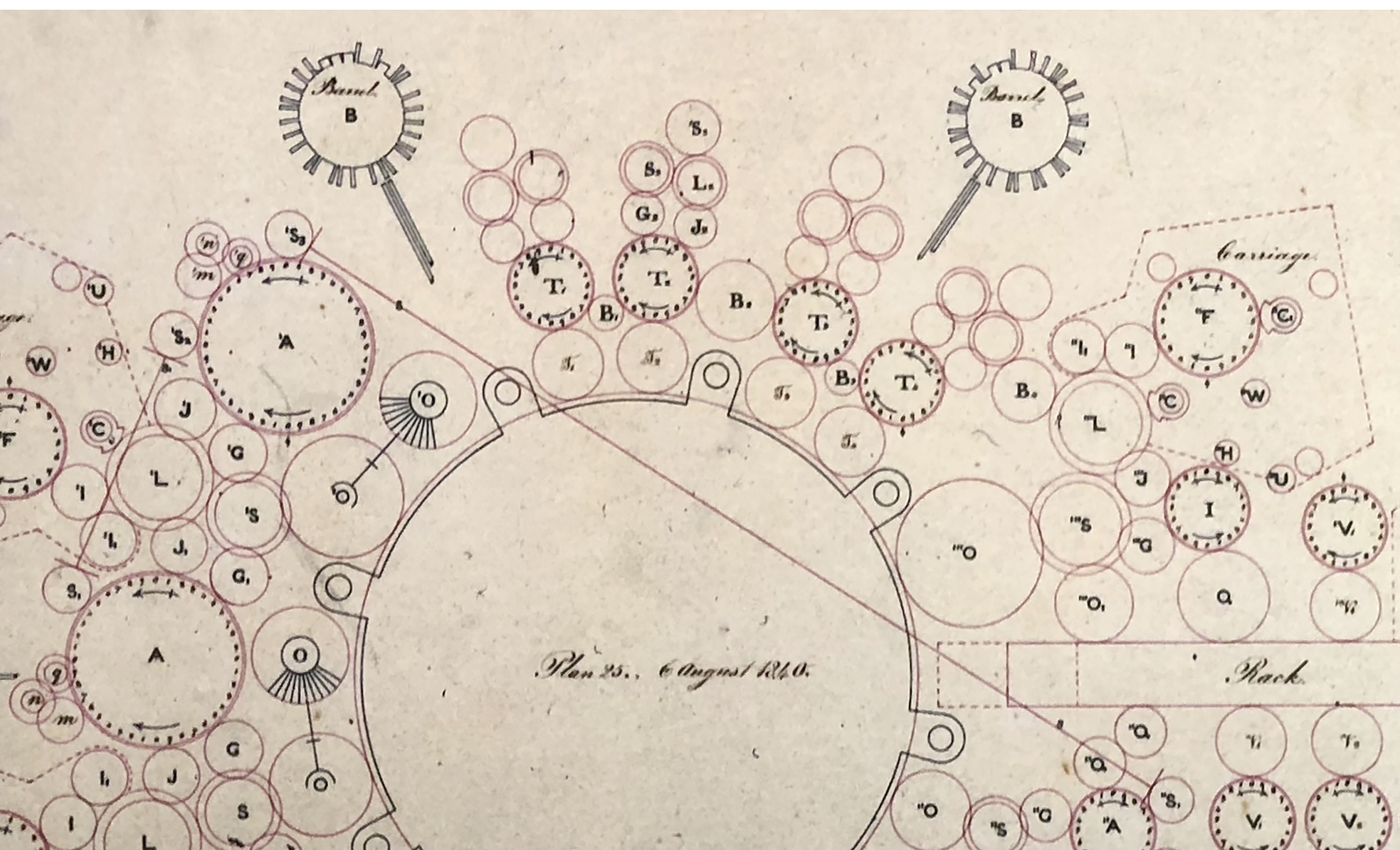


forallx@syr

An Introduction to Formal Logic

By P. D. Magnus, Tim Button, and Michael Rieppel

Spring 2021 Edition



© P. D. Magnus, Tim Button, and Michael Rieppel, 2005-2021

This book is based on *forall x: Cambridge*, by [Tim Button](#) (UCL), used under a [CC BY 4.0](#) license. Button's book is in turn based on the original *forall x*, by [P.D. Magnus](#) (SUNY Albany), used under a [CC BY 4.0](#) license.

The present text, *forallx@syr*, was revised and expanded by [Michael Rieppel](#) (Syracuse University). This work is licensed under a [CC BY-SA 4.0](#) (Creative Commons Attribution-ShareAlike 4.0) license. You are free to copy and redistribute the material in any medium or format, and remix, transform, and build upon the material for any purpose, even commercially, under the following terms:

- ▷ Attribution – You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
- ▷ ShareAlike – If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
- ▷ No additional restrictions – You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Cover Image: plan for Babbage's Analytical Engine (1840). Available at [Wikimedia Commons](#).

The \LaTeX source for this book is available on GitHub at <https://github.com/mrieppel>.

The present version was released on February 2, 2021.

Thanks to Yasmee Hembree, Scott Looney, and Ian York for spotting typos in earlier versions.

Preface

I first learned formal logic from Michael Byrd at UW-Madison, using Lemmon’s *Beginning Logic*, and first taught logic in 2008 as a teaching assistant for Branden Fitelson at UC Berkeley, using Graeme Forbes’ *Modern Logic*. When I started to develop my own logic course, I continued to use Forbes’ book, which I liked for its thorough treatment of the three central components of an introductory logic class: symbolization, semantics, and natural deduction.

However, over time I became dissatisfied with *Modern Logic* for two reasons. First, the Lemmon-style notation that it uses for natural deduction is much less accessible to beginners than the Fitch-style notation found in other texts. And second, Oxford University Press started printing fewer copies of the book, making it rather expensive — more expensive, at any rate, than I thought an introductory logic text should be. By the time I began teaching at Syracuse in 2015, I therefore started listing Forbes’ book only as a recommended text, and relied heavily on the detailed lecture notes I had put together over the years.

