mathbf{Q}$ , we have  $(b' + c)' = 0'$ . By axiom  $Q_1$ , we have  $(b' + c) = 0$ . But this means, by axiom  $Q_8$ , that  $c < 0$ , contradicting **Lemma 17.23**.

Now for the inductive step. We prove the case for  $n + 1$ , assuming the case for  $n$ . So suppose  $a < \overline{n+2}$ . Again using  $Q_3$  we can distinguish two cases:  $a = 0$  and for some  $b$ ,  $a = c'$ . In the first case,  $a = 0 \vee \dots \vee a = \overline{n+1}$  follows trivially. In the second case, we have  $c' < \overline{n+2}$ , i.e.,  $c' < \overline{n+1}'$ . By axiom  $Q_8$ , for some  $d$ ,  $(d' + c') = \overline{n+1}'$ . By axiom  $Q_5$ ,  $(d' + c)' = \overline{n+1}'$ . By axiom  $Q_1$ ,  $(d' + c) = \overline{n+1}$ , and so  $c < \overline{n+1}$  by axiom  $Q_8$ . By inductive hypothesis,  $c = 0 \vee \dots \vee c = \overline{n}$ . From this, we get  $c' = 0' \vee \dots \vee c' = \overline{n}'$  by logic, and so  $a = \overline{1} \vee \dots \vee a = \overline{n+1}$  since  $a = c'$ .  $\square$

**Lemma 17.25.** *For every natural number  $m$ ,*

$$\mathbf{Q} \vdash \forall y ((y < \overline{m} \vee \overline{m} < y) \vee y = \overline{m}).$$

*Proof.* By induction on  $m$ . First, consider the case  $m = 0$ .  $\mathbf{Q} \vdash \forall y (y = 0 \vee \exists z y = z')$  by  $Q_3$ . Let  $a$  be arbitrary. Then either  $a = 0$  or for some  $b$ ,  $a = b'$ . In the former case, we also have  $(a < 0 \vee 0 < a) \vee a = 0$ . But if  $a = b'$ , then  $(b' + 0) = (a + 0)$  by the logic of  $=$ . By  $Q_4$ ,  $(a + 0) = a$ , so we have  $(b' + 0) = a$ , and hence  $\exists z (z' + 0) = a$ . By the definition of  $<$  in  $Q_8$ ,  $0 < a$ . If  $0 < a$ , then also  $(0 < a \vee a < 0) \vee a = 0$ .

Now suppose we have

$$\mathbf{Q} \vdash \forall y ((y < \overline{m} \vee \overline{m} < y) \vee y = \overline{m})$$

and we want to show

$$\mathbf{Q} \vdash \forall y ((y < \overline{m+1} \vee \overline{m+1} < y) \vee y = \overline{m+1})$$

Let  $a$  be arbitrary. By  $Q_3$ , either  $a = 0$  or for some  $b$ ,  $a = \overline{b'}$ . In the first case, we have  $\overline{m'} + a = \overline{m+1}$  by  $Q_4$ , and so  $a < \overline{m+1}$  by  $Q_8$ .

Now consider the second case,  $a = b'$ . By the induction hypothesis,  $(b < \overline{m} \vee \overline{m} < b) \vee b = \overline{m}$ .

The first disjunct  $b < \overline{m}$  is equivalent (by  $Q_8$ ) to  $\exists z (z' + b) = \overline{m}$ . Suppose  $c$  has this property. If  $(c' + b) = \overline{m}$ , then also  $(c' + b)' = \overline{m'}$ . By  $Q_5$ ,  $(c' + b)' = (c' + b')$ . Hence,  $(c' + b') = \overline{m'}$ . We get  $\exists u (u' + b') = \overline{m+1}$  by existentially generalizing on  $c'$  and keeping in mind that  $\overline{m'} \equiv \overline{m+1}$ . Hence, if  $b < \overline{m}$  then  $b' < \overline{m+1}$  and so  $a < \overline{m+1}$ .

Now suppose  $\overline{m} < b$ , i.e.,  $\exists z (z' + \overline{m}) = b$ . Suppose  $c$  is such a  $z$ , i.e.,  $(c' + \overline{m}) = b$ . By logic,  $(c' + \overline{m})' = b'$ . By  $Q_5$ ,  $(c' + \overline{m}') = b'$ . Since  $a = b'$  and  $\overline{m'} \equiv \overline{m+1}$ ,  $(c' + \overline{m+1}) = a$ . By  $Q_8$ ,  $\overline{m+1} < a$ .

Finally, assume  $b = \overline{m}$ . Then, by logic,  $b' = \overline{m'}$ , and so  $a = \overline{m+1}$ .

Hence, from each disjunct of the case for  $m$  and  $b$ , we can obtain the corresponding disjunct for for  $m+1$  and  $a$ .  $\square$

**Proposition 17.26.** *If  $A_g(x, z, y)$  represents  $g(x, z)$  in  $\mathbf{Q}$ , then*

$$A_f(z, y) \equiv A_g(y, z, 0) \wedge \forall w (w < y \rightarrow \neg A_g(w, z, 0))$$

*represents  $f(z) = \mu x [g(x, z) = 0]$ .*

*Proof.* First we show that if  $f(n) = m$ , then  $\mathbf{Q} \vdash A_f(\overline{n}, \overline{m})$ , i.e.,

$$\mathbf{Q} \vdash A_g(\overline{m}, \overline{n}, 0) \wedge \forall w (w < \overline{m} \rightarrow \neg A_g(w, \overline{n}, 0)).$$



Since  $A_g(x, z, y)$  represents  $g(x, z)$  and  $g(m, n) = 0$  if  $f(n) = m$ , we have

$$\mathbf{Q} \vdash A_g(\bar{m}, \bar{n}, 0).$$

If  $f(n) = m$ , then for every  $k < m$ ,  $g(k, n) \neq 0$ . So

$$\mathbf{Q} \vdash \neg A_g(\bar{k}, \bar{n}, 0).$$

We get that

$$\mathbf{Q} \vdash \forall w (w < \bar{m} \rightarrow \neg A_g(w, \bar{n}, 0)). \quad (17.6)$$

by [Lemma 17.23](#) in case  $m = 0$  and by [Lemma 17.24](#) otherwise.

Now let's show that if  $f(n) = m$ , then  $\mathbf{Q} \vdash \forall y (A_f(\bar{n}, y) \rightarrow y = \bar{m})$ . We again sketch the argument informally, leaving the formalization to the reader.

Suppose  $A_f(\bar{n}, b)$ . From this we get (a)  $A_g(b, \bar{n}, 0)$  and (b)  $\forall w (w < b \rightarrow \neg A_g(w, \bar{n}, 0))$ . By [Lemma 17.25](#),  $(b < \bar{m} \vee \bar{m} < b) \vee b = \bar{m}$ . We'll show that both  $b < \bar{m}$  and  $\bar{m} < b$  leads to a contradiction.

If  $\bar{m} < b$ , then  $\neg A_g(\bar{m}, \bar{n}, 0)$  from (b). But  $m = f(n)$ , so  $g(m, n) = 0$ , and so  $\mathbf{Q} \vdash A_g(\bar{m}, \bar{n}, 0)$  since  $A_g$  represents  $g$ . So we have a contradiction.

Now suppose  $b < \bar{m}$ . Then since  $\mathbf{Q} \vdash \forall w (w < \bar{m} \rightarrow \neg A_g(w, \bar{n}, 0))$  by [eq. \(17.6\)](#), we get  $\neg A_g(b, \bar{n}, 0)$ . This again contradicts (a).  $\square$

## 17.8 Computable Functions are Representable in $\mathbf{Q}$

**Theorem 17.27.** *Every computable function is representable in  $\mathbf{Q}$ .*

*Proof.* For definiteness, and using the Church–Turing Thesis, let's say that a function is computable iff it is general recursive. The

general recursive functions are those which can be defined from the zero function  $\text{zero}$ , the successor function  $\text{succ}$ , and the projection function  $P_i^n$  using composition, primitive recursion, and regular minimization. By [Lemma 17.9](#), any function  $h$  that can be defined from  $f$  and  $g$  can also be defined using composition and regular minimization from  $f$ ,  $g$ , and  $\text{zero}$ ,  $\text{succ}$ ,  $P_i^n$ ,  $\text{add}$ ,  $\text{mult}$ ,  $\chi_=\text{}$ . Consequently, a function is general recursive iff it can be defined from  $\text{zero}$ ,  $\text{succ}$ ,  $P_i^n$ ,  $\text{add}$ ,  $\text{mult}$ ,  $\chi_=\text{}$  using composition and regular minimization.

We've furthermore shown that the basic functions in question are representable in  $\mathbf{Q}$  ([Propositions 17.10](#) to [17.13](#), [17.15](#) and [17.17](#)), and that any function defined from representable functions by composition or regular minimization ([Proposition 17.21](#), [Proposition 17.26](#)) is also representable. Thus every general recursive function is representable in  $\mathbf{Q}$ .  $\square$

We have shown that the set of computable functions can be characterized as the set of functions representable in  $\mathbf{Q}$ . In fact, the proof is more general. From the definition of representability, it is not hard to see that any theory extending  $\mathbf{Q}$  (or in which one can interpret  $\mathbf{Q}$ ) can represent the computable functions. But, conversely, in any derivation system in which the notion of derivation is computable, every representable function is computable. So, for example, the set of computable functions can be characterized as the set of functions representable in Peano arithmetic, or even Zermelo–Fraenkel set theory. As Gödel noted, this is somewhat surprising. We will see that when it comes to provability, questions are very sensitive to which theory you consider; roughly, the stronger the axioms, the more you can prove. But across a wide range of axiomatic theories, the representable functions are exactly the computable ones; stronger theories do not represent more functions as long as they are axiomatizable.

## 17.9 Representing Relations

Let us say what it means for a *relation* to be representable.

**Definition 17.28.** A relation  $R(x_0, \dots, x_k)$  on the natural numbers is *representable in  $\mathbf{Q}$*  if there is a formula  $A_R(x_0, \dots, x_k)$  such that whenever  $R(n_0, \dots, n_k)$  is true,  $\mathbf{Q}$  proves  $A_R(\bar{n}_0, \dots, \bar{n}_k)$ , and whenever  $R(n_0, \dots, n_k)$  is false,  $\mathbf{Q}$  proves  $\neg A_R(\bar{n}_0, \dots, \bar{n}_k)$ .

**Theorem 17.29.** *A relation is representable in  $\mathbf{Q}$  if and only if it is computable.*

*Proof.* For the forwards direction, suppose  $R(x_0, \dots, x_k)$  is represented by the formula  $A_R(x_0, \dots, x_k)$ . Here is an algorithm for computing  $R$ : on input  $n_0, \dots, n_k$ , simultaneously search for a proof of  $A_R(\bar{n}_0, \dots, \bar{n}_k)$  and a proof of  $\neg A_R(\bar{n}_0, \dots, \bar{n}_k)$ . By our hypothesis, the search is bound to find one or the other; if it is the first, report “yes,” and otherwise, report “no.”

In the other direction, suppose  $R(x_0, \dots, x_k)$  is computable. By definition, this means that the function  $\chi_R(x_0, \dots, x_k)$  is computable. By **Theorem 17.2**,  $\chi_R$  is represented by a formula, say  $A_{\chi_R}(x_0, \dots, x_k, y)$ . Let  $A_R(x_0, \dots, x_k)$  be the formula  $A_{\chi_R}(x_0, \dots, x_k, \bar{1})$ . Then for any  $n_0, \dots, n_k$ , if  $R(n_0, \dots, n_k)$  is true, then  $\chi_R(n_0, \dots, n_k) = 1$ , in which case  $\mathbf{Q}$  proves  $A_{\chi_R}(\bar{n}_0, \dots, \bar{n}_k, \bar{1})$ , and so  $\mathbf{Q}$  proves  $A_R(\bar{n}_0, \dots, \bar{n}_k)$ . On the other hand, if  $R(n_0, \dots, n_k)$  is false, then  $\chi_R(n_0, \dots, n_k) = 0$ . This means that  $\mathbf{Q}$  proves

$$\forall y (A_{\chi_R}(\bar{n}_0, \dots, \bar{n}_k, y) \rightarrow y = \bar{0}).$$

Since  $\mathbf{Q}$  proves  $\bar{0} \neq \bar{1}$ ,  $\mathbf{Q}$  proves  $\neg A_{\chi_R}(\bar{n}_0, \dots, \bar{n}_k, \bar{1})$ , and so it proves  $\neg A_R(\bar{n}_0, \dots, \bar{n}_k)$ .  $\square$

## 17.10 Undecidability

We call a theory  $\mathbf{T}$  *undecidable* if there is no computational procedure which, after finitely many steps and unfailingly, provides a correct answer to the question “does  $\mathbf{T}$  prove  $A$ ?” for any sentence  $A$  in the language of  $\mathbf{T}$ . So  $\mathbf{Q}$  would be decidable iff

there were a computational procedure which decides, given a sentence  $A$  in the language of arithmetic, whether  $\mathbf{Q} \vdash A$  or not. We can make this more precise by asking: Is the relation  $\text{Prov}_{\mathbf{Q}}(y)$ , which holds of  $y$  iff  $y$  is the Gödel number of a sentence provable in  $\mathbf{Q}$ , recursive? The answer is: no.

**Theorem 17.30.**  $\mathbf{Q}$  is undecidable, i.e., the relation

$$\text{Prov}_{\mathbf{Q}}(y) \Leftrightarrow \text{Sent}(y) \wedge \exists x \text{Prf}_{\mathbf{Q}}(x, y)$$

is not recursive.

*Proof.* Suppose it were. Then we could solve the halting problem as follows: Given  $e$  and  $n$ , we know that  $\varphi_e(n) \downarrow$  iff there is an  $s$  such that  $T(e, n, s)$ , where  $T$  is Kleene's predicate from [Theorem 15.28](#). Since  $T$  is primitive recursive it is representable in  $\mathbf{Q}$  by a formula  $B_T$ , that is,  $\mathbf{Q} \vdash B_T(\bar{e}, \bar{n}, \bar{s})$  iff  $T(e, n, s)$ . If  $\mathbf{Q} \vdash B_T(\bar{e}, \bar{n}, \bar{s})$  then also  $\mathbf{Q} \vdash \exists y B_T(\bar{e}, \bar{n}, y)$ . If no such  $s$  exists, then  $\mathbf{Q} \vdash \neg B_T(\bar{e}, \bar{n}, \bar{s})$  for every  $s$ . But  $\mathbf{Q}$  is  $\omega$ -consistent, i.e., if  $\mathbf{Q} \vdash \neg A(\bar{n})$  for every  $n \in \mathbb{N}$ , then  $\mathbf{Q} \not\vdash \exists y A(y)$ . We know this because the axioms of  $\mathbf{Q}$  are true in the standard model  $\mathbb{N}$ . So,  $\mathbf{Q} \not\vdash \exists y B_T(\bar{e}, \bar{n}, y)$ . In other words,  $\mathbf{Q} \vdash \exists y B_T(\bar{e}, \bar{n}, y)$  iff there is an  $s$  such that  $T(e, n, s)$ , i.e., iff  $\varphi_e(n) \downarrow$ . From  $e$  and  $n$  we can compute  $\# \exists y B_T(\bar{e}, \bar{n}, y) \#$ , let  $g(e, n)$  be the primitive recursive function which does that. So

$$h(e, n) = \begin{cases} 1 & \text{if } \text{Prov}_{\mathbf{Q}}(g(e, n)) \\ 0 & \text{otherwise.} \end{cases}$$

This would show that  $h$  is recursive if  $\text{Prov}_{\mathbf{Q}}$  is. But  $h$  is not recursive, by [Theorem 15.29](#), so  $\text{Prov}_{\mathbf{Q}}$  cannot be either.  $\square$

**Corollary 17.31.** *First-order logic is undecidable.*

*Proof.* If first-order logic were decidable, provability in  $\mathbf{Q}$  would be as well, since  $\mathbf{Q} \vdash A$  iff  $\vdash O \rightarrow A$ , where  $O$  is the conjunction of the axioms of  $\mathbf{Q}$ .  $\square$

## Summary

In order to show how theories like  $\mathbf{Q}$  can “talk” about computable functions—and especially about provability (via Gödel numbers)—we established that  $\mathbf{Q}$  **represents** all computable functions. By “ $\mathbf{Q}$  represents  $f(n)$ ” we mean that there is a formula  $A_f(x, y)$  in  $\mathcal{L}_A$  which expresses that  $f(x) = y$ , and  $\mathbf{Q}$  can prove that it does. This, in turn, means that whenever  $f(n) = m$ , then  $\mathbf{Q} \vdash A_f(\bar{n}, \bar{m})$  and  $\mathbf{Q} \vdash \forall y (A_f(\bar{n}, y) \rightarrow y = \bar{m})$ . (Here,  $\bar{n}$  is the **standard numeral** for  $n$ , i.e., the term  $0' \dots'$  with  $n$   $\prime$ s. The term  $\bar{n}$  picks out the number  $n$  in the standard model  $N$ , so it’s a convenient way of representing the number  $n$  in  $\mathcal{L}_A$ .) To prove that  $\mathbf{Q}$  represents all computable functions we go back to the characterization of computable functions as those that can be defined from zero, succ, and the projection functions, by composition, primitive recursion, and regular minimization. While it is relatively easy to prove that the basic functions are representable and that functions defined by composition and regular minimization from representable functions are also representable, primitive recursion is harder. We showed that we can actually avoid definition by primitive recursion, if we allow a few additional basic functions (namely, addition, multiplication, and the characteristic function of  $=$ ). This required a **beta function** which allows us to deal with sequences of numbers in a rudimentary way, and which can be defined without using primitive recursion.

## Problems

**Problem 17.1.** Show that the relations  $x < y$ ,  $x \mid y$ , and the function  $\text{rem}(x, y)$  can be defined without primitive recursion. You may use 0, successor, plus, times,  $\chi_=_$ , projections, and bounded minimization and quantification.

**Problem 17.2.** Prove that  $y = 0$ ,  $y = x'$ , and  $y = x_i$  represent zero, succ, and  $P_i^n$ , respectively.

**Problem 17.3.** Prove Lemma 17.18.

**Problem 17.4.** Use Lemma 17.18 to prove Proposition 17.17.

**Problem 17.5.** Using the proofs of Proposition 17.20 and Proposition 17.20 as a guide, carry out the proof of Proposition 17.21 in detail.

**Problem 17.6.** Show that if  $R$  is representable in  $\mathbf{Q}$ , so is  $\chi_R$ .

## CHAPTER 18

# *Incompleteness and Provability*

### 18.1 Introduction

Hilbert thought that a system of axioms for a mathematical structure, such as the natural numbers, is inadequate unless it allows one to derive all true statements about the structure. Combined with his later interest in formal systems of deduction, this suggests that he thought that we should guarantee that, say, the formal systems we are using to reason about the natural numbers is not only consistent, but also *complete*, i.e., every statement in its language is either derivable or its negation is. Gödel's first incompleteness theorem shows that no such system of axioms exists: there is no complete, consistent, axiomatizable formal system for arithmetic. In fact, no "sufficiently strong," consistent, axiomatizable mathematical theory is complete.

A more important goal of Hilbert's, the centerpiece of his program for the justification of modern ("classical") mathematics, was to find finitary consistency proofs for formal systems representing classical reasoning. With regard to Hilbert's program, then, Gödel's second incompleteness theorem was a much bigger blow. The second incompleteness theorem can be stated in vague terms, like the first incompleteness theorem. Roughly speaking,

it says that no sufficiently strong theory of arithmetic can prove its own consistency. We will have to take “sufficiently strong” to include a little bit more than  $\mathbf{Q}$ .

The idea behind Gödel’s original proof of the incompleteness theorem can be found in the Epimenides paradox. Epimenides, a Cretan, asserted that all Cretans are liars; a more direct form of the paradox is the assertion “this sentence is false.” Essentially, by replacing truth with derivability, Gödel was able to formalize a sentence which, in a roundabout way, asserts that it itself is not derivable. If that sentence were derivable, the theory would then be inconsistent. Gödel showed that the negation of that sentence is also not derivable from the system of axioms he was considering. (For this second part, Gödel had to assume that the theory  $\mathbf{T}$  is what’s called “ $\omega$ -consistent.”  $\omega$ -Consistency is related to consistency, but is a stronger property.<sup>1</sup> A few years after Gödel, Rosser showed that assuming simple consistency of  $\mathbf{T}$  is enough.)

The first challenge is to understand how one can construct a sentence that refers to itself. For every formula  $A$  in the language of  $\mathbf{Q}$ , let  $\ulcorner A \urcorner$  denote the numeral corresponding to  $\#A^\#$ . Think about what this means:  $A$  is a formula in the language of  $\mathbf{Q}$ ,  $\#A^\#$  is a natural number, and  $\ulcorner A \urcorner$  is a *term* in the language of  $\mathbf{Q}$ . So every formula  $A$  in the language of  $\mathbf{Q}$  has a *name*,  $\ulcorner A \urcorner$ , which is a term in the language of  $\mathbf{Q}$ ; this provides us with a conceptual framework in which formulas in the language of  $\mathbf{Q}$  can “say” things about other formulas. The following lemma is known as the fixed-point lemma.

**Lemma 18.1.** *Let  $\mathbf{T}$  be any theory extending  $\mathbf{Q}$ , and let  $B(x)$  be any formula with only the variable  $x$  free. Then there is a sentence  $A$  such that  $\mathbf{T} \vdash A \leftrightarrow B(\ulcorner A \urcorner)$ .*

The lemma asserts that given any property  $B(x)$ , there is a sentence  $A$  that asserts “ $B(x)$  is true of me,” and  $\mathbf{T}$  “knows” this.

<sup>1</sup>That is, any  $\omega$ -consistent theory is consistent, but not vice versa.



How can we construct such a sentence? Consider the following version of the Epimenides paradox, due to Quine:

“Yields falsehood when preceded by its quotation”  
yields falsehood when preceded by its quotation.

This sentence is not directly self-referential. It simply makes an assertion about the syntactic objects between quotes, and, in doing so, it is on par with sentences like

1. “Robert” is a nice name.
2. “I ran.” is a short sentence.
3. “Has three words” has three words.

But what happens when one takes the phrase “yields falsehood when preceded by its quotation,” and precedes it with a quoted version of itself? Then one has the original sentence! In short, the sentence asserts that it is false.

## 18.2 The Fixed-Point Lemma

The fixed-point lemma says that for any formula  $B(x)$ , there is a sentence  $A$  such that  $\mathbf{T} \vdash A \leftrightarrow B(\ulcorner A \urcorner)$ , provided  $\mathbf{T}$  extends  $\mathbf{Q}$ . In the case of the liar sentence, we’d want  $A$  to be equivalent (provably in  $\mathbf{T}$ ) to “ $\ulcorner A \urcorner$  is false,” i.e., the statement that  $\#A$  is the Gödel number of a false sentence. To understand the idea of the proof, it will be useful to compare it with Quine’s informal gloss of  $A$  as, “yields a falsehood when preceded by its own quotation’ yields a falsehood when preceded by its own quotation.” The operation of taking an expression, and then forming a sentence by preceding this expression by its own quotation may be called *diagonalizing* the expression, and the result its diagonalization. So, the diagonalization of ‘yields a falsehood when preceded by its own quotation’ is “‘yields a falsehood when preceded by its own quotation’ yields a falsehood when preceded by

its own quotation.” Now note that Quine’s liar sentence is not the diagonalization of ‘yields a falsehood’ but of ‘yields a falsehood when preceded by its own quotation.’ So the property being diagonalized to yield the liar sentence itself involves diagonalization!

In the language of arithmetic, we form quotations of a formula with one free variable by computing its Gödel numbers and then substituting the standard numeral for that Gödel number into the free variable. The diagonalization of  $E(x)$  is  $E(\bar{n})$ , where  $n = \#E(x)\#$ . (From now on, let’s abbreviate  $\#E(x)\#$  as  $\ulcorner E(x)\urcorner$ .) So if  $B(x)$  is “is a falsehood,” then “yields a falsehood if preceded by its own quotation,” would be “yields a falsehood when applied to the Gödel number of its diagonalization.” If we had a symbol *diag* for the function  $\text{diag}(n)$  which computes the Gödel number of the diagonalization of the formula with Gödel number  $n$ , we could write  $E(x)$  as  $B(\text{diag}(x))$ . And Quine’s version of the liar sentence would then be the diagonalization of it, i.e.,  $E(\ulcorner E(x)\urcorner)$  or  $B(\text{diag}(\ulcorner B(\text{diag}(x))\urcorner))$ . Of course,  $B(x)$  could now be any other property, and the same construction would work. For the incompleteness theorem, we’ll take  $B(x)$  to be “ $x$  is not derivable in  $\mathbf{T}$ .” Then  $E(x)$  would be “yields a sentence not derivable in  $\mathbf{T}$  when applied to the Gödel number of its diagonalization.”

To formalize this in  $\mathbf{T}$ , we have to find a way to formalize *diag*. The function  $\text{diag}(n)$  is computable, in fact, it is primitive recursive: if  $n$  is the Gödel number of a formula  $E(x)$ ,  $\text{diag}(n)$  returns the Gödel number of  $E(\ulcorner E(x)\urcorner)$ . (Recall,  $\ulcorner E(x)\urcorner$  is the standard numeral of the Gödel number of  $E(x)$ , i.e.,  $\#E(x)\#$ .) If *diag* were a function symbol in  $\mathbf{T}$  representing the function *diag*, we could take  $A$  to be the formula  $B(\text{diag}(\ulcorner B(\text{diag}(x))\urcorner))$ . Notice that

$$\begin{aligned} \text{diag}(\#B(\text{diag}(x))\#) &= \#B(\text{diag}(\ulcorner B(\text{diag}(x))\urcorner))\# \\ &= \#A\#. \end{aligned}$$

Assuming  $\mathbf{T}$  can derive

$$\text{diag}(\ulcorner B(\text{diag}(x))\urcorner) = \ulcorner A\urcorner,$$

it can derive  $B(\text{diag}(\ulcorner B(\text{diag}(x)) \urcorner)) \leftrightarrow B(\ulcorner A \urcorner)$ . But the left hand side is, by definition,  $A$ .

Of course,  $\text{diag}$  will in general not be a function symbol of  $\mathbf{T}$ , and certainly is not one of  $\mathbf{Q}$ . But, since  $\text{diag}$  is computable, it is *representable* in  $\mathbf{Q}$  by some formula  $D_{\text{diag}}(x, y)$ . So instead of writing  $B(\text{diag}(x))$  we can write  $\exists y (D_{\text{diag}}(x, y) \wedge B(y))$ . Otherwise, the proof sketched above goes through, and in fact, it goes through already in  $\mathbf{Q}$ .

**Lemma 18.2.** *Let  $B(x)$  be any formula with one free variable  $x$ . Then there is a sentence  $A$  such that  $\mathbf{Q} \vdash A \leftrightarrow B(\ulcorner A \urcorner)$ .*

*Proof.* Given  $B(x)$ , let  $E(x)$  be the formula  $\exists y (D_{\text{diag}}(x, y) \wedge B(y))$  and let  $A$  be its diagonalization, i.e., the formula  $E(\ulcorner E(x) \urcorner)$ .

Since  $D_{\text{diag}}$  represents  $\text{diag}$ , and  $\text{diag}(\ulcorner E(x) \urcorner) = \ulcorner A \urcorner$ ,  $\mathbf{Q}$  can derive

$$D_{\text{diag}}(\ulcorner E(x) \urcorner, \ulcorner A \urcorner) \quad (18.1)$$

$$\forall y (D_{\text{diag}}(\ulcorner E(x) \urcorner, y) \rightarrow y = \ulcorner A \urcorner). \quad (18.2)$$

Now we show that  $\mathbf{Q} \vdash A \leftrightarrow B(\ulcorner A \urcorner)$ . We argue informally, using just logic and facts derivable in  $\mathbf{Q}$ .

First, suppose  $A$ , i.e.,  $E(\ulcorner E(x) \urcorner)$ . Going back to the definition of  $E(x)$ , we see that  $E(\ulcorner E(x) \urcorner)$  just is

$$\exists y (D_{\text{diag}}(\ulcorner E(x) \urcorner, y) \wedge B(y)).$$

Consider such a  $y$ . Since  $D_{\text{diag}}(\ulcorner E(x) \urcorner, y)$ , by eq. (18.2),  $y = \ulcorner A \urcorner$ . So, from  $B(y)$  we have  $B(\ulcorner A \urcorner)$ .

Now suppose  $B(\ulcorner A \urcorner)$ . By eq. (18.1), we have

$$D_{\text{diag}}(\ulcorner E(x) \urcorner, \ulcorner A \urcorner) \wedge B(\ulcorner A \urcorner).$$

It follows that

$$\exists y (D_{\text{diag}}(\ulcorner E(x) \urcorner, y) \wedge B(y)).$$

But that's just  $E(\ulcorner E(x) \urcorner)$ , i.e.,  $A$ . □

You should compare this to the proof of the fixed-point lemma in computability theory. The difference is that here we want to define a *statement* in terms of itself, whereas there we wanted to define a *function* in terms of itself; this difference aside, it is really the same idea.

### 18.3 The First Incompleteness Theorem

We can now describe Gödel's original proof of the first incompleteness theorem. Let  $\mathbf{T}$  be any computably axiomatized theory in a language extending the language of arithmetic, such that  $\mathbf{T}$  includes the axioms of  $\mathbf{Q}$ . This means that, in particular,  $\mathbf{T}$  represents computable functions and relations.

We have argued that, given a reasonable coding of formulas and proofs as numbers, the relation  $\text{Prf}_T(x, y)$  is computable, where  $\text{Prf}_T(x, y)$  holds if and only if  $x$  is the Gödel number of a derivation of the formula with Gödel number  $y$  in  $\mathbf{T}$ . In fact, for the particular theory that Gödel had in mind, Gödel was able to show that this relation is primitive recursive, using the list of 45 functions and relations in his paper. The 45th relation,  $xBy$ , is just  $\text{Prf}_T(x, y)$  for his particular choice of  $\mathbf{T}$ . Remember that where Gödel uses the word “recursive” in his paper, we would now use the phrase “primitive recursive.”

Since  $\text{Prf}_T(x, y)$  is computable, it is representable in  $\mathbf{T}$ . We will use  $\text{Prf}_T(x, y)$  to refer to the formula that represents it. Let  $\text{Prov}_T(y)$  be the formula  $\exists x \text{Prf}_T(x, y)$ . This describes the 46th relation,  $\text{Bew}(y)$ , on Gödel's list. As Gödel notes, this is the only relation that “cannot be asserted to be recursive.” What he probably meant is this: from the definition, it is not clear that it is computable; and later developments, in fact, show that it isn't.

Let  $\mathbf{T}$  be an axiomatizable theory containing  $\mathbf{Q}$ . Then  $\text{Prf}_T(x, y)$  is decidable, hence representable in  $\mathbf{Q}$  by a formula  $\text{Prf}_T(x, y)$ . Let  $\text{Prov}_T(y)$  be the formula we described above. By the fixed-point lemma, there is a formula  $G_T$  such that  $\mathbf{Q}$  (and

hence  $\mathbf{T}$  derives

$$G_{\mathbf{T}} \leftrightarrow \neg \text{Prov}_{\mathbf{T}}(\ulcorner G_{\mathbf{T}} \urcorner). \quad (18.3)$$

Note that  $G_{\mathbf{T}}$  says, in essence, “ $G_{\mathbf{T}}$  is not derivable in  $\mathbf{T}$ .”

**Lemma 18.3.** *If  $\mathbf{T}$  is a consistent, axiomatizable theory extending  $\mathbf{Q}$ , then  $\mathbf{T} \not\vdash G_{\mathbf{T}}$ .*

*Proof.* Suppose  $\mathbf{T}$  derives  $G_{\mathbf{T}}$ . Then there is a derivation, and so, for some number  $m$ , the relation  $\text{Prf}_{\mathbf{T}}(m, \ulcorner G_{\mathbf{T}} \urcorner)$  holds. But then  $\mathbf{Q}$  derives the sentence  $\text{Prf}_{\mathbf{T}}(\bar{m}, \ulcorner G_{\mathbf{T}} \urcorner)$ . So  $\mathbf{Q}$  derives  $\exists x \text{Prf}_{\mathbf{T}}(x, \ulcorner G_{\mathbf{T}} \urcorner)$ , which is, by definition,  $\text{Prov}_{\mathbf{T}}(\ulcorner G_{\mathbf{T}} \urcorner)$ . By eq. (18.3),  $\mathbf{Q}$  derives  $\neg G_{\mathbf{T}}$ , and since  $\mathbf{T}$  extends  $\mathbf{Q}$ , so does  $\mathbf{T}$ . We have shown that if  $\mathbf{T}$  derives  $G_{\mathbf{T}}$ , then it also derives  $\neg G_{\mathbf{T}}$ , and hence it would be inconsistent.  $\square$

**Definition 18.4.** A theory  $\mathbf{T}$  is  $\omega$ -consistent if the following holds: if  $\exists x A(x)$  is any sentence and  $\mathbf{T}$  derives  $\neg A(\bar{0})$ ,  $\neg A(\bar{1})$ ,  $\neg A(\bar{2})$ , ... then  $\mathbf{T}$  does not prove  $\exists x A(x)$ .

Note that every  $\omega$ -consistent theory is also consistent. This follows simply from the fact that if  $\mathbf{T}$  is inconsistent, then  $\mathbf{T} \vdash A$  for every  $A$ . In particular, if  $\mathbf{T}$  is inconsistent, it derives both  $\neg A(\bar{n})$  for every  $n$  and also derives  $\exists x A(x)$ . So, if  $\mathbf{T}$  is inconsistent, it is  $\omega$ -inconsistent. By contraposition, if  $\mathbf{T}$  is  $\omega$ -consistent, it must be consistent.

**Lemma 18.5.** *If  $\mathbf{T}$  is an  $\omega$ -consistent, axiomatizable theory extending  $\mathbf{Q}$ , then  $\mathbf{T} \not\vdash \neg G_{\mathbf{T}}$ .*

*Proof.* We show that if  $\mathbf{T}$  derives  $\neg G_{\mathbf{T}}$ , then it is  $\omega$ -inconsistent. Suppose  $\mathbf{T}$  derives  $\neg G_{\mathbf{T}}$ . If  $\mathbf{T}$  is inconsistent, it is  $\omega$ -inconsistent, and we are done. Otherwise,  $\mathbf{T}$  is consistent, so it does not derive  $G_{\mathbf{T}}$  by Lemma 18.3. Since there is no derivation of  $G_{\mathbf{T}}$  in  $\mathbf{T}$ ,  $\mathbf{Q}$  derives

$$\neg \text{Prf}_{\mathbf{T}}(\bar{0}, \ulcorner G_{\mathbf{T}} \urcorner), \neg \text{Prf}_{\mathbf{T}}(\bar{1}, \ulcorner G_{\mathbf{T}} \urcorner), \neg \text{Prf}_{\mathbf{T}}(\bar{2}, \ulcorner G_{\mathbf{T}} \urcorner), \dots$$

and so does  $\mathbf{T}$ . On the other hand, by eq. (18.3),  $\neg G_{\mathbf{T}}$  is equivalent to  $\exists x \text{Prf}_{\mathbf{T}}(x, \ulcorner G_{\mathbf{T}} \urcorner)$ . So  $\mathbf{T}$  is  $\omega$ -inconsistent.  $\square$

**Theorem 18.6.** *Let  $\mathbf{T}$  be any  $\omega$ -consistent, axiomatizable theory extending  $\mathbf{Q}$ . Then  $\mathbf{T}$  is not complete.*

*Proof.* If  $\mathbf{T}$  is  $\omega$ -consistent, it is consistent, so  $\mathbf{T} \not\vdash G_{\mathbf{T}}$  by Lemma 18.3. By Lemma 18.5,  $\mathbf{T} \not\vdash \neg G_{\mathbf{T}}$ . This means that  $\mathbf{T}$  is incomplete, since it derives neither  $G_{\mathbf{T}}$  nor  $\neg G_{\mathbf{T}}$ .  $\square$

## 18.4 Rosser's Theorem

Can we modify Gödel's proof to get a stronger result, replacing " $\omega$ -consistent" with simply "consistent"? The answer is "yes," using a trick discovered by Rosser. Rosser's trick is to use a "modified" derivability predicate  $\text{RProv}_{\mathbf{T}}(y)$  instead of  $\text{Prov}_{\mathbf{T}}(y)$ .