But in the longer term, I faced a choice: either select a new textbook, or transform my own notes into a book. The open-source nature of Tim Button’s *forall x: Cambridge* offered me way to do both: I could take his already excellent text (which *inter alia* included a Fitch-style deduction system) and supplement it with material of my own. And so I’ve come to make my own addition to the “groaning shelves” of logic textbooks, to borrow Forbes’ description — though with electronic distribution, the groan has thankfully become more metaphorical. The main additions I’ve made to *forall x: Cambridge* are the following:

- ▷ I’ve changed the first chapter to more closely reflect my own introductory lecture, by e.g. elucidating the modal notion of validity using possible worlds, and emphasizing logic’s focus on formal validity a bit more.
- ▷ I’ve added more on the semantics of truth-functional and especially first-order logic, particularly as concerns the construction of countermodels.
- ▷ I’ve made some changes to set of natural deduction rules in both parts.
- ▷ I’ve reordered the presentation of various topics and revised the practice problems.

Besides the obvious debt the present text owes to the *forall x* editions that P.D. Magnus and Tim Button have so generously made available, and to Forbes’ *Modern Logic*, I also draw on ideas I’ve picked up from Barwise and Etchemendy’s *Language, Proof, and Logic*, Belnap’s *The Art of Logic*, Goldfarb’s *Deductive Logic*, *forall x: Calgary Remix*, and lecture notes by Branden Fitelson, Daniel Warren, and John MacFarlane.

Contents

1	What is Logic?	1
1.1	Arguments	1
1.2	Background Concepts	3
1.3	Good and Bad Arguments	5
1.4	Formal Validity	7
I	Truth Functional Logic	10
2	Symbolization in TFL	11
2.1	Atomic sentences	11
2.2	Negation	12
2.3	Conjunction	13
2.4	Disjunction	16
2.5	Conditional	18
2.6	Biconditional	20
2.7	‘Unless’	21
2.8	Symbolizing Whole Arguments	22
2.9	The Syntax of TFL	23
3	The Semantics of TFL	29
3.1	Meanings for TFL Connectives	29
3.2	Truth-Functionality	31
3.3	Conditionals in TFL and English	32
3.4	Complete Truth Tables	34
3.5	Semantic Concepts	37
3.6	Validity in TFL	40
3.7	Truth Table Shortcuts	45
3.8	Partial Truth Tables	47
4	Natural Deduction for TFL	53
4.1	The Idea Behind Natural Deduction	53
4.2	Setting up Natural Deduction Proofs	54
4.3	Conjunction Rules	56
4.4	Conditional Rules	58
4.5	Additional assumptions and subproofs	61
4.6	Proving Theorems and Reiterating	64
4.7	Biconditional Rules	66

4.8	Negation Rules	68
4.9	Disjunction Rules	72
4.10	Proof strategies	76
4.11	Derived Rules	78
II	First-order logic	84
5	Symbolization in FOL	85
5.1	Names and Predicates	86
5.2	Quantifiers and Quantifier Scope	88
5.3	Common Quantifier Phrases and Domains	91
5.4	Quantifiers and Negation	94
5.5	The Utility of Paraphrase	96
5.6	Many-Place Predicates	99
5.7	Multiple Generality	101
5.8	Intermediate Steps to Symbolization	103
5.9	Adding Identity	106
5.10	The Syntax of FOL	111
6	Natural deduction for FOL	116
6.1	Universal elimination	117
6.2	Existential introduction	119
6.3	Universal introduction	122
6.4	Existential elimination	125
6.5	Rules for identity	131
7	The Semantics of FOL	135
7.1	Predicates and their Extensions	135
7.2	FOL Interpretations	137
7.3	Truth in FOL	139
7.4	Truth-Rules for Quantified Sentences	141
7.5	Truth in an Interpretation: Examples	143
7.6	Semantic Concepts	147
7.7	Countermodels with One-Place Predicates	148
7.8	Countermodels with Many-Place Predicates	151
7.9	Validity and Decidability	154
7.10	Working with Other Semantic Concepts	156
7.11	Semantics for Identity	158
7.12	Appendix: Semantics with Variable Assignments	159
8	Quick Reference	162

“When you come to any passage you don’t understand, *read it again*: if you *still* don’t understand it, *read it again*: if you fail, even after *three* readings, very likely your brain is getting a little tired. In that case, put the book away, and take to other occupations, and next day, when you come to it fresh, you will very likely find that it is *quite* easy.

“If possible, find some genial friend, who will read the book along with you, and will talk over the difficulties with you. *Talking* is a wonderful smoother-over of difficulties. When I come upon anything — in Logic or in any other hard subject — that entirely puzzles me, I find it a capital plan to talk it over, *aloud*, even when I am all alone.”

Lewis Carroll, *Symbolic Logic* (1897)

What is Logic?

1

1.1 Arguments

This book provides an introduction to logic. But what is logic? This is a surprisingly difficult question, still debated by philosophers. But generally speaking, logic is about distinguishing *valid* from *invalid* arguments.

In everyday language, the word ‘argument’ is often used to describe an activity that people engage in. On twitter, or youtube, or the news, people often have heated debates, and you’ve probably had arguments like this with your family or friends. Logicians tend to be a pretty sedate crowd, and they mean something very different by ‘argument’. In logic, an argument is just a collection of statements. More specifically:

An ARGUMENT is a collection of one or more *statements*, exactly one of which is the argument’s *conclusion* and the rest of which are its *premises*.

Here is an example of an argument:

- (1) All rabbits are mammals.
Bugs Bunny is a rabbit.
∴ Bugs Bunny is a mammal.

This argument consists of three statements. One of them is the argument’s conclusion, which we indicate by the three dots ∴. These dots are read as “therefore.” The rest are the premises. This argument has two premises, but arguments can have any number of premises (though, again, only one conclusion).

Notice that our logician’s definition of an argument is very permissive. Consider the following:

- There is a bassoon-playing dragon in the *Cathedra Romana*.
∴ Salvador Dali was a poker player.

We have a premise and a conclusion, and so we have an argument. Admittedly, it’s a *terrible* argument, but it is still an argument.

Here’s another argument, one that’s not as obviously terrible:

- (2) All rabbits are mammals.
Winnie the Pooh is a mammal.
∴ Winnie the Pooh is a rabbit.

In this case the premises at least involve the same concepts as the conclusion. But this argument still isn't as good as (1) from earlier: unlike the earlier example, this argument isn't *valid* — its conclusion doesn't *follow from* its premises. But what exactly does validity, or “following from,” consist in? What's wrong with argument (2) as compared to (1)?

One thing that's worse about the second argument is that its conclusion is false: Pooh isn't a rabbit, he's a bear! But that isn't what distinguishes valid from invalid arguments in general, because there are valid arguments with false conclusions, and invalid arguments with true conclusions. For example:

(3) All rabbits are birds.
Winnie the Pooh is a rabbit.
∴ Winnie the Pooh is a bird.

(4) All rabbits are mammals
Bugs Bunny is a mammal.
∴ Bugs Bunny is a rabbit.

Argument (3) is valid, but has a false conclusion. And argument (1) has a true conclusion, as well as true premises, but it still isn't valid, because its conclusion doesn't follow from its premises: that all rabbits are mammals and that Bugs Bunny is a mammal doesn't yet guarantee that Bugs Bunny is a rabbit.

Validity isn't determined by whether the premises or the conclusion are as a matter of fact true. It rather has to do with the *relationship* between the premises and the conclusion. When we ask about validity we want to know whether, if all the premises *were true*, the conclusion would also *have to be true*. Put another way:

An argument is **VALID** if and only if it is *impossible* for all of its premises to be true but its conclusion false.

Let's unpack some of the concepts involved in the two definitions we've encountered a little bit more.

■ Exercises 1.1

As you've seen, we always put the conclusion at the end of an argument and indicate it using the three “therefore dots” ∴. Informally presented arguments don't always have the conclusion at the end, however — it can appear at the beginning, or even in the middle. In each of the following arguments, highlight the phrase which expresses the conclusion:

1. It is sunny. So I should take my sunglasses.
2. It must have been sunny. I did wear my sunglasses, after all.
3. No one but you has had their hands in the cookie-jar. And the scene of the crime is littered with cookie-crumbs. You're the culprit!
4. Miss Scarlett and Professor Plum were in the study at the time of the murder. And Reverend Green had the candlestick in the ballroom, and we know that there is no blood on his hands. Hence Colonel Mustard did it in the kitchen with the lead-piping. Recall, after all, that the gun had not been fired.

1.2 Background Concepts

First, we said that an argument is a collection of STATEMENTS. Statements are sentences that are either true or false. Truth and falsity are called TRUTH-VALUES. The truth-value of a statement is determined by what the world is like. A statement like ‘Syracuse is in New York State’ describes the world as being a certain way. This statement happens to be true because the world in fact is as the statement describes it. ‘Syracuse is in Alaska’, by contrast, describes the world incorrectly, and is therefore false. As the ancient philosopher (and logician!) Aristotle put it in his book *Metaphysics*:

“To say of what is that it is not, or of what is not that it is, is false, while to say of what is that it is, and of what is not that it is not, is true.” (1011b25)

It’s important to notice that not all English sentences count as statements in this sense. For example, none of the following sentences can be assessed as true or false:

- Welcome to the Syracuse Airport!
- Please have your ID ready.
- Are there any liquids in your bag?

A sentence like ‘please have your ID ready’ isn’t meant to describe the world, but to ask you to do something. Similarly, ‘Welcome to Syracuse Airport’ is just a greeting, and isn’t meant to offer an accurate or inaccurate description of the world. And although the answer to the last question on this list has a truth-value, the question itself doesn’t. In general, things like greetings, requests, orders, and questions don’t have truth-values, therefore don’t count as statements, and for that reason can’t be premises or conclusions of arguments. Though it’s important to keep this point in mind, moving forward we’ll generally use the words “statement” and “sentence” interchangeably.

Next, notice that because a statement’s truth value depends on what the world is like, its truth-value could have been different if the world had been different. For example, the sentence ‘Rieppel is a philosopher’ is in fact true, but if I had taken up a different career it would have been false. Conversely, ‘Rieppel is a professional juggler’ is false, but if I had gone to juggling school instead of continuing with philosophy, it would have been true.

Philosophers often invoke the notion of a POSSIBLE WORLD in this connection. The idea is that besides the actual world, there are various other possible worlds, other ways things could have been — alternative histories, or alternative universes, if you like. ‘Rieppel is a professional juggler’ is false in the actual world, but it is true in other possible worlds, ones where I went to juggling school or joined the circus. Similarly, ‘Rieppel is a philosopher’ is as a matter of fact true, but it is false in other possible worlds where I didn’t pursue philosophy. Sentences like this, which are true in some possible worlds and false in others, are said to be CONTINGENT.

Other sentences are not contingent. For example, ‘Syracuse either is or is not in New York State’ isn’t just true in the actual world, it’s true in *every* possible world, that is, it’s a NECESSARY TRUTH. Mathematical truths are another example: ‘ $2 + 2 = 4$ ’ is again true in every possible world, and therefore a necessary truth. At the other extreme, sentences like ‘Syracuse both is and is not in New York State’ and ‘ $2 + 2 = 5$ ’ are false in every possible world, or NECESSARILY FALSE.

Returning to arguments, you can think of the notion of validity in terms of possible worlds too. We said that an argument is valid just in case it's *impossible* for all of its premises to be true but its conclusion false. Phrased in terms of possible worlds, this becomes:

An argument is VALID if and only if *there is no possible world* where all of its premises are true but its conclusion is false.

Equivalently put: an argument is valid if its conclusion is true in *every possible world* in which all of its premises are true.

This gives us an informal way to test whether an argument is valid: we imagine a world where all the premises are true, and then ask ourselves whether the conclusion would have to be true as well at that world. If so, the argument is valid. On the other hand, if you can imagine a world where all the premises are true but the conclusion is still false, the argument isn't valid. So again, whether an argument is valid or not isn't determined by whether its premises and conclusion are *actually* true or false. It's about the *connection* between them — whether there's *any way* for the premises to be true but the conclusion false.

There are other logical concepts that we'll encounter in this class that involve the notions of necessity and possibility, besides validity. Some we've already mentioned:

- ▷ A sentence is CONTINGENT if and only if it is possible for it to be true, and also possible for it to be false.
- ▷ A sentence is a NECESSARY TRUTH if and only if it is not possible for it to be false.
- ▷ A sentence is a NECESSARY FALSEHOOD if and only if it is not possible for it to be true.
- ▷ Two sentences are CONTRADICTORY if and only if they necessarily have opposite truth values.
- ▷ Two sentence are EQUIVALENT if and only if they necessarily have the same truth value.
- ▷ A collection of sentences is JOINTLY CONSISTENT if and only if it is possible for all of them to be true together, and JOINTLY INCONSISTENT otherwise.