**Theorem 18.7.** *Let  $\mathbf{T}$  be any consistent, axiomatizable theory extending  $\mathbf{Q}$ . Then  $\mathbf{T}$  is not complete.*

*Proof.* Recall that  $\text{Prov}_{\mathbf{T}}(y)$  is defined as  $\exists x \text{Prf}_{\mathbf{T}}(x, y)$ , where  $\text{Prf}_{\mathbf{T}}(x, y)$  represents the decidable relation which holds iff  $x$  is the Gödel number of a derivation of the sentence with Gödel number  $y$ . The relation that holds between  $x$  and  $y$  if  $x$  is the Gödel number of a *refutation* of the sentence with Gödel number  $y$  is also decidable. Let  $\text{not}(x)$  be the primitive recursive function which does the following: if  $x$  is the code of a formula  $A$ ,  $\text{not}(x)$  is a code of  $\neg A$ . Then  $\text{Ref}_{\mathbf{T}}(x, y)$  holds iff  $\text{Prf}_{\mathbf{T}}(x, \text{not}(y))$ . Let  $\text{Ref}_{\mathbf{T}}(x, y)$  represent it. Then, if  $\mathbf{T} \vdash \neg A$  and  $\delta$  is a corresponding derivation,  $\mathbf{Q} \vdash \text{Ref}_{\mathbf{T}}(\ulcorner \delta \urcorner, \ulcorner A \urcorner)$ . We define  $\text{RProv}_{\mathbf{T}}(y)$  as

$$\exists x (\text{Prf}_{\mathbf{T}}(x, y) \wedge \forall z (z < x \rightarrow \neg \text{Ref}_{\mathbf{T}}(z, y))).$$

Roughly,  $\text{RProv}_{\mathbf{T}}(y)$  says "there is a proof of  $y$  in  $\mathbf{T}$ , and there is no shorter refutation of  $y$ ." Assuming  $\mathbf{T}$  is consistent,  $\text{RProv}_{\mathbf{T}}(y)$  is true of the same numbers as  $\text{Prov}_{\mathbf{T}}(y)$ ; but from the point of

view of *provability* in  $\mathbf{T}$  (and we now know that there is a difference between truth and provability!) the two have different properties. If  $\mathbf{T}$  is *inconsistent*, then the two do *not* hold of the same numbers! ( $\text{RProv}_T(y)$  is often read as “ $y$  is Rosser provable.” Since, as just discussed, Rosser provability is not some special kind of provability—in inconsistent theories, there are sentences that are provable but not Rosser provable—this may be confusing. To avoid the confusion, you could instead read it as “ $y$  is shmovable.”)

By the fixed-point lemma, there is a formula  $R_T$  such that

$$\mathbf{Q} \vdash R_T \leftrightarrow \neg \text{RProv}_T(\ulcorner R_T \urcorner). \quad (18.4)$$

In contrast to the proof of [Theorem 18.6](#), here we claim that if  $\mathbf{T}$  is consistent,  $\mathbf{T}$  doesn't derive  $R_T$ , and  $\mathbf{T}$  also doesn't derive  $\neg R_T$ . (In other words, we don't need the assumption of  $\omega$ -consistency.)

First, let's show that  $\mathbf{T} \not\vdash R_T$ . Suppose it did, so there is a derivation of  $R_T$  from  $T$ ; let  $n$  be its Gödel number. Then  $\mathbf{Q} \vdash \text{Prf}_T(\bar{n}, \ulcorner R_T \urcorner)$ , since  $\text{Prf}_T$  represents  $\text{Prf}_T$  in  $\mathbf{Q}$ . Also, for each  $k < n$ ,  $k$  is not the Gödel number of a derivation of  $\neg R_T$ , since  $\mathbf{T}$  is consistent. So for each  $k < n$ ,  $\mathbf{Q} \vdash \neg \text{Ref}_T(\bar{k}, \ulcorner R_T \urcorner)$ . By [Lemma 17.24](#),  $\mathbf{Q} \vdash \forall z (z < \bar{n} \rightarrow \neg \text{Ref}_T(z, \ulcorner R_T \urcorner))$ . Thus,

$$\mathbf{Q} \vdash \exists x (\text{Prf}_T(x, \ulcorner R_T \urcorner) \wedge \forall z (z < x \rightarrow \neg \text{Ref}_T(z, \ulcorner R_T \urcorner))),$$

but that's just  $\text{RProv}_T(\ulcorner R_T \urcorner)$ . By [eq. \(18.4\)](#),  $\mathbf{Q} \vdash \neg R_T$ . Since  $\mathbf{T}$  extends  $\mathbf{Q}$ , also  $\mathbf{T} \vdash \neg R_T$ . We've assumed that  $\mathbf{T} \vdash R_T$ , so  $\mathbf{T}$  would be inconsistent, contrary to the assumption of the theorem.

Now, let's show that  $\mathbf{T} \not\vdash \neg R_T$ . Again, suppose it did, and suppose  $n$  is the Gödel number of a derivation of  $\neg R_T$ . Then  $\text{Ref}_T(n, \ulcorner \neg R_T \urcorner)$  holds, and since  $\text{Ref}_T$  represents  $\text{Ref}_T$  in  $\mathbf{Q}$ ,  $\mathbf{Q} \vdash \text{Ref}_T(\bar{n}, \ulcorner \neg R_T \urcorner)$ . We'll again show that  $\mathbf{T}$  would then be inconsistent because it would also derive  $R_T$ . Since

$$\mathbf{Q} \vdash R_T \leftrightarrow \neg \text{RProv}_T(\ulcorner R_T \urcorner),$$

and since  $\mathbf{T}$  extends  $\mathbf{Q}$ , it suffices to show that

$$\mathbf{Q} \vdash \neg \text{RProv}_T(\ulcorner R_T \urcorner).$$

The sentence  $\neg \text{RProv}_T(\ulcorner R_T \urcorner)$ , i.e.,

$$\neg \exists x (\text{Prf}_T(x, \ulcorner R_T \urcorner) \wedge \forall z (z < x \rightarrow \neg \text{Ref}_T(z, \ulcorner R_T \urcorner))),$$

is logically equivalent to

$$\forall x (\text{Prf}_T(x, \ulcorner R_T \urcorner) \rightarrow \exists z (z < x \wedge \text{Ref}_T(z, \ulcorner R_T \urcorner))).$$

We argue informally using logic, making use of facts about what  $\mathbf{Q}$  derives. Suppose  $x$  is arbitrary and  $\text{Prf}_T(x, \ulcorner R_T \urcorner)$ . We already know that  $\mathbf{T} \not\vdash R_T$ , and so for every  $k$ ,  $\mathbf{Q} \vdash \neg \text{Prf}_T(\bar{k}, \ulcorner R_T \urcorner)$ . Thus, for every  $k$  it follows that  $x \neq \bar{k}$ . In particular, we have (a) that  $x \neq \bar{n}$ . We also have  $\neg(x = \bar{0} \vee x = \bar{1} \vee \dots \vee x = \overline{n-1})$  and so by [Lemma 17.24](#), (b)  $\neg(x < \bar{n})$ . By [Lemma 17.25](#),  $\bar{n} < x$ . Since  $\mathbf{Q} \vdash \text{Ref}_T(\bar{n}, \ulcorner R_T \urcorner)$ , we have  $\bar{n} < x \wedge \text{Ref}_T(\bar{n}, \ulcorner R_T \urcorner)$ , and from that  $\exists z (z < x \wedge \text{Ref}_T(z, \ulcorner R_T \urcorner))$ . Since  $x$  was arbitrary we get, as required, that

$$\forall x (\text{Prf}_T(x, \ulcorner R_T \urcorner) \rightarrow \exists z (z < x \wedge \text{Ref}_T(z, \ulcorner R_T \urcorner))). \quad \square$$

## 18.5 Comparison with Gödel's Original Paper

It is worthwhile to spend some time with Gödel's 1931 paper. The introduction sketches the ideas we have just discussed. Even if you just skim through the paper, it is easy to see what is going on at each stage: first Gödel describes the formal system  $P$  (syntax, axioms, proof rules); then he defines the primitive recursive functions and relations; then he shows that  $xBy$  is primitive recursive, and argues that the primitive recursive functions and relations are represented in  $\mathbf{P}$ . He then goes on to prove the incompleteness theorem, as above. In Section 3, he shows that one can take the unprovable assertion to be a sentence in the language of arithmetic. This is the origin of the  $\beta$ -lemma, which is



what we also used to handle sequences in showing that the recursive functions are representable in  $\mathbf{Q}$ . Gödel doesn't go so far to isolate a minimal set of axioms that suffice, but we now know that  $\mathbf{Q}$  will do the trick. Finally, in Section 4, he sketches a proof of the second incompleteness theorem.

## 18.6 The Derivability Conditions for PA

Peano arithmetic, or  $\mathbf{PA}$ , is the theory extending  $\mathbf{Q}$  with induction axioms for all formulas. In other words, one adds to  $\mathbf{Q}$  axioms of the form

$$(A(0) \wedge \forall x (A(x) \rightarrow A(x'))) \rightarrow \forall x A(x)$$

for every formula  $A$ . Notice that this is really a *schema*, which is to say, infinitely many axioms (and it turns out that  $\mathbf{PA}$  is *not* finitely axiomatizable). But since one can effectively determine whether or not a string of symbols is an instance of an induction axiom, the set of axioms for  $\mathbf{PA}$  is computable.  $\mathbf{PA}$  is a much more robust theory than  $\mathbf{Q}$ . For example, one can easily prove that addition and multiplication are commutative, using induction in the usual way. In fact, most finitary number-theoretic and combinatorial arguments can be carried out in  $\mathbf{PA}$ .

Since  $\mathbf{PA}$  is computably axiomatized, the derivability predicate  $\text{Prf}_{\mathbf{PA}}(x, y)$  is computable and hence represented in  $\mathbf{Q}$  (and so, in  $\mathbf{PA}$ ). As before, we will take  $\text{Prf}_{\mathbf{PA}}(x, y)$  to denote the formula representing the relation. Let  $\text{Prov}_{\mathbf{PA}}(y)$  be the formula  $\exists x \text{Prf}_{\mathbf{PA}}(x, y)$ , which, intuitively says, “ $y$  is derivable from the axioms of  $\mathbf{PA}$ .” The reason we need a little bit more than the axioms of  $\mathbf{Q}$  is we need to know that the theory we are using is strong enough to derive a few basic facts about this derivability predicate. In fact, what we need are the following facts:

P1. If  $\mathbf{PA} \vdash A$ , then  $\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner)$ .

P2. For all formulas  $A$  and  $B$ ,

$$\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{PA}}(\ulcorner B \urcorner)).$$

P<sub>3</sub>. For every formula  $A$ ,

$$\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{PA}}(\ulcorner \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \urcorner).$$

The only way to verify that these three properties hold is to describe the formula  $\text{Prov}_{\mathbf{PA}}(y)$  carefully and use the axioms of  $\mathbf{PA}$  to describe the relevant formal derivations. Conditions (1) and (2) are easy; it is really condition (3) that requires work. (Think about what kind of work it entails ...) Carrying out the details would be tedious and uninteresting, so here we will ask you to take it on faith that  $\mathbf{PA}$  has the three properties listed above. A reasonable choice of  $\text{Prov}_{\mathbf{PA}}(y)$  will also satisfy

P<sub>4</sub>. If  $\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner)$ , then  $\mathbf{PA} \vdash A$ .

But we will not need this fact.

Incidentally, Gödel was lazy in the same way we are being now. At the end of the 1931 paper, he sketches the proof of the second incompleteness theorem, and promises the details in a later paper. He never got around to it; since everyone who understood the argument believed that it could be carried out (he did not need to fill in the details.)

## 18.7 The Second Incompleteness Theorem

How can we express the assertion that  $\mathbf{PA}$  doesn't prove its own consistency? Saying  $\mathbf{PA}$  is inconsistent amounts to saying that  $\mathbf{PA} \vdash 0 = 1$ . So we can take the consistency statement  $\text{Con}_{\mathbf{PA}}$  to be the sentence  $\neg\text{Prov}_{\mathbf{PA}}(\ulcorner 0 = 1 \urcorner)$ , and then the following theorem does the job:

**Theorem 18.8.** *Assuming  $\mathbf{PA}$  is consistent, then  $\mathbf{PA}$  does not derive  $\text{Con}_{\mathbf{PA}}$ .*

It is important to note that the theorem depends on the particular representation of  $\text{Con}_{\mathbf{PA}}$  (i.e., the particular representation of  $\text{Prov}_{\mathbf{PA}}(y)$ ). All we will use is that the representation of

$\text{Prov}_{\mathbf{PA}}(y)$  satisfies the three derivability conditions, so the theorem generalizes to any theory with a derivability predicate having these properties.

It is informative to read Gödel's sketch of an argument, since the theorem follows like a good punch line. It goes like this. Let  $G_{\mathbf{PA}}$  be the Gödel sentence that we constructed in the proof of [Theorem 18.6](#). We have shown "If  $\mathbf{PA}$  is consistent, then  $\mathbf{PA}$  does not derive  $G_{\mathbf{PA}}$ ." If we formalize this *in*  $\mathbf{PA}$ , we have a proof of

$$\text{Con}_{\mathbf{PA}} \rightarrow \neg \text{Prov}_{\mathbf{PA}}(\ulcorner G_{\mathbf{PA}} \urcorner).$$

Now suppose  $\mathbf{PA}$  derives  $\text{Con}_{\mathbf{PA}}$ . Then it derives  $\neg \text{Prov}_{\mathbf{PA}}(\ulcorner G_{\mathbf{PA}} \urcorner)$ . But since  $G_{\mathbf{PA}}$  is a Gödel sentence, this is equivalent to  $G_{\mathbf{PA}}$ . So  $\mathbf{PA}$  derives  $G_{\mathbf{PA}}$ .

But: we know that if  $\mathbf{PA}$  is consistent, it doesn't derive  $G_{\mathbf{PA}}$ ! So if  $\mathbf{PA}$  is consistent, it can't derive  $\text{Con}_{\mathbf{PA}}$ .

To make the argument more precise, we will let  $G_{\mathbf{PA}}$  be the Gödel sentence for  $\mathbf{PA}$  and use the derivability conditions (P1)–(P3) to show that  $\mathbf{PA}$  derives  $\text{Con}_{\mathbf{PA}} \rightarrow G_{\mathbf{PA}}$ . This will show that  $\mathbf{PA}$  doesn't derive  $\text{Con}_{\mathbf{PA}}$ . Here is a sketch of the proof, *in*  $\mathbf{PA}$ . (For simplicity, we drop the  $\mathbf{PA}$  subscripts.)

$$G \leftrightarrow \neg \text{Prov}(\ulcorner G \urcorner) \tag{18.5}$$

$G$  is a Gödel sentence

$$G \rightarrow \neg \text{Prov}(\ulcorner G \urcorner) \tag{18.6}$$

from [eq. \(18.5\)](#)

$$G \rightarrow (\text{Prov}(\ulcorner G \urcorner) \rightarrow \perp) \tag{18.7}$$

from [eq. \(18.6\)](#) by logic

$$\text{Prov}(\ulcorner G \rightarrow (\text{Prov}(\ulcorner G \urcorner) \rightarrow \perp) \urcorner) \tag{18.8}$$

by from [eq. \(18.7\)](#) by condition P1

$$\text{Prov}(\ulcorner G \urcorner) \rightarrow \text{Prov}(\ulcorner (\text{Prov}(\ulcorner G \urcorner) \rightarrow \perp) \urcorner) \tag{18.9}$$

from [eq. \(18.8\)](#) by condition P2

$$\text{Prov}(\ulcorner G \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner G \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner \perp \urcorner)) \tag{18.10}$$

from [eq. \(18.9\)](#) by condition P2 and logic

$$\text{Prov}(\ulcorner G \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner G \urcorner) \urcorner) \quad (18.11)$$

by P<sub>3</sub>

$$\text{Prov}(\ulcorner G \urcorner) \rightarrow \text{Prov}(\ulcorner \perp \urcorner) \quad (18.12)$$

from eq. (18.10) and eq. (18.11) by logic

$$\text{Con} \rightarrow \neg \text{Prov}(\ulcorner G \urcorner) \quad (18.13)$$

contraposition of eq. (18.12) and  $\text{Con} \equiv \neg \text{Prov}(\ulcorner \perp \urcorner)$

$$\text{Con} \rightarrow G$$

from eq. (18.5) and eq. (18.13) by logic

The use of logic in the above just elementary facts from propositional logic, e.g., eq. (18.7) uses  $\vdash \neg A \leftrightarrow (A \rightarrow \perp)$  and eq. (18.12) uses  $A \rightarrow (B \rightarrow C), A \rightarrow B \vdash A \rightarrow C$ . The use of condition P<sub>2</sub> in eq. (18.9) and eq. (18.10) relies on instances of P<sub>2</sub>,  $\text{Prov}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}(\ulcorner A \urcorner) \rightarrow \text{Prov}(\ulcorner B \urcorner))$ . In the first one,  $A \equiv G$  and  $B \equiv \text{Prov}(\ulcorner G \urcorner) \rightarrow \perp$ ; in the second,  $A \equiv \text{Prov}(\ulcorner G \urcorner)$  and  $B \equiv \perp$ .

The more abstract version of the second incompleteness theorem is as follows:

**Theorem 18.9.** *Let  $\mathbf{T}$  be any consistent, axiomatized theory extending  $\mathbf{Q}$  and let  $\text{Prov}_{\mathbf{T}}(y)$  be any formula satisfying derivability conditions P<sub>1</sub>–P<sub>3</sub> for  $\mathbf{T}$ . Then  $\mathbf{T}$  does not derive  $\text{Con}_{\mathbf{T}}$ .*

The moral of the story is that no “reasonable” consistent theory for mathematics can derive its own consistency statement. Suppose  $\mathbf{T}$  is a theory of mathematics that includes  $\mathbf{Q}$  and Hilbert’s “finitary” reasoning (whatever that may be). Then, the whole of  $\mathbf{T}$  cannot derive the consistency statement of  $\mathbf{T}$ , and so, a fortiori, the finitary fragment can’t derive the consistency statement of  $\mathbf{T}$  either. In that sense, there cannot be a finitary consistency proof for “all of mathematics.”

There is some leeway in interpreting the term “finitary,” and Gödel, in the 1931 paper, grants the possibility that something we may consider “finitary” may lie outside the kinds of mathematics Hilbert wanted to formalize. But Gödel was being charitable;

today, it is hard to see how we might find something that can reasonably be called finitary but is not formalizable in, say, **ZFC**, Zermelo–Fraenkel set theory with the axiom of choice.

## 18.8 Löb’s Theorem

The Gödel sentence for a theory  $\mathbf{T}$  is a fixed point of  $\neg\text{Prov}_T(y)$ , i.e., a sentence  $G$  such that

$$\mathbf{T} \vdash \neg\text{Prov}_T(\ulcorner G \urcorner) \leftrightarrow G.$$

It is not derivable, because if  $\mathbf{T} \vdash G$ , (a) by derivability condition (1),  $\mathbf{T} \vdash \text{Prov}_T(\ulcorner G \urcorner)$ , and (b)  $\mathbf{T} \vdash G$  together with  $\mathbf{T} \vdash \neg\text{Prov}_T(\ulcorner G \urcorner) \leftrightarrow G$  gives  $\mathbf{T} \vdash \neg\text{Prov}_T(\ulcorner G \urcorner)$ , and so  $\mathbf{T}$  would be inconsistent. Now it is natural to ask about the status of a fixed point of  $\text{Prov}_T(y)$ , i.e., a sentence  $H$  such that

$$\mathbf{T} \vdash \text{Prov}_T(\ulcorner H \urcorner) \leftrightarrow H.$$

If it were derivable,  $\mathbf{T} \vdash \text{Prov}_T(\ulcorner H \urcorner)$  by condition (1), but the same conclusion follows if we apply modus ponens to the equivalence above. Hence, we don’t get that  $\mathbf{T}$  is inconsistent, at least not by the same argument as in the case of the Gödel sentence. This of course does not show that  $\mathbf{T}$  *does* derive  $H$ .

We can make headway on this question if we generalize it a bit. The left-to-right direction of the fixed point equivalence,  $\text{Prov}_T(\ulcorner H \urcorner) \rightarrow H$ , is an instance of a general schema called a *reflection principle*:  $\text{Prov}_T(\ulcorner A \urcorner) \rightarrow A$ . It is called that because it expresses, in a sense, that  $\mathbf{T}$  can “reflect” about what it can derive; basically it says, “If  $\mathbf{T}$  can derive  $A$ , then  $A$  is true,” for any  $A$ . This is true for sound theories only, of course, and this suggests that theories will in general not derive every instance of it. So which instances can a theory (strong enough, and satisfying the derivability conditions) derive? Certainly all those where  $A$  itself is derivable. And that’s it, as the next result shows.

**Theorem 18.10.** *Let  $\mathbf{T}$  be an axiomatizable theory extending  $\mathbf{Q}$ , and suppose  $\text{Prov}_{\mathbf{T}}(y)$  is a formula satisfying conditions  $P_1$ – $P_3$  from section 18.7. If  $\mathbf{T}$  derives  $\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$ , then in fact  $\mathbf{T}$  derives  $A$ .*

Put differently, if  $\mathbf{T} \not\vdash A$ , then  $\mathbf{T} \not\vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$ . This result is known as Löb’s theorem.

The heuristic for the proof of Löb’s theorem is a clever proof that Santa Claus exists. (If you don’t like that conclusion, you are free to substitute any other conclusion you would like.) Here it is:

1. Let  $X$  be the sentence, “If  $X$  is true, then Santa Claus exists.”
2. Suppose  $X$  is true.
3. Then what it says holds; i.e., we have: if  $X$  is true, then Santa Claus exists.
4. Since we are assuming  $X$  is true, we can conclude that Santa Claus exists, by modus ponens from (2) and (3).
5. We have succeeded in deriving (4), “Santa Claus exists,” from the assumption (2), “ $X$  is true.” By conditional proof, we have shown: “If  $X$  is true, then Santa Claus exists.”
6. But this is just the sentence  $X$ . So we have shown that  $X$  is true.
7. But then, by the argument (2)–(4) above, Santa Claus exists.

A formalization of this idea, replacing “is true” with “is derivable,” and “Santa Claus exists” with  $A$ , yields the proof of Löb’s theorem. The trick is to apply the fixed-point lemma to the formula  $\text{Prov}_{\mathbf{T}}(y) \rightarrow A$ . The fixed point of that corresponds to the sentence  $X$  in the preceding sketch.

*Proof of Theorem 18.10.* Suppose  $A$  is a sentence such that  $\mathbf{T}$  derives  $\text{Prov}_T(\ulcorner A \urcorner) \rightarrow A$ . Let  $B(y)$  be the formula  $\text{Prov}_T(y) \rightarrow A$ , and use the fixed-point lemma to find a sentence  $D$  such that  $\mathbf{T}$  derives  $D \leftrightarrow B(\ulcorner D \urcorner)$ . Then each of the following is derivable in  $\mathbf{T}$ :

$$D \leftrightarrow (\text{Prov}_T(\ulcorner D \urcorner) \rightarrow A) \quad (18.14)$$

$D$  is a fixed point of  $B(y)$

$$D \rightarrow (\text{Prov}_T(\ulcorner D \urcorner) \rightarrow A) \quad (18.15)$$

from eq. (18.14)

$$\text{Prov}_T(\ulcorner D \rightarrow (\text{Prov}_T(\ulcorner D \urcorner) \rightarrow A) \urcorner) \quad (18.16)$$

from eq. (18.15) by condition P1

$$\text{Prov}_T(\ulcorner D \urcorner) \rightarrow \text{Prov}_T(\ulcorner \text{Prov}_T(\ulcorner D \urcorner) \rightarrow A \urcorner) \quad (18.17)$$

from eq. (18.16) using condition P2

$$\text{Prov}_T(\ulcorner D \urcorner) \rightarrow (\text{Prov}_T(\ulcorner \text{Prov}_T(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}_T(\ulcorner A \urcorner)) \quad (18.18)$$

from eq. (18.17) using P2 again

$$\text{Prov}_T(\ulcorner D \urcorner) \rightarrow \text{Prov}_T(\ulcorner \text{Prov}_T(\ulcorner D \urcorner) \urcorner) \quad (18.19)$$

by derivability condition P3

$$\text{Prov}_T(\ulcorner D \urcorner) \rightarrow \text{Prov}_T(\ulcorner A \urcorner) \quad (18.20)$$

from eq. (18.18) and eq. (18.19)

$$\text{Prov}_T(\ulcorner A \urcorner) \rightarrow A \quad (18.21)$$

by assumption of the theorem

$$\text{Prov}_T(\ulcorner D \urcorner) \rightarrow A \quad (18.22)$$

from eq. (18.20) and eq. (18.21)

$$(\text{Prov}_T(\ulcorner D \urcorner) \rightarrow A) \rightarrow D \quad (18.23)$$

from eq. (18.14)

$$D \quad (18.24)$$

from eq. (18.22) and eq. (18.23)

$$\text{Prov}_T(\ulcorner D \urcorner) \quad (18.25)$$

from eq. (18.24) by condition P1

$A$  from eq. (18.21) and eq. (18.25) □

With Löb's theorem in hand, there is a short proof of the second incompleteness theorem (for theories having a derivability predicate satisfying conditions P1–P3): if  $\mathbf{T} \vdash \text{Prov}_T(\ulcorner \perp \urcorner) \rightarrow \perp$ , then  $\mathbf{T} \vdash \perp$ . If  $\mathbf{T}$  is consistent,  $\mathbf{T} \not\vdash \perp$ . So,  $\mathbf{T} \not\vdash \text{Prov}_T(\ulcorner \perp \urcorner) \rightarrow \perp$ , i.e.,  $\mathbf{T} \not\vdash \text{Con}_T$ . We can also apply it to show that  $H$ , the fixed point of  $\text{Prov}_T(x)$ , is derivable. For since

$$\mathbf{T} \vdash \text{Prov}_T(\ulcorner H \urcorner) \leftrightarrow H$$

in particular

$$\mathbf{T} \vdash \text{Prov}_T(\ulcorner H \urcorner) \rightarrow H$$

and so by Löb's theorem,  $\mathbf{T} \vdash H$ .

## 18.9 The Undefinability of Truth

The notion of *definability* depends on having a formal semantics for the language of arithmetic. We have described a set of formulas and sentences in the language of arithmetic. The “intended interpretation” is to read such sentences as making assertions about the natural numbers, and such an assertion can be true or false. Let  $N$  be the structure with domain  $\mathbb{N}$  and the standard interpretation for the symbols in the language of arithmetic. Then  $N \vDash A$  means “ $A$  is true in the standard interpretation.”

**Definition 18.11.** A relation  $R(x_1, \dots, x_k)$  of natural numbers is *definable* in  $N$  if and only if there is a formula  $A(x_1, \dots, x_k)$  in the language of arithmetic such that for every  $n_1, \dots, n_k$ ,  $R(n_1, \dots, n_k)$  if and only if  $N \vDash A(\bar{n}_1, \dots, \bar{n}_k)$ .

Put differently, a relation is definable in  $N$  if and only if it is representable in the theory  $\mathbf{TA}$ , where  $\mathbf{TA} = \{A : N \vDash A\}$  is the set of true sentences of arithmetic. (If this is not immediately clear to you, you should go back and check the definitions and convince yourself that this is the case.)



**Lemma 18.12.** *Every computable relation is definable in  $N$ .*

*Proof.* It is easy to check that the formula representing a relation in  $\mathbf{Q}$  defines the same relation in  $N$ .  $\square$

Now one can ask, is the converse also true? That is, is every relation definable in  $N$  computable? The answer is no. For example:

**Lemma 18.13.** *The halting relation is definable in  $N$ .*

*Proof.* Let  $H$  be the halting relation, i.e.,

$$H = \{\langle e, x \rangle : \exists s T(e, x, s)\}.$$

Let  $D_T$  define  $T$  in  $N$ . Then

$$H = \{\langle e, x \rangle : N \vDash \exists s D_T(\bar{e}, \bar{x}, s)\},$$

so  $\exists s D_T(z, x, s)$  defines  $H$  in  $N$ .  $\square$

What about **TA** itself? Is it definable in arithmetic? That is: is the set  $\{\ulcorner A \urcorner : N \vDash A\}$  definable in arithmetic? Tarski's theorem answers this in the negative.

**Theorem 18.14.** *The set of true sentences of arithmetic is not definable in arithmetic.*

*Proof.* Suppose  $D(x)$  defined it, i.e.,  $N \vDash A$  iff  $N \vDash D(\ulcorner A \urcorner)$ . By the fixed-point lemma, there is a formula  $A$  such that  $\mathbf{Q} \vdash A \leftrightarrow \neg D(\ulcorner A \urcorner)$ , and hence  $N \vDash A \leftrightarrow \neg D(\ulcorner A \urcorner)$ . But then  $N \vDash A$  if and only if  $N \vDash \neg D(\ulcorner A \urcorner)$ , which contradicts the fact that  $D(y)$  is supposed to define the set of true statements of arithmetic.  $\square$

Tarski applied this analysis to a more general philosophical notion of truth. Given any language  $L$ , Tarski argued that an adequate notion of truth for  $L$  would have to satisfy, for each sentence  $X$ ,

‘ $X$ ’ is true if and only if  $X$ .

Tarski’s oft-quoted example, for English, is the sentence

‘Snow is white’ is true if and only if snow is white.

However, for any language strong enough to represent the diagonal function, and any linguistic predicate  $T(x)$ , we can construct a sentence  $X$  satisfying “ $X$  if and only if not  $T('X')$ .” Given that we do not want a truth predicate to declare some sentences to be both true and false, Tarski concluded that one cannot specify a truth predicate for all sentences in a language without, somehow, stepping outside the bounds of the language. In other words, a the truth predicate for a language cannot be defined in the language itself.

## 18.10 Tarski’s Theorem and Löb’s Theorem

Tarski’s Theorem shows that there’s a gap between the notions of “*sentence provable in the theory PA*” (focusing on **PA** as the best theory of arithmetic we have and that of ‘sentence true in  $N$ , the standard structure of natural numbers’). The former notion is *definable* in  $\mathcal{L}_A$ , and the latter is *not definable*. Remember, ‘true in  $N$ ’ is really the same as ‘provable in **TA**’, so the difference, if you want to put it this way, is between these two notions of provable.

Tarski’s paper “Truth and proof” (Tarski, 1969) makes the following interesting remarks, which emphasize this gap, but which put an optimistic gloss on it:

Nothing is detracted from the significance of this result [these results] by the fact that its philosophical implications are essentially negative in character. *The result shows indeed that in no domain of mathematics is the notion of provability a perfect substitute for the notion*

*of truth. The belief that formal proof can serve as an adequate instrument for establishing truth of all mathematical statements has proved to be unfounded.* The original triumph of formal methods has been followed by a serious setback.

Whatever can be said to conclude this discussion is bound to be an anticlimax. The notion of truth for formalized theories can now be introduced by means of a precise and adequate definition. It can therefore be used without any restrictions and reservations in metalogical discussion. It has actually become a basic metalogical notion involved in important problems and results. On the other hand, the notion of proof has not lost its significance either. *Proof is still the only method used to ascertain the truth of sentences within any specific mathematical theory. We are now aware of the fact, however, that there are sentences formulated in the language of the theory which are true but not provable, and we cannot discount the possibility that some such sentences occur among those in which we are interested and which we attempt to prove.* Hence in some situations we may wish to explore the possibility of widening the set of provable sentences. To this end we enrich the given theory by including new sentences in its axiom system or by providing it with new rules of proof. In doing so we use the notion of truth as a guide; for we do not wish to add a new axiom or a new rule of proof if we have reason to believe that the new axiom is not a true sentence, or that the new rule of proof when applied to true sentences may yield a false sentence. The process of extending a theory may of course be repeated arbitrarily many times. The notion of a true sentence functions thus as an ideal limit which can never be reached but which we try to approximate by gradually widening the set of provable sentences. (It seems likely, although for

different reasons, that the notion of truth plays an analogous role in the realm of empirical knowledge.) There is no conflict between the notions of truth and proof in the development of mathematics; the two notions are not at war but live in peaceful coexistence. (Tarski, 1969, p. 77, emphasis added).

What Tarski brings out here is the gap between truth and proof, and the fact that there must be, in effect, a dialectical relationship between the two, perhaps more than what he calls here ‘peaceful coexistence’. This gap was dramatically emphasised by Gödel, in a lecture in 1951, in a comment he made about the Second Incompleteness Theorem, a comment which supports Tarski’s point that we should use “truth as a guide”:

It is *this* [the second] theorem which makes the incompleteness of mathematics particularly evident. For, *it makes it impossible that someone should set up a certain well-defined system of axioms and rules and consistently make the following assertion about it: All of these axioms and rules I perceive (with mathematical certitude) to be correct, and moreover I believe that they contain all of mathematics.* If someone makes such a statement he contradicts himself.<sup>2</sup> For if he perceives the axioms under consideration to be correct, he also perceives (with the same certainty) that they are consistent. Hence he has a mathematical insight not derivable from his axioms. However, one has to be careful in order to understand clearly the meaning of this state of affairs. Does it mean that no well-defined system of correct axioms can contain all of mathematics proper? It does, if by mathematics proper is understood the system of all true mathematical propositions; it does

---

<sup>2</sup>[Gödel’s footnote:] If he only says “I believe I shall be able to perceive one after the other to be true” (where their number is supposed to be infinite), he does not contradict himself.

not, however, if one understands by it the system of all demonstrable mathematical propositions. (Gödel, 1995, p. 309)

The juxtaposition of Tarski's Theorem and the Gödel results show us that the gap between truth and provability is large; what Löb's Theorem does is emphasise how large that gap really is, and in particular how unlike the notion of truth the notion of proof is. Above all it shows that proof can't really act as a surrogate for truth, that provability can't be just a form of "supertruth", like logical truth or necessary truth.

Let us show how Löb's Theorem brings out the difference.

First, if we proceeded with truth in much the same way as we proceed in the prove of Löb's Theorem we would get something quite absurd, a proof of the existence of Santa Claus. Note the fact that we have "Santa Claus exists" here plays no role whatsoever. It could be "The Dark Lord exists", or any one of your favourite nonsense claims, e.g., "The Moon is made of green cheese". We presented a version of this argument in the textbook. Here is the same argument given by George Boolos in his book *The Logic of Provability*:

Let Sam be the sentence "if Sam is true, SC" [where "SC" abbreviates the sentence "Santa Claus exists"]. Assume that Sam is true; then "if Sam is true, SC" is true; thus if Sam is true, SC; and so SC by modus ponens. Thus we have shown that SC on the assumption that Sam is true and have therefore shown outright that if Sam is true, SC. But then "If Sam is true, SC" is true, i.e., Sam is true, and by modus ponens again, SC. (Boolos, 1993, p. 56)

To reiterate, *any* statement can be proved this way. What this means is that we can easily get a contradiction, which means that we must be operating with inconsistent assumptions! But which inconsistent assumptions? In fact, the argument above is just an

elaborate version of the Liar Paradox, and we are in fact operating with the same inconsistent assumption as is behind that.