Notice that these concepts apply to different things. Whereas the first three concern properties of single sentences, the next two concern relations between two sentences, and the last ones concern properties of whole collections of sentences. Validity is again slightly different, because it is a property had (or lacked) by only those collections of sentences that also have a *designated conclusion*, i.e. by those collections that are *arguments*.

■ Exercises 1.2

A. For each of the following: is it necessarily true, necessarily false, or contingent?

1. Caesar crossed the Rubicon.
2. Someone once crossed the Rubicon.
3. No one has ever crossed the Rubicon.

4. If Caesar crossed the Rubicon, then someone has.
5. Even though Caesar crossed the Rubicon, no one has ever crossed the Rubicon.
6. If anyone has ever crossed the Rubicon, it was Caesar.

B. Consider the following sentences:

- G1. There are at least four giraffes at the zoo.
- G2. There are exactly seven gorillas at the zoo.
- G3. There are not more than two martians at the zoo.
- G4. Every giraffe at the zoo is a martian.

Now, for each of the following, determine if the sentences in question are jointly consistent or jointly inconsistent:

1. G2, G3, and G4
2. G1, G3, and G4
3. G1, G2, and G4
4. G1, G2, and G3

C. Could there be:

1. Jointly consistent sentences, one of which is necessarily false?
2. Jointly consistent sentences, one of which is a necessary truth?
3. Jointly inconsistent sentences, one of which is a necessary truth?

In each case: if so, give an example; if not, explain why not.

1.3 Good and Bad Arguments

Being valid is certainly one thing that makes for a good argument, intuitively speaking. But there's more to being a good argument than that.

First off, if an argument has an obviously false premise, then even if it is valid, it remains of limited interest because it doesn't establish its conclusion. By contrast, if an argument is valid and all of its premises are true, then we know that its conclusion has to be true too. Arguments like this are said to be *sound*:

An argument is **SOUND** if and only if (i) it is valid, and (ii) has premises that are in fact true.

Arguments are generally intended to be not just valid, but sound. So if you're faced with an argument, in a philosophy class or elsewhere, whose conclusion you want to resist, you have two options: you can either try to show that the argument is not valid, or you can try to show that one of its premises is false (and the argument therefore isn't sound). What you *can't* do is accept it as valid, and concede that its premises are true, but still reject the conclusion as false: if it's valid, and has true premises, its conclusion has to be true too.

Although it's important in practice to determine whether or not the premises of an argument are in fact true, it is (for the most part) not the job of logic to do this. The job of logic

is just to determine whether or not an argument is valid. The task of determining whether the argument's premises are in fact true (and the argument sound) is usually best left to experts in the relevant field: biologists, historians, philosophers, physicists, economists, or whomever.

A second way in which validity is not all there is to good argumentation comes out if you consider the following:

In January 2016, it snowed in Syracuse.

In January 2017, it snowed in Syracuse.

In January 2018, it snowed in Syracuse.

In January 2019, it snowed in Syracuse.

In January 2020, it snowed in Syracuse.

So: It snows every January in Syracuse.

This argument generalizes from observations about several past cases to a conclusion about all cases. The argument isn't valid in our sense, because even if it snowed in January in many recent years, that doesn't mean it's *impossible* for it not to snow in some future year. The argument could be made stronger by adding additional premises, about other snowy Syracuse Januaries in the past. But however many premises of this sort we add, the argument will remain invalid.

That doesn't mean that it's a bad argument. Arguments like this one are called **INDUCTIVE** arguments, and they are often used legitimately and with great success in science and everyday life. In this book, we will set aside the difficult question of what makes for a good inductive argument. What logic studies is the different notion of **DEDUCTIVE** validity — where the truth of the premises has to *guarantee* the truth of the conclusion — and this will be the focus of our concern.

■ Exercises 1.3

Here are some exercises to test your understanding of deductive validity and related concepts we've discussed. For these questions, you don't need to worry about the distinction between validity and "validity in virtue of logical form" to be discussed in §1.4 below. You should just use the definition of validity we gave in §1.1 and §1.2 above.

A. Which of the following arguments are valid? Which are invalid?

1. Socrates is a man.
 All men are carrots.
 ∴ Socrates is a carrot.
2. Abe Lincoln was either born in Illinois or he was once president.
 Abe Lincoln was never president.
 ∴ Abe Lincoln was born in Illinois.
3. If I pull the trigger, Abe Lincoln will die.
 I do not pull the trigger.
 ∴ Abe Lincoln will not die.

4. Abe Lincoln was either from France or from Luxemborg.
Abe Lincoln was not from Luxemborg.
∴ Abe Lincoln was from France.
5. If the world were to end today, then I would not need to get up tomorrow morning.
I will need to get up tomorrow morning.
∴ The world will not end today.
6. Joe is now 19 years old.
Joe is now 87 years old.
∴ Bob is now 20 years old.

B. Could there be:

1. A valid argument that has one false premise and one true premise?
2. A valid argument that has only false premises but a true conclusion?
3. A valid argument with only false premises and a false conclusion?
4. A valid argument with only true premise but a false conclusion?
5. An invalid argument with only true premises and a true conclusion?
6. An invalid argument with only false premises but a true conclusion?
7. A sound argument with a false conclusion?
8. A sound argument with at least one false premise?
9. An invalid argument that can be made valid by the addition of a new premise?
10. A valid argument that can be made invalid by the addition of a new premise?
11. A valid argument, the conclusion of which is necessarily false?
12. An invalid argument, the conclusion of which is necessarily true?
13. A valid argument whose premises are jointly inconsistent?
14. A valid argument with only one premise?

In each case: if so, give an example; if not, explain why not.

1.4 Formal Validity

There's one last complication we have to address before setting out on our investigation of logic. Consider the following arguments:

- (5) This beach ball is green all over.
∴ This beach ball is not red all over.
- (6) Reihan is a bachelor.
∴ Reihan is not married.

In both cases it is impossible for the premise to be true and the conclusion false: if something's green all over it can't be any other color, and being unmarried is part of what it is to be a bachelor. Both arguments are therefore valid.

But there's an important difference between valid arguments like these and one like the following:

- (7) Jenny is either happy or sad.
 Jenny is not happy.
 \therefore Jenny is sad.

This argument is also valid, but there's more. It has a special structure, or logical form, that we might represent as follows:

A or B
 not- A
 $\therefore B$

This is an excellent structure for an argument to have, because any argument of this form will be valid, no matter what sentences we put in place of A and B ! Or consider our Bugs Bunny argument, which has the structure represented to the right:

- | | |
|--------------------------------------|-----------------------|
| (1) All rabbits are mammals. | All F are G |
| Bugs Bunny is a rabbit. | a is F |
| \therefore Bugs Bunny is a mammal. | $\therefore a$ is G |

Again, this is a great structure, because any argument of this form will be valid, no matter what predicates we put in for F and G or what name we put in for a .

The general point is that arguments like (7) and (1) are valid simply *in virtue of their logical form*. They each exhibit a logical structure which renders *any* argument with that structure valid. By contrast, arguments (5) and (6), though valid, are not valid in virtue of their logical form. For example, the form of (6) could be represented as follows:

- | | |
|-------------------------------------|------------------------------|
| (6) Reihan is a bachelor | a is F |
| \therefore Reihan is not married. | $\therefore a$ is not- G . |

Here the premise ascribes a certain property (being a bachelor) to an individual, and the conclusion then denies another property (being married) of that individual. However, there are other arguments that share this same structure but aren't valid:

Reihan is a runner.
 \therefore Reihan is not married.

This isn't valid because it's trivial to imagine a world where Reihan is a runner but also married. What made argument (6) valid wasn't its logical form, but the *specific meanings* of the words 'bachelor' and 'married' that occur in its premise and its conclusion. Other arguments that have the same form but involve words with different meanings (e.g. 'runner') may no longer be valid. Logic is all about identifying *patterns* that make arguments valid. So it only cares about FORMALLY VALID arguments like (1) and (7), not arguments like (5) and (6) that are valid for reasons other than their logical form.

Due to logic's concern with form, we will approach the task of distinguishing valid from invalid arguments in an indirect way. We will first introduce a formal language in which we can symbolize English arguments. Doing this lets us represent the logical forms of those arguments. We will then give a precise definition of validity for arguments cast in this formal notation. And this will in turn give us our indirect means of distinguishing valid from invalid

arguments in English: if an English argument can be symbolized as a valid argument in our formal notation, then that English argument is formally valid.

In fact, we will study two systems of logic, involving two different formal languages. These systems will differ in what words they treat as *logical constants*, that is, which words they treat as indicative of logically significant structure, or form. The first system we will study is *Truth-Functional Logic* (or TFL). It will let us represent the structure of arguments like (7) via the symbolization to the right:

(7) Jenny is either happy or sad.	$(A \vee B)$
Jenny is not happy.	$\neg A$
\therefore Jenny is sad.	$\therefore B$

This language treats words like ‘either ... or’ and ‘not’ as logical constants (represented by ‘ \vee ’ and ‘ \neg ’ respectively), and will use upper-case letters to represent complete statements (like ‘A’ for ‘Jenny is happy’, and ‘B’ for ‘Jenny is sad’). TFL is the topic of Part 1 of this book.

In Part 2 of the book, we will turn to *First-Order Logic* (or FOL). It will let us represent the structure of things like the Bugs Bunny argument via the symbolization to the right:

(1) All rabbits are mammals.	$\forall x(Fx \rightarrow Gx)$
Bugs Bunny is a rabbit.	Fa
\therefore Bugs Bunny is a mammal.	$\therefore Ga$

This system extends Truth-Functional Logic by treating words like ‘all’ as logical constants (represented by ‘ \forall ’).¹ It also lets us represent some of the *internal* structure of a simple statement like ‘Bugs Bunny is a rabbit’, showing that it is formed by combining the name ‘Bugs Bunny’ (represented as ‘a’) with the predicate ‘is a rabbit’ (represented by ‘F’). With this very short preview out of the way, let’s get started with logic!

¹At this point you might be wondering how logicians decide which English words to treat as logical constants, and represent by special logical symbols that indicate “logical form.” This is a difficult philosophical question about logic that we won’t have a chance to delve into here. If you’re interested, check out the *Stanford Encyclopedia of Philosophy*’s entries on [Logical Constants](#) and [Logical Consequence](#).

I

Truth Functional Logic

Symbolization in TFL

2

2.1 Atomic sentences

In this chapter, we'll look at how to symbolize English arguments in the language of *truth-functional logic* (or TFL), thereby revealing their truth-functional logical structure, or form.¹ Here is an example of a symbolization in TFL:

- (1) It is raining outside.
If it is raining outside, then Jenny is miserable.
 \therefore Jenny is miserable.

A	A
If A , then C	$(A \rightarrow C)$
$\therefore C$	$\therefore C$

In symbolizing this argument, I began by replacing *subsences* of larger sentences with uppercase letters. For example, 'it is raining outside' is a subsentence of 'If it is raining outside, then Jenny is miserable', and we replaced this subsentence with the letter 'A'. Similarly, we used 'C' to symbolize the subsentence 'Jenny is miserable'.

Our formal language, TFL, uses uppercase letters as its ATOMIC SENTENCES. These will be the basic building blocks, or "atoms," out of which more complex TFL sentences are built. There are only twenty-six letters of the alphabet, but in principle there is no limit to the number of atomic sentences that we might want to consider. By adding numerical subscripts to letters, we obtain as many atomic sentences as we need. So the following all count as atomic sentences of TFL:

$$A, B, C, P, X, Z, A_1, A_2, P_{12}, P_{234}$$

To indicate which atomic sentence of TFL is being used to represent which English sentence, we provide a SYMBOLIZATION KEY like the following:

- A: It is raining outside
C: Jenny is miserable

Doing this does *not* fix this symbolization *once and for all*. We are just saying that, for the time being, we will use the atomic TFL sentence 'A' to symbolize the English sentence 'It

¹What does "truth-functional" mean? We'll get to that in the next chapter, on the semantics of TFL.

is raining outside’, and ‘*C*’ to symbolize ‘Jenny is miserable’. Later, when we are dealing with different sentences or different arguments, we can provide a new symbolization key that associates these atomic TFL sentences with different English sentences.

It is important to recognize that whatever internal structure an English sentence might have is lost if it is symbolized by an atomic sentence of TFL. From the point of view of TFL, an atomic sentence has no internal (or “subatomic”) structure. It can be used to build more complex sentences, but it cannot be taken apart.