We can show this as follows. If you look carefully at the argument Boolos presents, you will notice that several times we make use of the principle:

$$\text{“}X\text{” is true iff } X \quad (\text{T})$$

which we saw earlier illustrated through the particular instance “‘Snow is white’ is true iff snow is white”. We use the principle here in the shift (left-to-right) from “‘If Sam is true, SC’ is true” to “If Sam is true, SC”. Then later we have the shift (right-to-left here) from “If Sam is true, SC” to “‘If Sam is true, SC’ is true”. The principle (T) (often referred to as “convention (T)”) is the basic principle that a correct truth-definition has to satisfy. Moreover, it is the principle which in English (taken as the meta-language for discussion of English) allows us to deduce a contradiction from “This sentence is false”, i.e., the statement which gives us the Liar Paradox.

Let’s now go back to Tarski’s Theorem in the form that we proved it, i.e., showing that no 1-place predicate of the language can define the set of Gödel-numbers of sentences true in the standard model. Assume that there were such a predicate  $D(x)$ . We would have that  $n$  is the Gödel number of sentence  $A$  (so  $n = \ulcorner A \urcorner$ ) true in the standard model if, and only if,  $N \models D(\ulcorner A \urcorner)$ . From this follows  $N \models A$  iff  $N \models D(\ulcorner A \urcorner)$ . But this means that  $N \models A \leftrightarrow D(\ulcorner A \urcorner)$ . And remember that saying “ $N \models X$ ” is the same as saying “ $\mathbf{TA} \vdash X$ ”. Put all this together and we have:

$$\mathbf{TA} \vdash A \leftrightarrow D(\ulcorner A \urcorner), \quad (\text{T}_{\mathbf{TA}})$$

which, if we read  $D(x)$  as “is true”, really says that  $\mathbf{TA}$  can produce the principle eq. (T) above. We often say that principles like eq. (T<sub>TA</sub>) are “truth-definitions”, holding as they do for all sentences of the language concerned. But this immediately gives rise here to a contradiction, for we can apply the Diagonal Lemma to  $\neg D(x)$  and get a sentence  $L$  such that

$$\mathbf{TA} \vdash L \leftrightarrow \neg D(\ulcorner L \urcorner).$$

But eq.  $(T_{TA})$  above instantiated for  $L$  will give us:

$$\mathbf{TA} \vdash L \leftrightarrow D(\ulcorner L \urcorner),$$

thus contradiction.

Let's come back to the difference between truth and provability. First, as we've seen, what the principle eq.  $(T)$  allows is switching as we please between " $X$  is true" and " $X$ ". But the reasoning presented in the proof of Löb's Theorem "mimics" the Santa Claus argument when "is true" replaced by "is provable", and where the requisite switching between " $X$  is provable" and " $X$ " is now provided by the derivability conditions. We do not, of course, get a contradiction.

Second, we can use Löb's Theorem to give us Tarski's Theorem in the form that we cannot have a truth-definition, i.e., that "true in  $N$ " cannot be definable. The argument goes roughly like this:

Suppose  $\text{Tr}(x)$  is a truth-predicate for an appropriate theory  $\mathbf{T}$ , and the truth-definition

$$\mathbf{T} \vdash A \leftrightarrow \text{Tr}(\ulcorner A \urcorner) \tag{Tr}$$

holds. Then it's clear using this that all of the derivability conditions hold for  $\text{Tr}(x)$ , which means that we must have a Löb Theorem for  $\text{Tr}(x)$ . (Remember, in the statement of Löb's Theorem, all we really need to know about the provability predicate is that the derivability conditions hold for it. We know nothing else about the "inner workings" of the predicate. Thus, if there were a truth-predicate, it would satisfy the conditions of Löb's Theorem.) Now take *any* sentence  $A$ . Then by  $(Tr)$  we have  $\mathbf{T} \vdash \text{Tr}(\ulcorner A \urcorner) \leftrightarrow A$ , so in particular  $\mathbf{T} \vdash \text{Tr}(\ulcorner A \urcorner) \rightarrow A$ . By Löb's Theorem applied to  $\text{Tr}(x)$ , this means that we have  $\mathbf{T} \vdash A$ . So *any* sentence would be provable in  $\mathbf{T}$  if  $\mathbf{T}$  were equipped with a truth-definition, which of course would make  $\mathbf{T}$  inconsistent. So a sufficient condition for a sentence  $A$  to be provable is that  $\mathbf{T} \vdash \text{Tr}(\ulcorner A \urcorner) \rightarrow A$ , but  $(Tr)$  says that that must be the case for any

sentence. Hence, using Löb's Theorem, we've shown that if **T** is consistent, there can't be a truth-definition for **T**.

In short, if we could get Löb for truth, we would get contradiction, but Löb for provability *does not* give us contradiction. But if go further into the comparison, we can see how wide the gap is.

Boolos has some comments on the remarkable nature of the Löb Theorem which touch on what we've just observed, and which emphasize just how big the gap really is. In reading these comments, first think of "Prov( $x$ )" as a natural replacement for "truth( $x$ )"; that will highlight the surprise. Boolos comments as follows:

Löb's theorem is utterly astonishing for at least five reasons. In the first place, it is often hard to understand how vast the mathematical gap is between truth and provability. And to one who lacks that understanding and does not distinguish between truth and provability,  $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$ , which the hypothesis of Löb's theorem asserts to be provable, might appear to be trivially true in all cases, whether  $S$  is true or false, provable or unprovable. But if  $S$  is false,  $S$  had better not be provable. Thus it would seem that  $S$  ought not always to be provable provided merely that (the possibly trivial-seeming)  $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$  is provable.

Secondly, Prov seems here to be working like negation. After all, if  $\neg S \rightarrow S$  is provable, then so is  $S$ ; proving  $S$  by proving  $\neg S \rightarrow S$  is called *reductio ad absurdum* (or, sometimes, the law of Clavius). Moreover, inferring  $S$  solely on the ground that  $(S \rightarrow S)$  is demonstrable is known as begging the question, or reasoning in a circle. To one who conflates truth and provability, it may then seem that Löb's theorem asserts that begging the question is an admissible form of reasoning in PA.



Thirdly, one might have thought that *at least on occasion*, PA would claim to be sound with regard to an unprovable sentence  $S$ , i.e., claim that *if* it proves  $S$ , then  $S$  holds. But Löb's theorem tells us that it never does so: PA makes the claim  $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$  that it is sound with regard to  $S$  only when it obviously must, when the consequent  $S$  is actually provable. As Rohit Parikh once put it, "PA couldn't be more modest about its own veracity".

Fourthly, one might very naturally suppose that provability is a kind of necessity, and therefore, just as  $\Box(\Box p \rightarrow p)$  always expresses a truth if the box is interpreted as "it is necessary that" — for then  $\Box(\Box p \rightarrow p)$  says that it is necessarily true that if a statement is necessarily true, it is true —  $\text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \rightarrow S \urcorner)$  would also always be true or at least true in some cases in which  $S$  is false and not true only in the rather exceptional cases in which  $S$  is actually provable.

Finally, it seems wholly bizarre that the statement that if  $S$  is provable, then  $S$  is true is not itself provable, in general. For isn't it perfectly obvious, for any  $S$ , that  $S$  is true if provable? Why are we bothering with PA if its theorems are false? And how could any such (apparently) obvious truth not be provable? (Boolos, 1993, pp. 54–55)<sup>3</sup>

Let's elaborate a little on some of these points.

The first point. Take an example like " $2 + 2 = 5$ ". We certainly can't have it that **PA** can prove  $2 + 2 = 5$ , which means (following Löb's Theorem) that it can't be the case that  $\mathbf{PA} \vdash \text{Prov}(\ulcorner 2 + 2 = 5 \urcorner) \rightarrow 2 + 2 = 5$ ! In other words, **PA** can't tell us that whatever it proves is true, or that the notion of proof it recognizes is a good one! And haven't we designed it so that

---

<sup>3</sup>We've replaced Boolos' use of Gödel's *Bew* with our 'Prov'.

the formal provability predicate “ $\text{Prov}(x)$ ” (for **PA**) matches “is provable in **PA**”?

The third point extends this. Suppose we take something like the Goldbach Conjecture (*GC*) or the Twin Prime Conjecture (*TPC*), propositions that we don’t know to be provable in **PA**. We would hope that **PA** would be able to prove of such propositions, at least some of the time, that  $\text{Prov}(\ulcorner GC \urcorner) \rightarrow GC$  or  $\text{Prov}(\ulcorner TPC \urcorner) \rightarrow TPC$ , i.e., that it knows a **PA** proof of *TPC* would show that *TPC* is correct. But Löb’s Theorem tells us it can’t do that, or rather that it can only do that if there is *already* a **PA** proof of *GC* or *TPC*!

Think about how different this is from truth. For a truth-definition (if we have one), it has got to be the case that  $\text{Tr}(\ulcorner A \urcorner) \leftrightarrow A$  holds regardless of whether *A* is true, false, contradictory or ridiculous: and therefore  $\text{Tr}(\ulcorner A \urcorner) \rightarrow A$  holds also, regardless of what *A* asserts. In other words (switching to English) “‘The moon is made of green cheese’ is true  $\Leftrightarrow$  [or even just  $\Rightarrow$ ] the moon is made of green cheese”, even though the sentence involved here is not true, and is even faintly ludicrous, one we’re only prepared to contemplate by reason of its syntactic formulation. But Löb’s Theorem tells us that “ $\mathbf{T} \vdash S$ ” cannot be a surrogate for truth in this sense, for we could then, it seems, prove that the moon is made of green cheese (or any other kind of cheese, or marshmallow or sphagnum moss ...), or some arithmetic equivalents. That shows us (Tarski’s Theorem) that we can’t have a truth-definition for **PA** where  $(\mathbf{T}_{\mathbf{PA}})$  is provable in **PA**. But if “provable” is a surrogate for “is true”, then we would certainly expect **PA** to be able to prove “ $\text{Prov}(\ulcorner A \urcorner) \rightarrow A$ ” in general. This is stressed in Point 5.

Lastly, consider the second point. As Boolos points out, suppose we think that “ $\text{Prov}(\ulcorner S \urcorner)$ ” is really some sort of strong affirmation of *S* (“being true and more”). In this case, then “ $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$ ” would be something like “ $S \rightarrow S$ ” (or at least this should follow from “ $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$ ”). But then using this as a *justification* for “ $\mathbf{T} \vdash S$ ” is really like arguing by begging the question! However, if “ $\text{Prov}(\ulcorner S \urcorner) \rightarrow S$ ” really acts as a genuine

justification for  $S$ , then it looks as if “ $\text{Prov}(\ulcorner S \urcorner)$ ” is behaving, not like an affirmation at all, but really rather like a *negation*, as in the proof of “ $S$ ” given by proving “ $\neg S \rightarrow S$ ”! So it seems as if it can’t be the case that “ $\text{Prov}(\ulcorner S \urcorner)$ ” is like a strong *affirmation* of  $S$  (“being true and more”) at all, and is actually more like a *denial* of  $S$ .

## Summary

The **first incompleteness theorem** states that for any consistent, axiomatizable theory  $\mathbf{T}$  that extends  $\mathbf{Q}$ , there is a sentence  $G_{\mathbf{T}}$  such that  $\mathbf{T} \not\vdash G_{\mathbf{T}}$ .  $G_{\mathbf{T}}$  is constructed in such a way that  $G_{\mathbf{T}}$ , in a roundabout way, says “ $\mathbf{T}$  does not prove  $G_{\mathbf{T}}$ .” Since  $\mathbf{T}$  does not, in fact, prove it, what it says is true. If  $N \models \mathbf{T}$ , then  $\mathbf{T}$  does not prove any false claims, so  $\mathbf{T} \not\vdash \neg G_{\mathbf{T}}$ . Such a sentence is **independent** or **undecidable** in  $\mathbf{T}$ . Gödel’s original proof established that  $G_{\mathbf{T}}$  is independent on the assumption that  $\mathbf{T}$  is  $\omega$ -**consistent**. Rosser improved the result by finding a different sentence  $R_{\mathbf{T}}$  with is neither provable nor refutable in  $\mathbf{T}$  as long as  $\mathbf{T}$  is simply consistent.

The construction of  $G_{\mathbf{T}}$  is effective: given an axiomatization of  $\mathbf{T}$  we could, in principle, write down  $G_{\mathbf{T}}$ . The “roundabout way” in which  $G_{\mathbf{T}}$  states its own unprovability, is a special case of a general result, the **fixed-point lemma**. It states that for any formula  $B(y)$  in  $\mathcal{L}_A$ , there is a sentence  $A$  such that  $\mathbf{Q} \vdash A \leftrightarrow B(\ulcorner A \urcorner)$ . (Here,  $\ulcorner A \urcorner$  is the standard numeral for the Gödel number of  $A$ , i.e.,  $\overline{\ulcorner A \urcorner}$ .) To obtain  $G_{\mathbf{T}}$ , we use the formula  $\neg \text{Prov}_{\mathbf{T}}(y)$  as  $B(y)$ . We get  $\text{Prov}_{\mathbf{T}}$  as the culmination of our previous efforts: We know that  $\text{Prf}_{\mathbf{T}}(n, m)$ , which holds if  $n$  is the Gödel number of a derivation of the sentence with Gödel number  $m$  from  $\mathbf{T}$ , is primitive recursive. We also know that  $\mathbf{Q}$  represents all primitive recursive relations, and so there is some formula  $\text{Prf}_{\mathbf{T}}(x, y)$  that represents  $\text{Prf}_{\mathbf{T}}$  in  $\mathbf{Q}$ . The **provability predicate** for  $\mathbf{T}$  is  $\text{Prov}_{\mathbf{T}}(y)$  is  $\exists x \text{Prf}_{\mathbf{T}}(x, y)$  then expresses provability in  $\mathbf{T}$ . (It doesn’t represent it though: if  $\mathbf{T} \vdash A$ , then  $\mathbf{Q} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner)$ ; but if  $\mathbf{T} \not\vdash A$ ,

then  $\mathbf{Q}$  does not in general prove  $\neg\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner)$ .)

The **second incompleteness theorem** establishes that the sentence  $\text{Con}_{\mathbf{T}}$  that expresses that  $\mathbf{T}$  is consistent, i.e.,  $\mathbf{T}$  also does not prove  $\neg\text{Prov}_{\mathbf{T}}(\ulcorner \perp \urcorner)$ . The proof of the second incompleteness theorem requires some additional conditions on  $\mathbf{T}$ , the **provability conditions**.  $\mathbf{PA}$  satisfies them, although  $\mathbf{Q}$  does not. Theories that satisfy the provability conditions also satisfy **Löb's theorem**:  $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$  iff  $\mathbf{T} \vdash A$ .

The fixed-point theorem also has another important consequence. We say a relation  $R(n)$  is **definable** in  $\mathcal{L}_A$  if there is a formula  $A_R(x)$  such that  $N \models A_R(\bar{n})$  iff  $R(n)$  holds. For instance,  $\text{Prov}_{\mathbf{T}}$  is definable, since  $\text{Prov}_{\mathbf{T}}$  defines it. The property  $n$  has iff it is the Gödel number of a sentence true in  $N$ , however, is not definable. This is **Tarski's theorem** about the undefinability of truth.

## Problems

**Problem 18.1.** A formula  $A(x)$  is a *truth definition* if  $\mathbf{Q} \vdash B \leftrightarrow A(\ulcorner B \urcorner)$  for all sentences  $B$ . Show that no formula is a truth definition by using the fixed-point lemma.

**Problem 18.2.** Every  $\omega$ -consistent theory is consistent. Show that the converse does not hold, i.e., that there are consistent but  $\omega$ -inconsistent theories. Do this by showing that  $\mathbf{Q} \cup \{\neg G_{\mathbf{Q}}\}$  is consistent but  $\omega$ -inconsistent.

**Problem 18.3.** Two sets  $A$  and  $B$  of natural numbers are said to be *computably inseparable* if there is no decidable set  $X$  such that  $A \subseteq X$  and  $B \subseteq \bar{X}$  ( $\bar{X}$  is the complement,  $\mathbb{N} \setminus X$ , of  $X$ ). Let  $\mathbf{T}$  be a consistent axiomatizable extension of  $\mathbf{Q}$ . Suppose  $A$  is the set of Gödel numbers of sentences provable in  $\mathbf{T}$  and  $B$  the set of Gödel numbers of sentences refutable in  $\mathbf{T}$ . Prove that  $A$  and  $B$  are computably inseparable.

**Problem 18.4.** Show that  $\mathbf{PA}$  derives  $G_{\mathbf{PA}} \rightarrow \text{Con}_{\mathbf{PA}}$ .

**Problem 18.5.** Let  $\mathbf{T}$  be a computably axiomatized theory, and let  $\text{Prov}_T$  be a derivability predicate for  $\mathbf{T}$ . Consider the following four statements:

1. If  $T \vdash A$ , then  $T \vdash \text{Prov}_T(\ulcorner A \urcorner)$ .
2.  $T \vdash A \rightarrow \text{Prov}_T(\ulcorner A \urcorner)$ .
3. If  $T \vdash \text{Prov}_T(\ulcorner A \urcorner)$ , then  $T \vdash A$ .
4.  $T \vdash \text{Prov}_T(\ulcorner A \urcorner) \rightarrow A$

Under what conditions are each of these statements true?

**Problem 18.6.** Show that  $Q(n) \Leftrightarrow n \in \{^{\#}A^{\#} : \mathbf{Q} \vdash A\}$  is definable in arithmetic.

## CHAPTER 19

# *Models of Arithmetic*

### 19.1 Introduction

The *standard model* of arithmetic is the structure  $N$  with  $|N| = \mathbb{N}$  in which  $0$ ,  $\iota$ ,  $+$ ,  $\times$ , and  $<$  are interpreted as you would expect. That is,  $0$  is  $0$ ,  $\iota$  is the successor function,  $+$  is interpreted as addition and  $\times$  as multiplication of the numbers in  $\mathbb{N}$ . Specifically,

$$\begin{aligned}0^N &= 0 \\ \iota^N(n) &= n + 1 \\ +^N(n, m) &= n + m \\ \times^N(n, m) &= nm\end{aligned}$$

Of course, there are structures for  $\mathcal{L}_A$  that have domains other than  $\mathbb{N}$ . For instance, we can take  $M$  with domain  $|M| = \{a\}^*$  (the finite sequences of the single symbol  $a$ , i.e.,  $\emptyset$ ,  $a$ ,  $aa$ ,  $aaa$ , ...), and interpretations

$$\begin{aligned}0^M &= \emptyset \\ \iota^M(s) &= s \frown a \\ +^M(n, m) &= a^{n+m}\end{aligned}$$

$$\times^M(n, m) = a^{nm}$$

These two structures are “essentially the same” in the sense that the only difference is the elements of the domains but not how the elements of the domains are related among each other by the interpretation functions. We say that the two structures are *isomorphic*.

It is an easy consequence of the compactness theorem that any theory true in  $N$  also has models that are not isomorphic to  $N$ . Such structures are called *non-standard*. The interesting thing about them is that while the elements of a standard model (i.e.,  $N$ , but also all structures isomorphic to it) are exhausted by the values of the standard numerals  $\bar{n}$ , i.e.,

$$|N| = \{\text{Val}^N(\bar{n}) : n \in \mathbb{N}\}$$

that isn't the case in non-standard models: if  $M$  is non-standard, then there is at least one  $x \in |M|$  such that  $x \neq \text{Val}^M(\bar{n})$  for all  $n$ .

These non-standard elements are pretty neat: they are “infinite natural numbers.” But their existence also explains, in a sense, the incompleteness phenomena. Consider an example, e.g., the consistency statement for Peano arithmetic,  $\text{Con}_{\mathbf{PA}}$ , i.e.,  $\neg \exists x \text{Prf}_{\mathbf{PA}}(x, \ulcorner \perp \urcorner)$ . Since  $\mathbf{PA}$  neither proves  $\text{Con}_{\mathbf{PA}}$  nor  $\neg \text{Con}_{\mathbf{PA}}$ , either can be consistently added to  $\mathbf{PA}$ . Since  $\mathbf{PA}$  is consistent,  $N \models \text{Con}_{\mathbf{PA}}$ , and consequently  $N \not\models \neg \text{Con}_{\mathbf{PA}}$ . So  $N$  is *not* a model of  $\mathbf{PA} \cup \{\neg \text{Con}_{\mathbf{PA}}\}$ , and all its models must be nonstandard. Models of  $\mathbf{PA} \cup \{\neg \text{Con}_{\mathbf{PA}}\}$  must contain some element that serves as the witness that makes  $\exists x \text{Prf}_{\mathbf{PA}}(\ulcorner \perp \urcorner)$  true, i.e., a Gödel number of a derivation of a contradiction from  $\mathbf{PA}$ . Such an element can't be standard—since  $\mathbf{PA} \vdash \neg \text{Prf}_{\mathbf{PA}}(\bar{n}, \ulcorner \perp \urcorner)$  for every  $n$ .

## 19.2 Reducts and Expansions

Often it is useful or necessary to compare languages which have symbols in common, as well as structures for these languages. The most common case is when all the symbols in a language  $\mathcal{L}$

are also part of a language  $\mathcal{L}'$ , i.e.,  $\mathcal{L} \subseteq \mathcal{L}'$ . An  $\mathcal{L}$ -structure  $M$  can then always be expanded to an  $\mathcal{L}'$ -structure by adding interpretations of the additional symbols while leaving the interpretations of the common symbols the same. On the other hand, from an  $\mathcal{L}'$ -structure  $M'$  we can obtain an  $\mathcal{L}$ -structure simply by “forgetting” the interpretations of the symbols that do not occur in  $\mathcal{L}$ .

**Definition 19.1.** Suppose  $\mathcal{L} \subseteq \mathcal{L}'$ ,  $M$  is an  $\mathcal{L}$ -structure and  $M'$  is an  $\mathcal{L}'$ -structure.  $M$  is the *reduct* of  $M'$  to  $\mathcal{L}$ , and  $M'$  is an *expansion* of  $M$  to  $\mathcal{L}'$  iff

1.  $|M| = |M'|$
2. For every constant symbol  $c \in \mathcal{L}$ ,  $c^M = c^{M'}$ .
3. For every function symbol  $f \in \mathcal{L}$ ,  $f^M = f^{M'}$ .
4. For every predicate symbol  $P \in \mathcal{L}$ ,  $P^M = P^{M'}$ .

**Proposition 19.2.** *If an  $\mathcal{L}$ -structure  $M$  is a reduct of an  $\mathcal{L}'$ -structure  $M'$ , then for all  $\mathcal{L}$ -sentences  $A$ ,*

$$M \models A \text{ iff } M' \models A.$$

*Proof.* Exercise. □

**Definition 19.3.** When we have an  $\mathcal{L}$ -structure  $M$ , and  $\mathcal{L}' = \mathcal{L} \cup \{P\}$  is the expansion of  $\mathcal{L}$  obtained by adding a single  $n$ -place predicate symbol  $P$ , and  $R \subseteq |M|^n$  is an  $n$ -place relation, then we write  $(M, R)$  for the expansion  $M'$  of  $M$  with  $P^{M'} = R$ .

### 19.3 Isomorphic Structures

First-order structures can be alike in one of two ways. One way in which they can be alike is that they make the same sentences



true. We call such structures *elementarily equivalent*. But structures can be very different and still make the same sentences true—for instance, one can be countable and the other not. This is because there are lots of features of a structure that cannot be expressed in first-order languages, either because the language is not rich enough, or because of fundamental limitations of first-order logic such as the Löwenheim–Skolem theorem. So another, stricter, aspect in which structures can be alike is if they are fundamentally the same, in the sense that they only differ in the objects that make them up, but not in their structural features. A way of making this precise is by the notion of an *isomorphism*.

**Definition 19.4.** Given two structures  $M$  and  $M'$  for the same language  $\mathcal{L}$ , we say that  $M$  is *elementarily equivalent to  $M'$* , written  $M \equiv M'$ , if and only if for every sentence  $A$  of  $\mathcal{L}$ ,  $M \models A$  iff  $M' \models A$ .

**Definition 19.5.** Given two structures  $M$  and  $M'$  for the same language  $\mathcal{L}$ , we say that  $M$  is *isomorphic to  $M'$* , written  $M \simeq M'$ , if and only if there is a function  $h: |M| \rightarrow |M'|$  such that:

1.  $h$  is injective: if  $h(x) = h(y)$  then  $x = y$ ;
2.  $h$  is surjective: for every  $y \in |M'|$  there is  $x \in |M|$  such that  $h(x) = y$ ;
3. for every constant symbol  $c$ :  $h(c^M) = c^{M'}$ ;
4. for every  $n$ -place predicate symbol  $P$ :

$$\langle a_1, \dots, a_n \rangle \in P^M \quad \text{iff} \quad \langle h(a_1), \dots, h(a_n) \rangle \in P^{M'};$$

5. for every  $n$ -place function symbol  $f$ :

$$h(f^M(a_1, \dots, a_n)) = f^{M'}(h(a_1), \dots, h(a_n)).$$

**Theorem 19.6.** *If  $M \simeq M'$  then  $M \equiv M'$ .*

*Proof.* Let  $h$  be an isomorphism of  $M$  onto  $M'$ . For any assignment  $s$ ,  $h \circ s$  is the composition of  $h$  and  $s$ , i.e., the assignment in  $M'$  such that  $(h \circ s)(x) = h(s(x))$ . By induction on  $t$  and  $A$  one can prove the stronger claims:

- a.  $h(\text{Val}_s^M(t)) = \text{Val}_{h \circ s}^{M'}(t)$ .
- b.  $M, s \models A$  iff  $M', h \circ s \models A$ .

The first is proved by induction on the complexity of  $t$ .

1. If  $t \equiv c$ , then  $\text{Val}_s^M(c) = c^M$  and  $\text{Val}_{h \circ s}^{M'}(c) = c^{M'}$ . Thus,  $h(\text{Val}_s^M(t)) = h(c^M) = c^{M'}$  (by (3) of Definition 19.5) =  $\text{Val}_{h \circ s}^{M'}(t)$ .
2. If  $t \equiv x$ , then  $\text{Val}_s^M(x) = s(x)$  and  $\text{Val}_{h \circ s}^{M'}(x) = h(s(x))$ . Thus,  $h(\text{Val}_s^M(x)) = h(s(x)) = \text{Val}_{h \circ s}^{M'}(x)$ .
3. If  $t \equiv f(t_1, \dots, t_n)$ , then

$$\begin{aligned} \text{Val}_s^M(t) &= f^M(\text{Val}_s^M(t_1), \dots, \text{Val}_s^M(t_n)) \quad \text{and} \\ \text{Val}_{h \circ s}^{M'}(t) &= f^{M'}(\text{Val}_{h \circ s}^{M'}(t_1), \dots, \text{Val}_{h \circ s}^{M'}(t_n)). \end{aligned}$$

The induction hypothesis is that for each  $i$ ,  $h(\text{Val}_s^M(t_i)) = \text{Val}_{h \circ s}^{M'}(t_i)$ . So,

$$\begin{aligned} h(\text{Val}_s^M(t)) &= h(f^M(\text{Val}_s^M(t_1), \dots, \text{Val}_s^M(t_n))) \\ &= h(f^M(\text{Val}_{h \circ s}^{M'}(t_1), \dots, \text{Val}_{h \circ s}^{M'}(t_n))) \quad (19.1) \end{aligned}$$

$$\begin{aligned} &= f^{M'}(\text{Val}_{h \circ s}^{M'}(t_1), \dots, \text{Val}_{h \circ s}^{M'}(t_n)) \quad (19.2) \\ &= \text{Val}_{h \circ s}^{M'}(t) \end{aligned}$$

Here, eq. (19.1) follows by induction hypothesis and eq. (19.2) by (5) of Definition 19.5.

Part (b) is left as an exercise.

If  $A$  is a sentence, the assignments  $s$  and  $h \circ s$  are irrelevant, and we have  $M \models A$  iff  $M' \models A$ .  $\square$

**Definition 19.7.** An *automorphism* of a structure  $\mathfrak{M}$  is an isomorphism of  $\mathfrak{M}$  onto itself.

## 19.4 The Theory of a Structure

Every structure  $M$  makes some sentences true, and some false. The set of all the sentences it makes true is called its *theory*. That set is in fact a theory, since anything it entails must be true in all its models, including  $M$ .

**Definition 19.8.** Given a structure  $M$ , the *theory* of  $M$  is the set  $\text{Th}(M)$  of sentences that are true in  $M$ , i.e.,  $\text{Th}(M) = \{A : M \models A\}$ .

We also use the term “theory” informally to refer to sets of sentences having an intended interpretation, whether deductively closed or not.

**Proposition 19.9.** For any  $M$ ,  $\text{Th}(M)$  is complete.

*Proof.* For any sentence  $A$  either  $M \models A$  or  $M \models \neg A$ , so either  $A \in \text{Th}(M)$  or  $\neg A \in \text{Th}(M)$ .  $\square$

**Proposition 19.10.** If  $N \models A$  for every  $A \in \text{Th}(M)$ , then  $M \equiv N$ .

*Proof.* Since  $N \models A$  for all  $A \in \text{Th}(M)$ ,  $\text{Th}(M) \subseteq \text{Th}(N)$ . If  $N \models A$ , then  $N \not\models \neg A$ , so  $\neg A \notin \text{Th}(M)$ . Since  $\text{Th}(M)$  is complete,  $A \in \text{Th}(M)$ . So,  $\text{Th}(N) \subseteq \text{Th}(M)$ , and we have  $M \equiv N$ .  $\square$

*Remark 1.* Consider  $\mathbf{R} = \langle \mathbb{R}, < \rangle$ , the structure whose domain is the set  $\mathbb{R}$  of the real numbers, in the language comprising only a 2-place predicate symbol interpreted as the  $<$  relation over the reals. Clearly  $\mathbf{R}$  is uncountable; however, since  $\text{Th}(\mathbf{R})$  is obviously consistent, by the Löwenheim–Skolem theorem it has a countable model, say  $S$ , and by **Proposition 19.10**,  $\mathbf{R} \equiv S$ . Moreover, since  $\mathbf{R}$  and  $S$  are not isomorphic, this shows that the converse of **Theorem 19.6** fails in general.

## 19.5 Standard Models of Arithmetic

The language of arithmetic  $\mathcal{L}_A$  is obviously intended to be about numbers, specifically, about natural numbers. So, “the” standard model  $N$  is special: it is the model we want to talk about. But in logic, we are often just interested in structural properties, and any two structures that are isomorphic share those. So we can be a bit more liberal, and consider any structure that is isomorphic to  $N$  “standard.”

**Definition 19.11.** A structure for  $\mathcal{L}_A$  is *standard* if it is isomorphic to  $N$ .

**Proposition 19.12.** *If a structure  $M$  is standard, then its domain is the set of values of the standard numerals, i.e.,*

$$|M| = \{\text{Val}^M(\bar{n}) : n \in \mathbb{N}\}$$

*Proof.* Clearly, every  $\text{Val}^M(\bar{n}) \in |M|$ . We just have to show that every  $x \in |M|$  is equal to  $\text{Val}^M(\bar{n})$  for some  $n$ . Since  $M$  is standard, it is isomorphic to  $N$ . Suppose  $g: \mathbb{N} \rightarrow |M|$  is an isomorphism. Then  $g(n) = g(\text{Val}^N(\bar{n})) = \text{Val}^M(\bar{n})$ . But for every  $x \in |M|$ , there is an  $n \in \mathbb{N}$  such that  $g(n) = x$ , since  $g$  is surjective.  $\square$

If a structure  $M$  for  $\mathcal{L}_A$  is standard, the elements of its domain can all be named by the standard numerals  $\bar{0}, \bar{1}, \bar{2}, \dots$ , i.e., the terms  $0, 0', 0'', \dots$ . Of course, this does not mean that the elements of  $|M|$  are the numbers, just that we can pick them out the same way we can pick out the numbers in  $|N|$ .

**Proposition 19.13.** *If  $M \models \mathbf{Q}$ , and  $|M| = \{\text{Val}^M(\bar{n}) : n \in \mathbb{N}\}$ , then  $M$  is standard.*

*Proof.* We have to show that  $M$  is isomorphic to  $N$ . Consider the function  $g: \mathbb{N} \rightarrow |M|$  defined by  $g(n) = \text{Val}^M(\bar{n})$ . By the hypothesis,  $g$  is surjective. It is also injective:  $\mathbf{Q} \vdash \bar{n} \neq \bar{m}$  whenever

$n \neq m$ . Thus, since  $M \models \mathbf{Q}$ ,  $M \models \bar{n} \neq \bar{m}$ , whenever  $n \neq m$ . Thus, if  $n \neq m$ , then  $\text{Val}^M(\bar{n}) \neq \text{Val}^M(\bar{m})$ , i.e.,  $g(n) \neq g(m)$ .

We also have to verify that  $g$  is an isomorphism.

1. We have  $g(o^N) = g(0)$  since,  $o^N = 0$ . By definition of  $g$ ,  $g(0) = \text{Val}^M(\bar{0})$ . But  $\bar{0}$  is just  $o$ , and the value of a term which happens to be a constant symbol is given by what the structure assigns to that constant symbol, i.e.,  $\text{Val}^M(o) = o^M$ . So we have  $g(o^N) = o^M$  as required.
2.  $g(\iota^N(n)) = g(n+1)$ , since  $\iota$  in  $N$  is the successor function on  $\mathbb{N}$ . Then,  $g(n+1) = \text{Val}^M(\overline{n+1})$  by definition of  $g$ . But  $\overline{n+1}$  is the same term as  $\bar{n}'$ , so  $\text{Val}^M(\overline{n+1}) = \text{Val}^M(\bar{n}')$ . By the definition of the value function, this is  $\iota^M(\text{Val}^M(\bar{n}))$ . Since  $\text{Val}^M(\bar{n}) = g(n)$  we get  $g(\iota^N(n)) = \iota^M(g(n))$ .
3.  $g(+^N(n, m)) = g(n+m)$ , since  $+$  in  $N$  is the addition function on  $\mathbb{N}$ . Then,  $g(n+m) = \text{Val}^M(\overline{n+m})$  by definition of  $g$ . But  $\mathbf{Q} \vdash \overline{n+m} = \overline{(n+m)}$ , so  $\text{Val}^M(\overline{n+m}) = \text{Val}^M(\overline{n+m})$ . By the definition of the value function, this is  $+^M(\text{Val}^M(\bar{n}), \text{Val}^M(\bar{m}))$ . Since  $\text{Val}^M(\bar{n}) = g(n)$  and  $\text{Val}^M(\bar{m}) = g(m)$ , we get  $g(+^N(n, m)) = +^M(g(n), g(m))$ .
4.  $g(\times^N(n, m)) = \times^M(g(n), g(m))$ : Exercise.
5.  $\langle n, m \rangle \in <^N$  iff  $n < m$ . If  $n < m$ , then  $\mathbf{Q} \vdash \bar{n} < \bar{m}$ , and also  $M \models \bar{n} < \bar{m}$ . Thus  $\langle \text{Val}^M(\bar{n}), \text{Val}^M(\bar{m}) \rangle \in <^M$ , i.e.,  $\langle g(n), g(m) \rangle \in <^M$ . If  $n \not< m$ , then  $\mathbf{Q} \vdash \neg \bar{n} < \bar{m}$ , and consequently  $M \not\models \bar{n} < \bar{m}$ . Thus, as before,  $\langle g(n), g(m) \rangle \notin <^M$ . Together, we get:  $\langle n, m \rangle \in <^N$  iff  $\langle g(n), g(m) \rangle \in <^M$ .

□

The function  $g$  is the most obvious way of defining a mapping from  $\mathbb{N}$  to the domain of any other structure  $M$  for  $\mathcal{L}_A$ , since every such  $M$  contains elements named by  $\bar{0}, \bar{1}, \bar{2}$ , etc. So it isn't surprising that if  $M$  makes at least some basic statements about the  $\bar{n}$ 's true in the same way that  $N$  does, and  $g$  is also bijective,

then  $g$  will turn into an isomorphism. In fact, if  $|M|$  contains no elements other than what the  $\bar{n}$ 's name, it's the only one.

**Proposition 19.14.** *If  $M$  is standard, then  $g$  from the proof of Proposition 19.13 is the only isomorphism from  $N$  to  $M$ .*

*Proof.* Suppose  $h: \mathbb{N} \rightarrow |M|$  is an isomorphism between  $N$  and  $M$ . We show that  $g = h$  by induction on  $n$ . If  $n = 0$ , then  $g(0) = 0^M$  by definition of  $g$ . But since  $h$  is an isomorphism,  $h(0) = h(0^N) = 0^M$ , so  $g(0) = h(0)$ .

Now consider the case for  $n + 1$ . We have

$$\begin{aligned}
 g(n+1) &= \text{Val}^M(\overline{n+1}) \text{ by definition of } g \\
 &= \text{Val}^M(\bar{n}') \text{ since } \overline{n+1} \equiv \bar{n}' \\
 &= r^M(\text{Val}^M(\bar{n})) \text{ by definition of } \text{Val}^M(t') \\
 &= r^M(g(n)) \text{ by definition of } g \\
 &= r^M(h(n)) \text{ by induction hypothesis} \\
 &= h(r^N(n)) \text{ since } h \text{ is an isomorphism} \\
 &= h(n+1)
 \end{aligned}
 \quad \square$$

For any countably infinite set  $M$ , there's a bijection between  $\mathbb{N}$  and  $M$ , so every such set  $M$  is potentially the domain of a standard model  $M$ . In fact, once you pick an object  $z \in M$  and a suitable function  $s$  as  $0^M$  and  $r^M$ , the interpretations of  $+$ ,  $\times$ , and  $<$  is already fixed. Only functions  $s: M \rightarrow M \setminus \{z\}$  that are both injective and surjective are suitable in a standard model as  $r^M$ . The range of  $s$  cannot contain  $z$ , since otherwise  $\forall x 0 \neq x'$  would be false. That sentence is true in  $N$ , and so  $M$  also has to make it true. The function  $s$  has to be injective, since the successor function  $r^N$  in  $N$  is, and that  $r^N$  is injective is expressed by a sentence true in  $N$ . It has to be surjective because otherwise there would be some  $x \in M \setminus \{z\}$  not in the domain of  $s$ , i.e., the sentence  $\forall x (x = 0 \vee \exists y y' = x)$  would be false in  $M$ —but it is true in  $N$ .

## 19.6 Non-Standard Models

We call a structure for  $\mathcal{L}_A$  standard if it is isomorphic to  $N$ . If a structure isn't isomorphic to  $N$ , it is called non-standard.

**Definition 19.15.** A structure  $M$  for  $\mathcal{L}_A$  is *non-standard* if it is not isomorphic to  $N$ . The elements  $x \in |M|$  which are equal to  $\text{Val}^M(\bar{n})$  for some  $n \in \mathbb{N}$  are called *standard numbers* (of  $M$ ), and those not, *non-standard numbers*.

By **Proposition 19.12**, any standard structure for  $\mathcal{L}_A$  contains only standard elements. Consequently, a non-standard structure must contain at least one non-standard element. In fact, the existence of a non-standard element guarantees that the structure is non-standard.

**Proposition 19.16.** *If a structure  $M$  for  $\mathcal{L}_A$  contains a non-standard number,  $M$  is non-standard.*

*Proof.* Suppose not, i.e., suppose  $M$  standard but contains a non-standard number  $x$ . Let  $g: \mathbb{N} \rightarrow |M|$  be an isomorphism. It is easy to see (by induction on  $n$ ) that  $g(\text{Val}^N(\bar{n})) = \text{Val}^M(\bar{n})$ . In other words,  $g$  maps standard numbers of  $N$  to standard numbers of  $M$ . If  $M$  contains a non-standard number,  $g$  cannot be surjective, contrary to hypothesis.  $\square$

It is easy enough to specify non-standard structures for  $\mathcal{L}_A$ . For instance, take the structure with domain  $\mathbb{Z}$  and interpret all non-logical symbols as usual. Since negative numbers are not values of  $\bar{n}$  for any  $n$ , this structure is non-standard. Of course, it will not be a *model* of arithmetic in the sense that it makes the same sentences true as  $N$ . For instance,  $\forall x x' \neq 0$  is false. However, we can prove that non-standard models of arithmetic exist easily enough, using the compactness theorem.

**Proposition 19.17.** *Let  $\mathbf{TA} = \{A : N \vDash A\}$  be the theory of  $N$ .  $\mathbf{TA}$  has a countable non-standard model.*

*Proof.* Expand  $\mathcal{L}_A$  by a new constant symbol  $c$  and consider the set of sentences

$$\Gamma = \mathbf{TA} \cup \{c \neq \bar{0}, c \neq \bar{1}, c \neq \bar{2}, \dots\}$$

Any model  $M^c$  of  $\Gamma$  would contain an element  $x = c^M$  which is non-standard, since  $x \neq \text{Val}^M(\bar{n})$  for all  $n \in \mathbb{N}$ . Also, obviously,  $M^c \vDash \mathbf{TA}$ , since  $\mathbf{TA} \subseteq \Gamma$ . If we turn  $M^c$  into a structure  $M$  for  $\mathcal{L}_A$  simply by forgetting about  $c$ , its domain still contains the non-standard  $x$ , and also  $M \vDash \mathbf{TA}$ . The latter is guaranteed since  $c$  does not occur in  $\mathbf{TA}$ . So, it suffices to show that  $\Gamma$  has a model.

We use the compactness theorem to show that  $\Gamma$  has a model. If every finite subset of  $\Gamma$  is satisfiable, so is  $\Gamma$ . Consider any finite subset  $\Gamma_0 \subseteq \Gamma$ .  $\Gamma_0$  includes some sentences of  $\mathbf{TA}$  and some of the form  $c \neq \bar{n}$ , but only finitely many. Suppose  $k$  is the largest number so that  $c \neq \bar{k} \in \Gamma_0$ . Define  $N_k$  by expanding  $N$  to include the interpretation  $c^{N_k} = k + 1$ .  $N_k \vDash \Gamma_0$ : if  $A \in \mathbf{TA}$ ,  $N_k \vDash A$  since  $N_k$  is just like  $N$  in all respects except  $c$ , and  $c$  does not occur in  $A$ . And  $N_k \vDash c \neq \bar{n}$ , since  $n \leq k$ , and  $\text{Val}^{N_k}(c) = k + 1$ . Thus, every finite subset of  $\Gamma$  is satisfiable.  $\square$

## 19.7 Models of $\mathbf{Q}$

We know that there are non-standard structures that make the same sentences true as  $N$  does, i.e., is a model of  $\mathbf{TA}$ . Since  $N \vDash \mathbf{Q}$ , any model of  $\mathbf{TA}$  is also a model of  $\mathbf{Q}$ .  $\mathbf{Q}$  is much weaker than  $\mathbf{TA}$ , e.g.,  $\mathbf{Q} \not\vDash \forall x \forall y (x+y) = (y+x)$ . Weaker theories are easier to satisfy: they have more models. E.g.,  $\mathbf{Q}$  has models which make  $\forall x \forall y (x+y) = (y+x)$  false, but those cannot also be models of  $\mathbf{TA}$ , or  $\mathbf{PA}$  for that matter. Models of  $\mathbf{Q}$  are also relatively simple: we can specify them explicitly.



**Example 19.18.** Consider the structure  $K$  with domain  $|K| = \mathbb{N} \cup \{a\}$  and interpretations

$$\begin{aligned} 0^K &= 0 \\ \iota^K(x) &= \begin{cases} x+1 & \text{if } x \in \mathbb{N} \\ a & \text{if } x = a \end{cases} \\ +^K(x, y) &= \begin{cases} x+y & \text{if } x, y \in \mathbb{N} \\ a & \text{otherwise} \end{cases} \\ \times^K(x, y) &= \begin{cases} xy & \text{if } x, y \in \mathbb{N} \\ 0 & \text{if } x = 0 \text{ or } y = 0 \\ a & \text{otherwise} \end{cases} \\ <^K &= \{\langle x, y \rangle : x, y \in \mathbb{N} \text{ and } x < y\} \cup \{\langle x, a \rangle : x \in |K|\} \end{aligned}$$

To show that  $K \models \mathbf{Q}$  we have to verify that all axioms of  $\mathbf{Q}$  are true in  $K$ . For convenience, let's write  $x^*$  for  $\iota^K(x)$  (the “successor” of  $x$  in  $K$ ),  $x \oplus y$  for  $+^K(x, y)$  (the “sum” of  $x$  and  $y$  in  $K$ ),  $x \otimes y$  for  $\times^K(x, y)$  (the “product” of  $x$  and  $y$  in  $K$ ), and  $x \odot y$  for  $\langle x, y \rangle \in <^K$ . With these abbreviations, we can give the operations in  $K$  more perspicuously as

$x$	$x^*$	$x \oplus y$	$0$	$m$	$a$	$x \otimes y$	$0$	$m$	$a$
$n$	$n+1$	$0$	$0$	$m$	$a$	$0$	$0$	$0$	$0$
$a$	$a$	$n$	$n$	$n+m$	$a$	$n$	$0$	$nm$	$a$
		$a$	$a$	$a$	$a$	$a$	$0$	$a$	$a$

We have  $n \odot m$  iff  $n < m$  for  $n, m \in \mathbb{N}$  and  $x \odot a$  for all  $x \in |K|$ .

$K \models \forall x \forall y (x' = y' \rightarrow x = y)$  since  $*$  is injective.  $K \models \forall x 0 \neq x'$  since  $0$  is not a  $*$ -successor in  $K$ .  $K \models \forall x (x = 0 \vee \exists y x = y')$  since for every  $n > 0$ ,  $n = (n-1)^*$ , and  $a = a^*$ .

$K \models \forall x (x + 0) = x$  since  $n \oplus 0 = n + 0 = n$ , and  $a \oplus 0 = a$  by definition of  $\oplus$ .  $K \models \forall x \forall y (x + y') = (x + y)'$  is a bit trickier. If  $n, m$  are both standard, we have:

$$(n \oplus m^*) = (n + (m + 1)) = (n + m) + 1 = (n \oplus m)^*$$

since  $\oplus$  and  $*$  agree with  $+$  and  $\prime$  on standard numbers. Now suppose  $x \in |K|$ . Then

$$(x \oplus a^*) = (x \oplus a) = a = a^* = (x \oplus a)^*$$

The remaining case is if  $y \in |K|$  but  $x = a$ . Here we also have to distinguish cases according to whether  $y = n$  is standard or  $y = b$ :

$$(a \oplus n^*) = (a \oplus (n+1)) = a = a^* = (a \oplus n)^*$$

$$(a \oplus a^*) = (a \oplus a) = a = a^* = (a \oplus a)^*$$

This is of course a bit more detailed than needed. For instance, since  $a \oplus z = a$  whatever  $z$  is, we can immediately conclude  $a \oplus a^* = a$ . The remaining axioms can be verified the same way.

$K$  is thus a model of  $\mathbf{Q}$ . Its “addition”  $\oplus$  is also commutative. But there are other sentences true in  $N$  but false in  $K$ , and vice versa. For instance,  $a \otimes a$ , so  $K \models \exists x x < x$  and  $K \not\models \forall x \neg x < x$ . This shows that  $\mathbf{Q} \not\models \forall x \neg x < x$ .

**Example 19.19.** Consider the structure  $L$  with domain  $|L| = \mathbb{N} \cup \{a, b\}$  and interpretations  $\prime^L = *$ ,  $+^L = \oplus$  given by

$x$	$x^*$	$x \oplus y$	$m$	$a$	$b$
$n$	$n+1$	$n$	$n+m$	$b$	$a$
$a$	$a$	$a$	$a$	$b$	$a$
$b$	$b$	$b$	$b$	$b$	$a$

Since  $*$  is injective, 0 is not in its range, and every  $x \in |L|$  other than 0 is, axioms  $Q_1$ – $Q_3$  are true in  $L$ . For any  $x$ ,  $x \oplus 0 = x$ , so  $Q_4$  is true as well. For  $Q_5$ , consider  $x \oplus y^*$  and  $(x \oplus y)^*$ . They are equal if  $x$  and  $y$  are both standard, since then  $*$  and  $\oplus$  agree with  $\prime$  and  $+$ . If  $x$  is non-standard, and  $y$  is standard, we have  $x \oplus y^* = x = x^* = (x \oplus y)^*$ . If  $x$  and  $y$  are both non-standard, we have four cases:

$$a \oplus a^* = b = b^* = (a \oplus a)^*$$

$$b \oplus b^* = a = a^* = (b \oplus b)^*$$

$$b \oplus a^* = b = b^* = (b \oplus y)^*$$

$$a \oplus b^* = a = a^* = (a \oplus b)^*$$

If  $x$  is standard, but  $y$  is non-standard, we have

$$n \oplus a^* = n \oplus a = b = b^* = (n \oplus a)^*$$

$$n \oplus b^* = n \oplus b = a = a^* = (n \oplus b)^*$$

So,  $L \models Q_5$ . However,  $a \oplus 0 \neq 0 \oplus a$ , so  $L \not\models \forall x \forall y (x + y) = (y + x)$ .

We've explicitly constructed models of  $\mathbf{Q}$  in which the non-standard elements live “beyond” the standard elements. In fact, that much is required by the axioms. A non-standard element  $x$  cannot be  $\otimes 0$ , since  $\mathbf{Q} \vdash \forall x \neg x < 0$  (see [Lemma 17.23](#)). Also, for every  $n$ ,  $\mathbf{Q} \vdash \forall x (x < \bar{n}' \rightarrow (x = \bar{0} \vee x = \bar{1} \vee \dots \vee x = \bar{n}))$  ([Lemma 17.24](#)), so we can't have  $a \otimes n$  for any  $n > 0$ .

## 19.8 Models of PA

Any non-standard model of  $\mathbf{TA}$  is also one of  $\mathbf{PA}$ . We know that non-standard models of  $\mathbf{TA}$  and hence of  $\mathbf{PA}$  exist. We also know that such non-standard models contain non-standard “numbers,” i.e., elements of the domain that are “beyond” all the standard “numbers.” But how are they arranged? How many are there? We've seen that models of the weaker theory  $\mathbf{Q}$  can contain as few as a single non-standard number. But these simple structures are not models of  $\mathbf{PA}$  or  $\mathbf{TA}$ .

The key to understanding the structure of models of  $\mathbf{PA}$  or  $\mathbf{TA}$  is to see what facts are derivable in these theories. For instance, already  $\mathbf{PA}$  proves that  $\forall x x \neq x'$  and  $\forall x \forall y (x + y) = (y + x)$ , so this rules out simple structures (in which these sentences are false) as models of  $\mathbf{PA}$ .

Suppose  $M$  is a model of  $\mathbf{PA}$ . Then if  $\mathbf{PA} \vdash A$ ,  $M \models A$ . Let's again use  $\mathbf{z}$  for  $o^M$ ,  $*$  for  $\iota^M$ ,  $\oplus$  for  $+^M$ ,  $\otimes$  for  $\times^M$ , and  $\otimes$  for  $<^M$ . Any sentence  $A$  then states some condition about  $\mathbf{z}$ ,  $*$ ,  $\oplus$ ,  $\otimes$ , and

$\otimes$ , and if  $M \models A$  that condition must be satisfied. For instance, if  $M \models Q_1$ , i.e.,  $M \models \forall x \forall y (x' = y' \rightarrow x = y)$ , then  $*$  must be injective.

**Proposition 19.20.** *In  $M$ ,  $\otimes$  is a linear strict order, i.e., it satisfies:*

1. Not  $x \otimes x$  for any  $x \in |M|$ .
2. If  $x \otimes y$  and  $y \otimes z$  then  $x \otimes z$ .
3. For any  $x \neq y$ ,  $x \otimes y$  or  $y \otimes x$

*Proof.* **PA** proves:

1.  $\forall x \neg x < x$
2.  $\forall x \forall y \forall z ((x < y \wedge y < z) \rightarrow x < z)$
3.  $\forall x \forall y ((x < y \vee y < x) \vee x = y)$  □

**Proposition 19.21.**  *$\mathbf{z}$  is the least element of  $|M|$  in the  $\otimes$ -ordering. For any  $x$ ,  $x \otimes x^*$ , and  $x^*$  is the  $\otimes$ -least element with that property. For any  $x$ , there is a unique  $y$  such that  $y^* = x$ . (We call  $y$  the “predecessor” of  $x$  in  $M$ , and denote it by  ${}^*x$ .)*

*Proof.* Exercise. □

**Proposition 19.22.** *All standard elements of  $M$  are less than (according to  $\otimes$ ) all non-standard elements.*

*Proof.* We'll use  $n$  as short for  $\text{Val}^M(\bar{n})$ , a standard element of  $M$ . Already **Q** proves that, for any  $n \in \mathbb{N}$ ,  $\forall x (x < \bar{n}' \rightarrow (x = \bar{0} \vee x = \bar{1} \vee \dots \vee x = \bar{n}))$ . There are no elements that are  $\otimes \mathbf{z}$ . So if  $n$  is standard and  $x$  is non-standard, we cannot have  $x \otimes n$ . By definition, a non-standard element is one that isn't  $\text{Val}^M(\bar{n})$  for any  $n \in \mathbb{N}$ , so  $x \neq n$  as well. Since  $\otimes$  is a linear order, we must have  $n \otimes x$ . □

**Proposition 19.23.** *Every nonstandard element  $x$  of  $|M|$  is an element of the subset*

$$\dots^{***} x \ominus^{**} x \ominus^* x \ominus x \ominus x^* \ominus x^{**} \ominus x^{***} \ominus \dots$$

*We call this subset the block of  $x$  and write it as  $[x]$ . It has no least and no greatest element. It can be characterized as the set of those  $y \in |M|$  such that, for some standard  $n$ ,  $x \oplus n = y$  or  $y \oplus n = x$ .*

*Proof.* Clearly, such a set  $[x]$  always exists since every element  $y$  of  $|M|$  has a unique successor  $y^*$  and unique predecessor  $^*y$ . For successive elements  $y, y^*$  we have  $y \ominus y^*$  and  $y^*$  is the  $\ominus$ -least element of  $|M|$  such that  $y$  is  $\ominus$ -less than it. Since always  $^*y \ominus y$  and  $y \ominus y^*$ ,  $[x]$  has no least or greatest element. If  $y \in [x]$  then  $x \in [y]$ , for then either  $y^{*...*} = x$  or  $x^{*...*} = y$ . If  $y^{*...*} = x$  (with  $n$   $*$ 's), then  $y \oplus n = x$  and conversely, since  $\mathbf{PA} \vdash \forall x x^{*'\dots'} = (x + \bar{n})$  (if  $n$  is the number of  $'$ 's).  $\square$

**Proposition 19.24.** *If  $[x] \neq [y]$  and  $x \ominus y$ , then for any  $u \in [x]$  and any  $v \in [y]$ ,  $u \ominus v$ .*

*Proof.* Note that  $\mathbf{PA} \vdash \forall x \forall y (x < y \rightarrow (x' < y \vee x' = y))$ . Thus, if  $u \ominus v$ , we also have  $u \oplus n^* \ominus v$  for any  $n$  if  $[u] \neq [v]$ .

Any  $u \in [x]$  is  $\ominus y$ :  $x \ominus y$  by assumption. If  $u \ominus x$ ,  $u \ominus y$  by transitivity. And if  $x \ominus u$  but  $u \in [x]$ , we have  $u = x \oplus n^*$  for some  $n$ , and so  $u \ominus y$  by the fact just proved.

Now suppose that  $v \in [y]$  is  $\ominus y$ , i.e.,  $v \oplus m^* = y$  for some standard  $m$ . This rules out  $v \ominus x$ , otherwise  $y = v \oplus m^* \ominus x$ . Clearly also,  $x \neq v$ , otherwise  $x \oplus m^* = v \oplus m^* = y$  and we would have  $[x] = [y]$ . So,  $x \ominus v$ . But then also  $x \oplus n^* \ominus v$  for any  $n$ . Hence, if  $x \ominus u$  and  $u \in [x]$ , we have  $u \ominus v$ . If  $u \ominus x$  then  $u \ominus v$  by transitivity.

Lastly, if  $y \ominus v$ ,  $u \ominus v$  since, as we've shown,  $u \ominus y$  and  $y \ominus v$ .  $\square$

**Corollary 19.25.** *If  $[x] \neq [y]$ ,  $[x] \cap [y] = \emptyset$ .*

*Proof.* Suppose  $z \in [x]$  and  $x \otimes y$ . Then  $z \otimes u$  for all  $u \in [y]$ . If  $z \in [y]$ , we would have  $z \otimes z$ . Similarly if  $y \otimes x$ .  $\square$

This means that the blocks themselves can be ordered in a way that respects  $\otimes$ :  $[x] \otimes [y]$  iff  $x \otimes y$ , or, equivalently, if  $u \otimes v$  for any  $u \in [x]$  and  $v \in [y]$ . Clearly, the standard block  $[0]$  is the least block. It intersects with no non-standard block, and no two non-standard blocks intersect either. Specifically, you cannot “reach” a different block by taking repeated successors or predecessors.

**Proposition 19.26.** *If  $x$  and  $y$  are non-standard, then  $x \otimes x \oplus y$  and  $x \oplus y \notin [x]$ .*

*Proof.* If  $y$  is nonstandard, then  $y \neq \mathbf{z}$ .  $\mathbf{PA} \vdash \forall x (y \neq 0 \rightarrow x < (x+y))$ . Now suppose  $x \oplus y \in [x]$ . Since  $x \otimes x \oplus y$ , we would have  $x \oplus n^* = x \oplus y$ . But  $\mathbf{PA} \vdash \forall x \forall y \forall z ((x+y) = (x+z) \rightarrow y = z)$  (the cancellation law for addition). This would mean  $y = n^*$  for some standard  $n$ ; but  $y$  is assumed to be non-standard.  $\square$

**Proposition 19.27.** *There is no least non-standard block.*

*Proof.*  $\mathbf{PA} \vdash \forall x \exists y ((y+y) = x \vee (y+y)' = x)$ , i.e., that every  $x$  is divisible by 2 (possibly with remainder 1). If  $x$  is non-standard, so is  $y$ . By the preceding proposition,  $y \otimes y \oplus y$  and  $y \oplus y \notin [y]$ . Then also  $y \otimes (y \oplus y)^*$  and  $(y \oplus y)^* \notin [y]$ . But  $x = y \oplus y$  or  $x = (y \oplus y)^*$ , so  $y \otimes x$  and  $y \notin [x]$ .  $\square$

**Proposition 19.28.** *There is no largest block.*

*Proof.* Exercise.  $\square$

**Proposition 19.29.** *The ordering of the blocks is dense. That is, if  $x \otimes y$  and  $[x] \neq [y]$ , then there is a block  $[z]$  distinct from both that is between them.*

*Proof.* Suppose  $x \otimes y$ . As before,  $x \oplus y$  is divisible by two (possibly with remainder): there is a  $z \in |M|$  such that either  $x \oplus y = z \oplus z$  or  $x \oplus y = (z \oplus z)^*$ . The element  $z$  is the “average” of  $x$  and  $y$ , and  $x \otimes z$  and  $z \otimes y$ .  $\square$

The non-standard blocks are therefore ordered like the rationals: they form a countably infinite dense linear ordering without endpoints. One can show that any two such countably infinite orderings are isomorphic. It follows that for any two countable non-standard models  $M_1$  and  $M_2$  of true arithmetic, their reducts to the language containing  $<$  and  $=$  only are isomorphic. Indeed, an isomorphism  $h$  can be defined as follows: the standard parts of  $M_1$  and  $M_2$  are isomorphic to the standard model  $N$  and hence to each other. The blocks making up the non-standard part are themselves ordered like the rationals and therefore isomorphic; an isomorphism of the blocks can be extended to an isomorphism *within* the blocks by matching up arbitrary elements in each, and then taking the image of the successor of  $x$  in  $M_1$  to be the successor of the image of  $x$  in  $M_2$ . Note that it does *not* follow that  $\mathfrak{M}_1$  and  $\mathfrak{M}_2$  are isomorphic in the full language of arithmetic (indeed, isomorphism is always relative to a language), as there are non-isomorphic ways to define addition and multiplication over  $|M_1|$  and  $|M_2|$ . (This also follows from a famous theorem due to Vaught that the number of countable models of a complete theory cannot be 2.)

## 19.9 Computable Models of Arithmetic

The standard model  $N$  has two nice features. Its domain is the natural numbers  $\mathbb{N}$ , i.e., its elements are just the kinds of things we want to talk about using the language of arithmetic, and the standard numeral  $\bar{n}$  actually picks out  $n$ . The other nice feature

is that the interpretations of the non-logical symbols of  $\mathcal{L}_A$  are all *computable*. The successor, addition, and multiplication functions which serve as  $\iota^N$ ,  $+^N$ , and  $\times^N$  are computable functions of numbers. (Computable by Turing machines, or definable by primitive recursion, say.) And the less-than relation on  $N$ , i.e.,  $<^N$ , is decidable.

Non-standard models of arithmetical theories such as  $\mathbf{Q}$  and  $\mathbf{PA}$  must contain non-standard elements. Thus their domains typically include elements in addition to  $\mathbb{N}$ . However, any countable structure can be built on any countably infinite set, including  $\mathbb{N}$ . So there are also non-standard models with domain  $\mathbb{N}$ . In such models  $M$ , of course, at least some numbers cannot play the roles they usually play, since some  $k$  must be different from  $\text{Val}^M(\bar{n})$  for all  $n \in \mathbb{N}$ .

**Definition 19.30.** A structure  $M$  for  $\mathcal{L}_A$  is *computable* iff  $|M| = \mathbb{N}$  and  $\iota^M$ ,  $+^M$ ,  $\times^M$  are computable functions and  $<^M$  is a decidable relation.

**Example 19.31.** Recall the structure  $K$  from [Example 19.18](#). Its domain was  $|K| = \mathbb{N} \cup \{a\}$  and interpretations

$$\begin{aligned} 0^K &= 0 \\ \iota^K(x) &= \begin{cases} x+1 & \text{if } x \in \mathbb{N} \\ a & \text{if } x = a \end{cases} \\ +^K(x, y) &= \begin{cases} x+y & \text{if } x, y \in \mathbb{N} \\ a & \text{otherwise} \end{cases} \\ \times^K(x, y) &= \begin{cases} xy & \text{if } x, y \in \mathbb{N} \\ 0 & \text{if } x = 0 \text{ or } y = 0 \\ a & \text{otherwise} \end{cases} \\ <^K &= \{\langle x, y \rangle : x, y \in \mathbb{N} \text{ and } x < y\} \cup \{\langle x, a \rangle : x \in |K|\} \end{aligned}$$

But  $|K|$  is countably infinite and so is equinumerous with  $\mathbb{N}$ . For instance,  $g: \mathbb{N} \rightarrow |K|$  with  $g(0) = a$  and  $g(n) = n+1$  for  $n > 0$  is



a bijection. We can turn it into an isomorphism between a new model  $K'$  of  $\mathbf{Q}$  and  $K$ . In  $K'$ , we have to assign different functions and relations to the symbols of  $\mathcal{L}_A$ , since different elements of  $\mathbb{N}$  play the roles of standard and non-standard numbers.

Specifically, 0 now plays the role of  $a$ , not of the smallest standard number. The smallest standard number is now 1. So we assign  $0^{K'} = 1$ . The successor function is also different now: given a standard number, i.e., an  $n > 0$ , it still returns  $n+1$ . But 0 now plays the role of  $a$ , which is its own successor. So  $s^{K'}(0) = 0$ . For addition and multiplication we likewise have

$$+^{K'}(x,y) = \begin{cases} x+y-1 & \text{if } x, y > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\times^{K'}(x,y) = \begin{cases} 1 & \text{if } x = 1 \text{ or } y = 1 \\ xy - x - y + 2 & \text{if } x, y > 1 \\ 0 & \text{otherwise} \end{cases}$$

And we have  $\langle x, y \rangle \in <^{K'}$  iff  $x < y$  and  $x > 0$  and  $y > 0$ , or if  $y = 0$ .

All of these functions are computable functions of natural numbers and  $<^{K'}$  is a decidable relation on  $\mathbb{N}$ —but they are not the same functions as successor, addition, and multiplication on  $\mathbb{N}$ , and  $<^{K'}$  is not the same relation as  $<$  on  $\mathbb{N}$ .

**Example 19.31** shows that  $\mathbf{Q}$  has computable non-standard models with domain  $\mathbb{N}$ . However, the following result shows that this is not true for models of  $\mathbf{PA}$  (and thus also for models of  $\mathbf{TA}$ ).

**Theorem 19.32 (Tennenbaum's Theorem).**  *$N$  is the only computable model of  $\mathbf{PA}$ .*

## Summary

A **model of arithmetic** is a structure for the language  $\mathcal{L}_A$  of arithmetic. There is one distinguished such model, the **standard**

**model**  $N$ , with  $|N| = \mathbb{N}$  and interpretations of  $0$ ,  $\iota$ ,  $+$ ,  $\times$ , and  $<$  given by  $0$ , the successor, addition, and multiplication functions on  $\mathbb{N}$ , and the less-than relation.  $N$  is a model of the theories  $\mathbf{Q}$  and  $\mathbf{PA}$ .

More generally, a structure for  $\mathcal{L}_A$  is called **standard** iff it is isomorphic to  $N$ . Two structures are isomorphic if there is an **isomorphism** between them, i.e., a bijective function which preserves the interpretations of constant symbols, function symbols, and predicate symbols. By the **isomorphism theorem**, isomorphic structures are **elementarily equivalent**, i.e., they make the same sentences true. In standard models, the domain is just the set of values of all the numerals  $\bar{n}$ .

Models of  $\mathbf{Q}$  and  $\mathbf{PA}$  that are not isomorphic to  $N$  are called **non-standard**. In non-standard models, the domain is not exhausted by the values of the numerals. An element  $x \in |M|$  where  $x \neq \text{Val}^M(\bar{n})$  for all  $n \in \mathbb{N}$  is called a **non-standard element** of  $M$ . If  $M \models \mathbf{Q}$ , non-standard elements must obey the axioms of  $\mathbf{Q}$ , e.g., they have unique successors, they can be added and multiplied, and compared using  $<$ . The standard elements of  $M$  are all  $<^M$  all the non-standard elements. Non-standard models exist because of the compactness theorem, and for  $\mathbf{Q}$  they can relatively easily be given explicitly. Such models can be used to show that, e.g.,  $\mathbf{Q}$  is not strong enough to prove certain sentences, e.g.,  $\mathbf{Q} \not\models \forall x \forall y (x+y) = (y+x)$ . This is done by defining a non-standard  $M$  in which non-standard elements don't obey the law of commutativity.

Non-standard models of  $\mathbf{PA}$  cannot be so easily specified explicitly. By showing that  $\mathbf{PA}$  proves certain sentences, we can investigate the structure of the non-standard part of a non-standard model of  $\mathbf{PA}$ . If a non-standard model  $M$  of  $\mathbf{PA}$  is countable, every non-standard element is part of a "block" of non-standard elements which are ordered like  $\mathbb{Z}$  by  $<^M$ . These blocks themselves are arranged like  $\mathbb{Q}$ , i.e., there is no smallest or largest block, and there is always a block in between any two blocks.

Any countable model is isomorphic to one with domain  $\mathbb{N}$ . If the interpretations of  $\iota$ ,  $+$ ,  $\times$ , and  $<$  in such a model are com-

putable functions, we say it is a **computable model**. The standard model  $N$  is computable, since the successor, addition, and multiplication functions and the less-than relation on  $\mathbb{N}$  are computable. It is possible to define computable non-standard models of  $\mathcal{Q}$ , but  $N$  is the only computable model of **PA**. This is **Tennenbaum's Theorem**.

## Problems

**Problem 19.1.** Prove **Proposition 19.2**.

**Problem 19.2.** Carry out the proof of (b) of **Theorem 19.6** in detail. Make sure to note where each of the five properties characterizing isomorphisms of **Definition 19.5** is used.

**Problem 19.3.** Show that for any structure  $M$ , if  $X$  is a definable subset of  $M$ , and  $h$  is an automorphism of  $M$ , then  $X = \{h(x) : x \in X\}$  (i.e.,  $X$  is fixed under  $h$ ).

**Problem 19.4.** Show that the converse of **Proposition 19.12** is false, i.e., give an example of a structure  $M$  with  $|M| = \{\text{Val}^M(\bar{n}) : n \in \mathbb{N}\}$  that is not isomorphic to  $N$ .

**Problem 19.5.** Recall that  $\mathcal{Q}$  contains the axioms

$$\forall x \forall y (x' = y' \rightarrow x = y) \quad (Q_1)$$

$$\forall x 0 \neq x' \quad (Q_2)$$

$$\forall x (x = 0 \vee \exists y x = y') \quad (Q_3)$$

Give structures  $M_1, M_2, M_3$  such that

1.  $M_1 \models Q_1, M_1 \models Q_2, M_1 \not\models Q_3$ ;
2.  $M_2 \models Q_1, M_2 \not\models Q_2, M_2 \models Q_3$ ; and
3.  $M_3 \not\models Q_1, M_3 \models Q_2, M_3 \models Q_3$ ;

Obviously, you just have to specify  $0^{M_i}$  and  $'^{M_i}$  for each.

**Problem 19.6.** Prove that  $K$  from [Example 19.18](#) satisfies the remaining axioms of  $\mathbf{Q}$ ,

$$\forall x (x \times 0) = 0 \quad (Q_6)$$

$$\forall x \forall y (x \times y') = ((x \times y) + x) \quad (Q_7)$$

$$\forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y) \quad (Q_8)$$

Find a sentence only involving  $'$  true in  $N$  but false in  $K$ .

**Problem 19.7.** Expand  $L$  of [Example 19.19](#) to include  $\otimes$  and  $\odot$  that interpret  $\times$  and  $<$ . Show that your structure satisfies the remaining axioms of  $\mathbf{Q}$ ,

$$\forall x (x \otimes 0) = 0 \quad (Q_6)$$

$$\forall x \forall y (x \otimes y') = ((x \otimes y) + x) \quad (Q_7)$$

$$\forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y) \quad (Q_8)$$

**Problem 19.8.** In  $L$  of [Example 19.19](#),  $a^* = a$  and  $b^* = b$ . Is there a model of  $\mathbf{Q}$  in which  $a^* = b$  and  $b^* = a$ ?

**Problem 19.9.** Find sentences in  $\mathcal{L}_A$  derivable in  $\mathbf{PA}$  (and hence true in  $N$ ) which guarantee the properties of  $\mathbf{z}$ ,  $*$ , and  $\odot$  in [Proposition 19.21](#)

**Problem 19.10.** Show that in a non-standard model of  $\mathbf{PA}$ , there is no largest block.

**Problem 19.11.** Write out a detailed proof of [Proposition 19.29](#). Which sentence must  $\mathbf{PA}$  derive in order to guarantee the existence of  $z$ ? Why is  $x \odot z$  and  $z \odot y$ , and why is  $[x] \neq [z]$  and  $[z] \neq [y]$ ?