For this reason, English sentences that have an internal logical structure to them, like the conditional ‘If it’s raining outside, then Jenny is miserable’, should not be symbolized using atomic sentences of TFL. If we just symbolized this sentence as ‘*P*’, for example, our symbolization would obscure the fact that it has the form of an ‘if ... then ...’ statement, and that it contains subsentences that also occur on their own as premise and conclusion in our argument (1) above. We would therefore miss out on the logical form in virtue of which argument (1) is valid. For this reason we have to symbolize logically complex sentences of English via complex (i.e. non-atomic) sentence of TFL, in this case ‘ $(A \rightarrow C)$ ’.

Complex sentences of TFL are built up by combining atomic sentences with *connectives*. TFL connectives will be used to symbolize English connectives like ‘if ... then’, ‘or’ and ‘not’. Just as these English connectives can be applied to sentences to form new, bigger sentences, so TFL connectives can be applied to atomic sentences to build larger, complex sentences of TFL. In total, there are five connectives in TFL. This table summarizes them, and gives you a rough indication of their meaning:

symbol	name	rough meaning
\neg	negation	‘It is not the case that. ...’
\wedge	conjunction	‘Both... and ...’
\vee	disjunction	‘Either... or ...’
\rightarrow	conditional	‘If ... then ...’
\leftrightarrow	biconditional	‘... if and only if ...’

For the remainder of this chapter, we have two main objectives, a practical one and a technical one. The first, practical objective is to learn how to symbolize logically complex English sentences using our TFL connectives. The second, more technical objective will be to give a precise grammar, or syntax, for the language of TFL, which explains exactly how TFL connectives combine with atomic sentences to produce complex TFL sentences. We’ll turn to this technical task later, after we’ve gotten some practice using the language of TFL to symbolize English sentences.

2.2 Negation

Consider how we might symbolize these sentences:

- (2) Mary is in Barcelona.
- (3) It is not the case that Mary is in Barcelona.
- (4) Mary is not in Barcelona.

In order to symbolize (2), we will need an atomic sentence. We might offer this symbolization key:

B: Mary is in Barcelona.

Since (3) is obviously related to (2), we do not want to symbolize it with a completely different sentence, say '*A*'. Sentence (3) basically says 'It is not the case that *B*'. In order to symbolize this, we need the symbol for negation, ' \neg '. So we can symbolize (3) as ' $\neg B$ '.

Sentence (4) also contains the word 'not'. And it is obviously equivalent to (3), so we can also symbolize it with ' $\neg B$ '.

A sentence can be symbolized as $\neg\phi$ if it can be paraphrased in English as 'It is not the case that ϕ '.

Here's another example:

- (5) The widget can be replaced.
- (6) The widget is irreplaceable.
- (7) The widget is not irreplaceable.

Let's use the following symbolization key:

R: The widget is replaceable

Sentence (5) can then be symbolized by '*R*'. Next, (6) says the widget is irreplaceable, which means that it is not the case that the widget is replaceable. So even though (6) does not contain the word 'not', we can still symbolize it as ' $\neg R$ '. Sentence (7) can be paraphrased as 'It is not the case that the widget is irreplaceable.' Which can again be paraphrased as 'It is not the case that it is not the case that the widget is replaceable'. So we can symbolize this with the TFL sentence ' $\neg\neg R$ '.

At this point you might be wondering: don't double negatives cancel out, so that ' $\neg\neg R$ ' is equivalent to plain '*R*'? We'll get into meaning, or semantics, of TFL sentences in the next chapter; but the quick answer is that, yes, these are equivalent, and (7) *could* also be symbolized as '*R*'. Still, ' $\neg\neg R$ ' is the *preferred* symbolization, because it represents more of the logical structure implicit in the English sentence (7).

Some care is needed when handling negations. Consider:

- (8) Jane is happy.
- (9) Jane is unhappy.

If we let '*H*' symbolize 'Jane is happy', we can symbolize (8) as '*H*'. However, it would be a mistake to symbolize (9) with ' $\neg H$ '. Sentence (9) does not mean the same thing as 'It is not the case that Jane is happy'. Jane might be neither happy nor unhappy; she might be in a state of blank indifference. In order to symbolize (9), then, we would need a new atomic sentence of TFL.

2.3 Conjunction

Consider the sentence:

(10) Adam is athletic, and Barbara is athletic.

We will need separate atomic sentences to symbolize the two subsentences in (10):

A: Adam is athletic.

B: Barbara is athletic.

We will use ' \wedge ' to symbolize 'and', and thus symbolize (10) as ' $(A \wedge B)$ '. This connective is called CONJUNCTION. We also say that '*A*' and '*B*' are the two CONJUNCTS of the conjunction ' $(A \wedge B)$ '.

There are few things to notice about conjunction. First, in English the word 'and' doesn't always conjoin two sentences:

(11) Barbara is athletic and energetic.

(12) Barbara and Adam are both athletic.

In (11) the word 'and' conjoins two adjectives, rather than two sentences. But it can be paraphrased as 'Barbara is athletic and Barbara is energetic' where 'and' now does conjoin two sentences. So if we use '*E*' symbolize 'Barbara is energetic', we can symbolize the entire sentence as ' $(B \wedge E)$ '. In sentence (12), 'and' conjoins two names. Again, though, this can be paraphrased in terms of a conjunction of two sentences, as 'Barbara is athletic and Adam is athletic', and can therefore be symbolized as ' $(B \wedge A)$ '.²

Second, conjunction can be expressed in English using words other than 'and':

(13) Although Barbara is energetic, she is not athletic.

(14) Adam is athletic, but Barbara is not.

In (13), the word 'although' sets up a contrast between the first part of the sentence and the second part. Nevertheless, the sentence tells us both that Barbara is energetic and that she is not athletic. So we can symbolize it as a conjunction:

B: Barbara is athletic.

E: Barbara is energetic.

Symbolization of (13): $(E \wedge \neg B)$

Of course we have lost all sorts of nuance in this symbolization. There is a distinct difference in tone between the English sentence (13) and 'Both Barbara is energetic and it is not the case that Barbara is athletic'. TFL does not (and cannot) preserve these nuances, however, and so we do not attend to them when symbolizing English into TFL. Notice also that in the symbolization key we replaced the pronoun 'she' with 'Barbara', since it might otherwise be unclear who 'she' is meant to refer to. In general: always use names in place of pronouns in your symbolization key!