**Problem 19.12.** Give a structure  $L'$  with  $|L'| = \mathbb{N}$  isomorphic to  $L$  of [Example 19.19](#).

## APPENDIX A

# *Derivations in Arithmetic Theories*

When we showed that all general recursive functions are representable in  $\mathbf{Q}$ , and in the proofs of the incompleteness theorems, we claimed that various things are provable in  $\mathbf{Q}$  and  $\mathbf{PA}$ . The proofs of these claims, however, just gave the arguments informally without exhibiting actual derivations in natural deduction. We provide some of these derivations in this chapter.

For instance, in [Lemma 17.16](#) we proved that, for all  $n$  and  $m \in \mathbb{N}$ ,  $\mathbf{Q} \vdash (\overline{n} + \overline{m}) = \overline{n + m}$ . We did this by induction on  $m$ .

*Proof of Lemma 17.16.* Base case:  $m = 0$ . Then what has to be proved is that, for all  $n$ ,  $\mathbf{Q} \vdash \overline{n} + \overline{0} = \overline{n + 0}$ . Since  $\overline{0}$  is just  $0$  and  $\overline{n + 0}$  is  $\overline{n}$ , this amounts to showing that  $\mathbf{Q} \vdash (\overline{n} + 0) = \overline{n}$ . The derivation

$$\frac{\forall x (x + 0) = x}{(\overline{n} + 0) = \overline{n}} \forall\text{Elim}$$

is a natural deduction derivation of  $(\overline{n} + 0) = \overline{n}$  with one undischarged assumption, and that undischarged assumption is an ax-

iom of  $\mathbf{Q}$ .

Inductive step: Suppose that, for any  $n$ ,  $\mathbf{Q} \vdash (\bar{n} + \bar{m}) = \overline{n + m}$  (say, by a derivation  $\delta_{n,m}$ ). We have to show that also  $\mathbf{Q} \vdash (\bar{n} + \bar{m} + 1) = \overline{n + m + 1}$ . Note that  $\overline{m + 1} \equiv \bar{m}'$ , and that  $\overline{n + m + 1} \equiv \overline{n + m}'$ . So we are looking for a derivation of  $(\bar{n} + \bar{m}') = \overline{n + m}'$  from the axioms of  $\mathbf{Q}$ . Our derivation may use the derivation  $\delta_{n,m}$  which exists by inductive hypothesis.

$$\frac{\begin{array}{c} \vdots \\ \delta_{n,m} \\ \vdots \\ (\bar{n} + \bar{m}) = \overline{n + m} \end{array} \quad \frac{\frac{\forall x \forall y (x + y') = (x + y)'}{\forall y (\bar{n} + y') = (\bar{n} + y)'} \text{ } \forall\text{Elim}}{\frac{(\bar{n} + \bar{m}') = (\bar{n} + \bar{m})'}{\overline{n + m}' = \overline{n + m}} \text{ } \forall\text{Elim}} \text{ } =\text{Elim}$$

In the last =Elim inference, we replace the subterm  $\bar{n} + \bar{m}$  of the right side  $(\bar{n} + \bar{m})'$  of the right premise by the term  $\overline{n + m}$ .  $\square$

In **Lemma 17.23**, we showed that  $\mathbf{Q} \vdash \forall x \neg x < 0$ . What does an actual derivation look like?

*Proof of Lemma 17.23.* To prove a universal claim like this, we use  $\forall$ Intro, which requires a derivation of  $\neg a < 0$ . Looking at axiom  $Q_8$ , this means proving  $\neg \exists z (z' + a) = 0$ . Specifically, if we had a proof of the latter,  $Q_8$  would allow us to prove the former (recall that  $A \leftrightarrow B$  is short for  $(A \rightarrow B) \wedge (B \rightarrow A)$ ).

$$\frac{\frac{\frac{\forall x \forall y (x < y \leftrightarrow \exists z (z' + x) = y)}{\forall y (a < y \leftrightarrow \exists z (z' + a) = y)} \text{ } \forall\text{Elim}}{a < 0 \leftrightarrow \exists z (z' + a) = 0} \text{ } \forall\text{Elim}}{a < 0 \rightarrow \exists z (z' + a) = 0} \text{ } \wedge\text{Elim} \quad [a < 0]^1}{\frac{\neg \exists z (z' + a) = 0 \quad \exists z (z' + a) = 0}{1} \text{ } \rightarrow\text{Elim}} \text{ } \neg\text{Elim}$$

$$\frac{1}{\neg a < 0} \text{ } \neg\text{Intro}$$

This is a derivation of  $\neg a < 0$  from  $\neg \exists z (z' + a) = 0$  (and  $Q_8$ ); let's call it  $\delta_1$ .

Now how do we prove  $\neg \exists z (z' + a) = 0$  from the axioms of  $\mathbf{Q}$ ? To prove a negated claim like this, we'd need a derivation of the form

$$\begin{array}{c} [\exists z (z' + a) = 0]^2 \\ \vdots \\ \perp \\ \hline 2 \quad \neg \exists z (z' + a) = 0 \quad \neg\text{Intro} \end{array}$$

To get a contradiction from an existential claim, we introduce a constant  $b$  for the existentially quantified variable  $z$  and use  $\exists\text{Elim}$ :

$$\begin{array}{c} [(b' + a) = 0]^3 \\ \vdots \\ \delta_2 \\ \vdots \\ \perp \\ \hline 3 \quad \frac{[\exists z (z' + a) = 0]^2}{\neg \exists z (z' + a) = 0} \neg\text{Intro} \quad \exists\text{Elim} \end{array}$$

Now the task is to fill in  $\delta_2$ , i.e., prove  $\perp$  from  $(b' + a) = 0$  and the axioms of  $\mathbf{Q}$ .  $Q_2$  says that  $0$  can't be the successor of some number, so one way of doing that would be to show that  $(b' + a)$  is equal to the successor of some number. Since that expression itself is a sum, the axioms for addition must come into play. If  $a = 0$ ,  $Q_5$  would tell us that  $(b' + a) = b'$ , i.e.,  $b' + a$  is the successor of some number, namely of  $b$ . On the other hand, if  $a = c'$  for some  $c$ , then  $(b' + a) = (b' + c')$  by  $=\text{Elim}$ , and  $(b' + c') = (b' + c)'$  by  $Q_6$ . So again,  $b' + a$  is the successor of a number—in this case,  $b' + c$ . So the strategy is to divide the task into these two cases. We also have to verify that  $\mathbf{Q}$  proves that one of these cases holds, i.e.,  $\mathbf{Q} \vdash a = 0 \vee \exists y (a = y')$ , but this follows directly from  $Q_3$  by  $\vee\text{Elim}$ . Here are the two cases:

Case 1: Prove  $\perp$  from  $a = 0$  and  $(b' + a) = 0$  (and axioms  $Q_2$ ,  $Q_5$ ):

$$\frac{\frac{\forall x \neg 0 = x'}{\neg 0 = b'} \vee\text{Elim} \quad \frac{\frac{\forall x (x + 0) = x}{(b' + 0) = b'} \vee\text{Elim} \quad \frac{\frac{a = 0 \quad (b' + a) = 0}{(b' + 0) = 0} =\text{Elim}}{0 = (b' + 0)} =\text{Elim}}{0 = b'} \neg\text{Elim}}{\perp} \neg\text{Elim}$$

Call this derivation  $\delta_3$ . (We've abbreviated the derivation of  $0 = (b' + 0)$  from  $(b' + 0) = 0$  by a double inference line.)

Case 2: Prove  $\perp$  from  $\exists y a = y'$  and  $(b' + a) = 0$  (and axioms  $Q_2, Q_6$ ). We first show how to derive  $\perp$  from  $a = c'$  and  $(b' + a) = 0$ .

$$\frac{\frac{\forall x \neg 0 = x'}{\neg 0 = (b' + c)'} \forall\text{Elim} \quad \frac{\frac{a = c' \quad (b' + a) = 0}{(b' + c') = 0} =\text{Elim} \quad \frac{\frac{\forall x \forall y (x + y') = (x + y)'}{\forall y (b' + y') = (b' + y)'} \forall\text{Elim}}{(b' + c') = (b' + c)'} \forall\text{Elim}}{0 = (b' + c)'} =\text{Elim}}{\perp} \neg\text{Elim}$$

Call this  $\delta_4$ . We get the required derivation  $\delta_5$  by applying  $\exists\text{Elim}$  and discharging the assumption  $a = c'$ :

$$\frac{6 \quad \exists y a = y' \quad \frac{[a = c']^6 \quad (b' + a) = 0}{\vdots} \delta_4 \quad \perp}{\perp} \exists\text{Elim}$$

Putting everything together, the full proof looks like this:

$$\frac{3 \quad \frac{[\exists z (z' + a) = 0]^2 \quad \frac{7 \quad \frac{\forall x (x = 0 \vee \exists y (a = y'))}{a = 0 \vee \exists y (a = y')} \forall\text{Elim} \quad \frac{[a = 0]^7 \quad [\exists y a = y']^7}{[(b' + a) = 0]^3 \quad [(b' + a) = 0]^3} \quad \begin{matrix} \vdots \\ \delta_3 \\ \vdots \\ \perp \end{matrix} \quad \begin{matrix} \vdots \\ \delta_5 \\ \vdots \\ \perp \end{matrix}}{\perp} \forall\text{Elim}}{\perp} \exists\text{Elim}}{\perp} \neg\text{Intro} \quad \frac{\perp}{\neg \exists z (z' + a) = 0} \neg\text{Intro} \quad \frac{\perp}{\neg a < 0} \forall\text{Intro} \quad \frac{\neg a < 0}{\forall x \neg x < 0} \forall\text{Intro}}{\perp} \delta_1$$

□



In the proof of **Theorem 18.7**, we defined  $RProv(y)$  as

$$\exists x (\text{Prf}(x, y) \wedge \forall z (z < x \rightarrow \neg \text{Ref}(z, y))).$$

$\text{Prf}(x, y)$  is the formula representing the proof relation of  $\mathbf{T}$  (a consistent, axiomatizable extension of  $\mathbf{Q}$ ) in  $\mathbf{Q}$ , and  $\text{Ref}(z, y)$  is the formula representing the refutation relation. That means that if  $n$  is the Gödel number of a proof of  $A$ , then  $\mathbf{Q} \vdash \text{Prf}(\bar{n}, \ulcorner A \urcorner)$ , and otherwise  $\mathbf{Q} \vdash \neg \text{Prf}(\bar{n}, \ulcorner A \urcorner)$ . Similarly, if  $n$  is the Gödel number of a proof of  $\neg A$ , then  $\mathbf{Q} \vdash \text{Ref}(\bar{n}, \ulcorner A \urcorner)$ , and otherwise  $\mathbf{Q} \vdash \neg \text{Ref}(\bar{n}, \ulcorner A \urcorner)$ . We use the Diagonal Lemma to find a sentence  $R$  such that  $\mathbf{Q} \vdash R \leftrightarrow \neg RProv(\ulcorner R \urcorner)$ . Rosser's Theorem states that  $\mathbf{T} \not\vdash R$  and  $\mathbf{T} \not\vdash \neg R$ . Both claims were proved indirectly: we show that if  $\mathbf{T} \vdash R$ ,  $\mathbf{T}$  is inconsistent, i.e.,  $\mathbf{T} \vdash \perp$ , and the same if  $\mathbf{T} \vdash \neg R$ .

*Proof of Theorem 18.7.* First we prove something things about  $<$ . By **Lemma 17.24**, we know that  $\mathbf{Q} \vdash \forall x (x < \overline{n+1} \rightarrow (x = 0 \vee \dots \vee x = \bar{n}))$  for every  $n$ . So of course also (if  $n > 1$ ),  $\mathbf{Q} \vdash \forall x (x < \bar{n} \rightarrow (x = 0 \vee \dots \vee x = \overline{n-1}))$ . We can use this to derive  $a = 0 \vee \dots \vee a = \overline{n-1}$  from  $a < \bar{n}$ :

$$\frac{\begin{array}{c} \vdots \\ \vdots \\ \forall x (x < \bar{n} \rightarrow (x = \bar{0} \vee \dots \vee x = \overline{n-1})) \end{array}}{\frac{a < \bar{n} \quad \frac{\forall x (x < \bar{n} \rightarrow (x = \bar{0} \vee \dots \vee x = \overline{n-1}))}{a < \bar{n} \rightarrow (a = \bar{0} \vee \dots \vee a = \overline{n-1})} \vee \text{Elim}}{a = \bar{0} \vee \dots \vee a = \overline{n-1}} \rightarrow \text{Elim}}$$

Let's call this derivation  $\lambda_1$ .

Now, to show that  $\mathbf{T} \not\vdash R$ , we assume that  $\mathbf{T} \vdash R$  (with a derivation  $\delta$ ) and show that  $\mathbf{T}$  then would be inconsistent. Let  $n$  be the Gödel number of  $\delta$ . Since  $\text{Prf}$  represents the proof relation in  $\mathbf{Q}$ , there is a derivation  $\delta_1$  of  $\text{Prf}(\bar{n}, \ulcorner R \urcorner)$ . Furthermore, no  $k < n$  is the Gödel number of a refutation of  $R$  since  $\mathbf{T}$  is assumed to be consistent, so for each  $k < n$ ,  $\mathbf{Q} \vdash \neg \text{Ref}(\bar{k}, \ulcorner R \urcorner)$ ; let  $\rho_k$  be the corresponding derivation. We get a derivation of  $RProv(\ulcorner R \urcorner)$ :



$$\begin{array}{c}
 \vdots \lambda_2 \\
 a = \bar{0} \vee \dots \vee \\
 \hline
 a = \bar{m} \vee \bar{m} < a \quad \dots \quad \frac{\perp}{\bar{m} < a} \perp_I \quad \dots \quad [\bar{m} < a]^2 \quad \dots \quad \rho_1 \\
 \hline
 \bar{m} < a \quad \vee \text{Elim}^* \quad \text{Ref}(\bar{m}, \ulcorner R \urcorner) \\
 \hline
 \frac{\bar{m} < a \wedge \text{Ref}(\bar{m}, \ulcorner R \urcorner)}{\exists z(z < a \wedge \text{Ref}(z, \ulcorner R \urcorner))} \exists \text{Intro} \\
 \frac{1 \quad \frac{\text{Prf}(a, \ulcorner R \urcorner) \rightarrow \exists z(z < a \wedge \text{Ref}(z, \ulcorner R \urcorner))}{\forall x(\text{Prf}(x, \ulcorner R \urcorner) \rightarrow \exists z(z < x \wedge \text{Ref}(z, \ulcorner R \urcorner)))} \rightarrow \text{Intro}}{\forall x(\text{Prf}(x, \ulcorner R \urcorner) \rightarrow \exists z(z < x \wedge \text{Ref}(z, \ulcorner R \urcorner)))} \forall \text{Intro} \\
 \vdots \\
 \neg \exists x(\text{Prf}(x, \ulcorner R \urcorner) \wedge \forall z(z < x \rightarrow \neg \text{Ref}(z, \ulcorner R \urcorner))) \\
 \wedge \text{Intro}
 \end{array}$$

where  $\pi'_k$  is the derivation

$$\frac{\frac{\vdots \pi_k}{\neg \text{Prf}(\bar{k}, \ulcorner R \urcorner)} \quad \frac{a = \bar{k} \quad \text{Prf}(a, \ulcorner R \urcorner)}{\text{Prf}(\bar{k}, \ulcorner R \urcorner)} = \text{Elim}}{\perp} \neg \text{Elim}$$

and  $\lambda_2$  is

$$\begin{array}{c}
 \vdots \lambda_3 \\
 (a < \bar{m} \vee \\
 a = \bar{m}) \vee \\
 \bar{m} < a \\
 \hline
 a = \bar{0} \vee \dots \vee \\
 a = \overline{m-1} \\
 a = \bar{m} \vee \bar{m} < a \\
 \hline
 a = \bar{0} \vee \dots \vee \\
 a = \bar{m} \vee \bar{m} < a \\
 \hline
 [a < \bar{m}]^3 \\
 \vdots \lambda_1 \\
 a = \bar{0} \vee \dots \vee \\
 a = \overline{m-1} \\
 a = \bar{m} \vee \bar{m} < a \\
 \hline
 [a = \bar{m}]^3 \\
 a = \bar{0} \vee \dots \vee \\
 a = \bar{m} \vee \bar{m} < a \\
 \hline
 [\bar{m} < a]^3 \\
 a = \bar{0} \vee \dots \vee \\
 a = \bar{m} \vee \bar{m} < a \\
 \hline
 a = \bar{0} \vee \dots \vee a = \bar{m} \vee \bar{m} < a \\
 \hline
 \vee \text{Intro}^* \\
 \hline
 \vee \text{Elim}^2
 \end{array}$$

(The derivation  $\lambda_3$  exists by [Lemma 17.25](#). We abbreviate repeated use of  $\vee \text{Intro}$  by  $\vee \text{Intro}^*$  and the double use of  $\vee \text{Elim}$  to

derive  $a = \bar{0} \vee \dots \vee a = \bar{m} \vee \bar{m} < a$  from  $(a < \bar{m} \vee a = \bar{m}) \vee \bar{m} < a$   
as  $\vee\text{Elim}^2$ .) □

## APPENDIX B

# *Proofs*

### B.1 Introduction

Based on your experiences in introductory logic, you might be comfortable with a derivation system—probably a natural deduction or Fitch style derivation system, or perhaps a proof-tree system. You probably remember doing proofs in these systems, either proving a formula or show that a given argument is valid. In order to do this, you applied the rules of the system until you got the desired end result. In reasoning *about* logic, we also prove things, but in most cases we are not using a derivation system. In fact, most of the proofs we consider are done in English (perhaps, with some symbolic language thrown in) rather than entirely in the language of first-order logic. When constructing such proofs, you might at first be at a loss—how do I prove something without a derivation system? How do I start? How do I know if my proof is correct?

Before attempting a proof, it's important to know what a proof is and how to construct one. As implied by the name, a *proof* is meant to show that something is true. You might think of this in terms of a dialogue—someone asks you if something is true, say, if every prime other than two is an odd number. To answer “yes” is not enough; they might want to know *why*. In this case, you'd give them a proof.

In everyday discourse, it might be enough to gesture at an

answer, or give an incomplete answer. In logic and mathematics, however, we want rigorous proof—we want to show that something is true beyond *any* doubt. This means that every step in our proof must be justified, and the justification must be cogent (i.e., the assumption you’re using is actually assumed in the statement of the theorem you’re proving, the definitions you apply must be correctly applied, the justifications appealed to must be correct inferences, etc.).

Usually, we’re proving some statement. We call the statements we’re proving by various names: propositions, theorems, lemmas, or corollaries. A proposition is a basic proof-worthy statement: important enough to record, but perhaps not particularly deep nor applied often. A theorem is a significant, important proposition. Its proof often is broken into several steps, and sometimes it is named after the person who first proved it (e.g., Cantor’s Theorem, the Löwenheim–Skolem theorem) or after the fact it concerns (e.g., the completeness theorem). A lemma is a proposition or theorem that is used in the proof of a more important result. Confusingly, sometimes lemmas are important results in themselves, and also named after the person who introduced them (e.g., Zorn’s Lemma). A corollary is a result that easily follows from another one.

A statement to be proved often contains assumptions that clarify which kinds of things we’re proving something about. It might begin with “Let  $A$  be a formula of the form  $B \rightarrow C$ ” or “Suppose  $\Gamma \vdash A$ ” or something of the sort. These are *hypotheses* of the proposition, theorem, or lemma, and you may assume these to be true in your proof. They restrict what we’re proving, and also introduce some names for the objects we’re talking about. For instance, if your proposition begins with “Let  $A$  be a formula of the form  $B \rightarrow C$ ,” you’re proving something about all formulas of a certain sort only (namely, conditionals), and it’s understood that  $B \rightarrow C$  is an arbitrary conditional that your proof will talk about.

## B.2 Starting a Proof

But where do you even start?

You've been given something to prove, so this should be the last thing that is mentioned in the proof (you can, obviously, *announce* that you're going to prove it at the beginning, but you don't want to use it as an assumption). Write what you are trying to prove at the bottom of a fresh sheet of paper—this way you don't lose sight of your goal.

Next, you may have some assumptions that you are able to use (this will be made clearer when we talk about the *type* of proof you are doing in the next section). Write these at the top of the page and make sure to flag that they are assumptions (i.e., if you are assuming  $p$ , write “assume that  $p$ ,” or “suppose that  $p$ ”). Finally, there might be some definitions in the question that you need to know. You might be told to use a specific definition, or there might be various definitions in the assumptions or conclusion that you are working towards. *Write these down and ensure that you understand what they mean.*

How you set up your proof will also be dependent upon the form of the question. The next section provides details on how to set up your proof based on the type of sentence.

## B.3 Using Definitions

We mentioned that you must be familiar with all definitions that may be used in the proof, and that you can properly apply them. This is a really important point, and it is worth looking at in a bit more detail. Definitions are used to abbreviate properties and relations so we can talk about them more succinctly. The introduced abbreviation is called the *definiendum*, and what it abbreviates is the *definiens*. In proofs, we often have to go back to how the definiendum was introduced, because we have to exploit the logical structure of the definiens (the long version of which the defined term is the abbreviation) to get through our proof. By

unpacking definitions, you're ensuring that you're getting to the heart of where the logical action is.

We'll start with an example. Suppose you want to prove the following:

**Proposition B.1.** *For any sets  $A$  and  $B$ ,  $A \cup B = B \cup A$ .*

In order to even start the proof, we need to know what it means for two sets to be identical; i.e., we need to know what the “=” in that equation means for sets. Sets are defined to be identical whenever they have the same elements. So the definition we have to unpack is:

**Definition B.2.** Sets  $A$  and  $B$  are *identical*,  $A = B$ , iff every element of  $A$  is an element of  $B$ , and vice versa.

This definition uses  $A$  and  $B$  as placeholders for arbitrary sets. What it defines—the *definiendum*—is the expression “ $A = B$ ” by giving the condition under which  $A = B$  is true. This condition—“every element of  $A$  is an element of  $B$ , and vice versa”—is the *definiens*.<sup>1</sup> The definition specifies that  $A = B$  is true if, and only if (we abbreviate this to “iff”) the condition holds.

When you apply the definition, you have to match the  $A$  and  $B$  in the definition to the case you're dealing with. In our case, it means that in order for  $A \cup B = B \cup A$  to be true, each  $z \in A \cup B$  must also be in  $B \cup A$ , and vice versa. The expression  $A \cup B$  in the proposition plays the role of  $A$  in the definition, and  $B \cup A$  that of  $B$ . Since  $A$  and  $B$  are used both in the definition and in the statement of the proposition we're proving, but in different uses, you have to be careful to make sure you don't mix up the two. For instance, it would be a mistake to think that you could prove the proposition by showing that every element of  $A$  is an element

---

<sup>1</sup>In this particular case—and very confusingly!—when  $A = B$ , the sets  $A$  and  $B$  are just one and the same set, even though we use different letters for it on the left and the right side. But the ways in which that set is picked out may be different, and that makes the definition non-trivial.



of  $B$ , and vice versa—that would show that  $A = B$ , not that  $A \cup B = B \cup A$ . (Also, since  $A$  and  $B$  may be any two sets, you won't get very far, because if nothing is assumed about  $A$  and  $B$  they may well be different sets.)

Within the proof we are dealing with set-theoretic notions such as union, and so we must also know the meanings of the symbol  $\cup$  in order to understand how the proof should proceed. And sometimes, unpacking the definition gives rise to further definitions to unpack. For instance,  $A \cup B$  is defined as  $\{z : z \in A \text{ or } z \in B\}$ . So if you want to prove that  $x \in A \cup B$ , unpacking the definition of  $\cup$  tells you that you have to prove  $x \in \{z : z \in A \text{ or } z \in B\}$ . Now you also have to remember that  $x \in \{z : \dots z \dots\}$  iff  $\dots x \dots$ . So, further unpacking the definition of the  $\{z : \dots z \dots\}$  notation, what you have to show is:  $x \in A$  or  $x \in B$ . So, “every element of  $A \cup B$  is also an element of  $B \cup A$ ” really means: “for every  $x$ , if  $x \in A$  or  $x \in B$ , then  $x \in B$  or  $x \in A$ .” If we fully unpack the definitions in the proposition, we see that what we have to show is this:

**Proposition B.3.** *For any sets  $A$  and  $B$ : (a) for every  $x$ , if  $x \in A$  or  $x \in B$ , then  $x \in B$  or  $x \in A$ , and (b) for every  $x$ , if  $x \in B$  or  $x \in A$ , then  $x \in A$  or  $x \in B$ .*

What's important is that unpacking definitions is a necessary part of constructing a proof. Properly doing it is sometimes difficult: you must be careful to distinguish and match the variables in the definition and the terms in the claim you're proving. In order to be successful, you must know what the question is asking and what all the terms used in the question mean—you will often need to unpack more than one definition. In simple proofs such as the ones below, the solution follows almost immediately from the definitions themselves. Of course, it won't always be this simple.

## B.4 Inference Patterns

Proofs are composed of individual inferences. When we make an inference, we typically indicate that by using a word like “so,” “thus,” or “therefore.” The inference often relies on one or two facts we already have available in our proof—it may be something we have assumed, or something that we’ve concluded by an inference already. To be clear, we may label these things, and in the inference we indicate what other statements we’re using in the inference. An inference will often also contain an explanation of *why* our new conclusion follows from the things that come before it. There are some common patterns of inference that are used very often in proofs; we’ll go through some below. Some patterns of inference, like proofs by induction, are more involved (and will be discussed later).

We’ve already discussed one pattern of inference: unpacking, or applying, a definition. When we unpack a definition, we just restate something that involves the definiendum by using the definiens. For instance, suppose that we have already established in the course of a proof that  $D = E$  (a). Then we may apply the definition of  $=$  for sets and infer: “Thus, by definition from (a), every element of  $D$  is an element of  $E$  and vice versa.”

Somewhat confusingly, we often do not write the justification of an inference when we actually make it, but before. Suppose we haven’t already proved that  $D = E$ , but we want to. If  $D = E$  is the conclusion we aim for, then we can restate this aim also by applying the definition: to prove  $D = E$  we have to prove that every element of  $D$  is an element of  $E$  and vice versa. So our proof will have the form: (a) prove that every element of  $D$  is an element of  $E$ ; (b) every element of  $E$  is an element of  $D$ ; (c) therefore, from (a) and (b) by definition of  $=$ ,  $D = E$ . But we would usually not write it this way. Instead we might write something like,

We want to show  $D = E$ . By definition of  $=$ , this amounts to showing that every element of  $D$  is an el-

ement of  $E$  and vice versa.

(a) ... (a proof that every element of  $D$  is an element of  $E$ ) ...

(b) ... (a proof that every element of  $E$  is an element of  $D$ ) ...

## Using a Conjunction

Perhaps the simplest inference pattern is that of drawing as conclusion one of the conjuncts of a conjunction. In other words: if we have assumed or already proved that  $p$  and  $q$ , then we're entitled to infer that  $p$  (and also that  $q$ ). This is such a basic inference that it is often not mentioned. For instance, once we've unpacked the definition of  $D = E$  we've established that every element of  $D$  is an element of  $E$  and vice versa. From this we can conclude that every element of  $E$  is an element of  $D$  (that's the "vice versa" part).

## Proving a Conjunction

Sometimes what you'll be asked to prove will have the form of a conjunction; you will be asked to "prove  $p$  and  $q$ ." In this case, you simply have to do two things: prove  $p$ , and then prove  $q$ . You could divide your proof into two sections, and for clarity, label them. When you're making your first notes, you might write "(1) Prove  $p$ " at the top of the page, and "(2) Prove  $q$ " in the middle of the page. (Of course, you might not be explicitly asked to prove a conjunction but find that your proof requires that you prove a conjunction. For instance, if you're asked to prove that  $D = E$  you will find that, after unpacking the definition of  $=$ , you have to prove: every element of  $D$  is an element of  $E$  *and* every element of  $E$  is an element of  $D$ ).

## Proving a Disjunction

When what you are proving takes the form of a disjunction (i.e., it is an statement of the form “ $p$  or  $q$ ”), it is enough to show that one of the disjuncts is true. However, it basically never happens that either disjunct just follows from the assumptions of your theorem. More often, the assumptions of your theorem are themselves disjunctive, or you’re showing that all things of a certain kind have one of two properties, but some of the things have the one and others have the other property. This is where proof by cases is useful (see below).

## Conditional Proof

Many theorems you will encounter are in conditional form (i.e., show that if  $p$  holds, then  $q$  is also true). These cases are nice and easy to set up—simply assume the antecedent of the conditional (in this case,  $p$ ) and prove the conclusion  $q$  from it. So if your theorem reads, “If  $p$  then  $q$ ,” you start your proof with “assume  $p$ ” and at the end you should have proved  $q$ .

Conditionals may be stated in different ways. So instead of “If  $p$  then  $q$ ,” a theorem may state that “ $p$  only if  $q$ ,” “ $q$  if  $p$ ,” or “ $q$ , provided  $p$ .” These all mean the same and require assuming  $p$  and proving  $q$  from that assumption. Recall that a biconditional (“ $p$  if and only if (iff)  $q$ ”) is really two conditionals put together: if  $p$  then  $q$ , and if  $q$  then  $p$ . All you have to do, then, is two instances of conditional proof: one for the first conditional and another one for the second. Sometimes, however, it is possible to prove an “iff” statement by chaining together a bunch of other “iff” statements so that you start with “ $p$ ” an end with “ $q$ ”—but in that case you have to make sure that each step really is an “iff.”

## Universal Claims

Using a universal claim is simple: if something is true for anything, it’s true for each particular thing. So if, say, the hypothesis of your proof is  $A \subseteq B$ , that means (unpacking the definition

of  $\subseteq$ ), that, for every  $x \in A$ ,  $x \in B$ . Thus, if you already know that  $z \in A$ , you can conclude  $z \in B$ .

Proving a universal claim may seem a little bit tricky. Usually these statements take the following form: “If  $x$  has  $P$ , then it has  $Q$ ” or “All  $P$ s are  $Q$ s.” Of course, it might not fit this form perfectly, and it takes a bit of practice to figure out what you’re asked to prove exactly. But: we often have to prove that all objects with some property have a certain other property.

The way to prove a universal claim is to introduce names or variables, for the things that have the one property and then show that they also have the other property. We might put this by saying that to prove something for *all*  $P$ s you have to prove it for an *arbitrary*  $P$ . And the name introduced is a name for an arbitrary  $P$ . We typically use single letters as these names for arbitrary things, and the letters usually follow conventions: e.g., we use  $n$  for natural numbers,  $A$  for formulas,  $A$  for sets,  $f$  for functions, etc.

The trick is to maintain generality throughout the proof. You start by assuming that an arbitrary object (“ $x$ ”) has the property  $P$ , and show (based only on definitions or what you are allowed to assume) that  $x$  has the property  $Q$ . Because you have not stipulated what  $x$  is specifically, other than that it has the property  $P$ , then you can assert that all every  $P$  has the property  $Q$ . In short,  $x$  is a stand-in for *all* things with property  $P$ .

**Proposition B.4.** *For all sets  $A$  and  $B$ ,  $A \subseteq A \cup B$ .*

*Proof.* Let  $A$  and  $B$  be arbitrary sets. We want to show that  $A \subseteq A \cup B$ . By definition of  $\subseteq$ , this amounts to: for every  $x$ , if  $x \in A$  then  $x \in A \cup B$ . So let  $x \in A$  be an arbitrary element of  $A$ . We have to show that  $x \in A \cup B$ . Since  $x \in A$ ,  $x \in A$  or  $x \in B$ . Thus,  $x \in \{x : x \in A \vee x \in B\}$ . But that, by definition of  $\cup$ , means  $x \in A \cup B$ .  $\square$

## Proof by Cases

Suppose you have a disjunction as an assumption or as an already established conclusion—you have assumed or proved that  $p$  or  $q$  is true. You want to prove  $r$ . You do this in two steps: first you assume that  $p$  is true, and prove  $r$ , then you assume that  $q$  is true and prove  $r$  again. This works because we assume or know that one of the two alternatives holds. The two steps establish that either one is sufficient for the truth of  $r$ . (If both are true, we have not one but two reasons for why  $r$  is true. It is not necessary to separately prove that  $r$  is true assuming both  $p$  and  $q$ .) To indicate what we're doing, we announce that we “distinguish cases.” For instance, suppose we know that  $x \in B \cup C$ .  $B \cup C$  is defined as  $\{x : x \in B \text{ or } x \in C\}$ . In other words, by definition,  $x \in B$  or  $x \in C$ . We would prove that  $x \in A$  from this by first assuming that  $x \in B$ , and proving  $x \in A$  from this assumption, and then assume  $x \in C$ , and again prove  $x \in A$  from this. You would write “We distinguish cases” under the assumption, then “Case (1):  $x \in B$ ” underneath, and “Case (2):  $x \in C$ ” halfway down the page. Then you'd proceed to fill in the top half and the bottom half of the page.

Proof by cases is especially useful if what you're proving is itself disjunctive. Here's a simple example:

**Proposition B.5.** *Suppose  $B \subseteq D$  and  $C \subseteq E$ . Then  $B \cup C \subseteq D \cup E$ .*

*Proof.* Assume (a) that  $B \subseteq D$  and (b)  $C \subseteq E$ . By definition, any  $x \in B$  is also  $\in D$  (c) and any  $x \in C$  is also  $\in E$  (d). To show that  $B \cup C \subseteq D \cup E$ , we have to show that if  $x \in B \cup C$  then  $x \in D \cup E$  (by definition of  $\subseteq$ ).  $x \in B \cup C$  iff  $x \in B$  or  $x \in C$  (by definition of  $\cup$ ). Similarly,  $x \in D \cup E$  iff  $x \in D$  or  $x \in E$ . So, we have to show: for any  $x$ , if  $x \in B$  or  $x \in C$ , then  $x \in D$  or  $x \in E$ .

So far we've only unpacked definitions! We've reformulated our proposition without  $\subseteq$  and  $\cup$  and are left with trying to prove a universal conditional claim. By what we've discussed above, this is done by assuming

that  $x$  is something about which we assume the “if” part is true, and we’ll go on to show that the “then” part is true as well. In other words, we’ll assume that  $x \in B$  or  $x \in C$  and show that  $x \in D$  or  $x \in E$ .<sup>2</sup>

Suppose that  $x \in B$  or  $x \in C$ . We have to show that  $x \in D$  or  $x \in E$ . We distinguish cases.

Case 1:  $x \in B$ . By (c),  $x \in D$ . Thus,  $x \in D$  or  $x \in E$ . (Here we’ve made the inference discussed in the preceding subsection!)

Case 2:  $x \in C$ . By (d),  $x \in E$ . Thus,  $x \in D$  or  $x \in E$ .  $\square$

## Proving an Existence Claim

When asked to prove an existence claim, the question will usually be of the form “prove that there is an  $x$  such that  $\dots x \dots$ ”, i.e., that some object that has the property described by “ $\dots x \dots$ ”. In this case you’ll have to identify a suitable object show that it has the required property. This sounds straightforward, but a proof of this kind can be tricky. Typically it involves *constructing* or *defining* an object and proving that the object so defined has the required property. Finding the right object may be hard, proving that it has the required property may be hard, and sometimes it’s even tricky to show that you’ve succeeded in defining an object at all!

Generally, you’d write this out by specifying the object, e.g., “let  $x$  be  $\dots$ ” (where  $\dots$  specifies which object you have in mind), possibly proving that  $\dots$  in fact describes an object that exists, and then go on to show that  $x$  has the property  $Q$ . Here’s a simple example.

**Proposition B.6.** *Suppose that  $x \in B$ . Then there is an  $A$  such that  $A \subseteq B$  and  $A \neq \emptyset$ .*

*Proof.* Assume  $x \in B$ . Let  $A = \{x\}$ .

<sup>2</sup>This paragraph just explains what we’re doing—it’s not part of the proof, and you don’t have to go into all this detail when you write down your own proofs.

Here we've defined the set  $A$  by enumerating its elements. Since we assume that  $x$  is an object, and we can always form a set by enumerating its elements, we don't have to show that we've succeeded in defining a set  $A$  here. However, we still have to show that  $A$  has the properties required by the proposition. The proof isn't complete without that!

Since  $x \in A$ ,  $A \neq \emptyset$ .

This relies on the definition of  $A$  as  $\{x\}$  and the obvious facts that  $x \in \{x\}$  and  $x \notin \emptyset$ .

Since  $x$  is the only element of  $\{x\}$ , and  $x \in B$ , every element of  $A$  is also an element of  $B$ . By definition of  $\subseteq$ ,  $A \subseteq B$ .  $\square$

## Using Existence Claims

Suppose you know that some existence claim is true (you've proved it, or it's a hypothesis you can use), say, "for some  $x$ ,  $x \in A$ " or "there is an  $x \in A$ ." If you want to use it in your proof, you can just pretend that you have a name for one of the things which your hypothesis says exist. Since  $A$  contains at least one thing, there are things to which that name might refer. You might of course not be able to pick one out or describe it further (other than that it is  $\in A$ ). But for the purpose of the proof, you can pretend that you have picked it out and give a name to it. It's important to pick a name that you haven't already used (or that appears in your hypotheses), otherwise things can go wrong. In your proof, you indicate this by going from "for some  $x$ ,  $x \in A$ " to "Let  $a \in A$ ." Now you can reason about  $a$ , use some other hypotheses, etc., until you come to a conclusion,  $p$ . If  $p$  no longer mentions  $a$ ,  $p$  is independent of the assumption that  $a \in A$ , and you've shown that it follows just from the assumption "for some  $x$ ,  $x \in A$ ."



**Proposition B.7.** *If  $A \neq \emptyset$ , then  $A \cup B \neq \emptyset$ .*

*Proof.* Suppose  $A \neq \emptyset$ . So for some  $x$ ,  $x \in A$ .

Here we first just restated the hypothesis of the proposition. This hypothesis, i.e.,  $A \neq \emptyset$ , hides an existential claim, which you get to only by unpacking a few definitions. The definition of  $=$  tells us that  $A = \emptyset$  iff every  $x \in A$  is also  $\in \emptyset$  and every  $x \in \emptyset$  is also  $\in A$ . Negating both sides, we get:  $A \neq \emptyset$  iff either some  $x \in A$  is  $\notin \emptyset$  or some  $x \in \emptyset$  is  $\notin A$ . Since nothing is  $\in \emptyset$ , the second disjunct can never be true, and “ $x \in A$  and  $x \notin \emptyset$ ” reduces to just  $x \in A$ . So  $x \neq \emptyset$  iff for some  $x$ ,  $x \in A$ . That’s an existence claim. Now we use that existence claim by introducing a name for one of the elements of  $A$ :

Let  $a \in A$ .

Now we’ve introduced a name for one of the things  $\in A$ . We’ll continue to argue about  $a$ , but we’ll be careful to only assume that  $a \in A$  and nothing else:

Since  $a \in A$ ,  $a \in A \cup B$ , by definition of  $\cup$ . So for some  $x$ ,  $x \in A \cup B$ , i.e.,  $A \cup B \neq \emptyset$ .

In that last step, we went from “ $a \in A \cup B$ ” to “for some  $x$ ,  $x \in A \cup B$ .” That doesn’t mention  $a$  anymore, so we know that “for some  $x$ ,  $x \in A \cup B$ ” follows from “for some  $x$ ,  $x \in A$  alone.” But that means that  $A \cup B \neq \emptyset$ . □

It’s maybe good practice to keep bound variables like “ $x$ ” separate from hypothetical names like  $a$ , like we did. In practice, however, we often don’t and just use  $x$ , like so:

Suppose  $A \neq \emptyset$ , i.e., there is an  $x \in A$ . By definition of  $\cup$ ,  $x \in A \cup B$ . So  $A \cup B \neq \emptyset$ .

However, when you do this, you have to be extra careful that you use different  $x$ 's and  $y$ 's for different existential claims. For instance, the following is *not* a correct proof of “If  $A \neq \emptyset$  and  $B \neq \emptyset$  then  $A \cap B \neq \emptyset$ ” (which is not true).

Suppose  $A \neq \emptyset$  and  $B \neq \emptyset$ . So for some  $x$ ,  $x \in A$  and also for some  $x$ ,  $x \in B$ . Since  $x \in A$  and  $x \in B$ ,  $x \in A \cap B$ , by definition of  $\cap$ . So  $A \cap B \neq \emptyset$ .

Can you spot where the incorrect step occurs and explain why the result does not hold?

## B.5 An Example

Our first example is the following simple fact about unions and intersections of sets. It will illustrate unpacking definitions, proofs of conjunctions, of universal claims, and proof by cases.

**Proposition B.8.** *For any sets  $A$ ,  $B$ , and  $C$ ,  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$*

Let's prove it!

*Proof.* We want to show that for any sets  $A$ ,  $B$ , and  $C$ ,  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$

First we unpack the definition of “=” in the statement of the proposition. Recall that proving sets identical means showing that the sets have the same elements. That is, all elements of  $A \cup (B \cap C)$  are also elements of  $(A \cup B) \cap (A \cup C)$ , and vice versa. The “vice versa” means that also every element of  $(A \cup B) \cap (A \cup C)$  must be an element of  $A \cup (B \cap C)$ . So in unpacking the definition, we see that we have to prove a conjunction. Let's record this:

By definition,  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$  iff every element of  $A \cup (B \cap C)$  is also an element of  $(A \cup B) \cap (A \cup C)$ , and every element of  $(A \cup B) \cap (A \cup C)$  is an element of  $A \cup (B \cap C)$ .

Since this is a conjunction, we must prove each conjunct separately. Let's start with the first: let's prove that every element of  $A \cup (B \cap C)$  is also an element of  $(A \cup B) \cap (A \cup C)$ .

This is a universal claim, and so we consider an arbitrary element of  $A \cup (B \cap C)$  and show that it must also be an element of  $(A \cup B) \cap (A \cup C)$ . We'll pick a variable to call this arbitrary element by, say,  $z$ . Our proof continues:

First, we prove that every element of  $A \cup (B \cap C)$  is also an element of  $(A \cup B) \cap (A \cup C)$ . Let  $z \in A \cup (B \cap C)$ . We have to show that  $z \in (A \cup B) \cap (A \cup C)$ .

Now it is time to unpack the definition of  $\cup$  and  $\cap$ . For instance, the definition of  $\cup$  is:  $A \cup B = \{z : z \in A \text{ or } z \in B\}$ . When we apply the definition to " $A \cup (B \cap C)$ ," the role of the " $B$ " in the definition is now played by " $B \cap C$ ," so  $A \cup (B \cap C) = \{z : z \in A \text{ or } z \in B \cap C\}$ . So our assumption that  $z \in A \cup (B \cap C)$  amounts to:  $z \in \{z : z \in A \text{ or } z \in B \cap C\}$ . And  $z \in \{z : \dots z \dots\}$  iff  $\dots z \dots$ , i.e., in this case,  $z \in A$  or  $z \in B \cap C$ .

By the definition of  $\cup$ , either  $z \in A$  or  $z \in B \cap C$ .

Since this is a disjunction, it will be useful to apply proof by cases. We take the two cases, and show that in each one, the conclusion we're aiming for (namely, " $z \in (A \cup B) \cap (A \cup C)$ ") obtains.

Case 1: Suppose that  $z \in A$ .

There's not much more to work from based on our assumptions. So let's look at what we have to work with in the conclusion. We want to show that  $z \in (A \cup B) \cap (A \cup C)$ . Based on the definition of  $\cap$ , if we want to show that  $z \in (A \cup B) \cap (A \cup C)$ , we have to show that it's in both  $(A \cup B)$  and  $(A \cup C)$ . But  $z \in A \cup B$  iff  $z \in A$  or  $z \in B$ , and we already have (as the assumption of case 1) that  $z \in A$ . By the same reasoning—switching  $C$  for  $B$ — $z \in A \cup C$ . This argument went in the reverse direction, so let's record our reasoning in the direction needed in our proof.

Since  $z \in A$ ,  $z \in A$  or  $z \in B$ , and hence, by definition of  $\cup$ ,  $z \in A \cup B$ . Similarly,  $z \in A \cup C$ . But this means that  $z \in (A \cup B) \cap (A \cup C)$ , by definition of  $\cap$ .

This completes the first case of the proof by cases. Now we want to derive the conclusion in the second case, where  $z \in B \cap C$ .

Case 2: Suppose that  $z \in B \cap C$ .

Again, we are working with the intersection of two sets. Let's apply the definition of  $\cap$ :

Since  $z \in B \cap C$ ,  $z$  must be an element of both  $B$  and  $C$ , by definition of  $\cap$ .

It's time to look at our conclusion again. We have to show that  $z$  is in both  $(A \cup B)$  and  $(A \cup C)$ . And again, the solution is immediate.

Since  $z \in B$ ,  $z \in (A \cup B)$ . Since  $z \in C$ , also  $z \in (A \cup C)$ . So,  $z \in (A \cup B) \cap (A \cup C)$ .

Here we applied the definitions of  $\cup$  and  $\cap$  again, but since we've already recalled those definitions, and already showed that if  $z$  is in one of two sets it is in

their union, we don't have to be as explicit in what we've done.

We've completed the second case of the proof by cases, so now we can assert our first conclusion.

So, if  $z \in A \cup (B \cap C)$  then  $z \in (A \cup B) \cap (A \cup C)$ .

Now we just want to show the other direction, that every element of  $(A \cup B) \cap (A \cup C)$  is an element of  $A \cup (B \cap C)$ . As before, we prove this universal claim by assuming we have an arbitrary element of the first set and show it must be in the second set. Let's state what we're about to do.

Now, assume that  $z \in (A \cup B) \cap (A \cup C)$ . We want to show that  $z \in A \cup (B \cap C)$ .

We are now working from the hypothesis that  $z \in (A \cup B) \cap (A \cup C)$ . It hopefully isn't too confusing that we're using the same  $z$  here as in the first part of the proof. When we finished that part, all the assumptions we've made there are no longer in effect, so now we can make new assumptions about what  $z$  is. If that is confusing to you, just replace  $z$  with a different variable in what follows.

We know that  $z$  is in both  $A \cup B$  and  $A \cup C$ , by definition of  $\cap$ . And by the definition of  $\cup$ , we can further unpack this to: either  $z \in A$  or  $z \in B$ , and also either  $z \in A$  or  $z \in C$ . This looks like a proof by cases again—except the “and” makes it confusing. You might think that this amounts to there being three possibilities:  $z$  is either in  $A$ ,  $B$  or  $C$ . But that would be a mistake. We have to be careful, so let's consider each disjunction in turn.

By definition of  $\cap$ ,  $z \in A \cup B$  and  $z \in A \cup C$ . By definition of  $\cup$ ,  $z \in A$  or  $z \in B$ . We distinguish cases.

Since we're focusing on the first disjunction, we haven't gotten our second disjunction (from unpacking  $A \cup C$ ) yet. In fact, we don't need it yet. The first case is  $z \in A$ , and an element of a set is also an element of the union of that set with any other. So case 1 is easy:

Case 1: Suppose that  $z \in A$ . It follows that  $z \in A \cup (B \cap C)$ .

Now for the second case,  $z \in B$ . Here we'll unpack the second  $\cup$  and do another proof-by-cases:

Case 2: Suppose that  $z \in B$ . Since  $z \in A \cup C$ , either  $z \in A$  or  $z \in C$ . We distinguish cases further:

Case 2a:  $z \in A$ . Then, again,  $z \in A \cup (B \cap C)$ .

Ok, this was a bit weird. We didn't actually need the assumption that  $z \in B$  for this case, but that's ok.

Case 2b:  $z \in C$ . Then  $z \in B$  and  $z \in C$ , so  $z \in B \cap C$ , and consequently,  $z \in A \cup (B \cap C)$ .

This concludes both proofs-by-cases and so we're done with the second half.

So, if  $z \in (A \cup B) \cap (A \cup C)$  then  $z \in A \cup (B \cap C)$ . □

## B.6 Another Example

**Proposition B.9.** *If  $A \subseteq C$ , then  $A \cup (C \setminus A) = C$ .*

*Proof.* Suppose that  $A \subseteq C$ . We want to show that  $A \cup (C \setminus A) = C$ .

We begin by observing that this is a conditional statement. It is tacitly universally quantified: the proposition holds for all sets  $A$  and  $C$ . So  $A$  and  $C$  are variables for arbitrary sets. To prove such a statement, we assume the antecedent and prove the consequent.

We continue by using the assumption that  $A \subseteq C$ . Let's unpack the definition of  $\subseteq$ : the assumption means that all elements of  $A$  are also elements of  $C$ . Let's write this down—it's an important fact that we'll use throughout the proof.

By the definition of  $\subseteq$ , since  $A \subseteq C$ , for all  $z$ , if  $z \in A$ , then  $z \in C$ .

We've unpacked all the definitions that are given to us in the assumption. Now we can move onto the conclusion. We want to show that  $A \cup (C \setminus A) = C$ , and so we set up a proof similarly to the last example: we show that every element of  $A \cup (C \setminus A)$  is also an element of  $C$  and, conversely, every element of  $C$  is an element of  $A \cup (C \setminus A)$ . We can shorten this to:  $A \cup (C \setminus A) \subseteq C$  and  $C \subseteq A \cup (C \setminus A)$ . (Here we're doing the opposite of unpacking a definition, but it makes the proof a bit easier to read.) Since this is a conjunction, we have to prove both parts. To show the first part, i.e., that every element of  $A \cup (C \setminus A)$  is also an element of  $C$ , we assume that  $z \in A \cup (C \setminus A)$  for an arbitrary  $z$  and show that  $z \in C$ . By the definition of  $\cup$ , we can conclude that  $z \in A$  or  $z \in C \setminus A$  from  $z \in A \cup (C \setminus A)$ . You should now be getting the hang of this.

$A \cup (C \setminus A) = C$  iff  $A \cup (C \setminus A) \subseteq C$  and  $C \subseteq (A \cup (C \setminus A))$ . First we prove that  $A \cup (C \setminus A) \subseteq C$ . Let  $z \in A \cup (C \setminus A)$ . So, either  $z \in A$  or  $z \in (C \setminus A)$ .

We've arrived at a disjunction, and from it we want to prove that  $z \in C$ . We do this using proof by cases.

Case 1:  $z \in A$ . Since for all  $z$ , if  $z \in A$ ,  $z \in C$ , we have that  $z \in C$ .

Here we've used the fact recorded earlier which followed from the hypothesis of the proposition that  $A \subseteq C$ . The first case is complete, and we turn to

the second case,  $z \in (C \setminus A)$ . Recall that  $C \setminus A$  denotes the *difference* of the two sets, i.e., the set of all elements of  $C$  which are not elements of  $A$ . But any element of  $C$  not in  $A$  is in particular an element of  $C$ .

Case 2:  $z \in (C \setminus A)$ . This means that  $z \in C$  and  $z \notin A$ . So, in particular,  $z \in C$ .

Great, we've proved the first direction. Now for the second direction. Here we prove that  $C \subseteq A \cup (C \setminus A)$ . So we assume that  $z \in C$  and prove that  $z \in A \cup (C \setminus A)$ .

Now let  $z \in C$ . We want to show that  $z \in A$  or  $z \in C \setminus A$ .

Since all elements of  $A$  are also elements of  $C$ , and  $C \setminus A$  is the set of all things that are elements of  $C$  but not  $A$ , it follows that  $z$  is either in  $A$  or in  $C \setminus A$ . This may be a bit unclear if you don't already know why the result is true. It would be better to prove it step-by-step. It will help to use a simple fact which we can state without proof:  $z \in A$  or  $z \notin A$ . This is called the "principle of excluded middle:" for any statement  $p$ , either  $p$  is true or its negation is true. (Here,  $p$  is the statement that  $z \in A$ .) Since this is a disjunction, we can again use proof-by-cases.

Either  $z \in A$  or  $z \notin A$ . In the former case,  $z \in A \cup (C \setminus A)$ . In the latter case,  $z \in C$  and  $z \notin A$ , so  $z \in C \setminus A$ . But then  $z \in A \cup (C \setminus A)$ .

Our proof is complete: we have shown that  $A \cup (C \setminus A) = C$ . □

## B.7 Proof by Contradiction

In the first instance, proof by contradiction is an inference pattern that is used to prove negative claims. Suppose you want to



show that some claim  $p$  is *false*, i.e., you want to show  $\neg p$ . The most promising strategy is to (a) suppose that  $p$  is true, and (b) show that this assumption leads to something you know to be false. “Something known to be false” may be a result that conflicts with—contradicts— $p$  itself, or some other hypothesis of the overall claim you are considering. For instance, a proof of “if  $q$  then  $\neg p$ ” involves assuming that  $q$  is true and proving  $\neg p$  from it. If you prove  $\neg p$  by contradiction, that means assuming  $p$  in addition to  $q$ . If you can prove  $\neg q$  from  $p$ , you have shown that the assumption  $p$  leads to something that contradicts your other assumption  $q$ , since  $q$  and  $\neg q$  cannot both be true. Of course, you have to use other inference patterns in your proof of the contradiction, as well as unpacking definitions. Let’s consider an example.

**Proposition B.10.** *If  $A \subseteq B$  and  $B = \emptyset$ , then  $A$  has no elements.*

*Proof.* Suppose  $A \subseteq B$  and  $B = \emptyset$ . We want to show that  $A$  has no elements.

Since this is a conditional claim, we assume the antecedent and want to prove the consequent. The consequent is:  $A$  has no elements. We can make that a bit more explicit: it’s not the case that there is an  $x \in A$ .

$A$  has no elements iff it’s not the case that there is an  $x$  such that  $x \in A$ .

So we’ve determined that what we want to prove is really a negative claim  $\neg p$ , namely: it’s not the case that there is an  $x \in A$ . To use proof by contradiction, we have to assume the corresponding positive claim  $p$ , i.e., there is an  $x \in A$ , and prove a contradiction from it. We indicate that we’re doing a proof by contradiction by writing “by way of contradiction, assume” or even just “suppose not,” and then state the assumption  $p$ .

Suppose not: there is an  $x \in A$ .

This is now the new assumption we'll use to obtain a contradiction. We have two more assumptions: that  $A \subseteq B$  and that  $B = \emptyset$ . The first gives us that  $x \in B$ :

Since  $A \subseteq B$ ,  $x \in B$ .

But since  $B = \emptyset$ , every element of  $B$  (e.g.,  $x$ ) must also be an element of  $\emptyset$ .

Since  $B = \emptyset$ ,  $x \in \emptyset$ . This is a contradiction, since by definition  $\emptyset$  has no elements.

This already completes the proof: we've arrived at what we need (a contradiction) from the assumptions we've set up, and this means that the assumptions can't all be true. Since the first two assumptions ( $A \subseteq B$  and  $B = \emptyset$ ) are not contested, it must be the last assumption introduced (there is an  $x \in A$ ) that must be false. But if we want to be thorough, we can spell this out.

Thus, our assumption that there is an  $x \in A$  must be false, hence,  $A$  has no elements by proof by contradiction.  $\square$

Every positive claim is trivially equivalent to a negative claim:  $p$  iff  $\neg\neg p$ . So proofs by contradiction can also be used to establish positive claims "indirectly," as follows: To prove  $p$ , read it as the negative claim  $\neg\neg p$ . If we can prove a contradiction from  $\neg p$ , we've established  $\neg\neg p$  by proof by contradiction, and hence  $p$ .

In the last example, we aimed to prove a negative claim, namely that  $A$  has no elements, and so the assumption we made for the purpose of proof by contradiction (i.e., that there is an  $x \in A$ ) was a positive claim. It gave us something to work with, namely the hypothetical  $x \in A$  about which we continued to reason until we got to  $x \in \emptyset$ .

When proving a positive claim indirectly, the assumption you'd make for the purpose of proof by contradiction would be negative. But very often you can easily reformulate a positive claim as a negative claim, and a negative claim as a positive claim. Our previous proof would have been essentially the same had we proved " $A = \emptyset$ " instead of the negative consequent " $A$  has no elements." (By definition of  $=$ , " $A = \emptyset$ " is a general claim, since it unpacks to "every element of  $A$  is an element of  $\emptyset$  and vice versa".) But it is easily seen to be equivalent to the negative claim "not: there is an  $x \in A$ ."

So it is sometimes easier to work with  $\neg p$  as an assumption than it is to prove  $p$  directly. Even when a direct proof is just as simple or even simpler (as in the next examples), some people prefer to proceed indirectly. If the double negation confuses you, think of a proof by contradiction of some claim as a proof of a contradiction from the *opposite* claim. So, a proof by contradiction of  $\neg p$  is a proof of a contradiction from the assumption  $p$ ; and proof by contradiction of  $p$  is a proof of a contradiction from  $\neg p$ .

**Proposition B.11.**  $A \subseteq A \cup B$ .

*Proof.* We want to show that  $A \subseteq A \cup B$ .

On the face of it, this is a positive claim: every  $x \in A$  is also in  $A \cup B$ . The negation of that is: some  $x \in A$  is  $\notin A \cup B$ . So we can prove the claim indirectly by assuming this negated claim, and showing that it leads to a contradiction.

Suppose not, i.e.,  $A \not\subseteq A \cup B$ .

We have a definition of  $A \subseteq A \cup B$ : every  $x \in A$  is also  $\in A \cup B$ . To understand what  $A \not\subseteq A \cup B$  means, we have to use some elementary logical manipulation on the unpacked definition: it's false that every  $x \in A$  is also  $\in A \cup B$  iff there is *some*  $x \in A$  that is  $\notin C$ . (This is a place where you want to be very careful:

many students' attempted proofs by contradiction fail because they analyze the negation of a claim like "all  $A$ s are  $B$ s" incorrectly.) In other words,  $A \not\subseteq A \cup B$  iff there is an  $x$  such that  $x \in A$  and  $x \notin A \cup B$ . From then on, it's easy.

So, there is an  $x \in A$  such that  $x \notin A \cup B$ . By definition of  $\cup$ ,  $x \in A \cup B$  iff  $x \in A$  or  $x \in B$ . Since  $x \in A$ , we have  $x \in A \cup B$ . This contradicts the assumption that  $x \notin A \cup B$ .  $\square$

**Proposition B.12.** *If  $A \subseteq B$  and  $B \subseteq C$  then  $A \subseteq C$ .*

*Proof.* Suppose  $A \subseteq B$  and  $B \subseteq C$ . We want to show  $A \subseteq C$ .

Let's proceed indirectly: we assume the negation of what we want to establish.

Suppose not, i.e.,  $A \not\subseteq C$ .

As before, we reason that  $A \not\subseteq C$  iff not every  $x \in A$  is also  $\in C$ , i.e., some  $x \in A$  is  $\notin C$ . Don't worry, with practice you won't have to think hard anymore to unpack negations like this.

In other words, there is an  $x$  such that  $x \in A$  and  $x \notin C$ .

Now we can use this to get to our contradiction. Of course, we'll have to use the other two assumptions to do it.

Since  $A \subseteq B$ ,  $x \in B$ . Since  $B \subseteq C$ ,  $x \in C$ . But this contradicts  $x \notin C$ .  $\square$

**Proposition B.13.** *If  $A \cup B = A \cap B$  then  $A = B$ .*

*Proof.* Suppose  $A \cup B = A \cap B$ . We want to show that  $A = B$ .

The beginning is now routine:

Assume, by way of contradiction, that  $A \neq B$ .

Our assumption for the proof by contradiction is that  $A \neq B$ . Since  $A = B$  iff  $A \subseteq B$  and  $B \subseteq A$ , we get that  $A \neq B$  iff  $A \not\subseteq B$  or  $B \not\subseteq A$ . (Note how important it is to be careful when manipulating negations!) To prove a contradiction from this disjunction, we use a proof by cases and show that in each case, a contradiction follows.

$A \neq B$  iff  $A \not\subseteq B$  or  $B \not\subseteq A$ . We distinguish cases.

In the first case, we assume  $A \not\subseteq B$ , i.e., for some  $x$ ,  $x \in A$  but  $x \notin B$ .  $A \cap B$  is defined as those elements that  $A$  and  $B$  have in common, so if something isn't in one of them, it's not in the intersection.  $A \cup B$  is  $A$  together with  $B$ , so anything in either is also in the union. This tells us that  $x \in A \cup B$  but  $x \notin A \cap B$ , and hence that  $A \cap B \neq A \cup B$ .

Case 1:  $A \not\subseteq B$ . Then for some  $x$ ,  $x \in A$  but  $x \notin B$ . Since  $x \notin B$ , then  $x \notin A \cap B$ . Since  $x \in A$ ,  $x \in A \cup B$ . So,  $A \cap B \neq A \cup B$ , contradicting the assumption that  $A \cap B = A \cup B$ .

Case 2:  $B \not\subseteq A$ . Then for some  $y$ ,  $y \in B$  but  $y \notin A$ . As before, we have  $y \in A \cup B$  but  $y \notin A \cap B$ , and so  $A \cap B \neq A \cup B$ , again contradicting  $A \cap B = A \cup B$ .  $\square$

## B.8 Reading Proofs

Proofs you find in textbooks and articles very seldom give all the details we have so far included in our examples. Authors often

do not draw attention to when they distinguish cases, when they give an indirect proof, or don't mention that they use a definition. So when you read a proof in a textbook, you will often have to fill in those details for yourself in order to understand the proof. Doing this is also good practice to get the hang of the various moves you have to make in a proof. Let's look at an example.

**Proposition B.14 (Absorption).** *For all sets  $A, B$ ,*

$$A \cap (A \cup B) = A$$

*Proof.* If  $z \in A \cap (A \cup B)$ , then  $z \in A$ , so  $A \cap (A \cup B) \subseteq A$ . Now suppose  $z \in A$ . Then also  $z \in A \cup B$ , and therefore also  $z \in A \cap (A \cup B)$ .  $\square$

The preceding proof of the absorption law is very condensed. There is no mention of any definitions used, no “we have to prove that” before we prove it, etc. Let's unpack it. The proposition proved is a general claim about any sets  $A$  and  $B$ , and when the proof mentions  $A$  or  $B$ , these are variables for arbitrary sets. The general claims the proof establishes is what's required to prove identity of sets, i.e., that every element of the left side of the identity is an element of the right and vice versa.

“If  $z \in A \cap (A \cup B)$ , then  $z \in A$ , so  $A \cap (A \cup B) \subseteq A$ .”

This is the first half of the proof of the identity: it establishes that if an arbitrary  $z$  is an element of the left side, it is also an element of the right, i.e.,  $A \cap (A \cup B) \subseteq A$ . Assume that  $z \in A \cap (A \cup B)$ . Since  $z$  is an element of the intersection of two sets iff it is an element of both sets, we can conclude that  $z \in A$  and also  $z \in A \cup B$ . In particular,  $z \in A$ , which is what we wanted to show. Since that's all that has to be done for the first half, we know that the rest of the proof must be a proof of the second half, i.e., a proof that  $A \subseteq A \cap (A \cup B)$ .

“Now suppose  $z \in A$ . Then also  $z \in A \cup B$ , and therefore also  $z \in A \cap (A \cup B)$ .”

We start by assuming that  $z \in A$ , since we are showing that, for any  $z$ , if  $z \in A$  then  $z \in A \cap (A \cup B)$ . To show that  $z \in A \cap (A \cup B)$ , we have to show (by definition of “ $\cap$ ”) that (i)  $z \in A$  and also (ii)  $z \in A \cup B$ . Here (i) is just our assumption, so there is nothing further to prove, and that’s why the proof does not mention it again. For (ii), recall that  $z$  is an element of a union of sets iff it is an element of at least one of those sets. Since  $z \in A$ , and  $A \cup B$  is the union of  $A$  and  $B$ , this is the case here. So  $z \in A \cup B$ . We’ve shown both (i)  $z \in A$  and (ii)  $z \in A \cup B$ , hence, by definition of “ $\cap$ ,”  $z \in A \cap (A \cup B)$ . The proof doesn’t mention those definitions; it’s assumed the reader has already internalized them. If you haven’t, you’ll have to go back and remind yourself what they are. Then you’ll also have to recognize why it follows from  $z \in A$  that  $z \in A \cup B$ , and from  $z \in A$  and  $z \in A \cup B$  that  $z \in A \cap (A \cup B)$ .

Here’s another version of the proof above, with everything made explicit:

*Proof.* [By definition of  $=$  for sets,  $A \cap (A \cup B) = A$  we have to show (a)  $A \cap (A \cup B) \subseteq A$  and (b)  $A \cap (A \cup B) \subseteq A$ . (a): By definition of  $\subseteq$ , we have to show that if  $z \in A \cap (A \cup B)$ , then  $z \in A$ .] If  $z \in A \cap (A \cup B)$ , then  $z \in A$  [since by definition of  $\cap$ ,  $z \in A \cap (A \cup B)$  iff  $z \in A$  and  $z \in A \cup B$ ], so  $A \cap (A \cup B) \subseteq A$ . [(b): By definition of  $\subseteq$ , we have to show that if  $z \in A$ , then  $z \in A \cap (A \cup B)$ .] Now suppose [(1)]  $z \in A$ . Then also [(2)]  $z \in A \cup B$  [since by (1)  $z \in A$  or  $z \in B$ , which by definition of  $\cup$  means  $z \in A \cup B$ ], and therefore also  $z \in A \cap (A \cup B)$  [since the definition of  $\cap$  requires that  $z \in A$ , i.e., (1), and  $z \in A \cup B$ , i.e., (2)].  $\square$

## B.9 I Can’t Do It!

We all get to a point where we feel like giving up. But you *can* do it. Your instructor and teaching assistant, as well as your fellow students, can help. Ask them for help! Here are a few tips to help you avoid a crisis, and what to do if you feel like giving up.

To make sure you can solve problems successfully, do the following:

1. *Start as far in advance as possible.* We get busy throughout the semester and many of us struggle with procrastination, one of the best things you can do is to start your homework assignments early. That way, if you're stuck, you have time to look for a solution (that isn't crying).
2. *Talk to your classmates.* You are not alone. Others in the class may also struggle—but they may struggle with different things. Talking it out with your peers can give you a different perspective on the problem that might lead to a breakthrough. Of course, don't just copy their solution: ask them for a hint, or explain where you get stuck and ask them for the next step. And when you do get it, reciprocate. Helping someone else along, and explaining things will help you understand better, too.
3. *Ask for help.* You have many resources available to you—your instructor and teaching assistant are there for you and *want* you to succeed. They should be able to help you work out a problem and identify where in the process you're struggling.
4. *Take a break.* If you're stuck, it *might* be because you've been staring at the problem for too long. Take a short break, have a cup of tea, or work on a different problem for a while, then return to the problem with a fresh mind. Sleep on it.

Notice how these strategies require that you've started to work on the proof well in advance? If you've started the proof at 2am the day before it's due, these might not be so helpful.

This might sound like doom and gloom, but solving a proof is a challenge that pays off in the end. Some people do this as a career—so there must be something to enjoy about it. Like



basically everything, solving problems and doing proofs is something that requires practice. You might see classmates who find this easy: they've probably just had lots of practice already. Try not to give in too easily.

If you do run out of time (or patience) on a particular problem: that's ok. It doesn't mean you're stupid or that you will never get it. Find out (from your instructor or another student) how it is done, and identify where you went wrong or got stuck, so you can avoid doing that the next time you encounter a similar issue. Then try to do it without looking at the solution. And next time, start (and ask for help) earlier.

## B.10 Other Resources

There are many books on how to do proofs in mathematics which may be useful. Check out *How to Read and do Proofs: An Introduction to Mathematical Thought Processes* (Solow, 2013) and *How to Prove It: A Structured Approach* (Velleman, 2019) in particular. The *Book of Proof* (Hammack, 2013) and *Mathematical Reasoning* (Sandstrum, 2019) are books on proof that are freely available online. Philosophers might find *More Precisely: The Math you need to do Philosophy* (Steinhart, 2018) to be a good primer on mathematical reasoning.

There are also various shorter guides to proofs available on the internet; e.g., “Introduction to Mathematical Arguments” (Hutchings, 2003) and “How to write proofs” (Cheng, 2004).

### Motivational Videos

Feel like you have no motivation to do your homework? Feeling down? These videos might help!

- [https://www.youtube.com/watch?v=ZXsQAXx\\_ao0](https://www.youtube.com/watch?v=ZXsQAXx_ao0)
- <https://www.youtube.com/watch?v=BQ4yd2W50No>
- <https://www.youtube.com/watch?v=StTqXEQ2l-Y>

## Problems

**Problem B.1.** Suppose you are asked to prove that  $A \cap B \neq \emptyset$ . Unpack all the definitions occurring here, i.e., restate this in a way that does not mention “ $\cap$ ”, “ $=$ ”, or “ $\emptyset$ ”.

**Problem B.2.** Prove *indirectly* that  $A \cap B \subseteq A$ .

**Problem B.3.** Expand the following proof of  $A \cup (A \cap B) = A$ , where you mention all the inference patterns used, why each step follows from assumptions or claims established before it, and where we have to appeal to which definitions.

*Proof.* If  $z \in A \cup (A \cap B)$  then  $z \in A$  or  $z \in A \cap B$ . If  $z \in A \cap B$ ,  $z \in A$ . Any  $z \in A$  is also  $\in A \cup (A \cap B)$ .  $\square$

## APPENDIX C

# *Induction*

### C.1 Introduction

Induction is an important proof technique which is used, in different forms, in almost all areas of logic, theoretical computer science, and mathematics. It is needed to prove many of the results in logic.

Induction is often contrasted with deduction, and characterized as the inference from the particular to the general. For instance, if we observe many green emeralds, and nothing that we would call an emerald that's not green, we might conclude that all emeralds are green. This is an inductive inference, in that it proceeds from many particular cases (this emerald is green, that emerald is green, etc.) to a general claim (all emeralds are green). *Mathematical* induction is also an inference that concludes a general claim, but it is of a very different kind than this "simple induction."

Very roughly, an inductive proof in mathematics concludes that all mathematical objects of a certain sort have a certain property. In the simplest case, the mathematical objects an inductive proof is concerned with are natural numbers. In that case an inductive proof is used to establish that all natural numbers have some property, and it does this by showing that

1. 0 has the property, and

2. whenever a number  $k$  has the property, so does  $k + 1$ .

Induction on natural numbers can then also often be used to prove general claims about mathematical objects that can be assigned numbers. For instance, finite sets each have a finite number  $n$  of elements, and if we can use induction to show that every number  $n$  has the property “all finite sets of size  $n$  are ...” then we will have shown something about all finite sets.

Induction can also be generalized to mathematical objects that are *inductively defined*. For instance, expressions of a formal language such as those of first-order logic are defined inductively. *Structural induction* is a way to prove results about all such expressions. Structural induction, in particular, is very useful—and widely used—in logic.

## C.2 Induction on $\mathbb{N}$

In its simplest form, induction is a technique used to prove results for all natural numbers. It uses the fact that by starting from 0 and repeatedly adding 1 we eventually reach every natural number. So to prove that something is true for every number, we can (1) establish that it is true for 0 and (2) show that whenever it is true for a number  $n$ , it is also true for the next number  $n+1$ . If we abbreviate “number  $n$  has property  $P$ ” by  $P(n)$  (and “number  $k$  has property  $P$ ” by  $P(k)$ , etc.), then a proof by induction that  $P(n)$  for all  $n \in \mathbb{N}$  consists of:

1. a proof of  $P(0)$ , and
2. a proof that, for any  $k$ , if  $P(k)$  then  $P(k + 1)$ .