²Some care is needed with this. Not *all* sentences where 'and' conjoins two names can be paraphrased in a way where 'and' conjoins two sentences. For example, 'Barbara and Adam carried the piano upstairs' may not mean the same as 'Barbara carried the piano upstairs and Adam carried the piano upstairs', since the latter (but not the former) is compatible with them each carrying it individually rather than together.

Sentence (14) raises similar issues. The word ‘but’ sets up a contrast between the two parts of the sentence, but this is not something that TFL can deal with. We can paraphrase the sentence as ‘Both Adam is athletic, and Barbara is not athletic’. Notice that the second conjunct involves a negation as well. Using the sentence letters already introduced, we can symbolize (14) as ‘ $(A \wedge \neg B)$ ’.

There are other words besides ‘although’ and ‘but’ that can be used to express conjunction. For example, ‘Barbara is energetic, however she is not athletic’ and ‘Barbara is energetic despite not being athletic’ expresses the same conjunctive claim as sentence (13) from above, and get symbolized using the same TFL sentence. In general, the symbolization guideline for conjunction is:

A sentence can be symbolized as $(\phi \wedge \psi)$ if (nuance aside) it can be paraphrased in English as ‘Both ϕ and ψ ’.

You might be wondering why we always put *brackets* around the conjunctions. The reason can be brought out by thinking about negation interacts with conjunction. Consider:

- (15) You will not get both soup and salad.
 (16) You will not get soup but you will get salad.

Sentence (15) can be paraphrased as ‘It is not the case that: both you will get soup and you will get salad’. Using the following symbolization key:

- S_1 : You get soup.
 S_2 : You get salad.

we would symbolize ‘both you will get soup and you will get salad’ as ‘ $(S_1 \wedge S_2)$ ’. To symbolize the whole of sentence (15), then, we negate this, giving us ‘ $\neg(S_1 \wedge S_2)$ ’. Sentence (16), on the other hand, gets symbolized as a conjunction whose first conjunct is negated ‘ $(\neg S_1 \wedge S_2)$ ’.

These English sentences mean different things, and their symbolizations differ accordingly. In one of them, the entire conjunction is negated. In the other, just one conjunct is negated. Brackets help us to keep track of things like the *scope* of the negation: whether it applies to the entire conjunction, or just to the first conjunct.

■ Exercises 2.3

A. Symbolize each English sentence in TFL.

1. Simon and his sister went shopping.
2. Melissa gave her dog a biscuit but he didn’t like it.
3. Simon and Melissa did not both go shopping.
4. Bill’s coat was not expensive despite being of good quality.

2.4 Disjunction

Consider these sentences:

- (17) Either Denison will play golf with me, or he will watch movies.
 (18) Either Denison or Ellery will play golf with me.

We can use the following symbolization key for these sentences (notice that we, as always, replace pronouns with names):

D : Denison will play golf with me.
 E : Ellery will play golf with me.
 M : Denison will watch movies.

Sentence (17) is an ‘either ... or’ statement, and gets symbolized as ‘ $(D \vee M)$ ’. The connective is called DISJUNCTION. We also say that ‘ D ’ and ‘ M ’ are the DISJUNCTS of the disjunction ‘ $(D \vee M)$ ’.

Sentence (18) is only slightly more complicated. Here ‘or’ occurs between two names rather than two complete sentences. However, we can paraphrase sentence (18) as ‘Either Denison will play golf with me, or Ellery will play golf with me’ where ‘or’ now connects two complete sentences. So we can symbolize it as ‘ $(D \vee E)$ ’.

A sentence can be symbolized as $(\phi \vee \psi)$ if it can be paraphrased in English as ‘Either ϕ or ψ .’

Sometimes in English, the word ‘or’ excludes the possibility that both disjuncts are true. This is called an EXCLUSIVE OR. An *exclusive or* is intended when it says, on a restaurant menu, ‘Entrees come with either soup or salad’: you may have soup; you may have salad; but, if you want *both* soup *and* salad, then you have to pay extra. At other times, the word ‘or’ allows for the possibility that both disjuncts might be true. This is probably the case with sentence (18), above: I might play golf with Denison, with Ellery, or with both Denison and Ellery. Sentence (18) merely says that I will play with *at least* one of them. This is an INCLUSIVE OR.

Importantly, the TFL symbol ‘ \vee ’ expresses *inclusive or*. Whenever you see the English words ‘either ... or’ in this book, you can assume that the inclusive sense is intended, and symbolize such sentences using ‘ \vee ’. When an exclusive sense is intended, we will always add an explicit ‘but not both’, as in:

- (19) The entree will come with either soup or salad, but not both.

Using ‘ S_1 ’ for ‘the entree will come with soup’ and ‘ S_2 ’ for ‘the entree will come with salad’, we can symbolize (19) as ‘ $((S_1 \vee S_2) \wedge \neg(S_1 \wedge S_2))$ ’. So although the TFL symbol ‘ \vee ’ always symbolizes *inclusive or*, we can symbolize an *exclusive or* in TFL. We just have to use a few other TFL symbols in addition to ‘ \vee ’!

There are some further interesting interactions between disjunction and negation that we should attend to. Consider the following:

- (20) Either Barbara will not get soup, or she will not get salad.

(21) Barbara will get neither soup nor salad.

Sentence (20) can be paraphrased as: ‘*Either* it is not the case that Barbara will get soup, *or* it is not the case that Barbara will get salad’. Using the following symbolization key:

P : Barbara will get soup.
 D : Barbara will get salad.

we can symbolize (20) as $(\neg P \vee \neg D)$. This has the form of a disjunction both disjuncts of which are negated. Sentence (21) has a different structure. It can be paraphrased as, ‘*It is not the case that*: either Barbara will get soup or Barbara will get salad’. Since this negates the entire disjunction, we symbolize sentence (21) as $\neg(P \vee D)$. This differs from our symbolization of (20), as it should given that the two English sentences mean different things.

You may have noticed that (20) means the same thing as ‘Barbara will not get both soup and salad.’ The latter can be symbolized as a negated conjunction, $\neg(P \wedge D)$. Since $\neg(P \wedge D)$ symbolizes an English sentence that’s equivalent to (20), and (20) is symbolized as $(\neg P \vee \neg D)$, we can conclude that the TFL sentences $\neg(P \wedge D)$ and $(\neg P \vee \neg D)$ are themselves equivalent.