To make this crystal clear, suppose we have both (1) and (2). Then (1) tells us that  $P(0)$  is true. If we also have (2), we know in particular that if  $P(0)$  then  $P(0 + 1)$ , i.e.,  $P(1)$ . This follows from the general statement “for any  $k$ , if  $P(k)$  then  $P(k + 1)$ ” by putting 0 for  $k$ . So by modus ponens, we have that  $P(1)$ . From (2) again, now taking 1 for  $n$ , we have: if  $P(1)$  then  $P(2)$ . Since we’ve

just established  $P(1)$ , by modus ponens, we have  $P(2)$ . And so on. For any number  $n$ , after doing this  $n$  times, we eventually arrive at  $P(n)$ . So (1) and (2) together establish  $P(n)$  for any  $n \in \mathbb{N}$ .

Let's look at an example. Suppose we want to find out how many different sums we can throw with  $n$  dice. Although it might seem silly, let's start with 0 dice. If you have no dice there's only one possible sum you can "throw": no dots at all, which sums to 0. So the number of different possible throws is 1. If you have only one die, i.e.,  $n = 1$ , there are six possible values, 1 through 6. With two dice, we can throw any sum from 2 through 12, that's 11 possibilities. With three dice, we can throw any number from 3 to 18, i.e., 16 different possibilities. 1, 6, 11, 16: looks like a pattern: maybe the answer is  $5n + 1$ ? Of course,  $5n + 1$  is the maximum possible, because there are only  $5n + 1$  numbers between  $n$ , the lowest value you can throw with  $n$  dice (all 1's) and  $6n$ , the highest you can throw (all 6's).

**Theorem C.1.** *With  $n$  dice one can throw all  $5n + 1$  possible values between  $n$  and  $6n$ .*

*Proof.* Let  $P(n)$  be the claim: "It is possible to throw any number between  $n$  and  $6n$  using  $n$  dice." To use induction, we prove:

1. The *induction basis*  $P(1)$ , i.e., with just one die, you can throw any number between 1 and 6.
2. The *induction step*, for all  $k$ , if  $P(k)$  then  $P(k + 1)$ .

(1) Is proved by inspecting a 6-sided die. It has all 6 sides, and every number between 1 and 6 shows up one on of the sides. So it is possible to throw any number between 1 and 6 using a single die.

To prove (2), we assume the antecedent of the conditional, i.e.,  $P(k)$ . This assumption is called the *inductive hypothesis*. We use it to prove  $P(k + 1)$ . The hard part is to find a way of thinking about the possible values of a throw of  $k + 1$  dice in terms of the

possible values of throws of  $k$  dice plus of throws of the extra  $k + 1$ -st die—this is what we have to do, though, if we want to use the inductive hypothesis.

The inductive hypothesis says we can get any number between  $k$  and  $6k$  using  $k$  dice. If we throw a 1 with our  $(k + 1)$ -st die, this adds 1 to the total. So we can throw any value between  $k + 1$  and  $6k + 1$  by throwing  $k$  dice and then rolling a 1 with the  $(k + 1)$ -st die. What's left? The values  $6k + 2$  through  $6k + 6$ . We can get these by rolling  $k$  6s and then a number between 2 and 6 with our  $(k + 1)$ -st die. Together, this means that with  $k + 1$  dice we can throw any of the numbers between  $k + 1$  and  $6(k + 1)$ , i.e., we've proved  $P(k + 1)$  using the assumption  $P(k)$ , the inductive hypothesis.  $\square$

Very often we use induction when we want to prove something about a series of objects (numbers, sets, etc.) that is itself defined “inductively,” i.e., by defining the  $(n + 1)$ -st object in terms of the  $n$ -th. For instance, we can define the sum  $s_n$  of the natural numbers up to  $n$  by

$$\begin{aligned}s_0 &= 0 \\ s_{n+1} &= s_n + (n + 1)\end{aligned}$$

This definition gives:

$$\begin{aligned}s_0 &= 0, \\ s_1 &= s_0 + 1 &&= 1, \\ s_2 &= s_1 + 2 &&= 1 + 2 = 3 \\ s_3 &= s_2 + 3 &&= 1 + 2 + 3 = 6, \text{ etc.}\end{aligned}$$

Now we can prove, by induction, that  $s_n = n(n + 1)/2$ .

**Proposition C.2.**  $s_n = n(n + 1)/2$ .

*Proof.* We have to prove (1) that  $s_0 = 0 \cdot (0 + 1)/2$  and (2) if  $s_k = k(k + 1)/2$  then  $s_{k+1} = (k + 1)(k + 2)/2$ . (1) is obvious. To

prove (2), we assume the inductive hypothesis:  $s_k = k(k+1)/2$ . Using it, we have to show that  $s_{k+1} = (k+1)(k+2)/2$ .

What is  $s_{k+1}$ ? By the definition,  $s_{k+1} = s_k + (k+1)$ . By inductive hypothesis,  $s_k = k(k+1)/2$ . We can substitute this into the previous equation, and then just need a bit of arithmetic of fractions:

$$\begin{aligned} s_{k+1} &= \frac{k(k+1)}{2} + (k+1) = \\ &= \frac{k(k+1)}{2} + \frac{2(k+1)}{2} = \\ &= \frac{k(k+1) + 2(k+1)}{2} = \\ &= \frac{(k+2)(k+1)}{2}. \quad \square \end{aligned}$$

The important lesson here is that if you're proving something about some inductively defined sequence  $a_n$ , induction is the obvious way to go. And even if it isn't (as in the case of the possibilities of dice throws), you can use induction if you can somehow relate the case for  $k+1$  to the case for  $k$ .

### C.3 Strong Induction

In the principle of induction discussed above, we prove  $P(0)$  and also if  $P(k)$ , then  $P(k+1)$ . In the second part, we assume that  $P(k)$  is true and use this assumption to prove  $P(k+1)$ . Equivalently, of course, we could assume  $P(k-1)$  and use it to prove  $P(k)$ —the important part is that we be able to carry out the inference from any number to its successor; that we can prove the claim in question for any number under the assumption it holds for its predecessor.

There is a variant of the principle of induction in which we don't just assume that the claim holds for the predecessor  $k-1$  of  $k$ , but for all numbers smaller than  $k$ , and use this assumption to establish the claim for  $k$ . This also gives us the claim  $P(n)$  for all  $n \in \mathbb{N}$ . For once we have established  $P(0)$ , we have

thereby established that  $P$  holds for all numbers less than 1. And if we know that if  $P(l)$  for all  $l < k$ , then  $P(k)$ , we know this in particular for  $k = 1$ . So we can conclude  $P(1)$ . With this we have proved  $P(0)$  and  $P(1)$ , i.e.,  $P(l)$  for all  $l < 2$ , and since we have also the conditional, if  $P(l)$  for all  $l < 2$ , then  $P(2)$ , we can conclude  $P(2)$ , and so on.

In fact, if we can establish the general conditional “for all  $k$ , if  $P(l)$  for all  $l < k$ , then  $P(k)$ ,” we do not have to establish  $P(0)$  anymore, since it follows from it. For remember that a general claim like “for all  $l < k$ ,  $P(l)$ ” is true if there are no  $l < k$ . This is a case of vacuous quantification: “all  $As$  are  $Bs$ ” is true if there are no  $As$ ,  $\forall x (A(x) \rightarrow B(x))$  is true if no  $x$  satisfies  $A(x)$ . In this case, the formalized version would be “ $\forall l (l < k \rightarrow P(l))$ ”—and that is true if there are no  $l < k$ . And if  $k = 0$  that’s exactly the case: no  $l < 0$ , hence “for all  $l < 0$ ,  $P(0)$ ” is true, whatever  $P$  is. A proof of “if  $P(l)$  for all  $l < k$ , then  $P(k)$ ” thus automatically establishes  $P(0)$ .

This variant is useful if establishing the claim for  $k$  can’t be made to just rely on the claim for  $k - 1$  but may require the assumption that it is true for one or more  $l < k$ .

## C.4 Inductive Definitions

In logic we very often define kinds of objects *inductively*, i.e., by specifying rules for what counts as an object of the kind to be defined which explain how to get new objects of that kind from old objects of that kind. For instance, we often define special kinds of sequences of symbols, such as the terms and formulas of a language, by induction. For a simple example, consider strings of consisting of letters  $a$ ,  $b$ ,  $c$ ,  $d$ , the symbol  $\circ$ , and brackets  $[$  and  $]$ , such as “ $[[c\circ d][$ ”, “ $[a[]\circ]$ ”, “ $a$ ” or “ $[[a\circ b]\circ d]$ ”. You probably feel that there’s something “wrong” with the first two strings: the brackets don’t “balance” at all in the first, and you might feel that the “ $\circ$ ” should “connect” expressions that themselves make sense. The third and fourth string look better: for every “[” there’s a



closing “]” (if there are any at all), and for any  $\circ$  we can find “nice” expressions on either side, surrounded by a pair of parentheses.

We would like to precisely specify what counts as a “nice term.” First of all, every letter by itself is nice. Anything that’s not just a letter by itself should be of the form “[ $t \circ s$ ]” where  $s$  and  $t$  are themselves nice. Conversely, if  $t$  and  $s$  are nice, then we can form a new nice term by putting a  $\circ$  between them and surround them by a pair of brackets. We might use these operations to *define* the set of nice terms. This is an *inductive definition*.

**Definition C.3 (Nice terms).** The set of *nice terms* is inductively defined as follows:

1. Any letter  $a, b, c, d$  is a nice term.
2. If  $s_1$  and  $s_2$  are nice terms, then so is [ $s_1 \circ s_2$ ].
3. Nothing else is a nice term.

This definition tells us that something counts as a nice term iff it can be constructed according to the two conditions (1) and (2) in some finite number of steps. In the first step, we construct all nice terms just consisting of letters by themselves, i.e.,

$$a, b, c, d$$

In the second step, we apply (2) to the terms we’ve constructed. We’ll get

$$[a \circ a], [a \circ b], [b \circ a], \dots, [d \circ d]$$

for all combinations of two letters. In the third step, we apply (2) again, to any two nice terms we’ve constructed so far. We get new nice term such as [ $a \circ [a \circ a]$ ]—where  $t$  is  $a$  from step 1 and  $s$  is [ $a \circ a$ ] from step 2—and [[ $b \circ c$ ]  $\circ$  [ $d \circ b$ ]] constructed out of the two terms [ $b \circ c$ ] and [ $d \circ b$ ] from step 2. And so on. Clause (3) rules out that anything not constructed in this way sneaks into the set of nice terms.

Note that we have not yet proved that every sequence of symbols that “feels” nice is nice according to this definition. However, it should be clear that everything we can construct does in fact “feel nice”: brackets are balanced, and  $\circ$  connects parts that are themselves nice.

The key feature of inductive definitions is that if you want to prove something about all nice terms, the definition tells you which cases you must consider. For instance, if you are told that  $t$  is a nice term, the inductive definition tells you what  $t$  can look like:  $t$  can be a letter, or it can be  $[s_1 \circ s_2]$  for some pair of nice terms  $s_1$  and  $s_2$ . Because of clause (3), those are the only possibilities.

When proving claims about all of an inductively defined set, the strong form of induction becomes particularly important. For instance, suppose we want to prove that for every nice term of length  $n$ , the number of  $[$  in it is  $< n/2$ . This can be seen as a claim about all  $n$ : for every  $n$ , the number of  $[$  in any nice term of length  $n$  is  $< n/2$ .

**Proposition C.4.** *For any  $n$ , the number of  $[$  in a nice term of length  $n$  is  $< n/2$ .*

*Proof.* To prove this result by (strong) induction, we have to show that the following conditional claim is true:

If for every  $l < k$ , any nice term of length  $l$  has  $< l/2$   $[$ 's, then any nice term of length  $k$  has  $< k/2$   $[$ 's.

To show this conditional, assume that its antecedent is true, i.e., assume that for any  $l < k$ , nice terms of length  $l$  contain  $< l/2$   $[$ 's. We call this assumption the inductive hypothesis. We want to show the same is true for nice terms of length  $k$ .

So suppose  $t$  is a nice term of length  $k$ . Because nice terms are inductively defined, we have two cases: (1)  $t$  is a letter by itself, or (2)  $t$  is  $[s_1 \circ s_2]$  for some nice terms  $s_1$  and  $s_2$ .

1.  $t$  is a letter. Then  $k = 1$ , and the number of  $[$  in  $t$  is 0. Since  $0 < 1/2$ , the claim holds.

2.  $t$  is  $[s_1 \circ s_2]$  for some nice terms  $s_1$  and  $s_2$ . Let's let  $l_1$  be the length of  $s_1$  and  $l_2$  be the length of  $s_2$ . Then the length  $k$  of  $t$  is  $l_1 + l_2 + 3$  (the lengths of  $s_1$  and  $s_2$  plus three symbols  $[, \circ, ]$ ). Since  $l_1 + l_2 + 3$  is always greater than  $l_1$ ,  $l_1 < k$ . Similarly,  $l_2 < k$ . That means that the induction hypothesis applies to the terms  $s_1$  and  $s_2$ : the number  $m_1$  of  $[$  in  $s_1$  is  $< l_1/2$ , and the number  $m_2$  of  $[$  in  $s_2$  is  $< l_2/2$ .

The number of  $[$  in  $t$  is the number of  $[$  in  $s_1$ , plus the number of  $[$  in  $s_2$ , plus 1, i.e., it is  $m_1 + m_2 + 1$ . Since  $m_1 < l_1/2$  and  $m_2 < l_2/2$  we have:

$$m_1 + m_2 + 1 < \frac{l_1}{2} + \frac{l_2}{2} + 1 = \frac{l_1 + l_2 + 2}{2} < \frac{l_1 + l_2 + 3}{2} = k/2.$$

In each case, we've shown that the number of  $[$  in  $t$  is  $< k/2$  (on the basis of the inductive hypothesis). By strong induction, the proposition follows.  $\square$

## C.5 Structural Induction

So far we have used induction to establish results about all natural numbers. But a corresponding principle can be used directly to prove results about all elements of an inductively defined set. This is often called *structural* induction, because it depends on the structure of the inductively defined objects.

Generally, an inductive definition is given by (a) a list of "initial" elements of the set and (b) a list of operations which produce new elements of the set from old ones. In the case of nice terms, for instance, the initial objects are the letters. We only have one operation: the operations are

$$o(s_1, s_2) = [s_1 \circ s_2]$$

You can even think of the natural numbers  $\mathbb{N}$  themselves as being given by an inductive definition: the initial object is 0, and the operation is the successor function  $x + 1$ .

In order to prove something about all elements of an inductively defined set, i.e., that every element of the set has a property  $P$ , we must:

1. Prove that the initial objects have  $P$
2. Prove that for each operation  $o$ , if the arguments have  $P$ , so does the result.

For instance, in order to prove something about all nice terms, we would prove that it is true about all letters, and that it is true about  $[s_1 \circ s_2]$  provided it is true of  $s_1$  and  $s_2$  individually.

**Proposition C.5.** *The number of [ equals the number of ] in any nice term  $t$ .*

*Proof.* We use structural induction. Nice terms are inductively defined, with letters as initial objects and the operation  $o$  for constructing new nice terms out of old ones.

1. The claim is true for every letter, since the number of [ in a letter by itself is 0 and the number of ] in it is also 0.
2. Suppose the number of [ in  $s_1$  equals the number of ], and the same is true for  $s_2$ . The number of [ in  $o(s_1, s_2)$ , i.e., in  $[s_1 \circ s_2]$ , is the sum of the number of [ in  $s_1$  and  $s_2$  plus one. The number of ] in  $o(s_1, s_2)$  is the sum of the number of ] in  $s_1$  and  $s_2$  plus one. Thus, the number of [ in  $o(s_1, s_2)$  equals the number of ] in  $o(s_1, s_2)$ .  $\square$

Let's give another proof by structural induction: a proper initial segment of a string  $t$  of symbols is any string  $s$  that agrees with  $t$  symbol by symbol, read from the left, but  $t$  is longer. So, e.g.,  $[a \circ$  is a proper initial segment of  $[a \circ b]$ , but neither are  $[b \circ$  (they disagree at the second symbol) nor  $[a \circ b]$  (they are the same length).

**Proposition C.6.** *Every proper initial segment of a nice term  $t$  has more [ 's than ] 's.*

*Proof.* By induction on  $t$ :

1.  $t$  is a letter by itself: Then  $t$  has no proper initial segments.
2.  $t = [s_1 \circ s_2]$  for some nice terms  $s_1$  and  $s_2$ . If  $r$  is a proper initial segment of  $t$ , there are a number of possibilities:
  - a)  $r$  is just [ : Then  $r$  has one more [ than it does ] .
  - b)  $r$  is  $[r_1$  where  $r_1$  is a proper initial segment of  $s_1$ : Since  $s_1$  is a nice term, by induction hypothesis,  $r_1$  has more [ than ] and the same is true for  $[r_1$ .
  - c)  $r$  is  $[s_1$  or  $[s_1 \circ$  : By the previous result, the number of [ and ] in  $s_1$  are equal; so the number of [ in  $[s_1$  or  $[s_1 \circ$  is one more than the number of ] .
  - d)  $r$  is  $[s_1 \circ r_2$  where  $r_2$  is a proper initial segment of  $s_2$ : By induction hypothesis,  $r_2$  contains more [ than ] . By the previous result, the number of [ and of ] in  $s_1$  are equal. So the number of [ in  $[s_1 \circ r_2$  is greater than the number of ] .
  - e)  $r$  is  $[s_1 \circ s_2$ : By the previous result, the number of [ and ] in  $s_1$  are equal, and the same for  $s_2$ . So there is one more [ in  $[s_1 \circ s_2$  than there are ] . □

## C.6 Relations and Functions

When we have defined a set of objects (such as the natural numbers or the nice terms) inductively, we can also define *relations on* these objects by induction. For instance, consider the following idea: a nice term  $t_1$  is a subterm of a nice term  $t_2$  if it occurs as a part of it. Let's use a symbol for it:  $t_1 \sqsubseteq t_2$ . Every nice term is a subterm of itself, of course:  $t \sqsubseteq t$ . We can give an inductive definition of this relation as follows:

**Definition C.7.** The relation of a nice term  $t_1$  being a subterm of  $t_2$ ,  $t_1 \sqsubseteq t_2$ , is defined by induction on  $t_2$  as follows:

1. If  $t_2$  is a letter, then  $t_1 \sqsubseteq t_2$  iff  $t_1 = t_2$ .
2. If  $t_2$  is  $[s_1 \circ s_2]$ , then  $t_1 \sqsubseteq t_2$  iff  $t_1 = t_2$ ,  $t_1 \sqsubseteq s_1$ , or  $t_1 \sqsubseteq s_2$ .

This definition, for instance, will tell us that  $a \sqsubseteq [b \circ a]$ . For (2) says that  $a \sqsubseteq [b \circ a]$  iff  $a = [b \circ a]$ , or  $a \sqsubseteq b$ , or  $a \sqsubseteq a$ . The first two are false:  $a$  clearly isn't identical to  $[b \circ a]$ , and by (1),  $a \sqsubseteq b$  iff  $a = b$ , which is also false. However, also by (1),  $a \sqsubseteq a$  iff  $a = a$ , which is true.

It's important to note that the success of this definition depends on a fact that we haven't proved yet: every nice term  $t$  is either a letter by itself, or there are *uniquely determined* nice terms  $s_1$  and  $s_2$  such that  $t = [s_1 \circ s_2]$ . "Uniquely determined" here means that if  $t = [s_1 \circ s_2]$  it isn't *also*  $= [r_1 \circ r_2]$  with  $s_1 \neq r_1$  or  $s_2 \neq r_2$ . If this were the case, then clause (2) may come in conflict with itself: reading  $t_2$  as  $[s_1 \circ s_2]$  we might get  $t_1 \sqsubseteq t_2$ , but if we read  $t_2$  as  $[r_1 \circ r_2]$  we might get not  $t_1 \sqsubseteq t_2$ . Before we prove that this can't happen, let's look at an example where it *can* happen.

**Definition C.8.** Define *bracketless terms* inductively by

1. Every letter is a bracketless term.
2. If  $s_1$  and  $s_2$  are bracketless terms, then  $s_1 \circ s_2$  is a bracketless term.
3. Nothing else is a bracketless term.

Bracketless terms are, e.g.,  $a$ ,  $b \circ d$ ,  $b \circ a \circ b$ . Now if we defined "subterm" for bracketless terms the way we did above, the second clause would read

If  $t_2 = s_1 \circ s_2$ , then  $t_1 \sqsubseteq t_2$  iff  $t_1 = t_2$ ,  $t_1 \sqsubseteq s_1$ , or  $t_1 \sqsubseteq s_2$ .

Now  $b \circ a \circ b$  is of the form  $s_1 \circ s_2$  with

$$s_1 = b \text{ and } s_2 = a \circ b.$$

It is also of the form  $r_1 \circ r_2$  with

$$r_1 = b \circ a \text{ and } r_2 = b.$$

Now is  $a \circ b$  a subterm of  $b \circ a \circ b$ ? The answer is yes if we go by the first reading, and no if we go by the second.

The property that the way a nice term is built up from other nice terms is unique is called *unique readability*. Since inductive definitions of relations for such inductively defined objects are important, we have to prove that it holds.

**Proposition C.9.** *Suppose  $t$  is a nice term. Then either  $t$  is a letter by itself, or there are uniquely determined nice terms  $s_1, s_2$  such that  $t = [s_1 \circ s_2]$ .*

*Proof.* If  $t$  is a letter by itself, the condition is satisfied. So assume  $t$  isn't a letter by itself. We can tell from the inductive definition that then  $t$  must be of the form  $[s_1 \circ s_2]$  for some nice terms  $s_1$  and  $s_2$ . It remains to show that these are uniquely determined, i.e., if  $t = [r_1 \circ r_2]$ , then  $s_1 = r_1$  and  $s_2 = r_2$ .

So suppose  $t = [s_1 \circ s_2]$  and also  $t = [r_1 \circ r_2]$  for nice terms  $s_1, s_2, r_1, r_2$ . We have to show that  $s_1 = r_1$  and  $s_2 = r_2$ . First,  $s_1$  and  $r_1$  must be identical, for otherwise one is a proper initial segment of the other. But by **Proposition C.6**, that is impossible if  $s_1$  and  $r_1$  are both nice terms. But if  $s_1 = r_1$ , then clearly also  $s_2 = r_2$ .  $\square$

We can also define functions inductively: e.g., we can define the function  $f$  that maps any nice term to the maximum depth of nested  $[ \dots ]$  in it as follows:

**Definition C.10.** The *depth* of a nice term,  $f(t)$ , is defined in-

ductively as follows:

$$f(t) = \begin{cases} 0 & \text{if } t \text{ is a letter} \\ \max(f(s_1), f(s_2)) + 1 & \text{if } t = [s_1 \circ s_2]. \end{cases}$$

For instance

$$\begin{aligned} f([a \circ b]) &= \max(f(a), f(b)) + 1 = \\ &= \max(0, 0) + 1 = 1, \text{ and} \\ f([[a \circ b] \circ c]) &= \max(f([a \circ b]), f(c)) + 1 = \\ &= \max(1, 0) + 1 = 2. \end{aligned}$$

Here, of course, we assume that  $s_1$  and  $s_2$  are nice terms, and make use of the fact that every nice term is either a letter or of the form  $[s_1 \circ s_2]$ . It is again important that it can be of this form in only one way. To see why, consider again the bracketless terms we defined earlier. The corresponding “definition” would be:

$$g(t) = \begin{cases} 0 & \text{if } t \text{ is a letter} \\ \max(g(s_1), g(s_2)) + 1 & \text{if } t = s_1 \circ s_2. \end{cases}$$

Now consider the bracketless term  $a \circ b \circ c \circ d$ . It can be read in more than one way, e.g., as  $s_1 \circ s_2$  with

$$s_1 = a \text{ and } s_2 = b \circ c \circ d,$$

or as  $r_1 \circ r_2$  with

$$r_1 = a \circ b \text{ and } r_2 = c \circ d.$$

Calculating  $g$  according to the first way of reading it would give

$$\begin{aligned} g(s_1 \circ s_2) &= \max(g(a), g(b \circ c \circ d)) + 1 = \\ &= \max(0, 2) + 1 = 3 \end{aligned}$$

while according to the other reading we get

$$g(r_1 \circ r_2) = \max(g(a \circ b), g(c \circ d)) + 1 =$$



$$= \max(1,1) + 1 = 2$$

But a function must always yield a unique value; so our “definition” of  $g$  doesn’t define a function at all.

## Problems

**Problem C.1.** Define the set of supernice terms by

1. Any letter  $a, b, c, d$  is a supernice term.
2. If  $s$  is a supernice term, then so is  $[s]$ .
3. If  $s_1$  and  $s_2$  are supernice terms, then so is  $[s_1 \circ s_2]$ .
4. Nothing else is a supernice term.

Show that the number of  $[$  in a supernice term  $t$  of length  $n$  is  $\leq n/2 + 1$ .

**Problem C.2.** Prove by structural induction that no nice term starts with  $]$ .

**Problem C.3.** Give an inductive definition of the function  $l$ , where  $l(t)$  is the number of symbols in the nice term  $t$ .

**Problem C.4.** Prove by structural induction on nice terms  $t$  that  $f(t) < l(t)$  (where  $l(t)$  is the number of symbols in  $t$  and  $f(t)$  is the depth of  $t$  as defined in [Definition C.10](#)).

## APPENDIX D

# *Biographies*

### D.1 Georg Cantor

An early biography of Georg Cantor (GAY-org KAHN-tor) claimed that he was born and found on a ship that was sailing for Saint Petersburg, Russia, and that his parents were unknown. This, however, is not true; although he was born in Saint Petersburg in 1845.

Cantor received his doctorate in mathematics at the University of Berlin in 1867. He is known for his work in set theory, and is credited with founding set theory as a distinctive research discipline. He was the first to prove that

there are infinite sets of different sizes. His theories, and especially his theory of infinities, caused much debate among mathematicians at the time, and his work was controversial.

Cantor's religious beliefs and his mathematical work were in-



*Fig. D.1:* Georg Cantor

extricably tied; he even claimed that the theory of transfinite numbers had been communicated to him directly by God. In later life, Cantor suffered from mental illness. Beginning in 1894, and more frequently towards his later years, Cantor was hospitalized. The heavy criticism of his work, including a falling out with the mathematician Leopold Kronecker, led to depression and a lack of interest in mathematics. During depressive episodes, Cantor would turn to philosophy and literature, and even published a theory that Francis Bacon was the author of Shakespeare's plays.

Cantor died on January 6, 1918, in a sanatorium in Halle.

**Further Reading** For full biographies of Cantor, see [Dauben \(1990\)](#) and [Grattan-Guinness \(1971\)](#). Cantor's radical views are also described in the BBC Radio 4 program *A Brief History of Mathematics* ([du Sautoy, 2014](#)). If you'd like to hear about Cantor's theories in rap form, see [Rose \(2012\)](#).

## D.2 Alonzo Church

Alonzo Church was born in Washington, DC on June 14, 1903. In early childhood, an air gun incident left Church blind in one eye. He finished preparatory school in Connecticut in 1920 and began his university education at Princeton that same year. He completed his doctoral studies in 1927. After a couple years abroad, Church returned to Princeton. Church was known exceedingly polite and careful. His blackboard writing was immaculate, and he would preserve important pa-



*Fig. D.2:* Alonzo Church

pers by carefully covering them in Duco cement (a clear glue). Outside of his academic pursuits, he enjoyed reading science fiction magazines and was not afraid to write to the editors if he spotted any inaccuracies in the writing.

Church's academic achievements were great. Together with his students Stephen Kleene and Barkley Rosser, he developed a theory of effective calculability, the lambda calculus, independently of Alan Turing's development of the Turing machine. The two definitions of computability are equivalent, and give rise to what is now known as the *Church–Turing Thesis*, that a function of the natural numbers is effectively computable if and only if it is computable via Turing machine (or lambda calculus). He also proved what is now known as *Church's Theorem*: The decision problem for the validity of first-order formulas is unsolvable.

Church continued his work into old age. In 1967 he left Princeton for UCLA, where he was professor until his retirement in 1990. Church passed away on August 1, 1995 at the age of 92.

**Further Reading** For a brief biography of Church, see [Enderton \(2019\)](#). Church's original writings on the lambda calculus and the Entscheidungsproblem (Church's Thesis) are [Church \(1936a,b\)](#). [Aspray \(1984\)](#) records an interview with Church about the Princeton mathematics community in the 1930s. Church wrote a series of book reviews of the *Journal of Symbolic Logic* from 1936 until 1979. They are all archived on John MacFarlane's website ([MacFarlane, 2015](#)).

### D.3 Gerhard Gentzen

Gerhard Gentzen is known primarily as the creator of structural proof theory, and specifically the creation of the natural deduction and sequent calculus derivation systems. He was born on November 24, 1909 in Greifswald, Germany. Gerhard was home-schooled for three years before attending preparatory school, where he was behind most of his classmates in terms of educa-

tion. Despite this, he was a brilliant student and showed a strong aptitude for mathematics. His interests were varied, and he, for instance, also wrote poems for his mother and plays for the school theatre.

Gentzen began his university studies at the University of Greifswald, but moved around to Göttingen, Munich, and Berlin. He received his doctorate in 1933 from the University of Göttingen under Hermann Weyl. (Paul Bernays supervised most of his work, but was dismissed from the university by the



*Fig. D.3:* Gerhard Gentzen

Nazis.) In 1934, Gentzen began work as an assistant to David Hilbert. That same year he developed the sequent calculus and natural deduction derivation systems, in his papers *Untersuchungen über das logische Schließen I–II* [*Investigations Into Logical Deduction I–II*]. He proved the consistency of the Peano axioms in 1936.

Gentzen's relationship with the Nazis is complicated. At the same time his mentor Bernays was forced to leave Germany, Gentzen joined the university branch of the SA, the Nazi paramilitary organization. Like many Germans, he was a member of the Nazi party. During the war, he served as a telecommunications officer for the air intelligence unit. However, in 1942 he was released from duty due to a nervous breakdown. It is unclear whether or not Gentzen's loyalties lay with the Nazi party, or whether he joined the party in order to ensure academic success.

In 1943, Gentzen was offered an academic position at the Mathematical Institute of the German University of Prague, which he accepted. However, in 1945 the citizens of Prague revolted against German occupation. Soviet forces arrived in the city and arrested all the professors at the university. Because of his membership in Nazi organizations, Gentzen was taken to a

forced labour camp. He died of malnutrition while in his cell on August 4, 1945 at the age of 35.

**Further Reading** For a full biography of Gentzen, see [Menzler-Trott \(2007\)](#). An interesting read about mathematicians under Nazi rule, which gives a brief note about Gentzen's life, is given by [Segal \(2014\)](#). Gentzen's papers on logical deduction are available in the original German ([Gentzen, 1935a,b](#)). English translations of Gentzen's papers have been collected in a single volume by [Szabo \(1969\)](#), which also includes a biographical sketch.

## D.4 Kurt Gödel

Kurt Gödel (GER-dle) was born on April 28, 1906 in Brünn in the Austro-Hungarian empire (now Brno in the Czech Republic). Due to his inquisitive and bright nature, young Kurt was often called “Der kleine Herr Warum” (Little Mr. Why) by his family. He excelled in academics from primary school onward, where he got less than the highest grade only in mathematics. Gödel was often absent from school due to poor health and was exempt from physical education. He was diagnosed with rheumatic fever during his childhood. Throughout his life, he believed this permanently affected his heart despite medical assessment saying otherwise.



*Fig. D.4:* Kurt Gödel

Gödel began studying at the University of Vienna in 1924 and completed his doctoral studies in 1929. He first intended to study physics, but his interests soon moved to mathematics and especially logic, in part due to the influence of the philosopher Rudolf Carnap. His dissertation, written under the supervision of Hans Hahn, proved the completeness theorem of first-order predicate logic with identity (Gödel, 1929). Only a year later, he obtained his most famous results—the first and second incompleteness theorems (published in Gödel 1931). During his time in Vienna, Gödel was heavily involved with the Vienna Circle, a group of scientifically-minded philosophers that included Carnap, whose work was especially influenced by Gödel's results.

In 1938, Gödel married Adele Nimbursky. His parents were not pleased: not only was she six years older than him and already divorced, but she worked as a dancer in a nightclub. Social pressures did not affect Gödel, however, and they remained happily married until his death.

After Nazi Germany annexed Austria in 1938, Gödel and Adele emigrated to the United States, where he took up a position at the Institute for Advanced Study in Princeton, New Jersey. Despite his introversion and eccentric nature, Gödel's time at Princeton was collaborative and fruitful. He published essays in set theory, philosophy and physics. Notably, he struck up a particularly strong friendship with his colleague at the IAS, Albert Einstein.

In his later years, Gödel's mental health deteriorated. His wife's hospitalization in 1977 meant she was no longer able to cook his meals for him. Having suffered from mental health issues throughout his life, he succumbed to paranoia. Deathly afraid of being poisoned, Gödel refused to eat. He died of starvation on January 14, 1978, in Princeton.

**Further Reading** For a complete biography of Gödel's life is available, see [John Dawson \(1997\)](#). For further biographical pieces, as well as essays about Gödel's contributions to logic and

philosophy, see Wang (1990), Baaz et al. (2011), Takeuti et al. (2003), and Sigmund et al. (2007).

Gödel's PhD thesis is available in the original German (Gödel, 1929). The original text of the incompleteness theorems is (Gödel, 1931). All of Gödel's published and unpublished writings, as well as a selection of correspondence, are available in English in his *Collected Papers* Feferman et al. (1986, 1990).

For a detailed treatment of Gödel's incompleteness theorems, see Smith (2013). For an informal, philosophical discussion of Gödel's theorems, see Mark Linsenmayer's podcast (Linsenmayer, 2014).

## D.5 Emmy Noether

Emmy Noether (NER-ter) was born in Erlangen, Germany, on March 23, 1882, to an upper-middle class scholarly family. Hailed as the “mother of modern algebra,” Noether made groundbreaking contributions to both mathematics and physics, despite significant barriers to women's education. In Germany at the time, young girls were meant to be educated in arts and were not allowed to attend college preparatory schools. However, after auditing classes at the Universities of Göttingen and Erlan-



Fig. D.5: Emmy Noether

gen (where her father was professor of mathematics), Noether was eventually able to enroll as a student at Erlangen in 1904, when their policy was updated to allow female students. She re-



ceived her doctorate in mathematics in 1907.

Despite her qualifications, Noether experienced much resistance during her career. From 1908–1915, she taught at Erlangen without pay. During this time, she caught the attention of David Hilbert, one of the world’s foremost mathematicians of the time, who invited her to Göttingen. However, women were prohibited from obtaining professorships, and she was only able to lecture under Hilbert’s name, again without pay. During this time she proved what is now known as Noether’s theorem, which is still used in theoretical physics today. Noether was finally granted the right to teach in 1919. Hilbert’s response to continued resistance of his university colleagues reportedly was: “Gentlemen, the faculty senate is not a bathhouse.”

In the later 1920s, she concentrated on work in abstract algebra, and her contributions revolutionized the field. In her proofs she often made use of the so-called ascending chain condition, which states that there is no infinite strictly increasing chain of certain sets. For instance, certain algebraic structures now known as Noetherian rings have the property that there are no infinite sequences of ideals  $I_1 \subseteq I_2 \subseteq \dots$ . The condition can be generalized to any partial order (in algebra, it concerns the special case of ideals ordered by the subset relation), and we can also consider the dual descending chain condition, where every strictly *decreasing* sequence in a partial order eventually ends. If a partial order satisfies the descending chain condition, it is possible to use induction along this order in a similar way in which we can use induction along the  $<$  order on  $\mathbb{N}$ . Such orders are called *well-founded* or *Noetherian*, and the corresponding proof principle *Noetherian induction*.

Noether was Jewish, and when the Nazis came to power in 1933, she was dismissed from her position. Luckily, Noether was able to emigrate to the United States for a temporary position at Bryn Mawr, Pennsylvania. During her time there she also lectured at Princeton, although she found the university to be unwelcoming to women (Dick, 1981, 81). In 1935, Noether underwent an operation to remove a uterine tumour. She died from an infection

as a result of the surgery, and was buried at Bryn Mawr.

**Further Reading** For a biography of Noether, see [Dick \(1981\)](#). The Perimeter Institute for Theoretical Physics has their lectures on Noether's life and influence available online ([Institute, 2015](#)). If you're tired of reading, *Stuff You Missed in History Class* has a podcast on Noether's life and influence ([Frey and Wilson, 2015](#)). The collected works of Noether are available in the original German ([Jacobson, 1983](#)).

## D.6 Rózsa Péter

Rózsa Péter was born Rózsa Politzer, in Budapest, Hungary, on February 17, 1905. She is best known for her work on recursive functions, which was essential for the creation of the field of recursion theory.

Péter was raised during harsh political times—WWI raged when she was a teenager—but was able to attend the affluent Maria Terezia Girls' School in Budapest, from where she graduated in 1922. She then studied



*Fig. D.6:* Rózsa Péter

at Pázmány Péter University (later renamed Loránd Eötvös University) in Budapest. She began studying chemistry at the insistence of her father, but later switched to mathematics, and graduated in 1927. Although she had the credentials to teach high school mathematics, the economic situation at the time was dire as the Great Depression affected the world economy. During this time, Péter took odd jobs as a tutor and private teacher of

mathematics. She eventually returned to university to take up graduate studies in mathematics. She had originally planned to work in number theory, but after finding out that her results had already been proven, she almost gave up on mathematics altogether. She was encouraged to work on Gödel's incompleteness theorems, and unknowingly proved several of his results in different ways. This restored her confidence, and Péter went on to write her first papers on recursion theory, inspired by David Hilbert's foundational program. She received her PhD in 1935, and in 1937 she became an editor for the *Journal of Symbolic Logic*.

Péter's early papers are widely credited as founding contributions to the field of recursive function theory. In Péter (1935a), she investigated the relationship between different kinds of recursion. In Péter (1935b), she showed that a certain recursively defined function is not primitive recursive. This simplified an earlier result due to Wilhelm Ackermann. Péter's simplified function is what's now often called the Ackermann function—and sometimes, more properly, the Ackermann–Péter function. She wrote the first book on recursive function theory (Péter, 1951).

Despite the importance and influence of her work, Péter did not obtain a full-time teaching position until 1945. During the Nazi occupation of Hungary during World War II, Péter was not allowed to teach due to anti-Semitic laws. In 1944 the government created a Jewish ghetto in Budapest; the ghetto was cut off from the rest of the city and attended by armed guards. Péter was forced to live in the ghetto until 1945 when it was liberated. She then went on to teach at the Budapest Teachers Training College, and from 1955 onward at Eötvös Loránd University. She was the first female Hungarian mathematician to become an Academic Doctor of Mathematics, and the first woman to be elected to the Hungarian Academy of Sciences.

Péter was known as a passionate teacher of mathematics, who preferred to explore the nature and beauty of mathematical problems with her students rather than to merely lecture. As a result, she was affectionately called “Aunt Rosa” by her students. Péter died in 1977 at the age of 71.

**Further Reading** For more biographical reading, see (O'Connor and Robertson, 2014) and (Andrásfai, 1986). Tamassy (1994) conducted a brief interview with Péter. For a fun read about mathematics, see Péter's book *Playing With Infinity* (Péter, 2010).

## D.7 Julia Robinson

Julia Bowman Robinson was an American mathematician. She is known mainly for her work on decision problems, and most famously for her contributions to the solution of Hilbert's tenth problem. Robinson was born in St. Louis, Missouri, on December 8, 1919. Robinson recalls being intrigued by numbers already as a child (Reid, 1986, 4). At age nine she contracted scarlet fever and suffered from several recurrent bouts of rheumatic fever. This forced her to spend much of



*Fig. D.7:* Julia Robinson

her time in bed, putting her behind in her education. Although she was able to catch up with the help of private tutors, the physical effects of her illness had a lasting impact on her life.

Despite her childhood struggles, Robinson graduated high school with several awards in mathematics and the sciences. She started her university career at San Diego State College, and transferred to the University of California, Berkeley, as a senior. There she was influenced by the mathematician Raphael Robinson. They became good friends, and married in 1941. As a spouse of a faculty member, Robinson was barred from teaching

in the mathematics department at Berkeley. Although she continued to audit mathematics classes, she hoped to leave university and start a family. Not long after her wedding, however, Robinson contracted pneumonia. She was told that there was substantial scar tissue build up on her heart due to the rheumatic fever she suffered as a child. Due to the severity of the scar tissue, the doctor predicted that she would not live past forty and she was advised not to have children (Reid, 1986, 13).

Robinson was depressed for a long time, but eventually decided to continue studying mathematics. She returned to Berkeley and completed her PhD in 1948 under the supervision of Alfred Tarski. The first-order theory of the real numbers had been shown to be decidable by Tarski, and from Gödel's work it followed that the first-order theory of the natural numbers is undecidable. It was a major open problem whether the first-order theory of the rationals is decidable or not. In her thesis (1949), Robinson proved that it was not.

Interested in decision problems, Robinson next attempted to find a solution to Hilbert's tenth problem. This problem was one of a famous list of 23 mathematical problems posed by David Hilbert in 1900. The tenth problem asks whether there is an algorithm that will answer, in a finite amount of time, whether or not a polynomial equation with integer coefficients, such as  $3x^2 - 2y + 3 = 0$ , has a solution in the integers. Such questions are known as *Diophantine problems*. After some initial successes, Robinson joined forces with Martin Davis and Hilary Putnam, who were also working on the problem. They succeeded in showing that exponential Diophantine problems (where the unknowns may also appear as exponents) are undecidable, and showed that a certain conjecture (later called "J.R.") implies that Hilbert's tenth problem is undecidable (Davis et al., 1961). Robinson continued to work on the problem throughout the 1960s. In 1970, the young Russian mathematician Yuri Matijasevich finally proved the J.R. hypothesis. The combined result is now called the Matijasevich–Robinson–Davis–Putnam theorem, or MRDP theorem for short. Matijasevich and Robinson became friends

and collaborated on several papers. In a letter to Matijasevich, Robinson once wrote that “actually I am very pleased that working together (thousands of miles apart) we are obviously making more progress than either one of us could alone” (Matijasevich, 1992, 45).

Robinson was the first female president of the American Mathematical Society, and the first woman to be elected to the National Academy of Science. She died on July 30, 1985 at the age of 65 after being diagnosed with leukemia.

**Further Reading** Robinson’s mathematical papers are available in her *Collected Works* (Robinson, 1996), which also includes a reprint of her National Academy of Sciences biographical memoir (Feferman, 1994). Robinson’s older sister Constance Reid published an “Autobiography of Julia,” based on interviews (Reid, 1986), as well as a full memoir (Reid, 1996). A short documentary about Robinson and Hilbert’s tenth problem was directed by George Csicsery (Csicsery, 2016). For a brief memoir about Yuri Matijasevich’s collaborations with Robinson, and her influence on his work, see (Matijasevich, 1992).

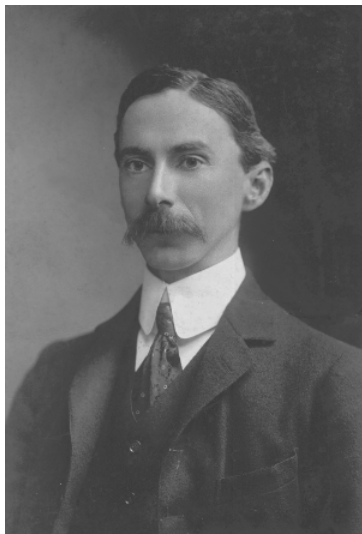
## D.8 Bertrand Russell

Bertrand Russell is hailed as one of the founders of modern analytic philosophy. Born May 18, 1872, Russell was not only known for his work in philosophy and logic, but wrote many popular books in various subject areas. He was also an ardent political activist throughout his life.

Russell was born in Trellech, Monmouthshire, Wales. His parents were members of the British nobility. They were free-thinkers, and even made friends with the radicals in Boston at the time. Unfortunately, Russell’s parents died when he was young, and Russell was sent to live with his grandparents. There, he was given a religious upbringing (something his parents had wanted to avoid at all costs). His grandmother was very strict in all

matters of morality. During adolescence he was mostly home-schooled by private tutors.

Russell's influence in analytic philosophy, and especially logic, is tremendous. He studied mathematics and philosophy at Trinity College, Cambridge, where he was influenced by the mathematician and philosopher Alfred North Whitehead. In 1910, Russell and Whitehead published the first volume of *Principia Mathematica*, where they championed the view that mathematics is reducible to logic. He went on to publish hundreds of books, essays and political pamphlets. In 1950, he won the Nobel Prize for literature.



*Fig. D.8:* Bertrand Russell

Russell's was deeply entrenched in politics and social activism. During World War I he was arrested and sent to prison for six months due to pacifist activities and protest. While in prison, he was able to write and read, and claims to have found the experience "quite agreeable." He remained a pacifist throughout his life, and was again incarcerated for attending a nuclear disarmament rally in 1961. He also survived a plane crash in 1948, where the only survivors were those sitting in the smoking section. As such, Russell claimed that he owed his life to smoking. Russell was married four times, but had a reputation for carrying on extra-marital affairs. He died on February 2, 1970 at the age of 97 in Penrhyndeudraeth, Wales.

**Further Reading** Russell wrote an autobiography in three parts, spanning his life from 1872–1967 (Russell, 1967, 1968,

1969). The Bertrand Russell Research Centre at McMaster University is home of the Bertrand Russell archives. See their website at [Duncan \(2015\)](#), for information on the volumes of his collected works (including searchable indexes), and archival projects. Russell’s paper *On Denoting* ([Russell, 1905](#)) is a classic of 20th century analytic philosophy.

The Stanford Encyclopedia of Philosophy entry on Russell ([Irvine, 2015](#)) has sound clips of Russell speaking on Desire and Political theory. Many video interviews with Russell are available online. To see him talk about smoking and being involved in a plane crash, e.g., see [Russell \(n.d.\)](#). Some of Russell’s works, including his *Introduction to Mathematical Philosophy* are available as free audiobooks on [LibriVox \(n.d.\)](#).

## D.9 Alfred Tarski

Alfred Tarski was born on January 14, 1901 in Warsaw, Poland (then part of the Russian Empire). Described as “Napoleonic,” Tarski was boisterous, talkative, and intense. His energy was often reflected in his lectures—he once set fire to a wastebasket while disposing of a cigarette during a lecture, and was forbidden from lecturing in that building again.

Tarski had a thirst for knowledge from a young age. Although later in life he would tell students that he studied



*Fig. D.9:* Alfred Tarski

logic because it was the only class in which he got a B, his high school records show that he got A’s across the board—even in



logic. He studied at the University of Warsaw from 1918 to 1924. Tarski first intended to study biology, but became interested in mathematics, philosophy, and logic, as the university was the center of the Warsaw School of Logic and Philosophy. Tarski earned his doctorate in 1924 under the supervision of Stanisław Leśniewski.

Before emigrating to the United States in 1939, Tarski completed some of his most important work while working as a secondary school teacher in Warsaw. His work on logical consequence and logical truth were written during this time. In 1939, Tarski was visiting the United States for a lecture tour. During his visit, Germany invaded Poland, and because of his Jewish heritage, Tarski could not return. His wife and children remained in Poland until the end of the war, but were then able to emigrate to the United States as well. Tarski taught at Harvard, the College of the City of New York, and the Institute for Advanced Study at Princeton, and finally the University of California, Berkeley. There he founded the multidisciplinary program in Logic and the Methodology of Science. Tarski died on October 26, 1983 at the age of 82.

**Further Reading** For more on Tarski's life, see the biography *Alfred Tarski: Life and Logic* (Feferman and Feferman, 2004). Tarski's seminal works on logical consequence and truth are available in English in (Corcoran, 1983). All of Tarski's original works have been collected into a four volume series, (Tarski, 1981).

## D.10 Alan Turing

Alan Turing was born in Maida Vale, London, on June 23, 1912. He is considered the father of theoretical computer science. Turing's interest in the physical sciences and mathematics started at a young age. However, as a boy his interests were not represented well in his schools, where emphasis was placed on literature and

classics. Consequently, he did poorly in school and was reprimanded by many of his teachers.

Turing attended King's College, Cambridge as an undergraduate, where he studied mathematics. In 1936 Turing developed (what is now called) the Turing machine as an attempt to precisely define the notion of a computable function and to prove the undecidability of the decision problem. He was beaten to the result by Alonzo Church, who proved the result via his own lambda calculus. Turing's paper was still published with reference to Church's re-



*Fig. D.10:* Alan Turing

sult. Church invited Turing to Princeton, where he spent 1936–1938, and obtained a doctorate under Church.

Despite his interest in logic, Turing's earlier interests in physical sciences remained prevalent. His practical skills were put to work during his service with the British cryptanalytic department at Bletchley Park during World War II. Turing was a central figure in cracking the cypher used by German Naval communications—the Enigma code. Turing's expertise in statistics and cryptography, together with the introduction of electronic machinery, gave the team the ability to crack the code by creating a de-crypting machine called a “bombe.” His ideas also helped in the creation of the world's first programmable electronic computer, the Colossus, also used at Bletchley park to break the German Lorenz cypher.

Turing was gay. Nevertheless, in 1942 he proposed to Joan Clarke, one of his teammates at Bletchley Park, but later broke off the engagement and confessed to her that he was homosexual. He had several lovers throughout his lifetime, although homosexual

acts were then criminal offences in the UK. In 1952, Turing's house was burgled by a friend of his lover at the time, and when filing a police report, Turing admitted to having a homosexual relationship, under the impression that the government was on their way to legalizing homosexual acts. This was not true, and he was charged with gross indecency. Instead of going to prison, Turing opted for a hormone treatment that reduced libido. Turing was found dead on June 8, 1954, of a cyanide overdose—most likely suicide. He was given a royal pardon by Queen Elizabeth II in 2013.

**Further Reading** For a comprehensive biography of Alan Turing, see [Hodges \(2014\)](#). Turing's life and work inspired a play, *Breaking the Code*, which was produced in 1996 for TV starring Derek Jacobi as Turing. *The Imitation Game*, an Academy Award nominated film starring Benedict Cumberbatch and Kiera Knightley, is also loosely based on Alan Turing's life and time at Bletchley Park ([Tyldum, 2014](#)).

[Radiolab \(2012\)](#) has several podcasts on Turing's life and work. BBC Horizon's documentary *The Strange Life and Death of Dr. Turing* is available to watch online ([Sykes, 1992](#)). ([Theelen, 2012](#)) is a short video of a working LEGO Turing Machine—made to honour Turing's centenary in 2012.

Turing's original paper on Turing machines and the decision problem is [Turing \(1937\)](#).

## D.11 Ernst Zermelo

Ernst Zermelo was born on July 27, 1871 in Berlin, Germany. He had five sisters, though his family suffered from poor health and only three survived to adulthood. His parents also passed away when he was young, leaving him and his siblings orphans when he was seventeen. Zermelo had a deep interest in the arts, and especially in poetry. He was known for being sharp, witty, and critical. His most celebrated mathematical achievements in-

clude the introduction of the axiom of choice (in 1904), and his axiomatization of set theory (in 1908).

Zermelo's interests at university were varied. He took courses in physics, mathematics, and philosophy. Under the supervision of Hermann Schwarz, Zermelo completed his dissertation *Investigations in the Calculus of Variations* in 1894 at the University of Berlin. In 1897, he decided to pursue more studies at the University of Göttingen, where he was heavily influenced by the foundational work of David Hilbert. In 1899 he became eligible for professorship, but did not get one until eleven years later—possibly due to his strange demeanour and “nervous haste.”



Fig. D.11: Ernst Zermelo

Zermelo finally received a paid professorship at the University of Zurich in 1910, but was forced to retire in 1916 due to tuberculosis. After his recovery, he was given an honorary professorship at the University of Freiburg in 1921. During this time he worked on foundational mathematics. He became irritated with the works of Thoralf Skolem and Kurt Gödel, and publicly criticized their approaches in his papers. He was dismissed from his position at Freiburg in 1935, due to his unpopularity and his opposition to Hitler's rise to power in Germany.

The later years of Zermelo's life were marked by isolation. After his dismissal in 1935, he abandoned mathematics. He moved to the country where he lived modestly. He married in 1944, and became completely dependent on his wife as he was going blind. Zermelo lost his sight completely by 1951. He passed away in Günterstal, Germany, on May 21, 1953.

**Further Reading** For a full biography of Zermelo, see [Ebbinghaus \(2015\)](#). Zermelo's seminal 1904 and 1908 papers are available to read in the original German ([Zermelo, 1904, 1908](#)). Zermelo's collected works, including his writing on physics, are available in English translation in ([Ebbinghaus et al., 2010](#); [Ebbinghaus and Kanamori, 2013](#)).

# *Photo Credits*

Georg Cantor, p. 483: Portrait of Georg Cantor by Otto Zeth courtesy of the [Universitätsarchiv, Martin-Luther Universität Halle–Wittenberg](#). UAHW Rep. 40-VI, Nr. 3 Bild 102.

Alonzo Church, p. 484: Portrait of Alonzo Church, undated, photographer unknown. Alonzo Church Papers; 1924–1995, (C0948) Box 60, Folder 3. [Manuscripts Division, Department of Rare Books and Special Collections, Princeton University Library](#). © Princeton University. The Open Logic Project has obtained permission to use this image for inclusion in non-commercial OLP-derived materials. Permission from Princeton University is required for any other use.

Gerhard Gentzen, p. 486: Portrait of Gerhard Gentzen playing ping-pong courtesy of Ekhart Mentzler-Trott.

Kurt Gödel, p. 487: Portrait of Kurt Gödel, ca. 1925, photographer unknown. From the [Shelby White and Leon Levy Archives Center, Institute for Advanced Study](#), Princeton, NJ, USA, on deposit at Princeton University Library, [Manuscript Division, Department of Rare Books and Special Collections](#), Kurt Gödel Papers, (C0282), Box 14b, #110000. The Open Logic Project has obtained permission from the Institute's Archives Center to use this image for inclusion in non-commercial OLP-derived materials. Permission from the Archives Center is required for any other use.

Emmy Noether, p. 489: Portrait of Emmy Noether, ca. 1922, courtesy of the [Abteilung für Handschriften und Seltene Drucke](#),

Niedersächsische Staats- und Universitätsbibliothek Göttingen, Cod. Ms. D. Hilbert 754, Bl. 14 Nr. 73. Restored from an original scan by Joel Fuller.

Rózsa Péter, p. 491: Portrait of Rózsa Péter, undated, photographer unknown. Courtesy of Béla Andrásfai.

Julia Robinson, p. 493: Portrait of Julia Robinson, unknown photographer, courtesy of Neil D. Reid. The Open Logic Project has obtained permission to use this image for inclusion in non-commercial OLP-derived materials. Permission is required for any other use.

Bertrand Russell, p. 496: Portrait of Bertrand Russell, ca. 1907, courtesy of the William Ready Division of Archives and Research Collections, McMaster University Library. **Bertrand Russell Archives**, Box 2, f. 4.

Alfred Tarski, p. 497: Passport photo of Alfred Tarski, 1939. Cropped and restored from a scan of Tarski's passport by Joel Fuller. Original courtesy of **Bancroft Library, University of California, Berkeley**. Alfred Tarski Papers, Banc MSS 84/49. The Open Logic Project has obtained permission to use this image for inclusion in non-commercial OLP-derived materials. Permission from Bancroft Library is required for any other use.

Alan Turing, p. 499: Portrait of **Alan Mathison Turing** by Elliott & Fry, 29 March 1951, NPG x82217, © National Portrait Gallery, London. Used under a **Creative Commons BY-NC-ND 3.0 license**.

Ernst Zermelo, p. 501: Portrait of Ernst Zermelo, ca. 1922, courtesy of the **Abteilung für Handschriften und Seltene Drucke, Niedersächsische Staats- und Universitätsbibliothek Göttingen**, Cod. Ms. D. Hilbert 754, Bl. 6 Nr. 25.

# Bibliography

- Andrásfai, Béla. 1986. Rózsa (Rosa) Péter. *Periodica Polytechnica Electrical Engineering* 30(2-3): 139–145. URL <http://www.pp.bme.hu/ee/article/view/4651>.
- Aspray, William. 1984. The Princeton mathematics community in the 1930s: Alonzo Church. URL [http://www.princeton.edu/mudd/finding\\_aids/mathoral/pmc05.htm](http://www.princeton.edu/mudd/finding_aids/mathoral/pmc05.htm). Interview.
- Baaz, Matthias, Christos H. Papadimitriou, Hilary W. Putnam, Dana S. Scott, and Charles L. Harper Jr. 2011. *Kurt Gödel and the Foundations of Mathematics: Horizons of Truth*. Cambridge: Cambridge University Press.
- Boolos, George. 1993. *The Logic of Provability*. Cambridge: Cambridge University Press.
- Cantor, Georg. 1892. Über eine elementare Frage der Mannigfaltigkeitslehre. *Jahresbericht der deutschen Mathematiker-Vereinigung* 1: 75–8.
- Cheng, Eugenia. 2004. How to write proofs: A quick guide. URL <http://http://eugeniacheng.com/wp-content/uploads/2017/02/cheng-proofguide.pdf>.
- Church, Alonzo. 1936a. A note on the Entscheidungsproblem. *The Journal of Symbolic Logic* 1: 40–41.



- Church, Alonzo. 1936b. An unsolvable problem of elementary number theory. *American Journal of Mathematics* 58: 345–363.
- Corcoran, John. 1983. *Logic, Semantics, Metamathematics*. Indianapolis: Hackett, 2nd ed.
- Csicsery, George. 2016. Zala films: Julia Robinson and Hilbert's tenth problem. URL <http://www.zalafilms.com/films/juliarobinson.html>.
- Dauben, Joseph. 1990. *Georg Cantor: His Mathematics and Philosophy of the Infinite*. Princeton: Princeton University Press.
- Davis, Martin, Hilary Putnam, and Julia Robinson. 1961. The decision problem for exponential Diophantine equations. *Annals of Mathematics* 74(3): 425–436. URL <http://www.jstor.org/stable/1970289>.
- Dick, Auguste. 1981. *Emmy Noether 1882–1935*. Boston: Birkhäuser.
- du Sautoy, Marcus. 2014. A brief history of mathematics: Georg Cantor. URL <http://www.bbc.co.uk/programmes/b00ss1j0>. Audio Recording.
- Duncan, Arlene. 2015. The Bertrand Russell Research Centre. URL <http://russell.mcmaster.ca/>.
- Ebbinghaus, Heinz-Dieter. 2015. *Ernst Zermelo: An Approach to his Life and Work*. Berlin: Springer-Verlag.
- Ebbinghaus, Heinz-Dieter, Craig G. Fraser, and Akihiro Kanamori. 2010. *Ernst Zermelo. Collected Works*, vol. 1. Berlin: Springer-Verlag.
- Ebbinghaus, Heinz-Dieter and Akihiro Kanamori. 2013. *Ernst Zermelo: Collected Works*, vol. 2. Berlin: Springer-Verlag.

- Enderton, Herbert B. 2019. Alonzo Church: Life and Work. In *The Collected Works of Alonzo Church*, eds. Tyler Burge and Herbert B. Enderton. Cambridge, MA: MIT Press.
- Feferman, Anita and Solomon Feferman. 2004. *Alfred Tarski: Life and Logic*. Cambridge: Cambridge University Press.
- Feferman, Solomon. 1994. Julia Bowman Robinson 1919–1985. *Biographical Memoirs of the National Academy of Sciences* 63: 1–28. URL <http://www.nasonline.org/publications/biographical-memoirs/memoir-pdfs/robinson-julia.pdf>.
- Feferman, Solomon, John W. Dawson Jr., Stephen C. Kleene, Gregory H. Moore, Robert M. Solovay, and Jean van Heijenoort. 1986. *Kurt Gödel: Collected Works. Vol. 1: Publications 1929–1936*. Oxford: Oxford University Press.
- Feferman, Solomon, John W. Dawson Jr., Stephen C. Kleene, Gregory H. Moore, Robert M. Solovay, and Jean van Heijenoort. 1990. *Kurt Gödel: Collected Works. Vol. 2: Publications 1938–1974*. Oxford: Oxford University Press.
- Frege, Gottlob. 1884. *Die Grundlagen der Arithmetik: Eine logisch mathematische Untersuchung über den Begriff der Zahl*. Breslau: Wilhelm Koebner. Translation in Frege (1953).
- Frege, Gottlob. 1953. *Foundations of Arithmetic*, ed. J. L. Austin. Oxford: Basil Blackwell & Mott, 2nd ed.
- Frey, Holly and Tracy V. Wilson. 2015. Stuff you missed in history class: Emmy Noether, mathematics trailblazer. URL <https://www.iheart.com/podcast/stuff-you-missed-in-history-cl-21124503/episode/emmy-noether-mathematics-trailblazer-30207491/>. Podcast audio.

- Gentzen, Gerhard. 1935a. Untersuchungen über das logische Schließen I. *Mathematische Zeitschrift* 39: 176–210. English translation in Szabo (1969), pp. 68–131.
- Gentzen, Gerhard. 1935b. Untersuchungen über das logische Schließen II. *Mathematische Zeitschrift* 39: 176–210, 405–431. English translation in Szabo (1969), pp. 68–131.
- Gödel, Kurt. 1929. Über die Vollständigkeit des Logikkalküls [On the completeness of the calculus of logic]. Dissertation, Universität Wien. Reprinted and translated in Feferman et al. (1986), pp. 60–101.
- Gödel, Kurt. 1931. über formal unentscheidbare Sätze der *Principia Mathematica* und verwandter Systeme I [On formally undecidable propositions of *Principia Mathematica* and related systems I]. *Monatshefte für Mathematik und Physik* 38: 173–198. Reprinted and translated in Feferman et al. (1986), pp. 144–195.
- Gödel, Kurt. 1995. Some basic theorems on the foundations of mathematics and their implications. In *Kurt Gödel: Collected Works*, eds. Solomon Feferman et al., vol. 3, 304–323. New York and Oxford: Oxford University Press.
- Grattan-Guinness, Ivor. 1971. Towards a biography of Georg Cantor. *Annals of Science* 27(4): 345–391.
- Hammack, Richard. 2013. *Book of Proof*. Richmond, VA: Virginia Commonwealth University. URL <http://www.people.vcu.edu/~rhammack/BookOfProof/BookOfProof.pdf>.
- Hodges, Andrew. 2014. *Alan Turing: The Enigma*. London: Vintage.
- Hutchings, Michael. 2003. Introduction to mathematical arguments. URL <https://math.berkeley.edu/~hutching/teach/proofs.pdf>.

- Institute, Perimeter. 2015. Emmy Noether: Her life, work, and influence. URL <https://www.youtube.com/watch?v=tNNyAyMRsgE>. Video Lecture.
- Irvine, Andrew David. 2015. Sound clips of Bertrand Russell speaking. URL <http://plato.stanford.edu/entries/russell/russell-soundclips.html>.
- Jacobson, Nathan. 1983. *Emmy Noether: Gesammelte Abhandlungen—Collected Papers*. Berlin: Springer-Verlag.
- John Dawson, Jr. 1997. *Logical Dilemmas: The Life and Work of Kurt Gödel*. Boca Raton: CRC Press.
- LibriVox. n.d. Bertrand Russell. URL [https://librivox.org/author/1508?primary\\_key=1508&search\\_category=author&search\\_page=1&search\\_form=get\\_results](https://librivox.org/author/1508?primary_key=1508&search_category=author&search_page=1&search_form=get_results). Collection of public domain audiobooks.
- Linsenmayer, Mark. 2014. The partially examined life: Gödel on math. URL <http://www.partiallyexaminedlife.com/2014/06/16/ep95-godel/>. Podcast audio.
- MacFarlane, John. 2015. Alonzo Church's JSL reviews. URL <http://johnmacfarlane.net/church.html>.
- Magnus, P. D., Tim Button, J. Robert Loftis, Aaron Thomas-Bolduc, Robert Trueman, and Richard Zach. 2021. *Forall x: Calgary. An Introduction to Formal Logic*. Calgary: Open Logic Project, f21 ed. URL <https://forallx.openlogicproject.org/>.
- Matijasevich, Yuri. 1992. My collaboration with Julia Robinson. *The Mathematical Intelligencer* 14(4): 38–45.
- Menzler-Trott, Eckart. 2007. *Logic's Lost Genius: The Life of Gerhard Gentzen*. Providence: American Mathematical Society.

- O'Connor, John J. and Edmund F. Robertson. 2014. Rózsa Péter. URL <http://www-groups.dcs.st-and.ac.uk/~history/Biographies/Peter.html>.
- Péter, Rózsa. 1935a. Über den Zusammenhang der verschiedenen Begriffe der rekursiven Funktion. *Mathematische Annalen* 110: 612–632.
- Péter, Rózsa. 1935b. Konstruktion nichtrekursiver Funktionen. *Mathematische Annalen* 111: 42–60.
- Péter, Rózsa. 1951. *Rekursive Funktionen*. Budapest: Akadémiai Kiado. English translation in (Péter, 1967).
- Péter, Rózsa. 1967. *Recursive Functions*. New York: Academic Press.
- Péter, Rózsa. 2010. *Playing with Infinity*. New York: Dover. URL [https://books.google.ca/books?id=6V3wNs4uv\\_4C&lpg=PP1&ots=BkQZaHcR99&lr&pg=PP1#v=onepage&q&f=false](https://books.google.ca/books?id=6V3wNs4uv_4C&lpg=PP1&ots=BkQZaHcR99&lr&pg=PP1#v=onepage&q&f=false).
- Potter, Michael. 2004. *Set Theory and its Philosophy*. Oxford: Oxford University Press.
- Radiolab. 2012. The Turing problem. URL <http://www.radiolab.org/story/193037-turing-problem/>. Podcast audio.
- Reid, Constance. 1986. The autobiography of Julia Robinson. *The College Mathematics Journal* 17: 3–21.
- Reid, Constance. 1996. *Julia: A Life in Mathematics*. Cambridge: Cambridge University Press. URL <https://books.google.ca/books?id=1RtSzQyHf9UC&lpg=PP1&pg=PP1#v=onepage&q&f=false>.
- Robinson, Julia. 1949. Definability and decision problems in arithmetic. *The Journal of Symbolic Logic* 14(2): 98–114. URL <http://www.jstor.org/stable/2266510>.

- Robinson, Julia. 1996. *The Collected Works of Julia Robinson*. Providence: American Mathematical Society.
- Rose, Daniel. 2012. A song about Georg Cantor. URL <https://www.youtube.com/watch?v=QUP5Z4Fb5k4>. Audio Recording.
- Russell, Bertrand. 1905. On denoting. *Mind* 14: 479–493.
- Russell, Bertrand. 1967. *The Autobiography of Bertrand Russell*, vol. 1. London: Allen and Unwin.
- Russell, Bertrand. 1968. *The Autobiography of Bertrand Russell*, vol. 2. London: Allen and Unwin.
- Russell, Bertrand. 1969. *The Autobiography of Bertrand Russell*, vol. 3. London: Allen and Unwin.
- Russell, Bertrand. n.d. Bertrand Russell on smoking. URL [https://www.youtube.com/watch?v=80oLTiVW\\_1c](https://www.youtube.com/watch?v=80oLTiVW_1c). Video Interview.
- Sandstrum, Ted. 2019. *Mathematical Reasoning: Writing and Proof*. Allendale, MI: Grand Valley State University. URL <https://scholarworks.gvsu.edu/books/7/>.
- Segal, Sanford L. 2014. *Mathematicians under the Nazis*. Princeton: Princeton University Press.
- Sigmund, Karl, John Dawson, Kurt Mühlberger, Hans Magnus Enzensberger, and Juliette Kennedy. 2007. Kurt Gödel: Das Album—The Album. *The Mathematical Intelligencer* 29(3): 73–76.
- Smith, Peter. 2013. *An Introduction to Gödel's Theorems*. Cambridge: Cambridge University Press.
- Solow, Daniel. 2013. *How to Read and Do Proofs*. Hoboken, NJ: Wiley.

- Steinhart, Eric. 2018. *More Precisely: The Math You Need to Do Philosophy*. Peterborough, ON: Broadview, 2nd ed.
- Sykes, Christopher. 1992. BBC Horizon: The strange life and death of Dr. Turing. URL <https://www.youtube.com/watch?v=gyusnGbBSHE>.
- Szabo, Manfred E. 1969. *The Collected Papers of Gerhard Gentzen*. Amsterdam: North-Holland.
- Takeuti, Gaisi, Nicholas Passell, and Mariko Yasugi. 2003. *Memoirs of a Proof Theorist: Gödel and Other Logicians*. Singapore: World Scientific.
- Tamassy, Istvan. 1994. Interview with Róza Péter. *Modern Logic* 4(3): 277–280.
- Tarski, Alfred. 1969. Truth and proof. *Scientific American* 220(6): 63–77. URL <https://www.jstor.org/stable/24926385>.
- Tarski, Alfred. 1981. *The Collected Works of Alfred Tarski*, vol. I–IV. Basel: Birkhäuser.
- Theelen, Andre. 2012. Lego turing machine. URL <https://www.youtube.com/watch?v=FTSAiF9AHN4>.
- Turing, Alan M. 1937. On computable numbers, with an application to the “Entscheidungsproblem”. *Proceedings of the London Mathematical Society, 2nd Series* 42: 230–265.
- Tyldum, Morten. 2014. The imitation game. Motion picture.
- Velleman, Daniel J. 2019. *How to Prove It: A Structured Approach*. Cambridge: Cambridge University Press, 3rd ed.
- Wang, Hao. 1990. *Reflections on Kurt Gödel*. Cambridge: MIT Press.
- Zermelo, Ernst. 1904. Beweis, daß jede Menge wohlgeordnet werden kann. *Mathematische Annalen* 59: 514–516. English translation in (Ebbinghaus et al., 2010, pp. 115–119).

Zermelo, Ernst. 1908. Untersuchungen über die Grundlagen der Mengenlehre I. *Mathematische Annalen* 65(2): 261–281. English translation in (Ebbinghaus et al., 2010, pp. 189–229).



# *About the Open Logic Project*

The *Open Logic Text* is an open-source, collaborative textbook of formal meta-logic and formal methods, starting at an intermediate level (i.e., after an introductory formal logic course). Though aimed at a non-mathematical audience (in particular, students of philosophy and computer science), it is rigorous.

Coverage of some topics currently included may not yet be complete, and many sections still require substantial revision. We plan to expand the text to cover more topics in the future. We also plan to add features to the text, such as a glossary, a list of further reading, historical notes, pictures, better explanations, sections explaining the relevance of results to philosophy, computer science, and mathematics, and more problems and examples. If you find an error, or have a suggestion, **please let the project team know**.

The project operates in the spirit of open source. Not only is the text freely available, we provide the LaTeX source under the Creative Commons Attribution license, which gives anyone the right to download, use, modify, re-arrange, convert, and re-distribute our work, as long as they give appropriate credit. Please see the Open Logic Project website at [openlogicproject.org](http://openlogicproject.org) for additional information.