Similarly, (21) means the same as ‘Barbara will not get soup and Barbara will not get salad’, which can be symbolized as a conjunction of two negations $(\neg P \wedge \neg D)$. So again, since $(\neg P \wedge \neg D)$ symbolizes an English sentence that’s equivalent to (21), and (21) is symbolized as $\neg(P \vee D)$, we can conclude that these two TFL sentences are themselves equivalent. What we’ve discovered are:

DEMORGAN’S LAWS:

$\neg(\varphi \vee \psi)$ is equivalent to $(\neg\varphi \wedge \neg\psi)$
 $\neg(\varphi \wedge \psi)$ is equivalent to $(\neg\varphi \vee \neg\psi)$

These laws are named after Augustus DeMorgan who first explicitly formulated them in the nineteenth century. These laws will stay with us throughout our study of logic, and we will see how to prove that these equivalences hold in the next chapter.³

■ Exercises 2.4

A. Symbolize each English sentence in TFL.

1. Socrates was neither tall nor handsome.
2. Liz isn’t flying to both San Francisco and New York today.
3. The package ended up in either Korea or Japan, but not both.
4. The book was both intelligent and funny, but the movie was neither.
5. Li and Virag did not both go to the party.

³It is noteworthy that these equivalences only hold in TFL given that ‘ \vee ’ expresses inclusive disjunction. The fact that the corresponding English sentences are also intuitively equivalent again shows that English ‘or’ often expresses inclusive disjunction.

6. Li and Virag both did not go to the party.

B. We symbolized *exclusive or*, or *xor*, using ' \vee ', ' \wedge ', and ' \neg '. How could you symbolize an *xor* using only two connectives? Hint: think about how you might use DeMorgan's Laws to eliminate one of the connectives in your symbolization.

2.5 Conditional

Consider these sentences:

- (22) If Jean is in Paris, then Jean is in France.
 (23) Jean is in France if Jean is in Paris.
 (24) Jean is in France only if Jean is in Paris.

Let's use the following symbolization key:

P : Jean is in Paris.
 F : Jean is in France

Sentence (22) is roughly of this form: 'if P , then F '. We will use the symbol ' \rightarrow ' to symbolize the English 'if... then...' structure. So we symbolize sentence 23 by ' $(P \rightarrow F)$ '. The connective is called THE CONDITIONAL. Here, ' P ' is called the ANTECEDENT of the conditional ' $(P \rightarrow F)$ ', and ' F ' is called the CONSEQUENT.

Sentence (23) looks different from (22) since the word 'if' occurs in the middle of the sentence rather than at the beginning. But clearly (23) it is equivalent to (22), so we can also symbolize it as ' $(P \rightarrow F)$ '. In general, the 'if'-clause of an English conditional always introduces the *antecedent*, whether it occurs first or second in the English sentence, and the rest of the sentence then functions as the *consequent*.

Sentence (24) is also a conditional. Since the word 'if' appears in the second half of the sentence, it might be tempting to symbolize this in the same way as sentence (23), as ' $(P \rightarrow F)$ '. But that would be a mistake. My knowledge of geography tells me that sentence (23) is unproblematically true: there is no way for Jean to be in Paris that doesn't involve Jean being in France. But sentence (24) is not so straightforward: were Jean in Marseilles, Lyon, or Toulouse, Jean would be in France without being in Paris, thereby rendering sentence (24) false. Since geography alone dictates the truth of sentence (23), whereas travel plans (say) are needed to know the truth of sentence (24), they must mean different things, and (24) can't be symbolized as ' $(P \rightarrow F)$ '.

The moral is that 'only if' means something very different from plain 'if'. The 'if'-clause of a conditional introduces a SUFFICIENT CONDITION: (22) and (23) say that Jean's being in Paris is *sufficient for* his being in France (which is unproblematically true). 'Only if', by contrast, introduces a NECESSARY CONDITION: sentence (24) claims that Jean's being in Paris is *necessary for* his being in France (which is likely false, since there are other ways to be in France). In TFL, the antecedent ϕ of a conditional ($\phi \rightarrow \psi$) always indicates the sufficient condition, whereas the consequent ψ indicates the necessary condition. Since Jean's being in Paris is claimed as a necessary condition in (24), this sentence is symbolized as ' $(F \rightarrow P)$ '. In general, whereas 'if' introduces the antecedent of a conditional

(the sufficient condition), ‘only if’ introduces the consequent (the necessary condition). So our symbolization guidelines for conditionals are:

A sentence can be symbolized as $(\phi \rightarrow \psi)$ if it can be paraphrased in English as ‘If ϕ , then ψ ’, or as ‘ ψ if ϕ ’, or as ‘ ϕ only if ψ ’.

The fact that (22) is symbolized as ‘ $(P \rightarrow F)$ ’ means that this ‘if...then’ statement can also be paraphrased as an ‘only if’ statement: ‘Jean is in Paris only if Jean is in France’. That’s intuitively correct: since his being in Paris is sufficient for his being in France, it’s also true that his being in France is necessary for his being in Paris.

This connection between conditionals and necessary and sufficient conditions also means that our connective ‘ \rightarrow ’ can represent other English constructions that don’t involve the word ‘if’ at all. The following are all ways of saying that the truth of ϕ is *sufficient* for the truth of ψ :

- ▷ If ϕ then ψ
- ▷ ψ if ϕ
- ▷ ψ provided that ϕ
- ▷ ψ whenever ϕ
- ▷ ψ is the case as long as ϕ is the case

So sentences of this form would all be symbolized $(\phi \rightarrow \psi)$. On the other hand, the following are ways of saying that the truth of ϕ is *necessary* for the truth of ψ :

- ▷ ψ only if ϕ
- ▷ ψ is contingent on ϕ
- ▷ For ψ to be the case it is *necessary* that ϕ be the case

So sentences of this form would all be symbolized as $(\psi \rightarrow \phi)$, with the arrow running in the opposite direction compared to the first set of examples.

One final point: it is important to bear in mind that the connective ‘ \rightarrow ’ just says that the truth of the antecedent is sufficient for the truth of the consequent (or that the truth of the consequent is necessary for the truth of the antecedent). It says nothing about a *causal* or *explanatory* connection between two events, though English conditionals often carry such a suggestion. So something of the form $(\phi \rightarrow \psi)$ doesn’t mean that ϕ explains ψ , or that the truth of ϕ caused, or brought about, the truth of ψ . It only represents a *logical* relationship between the two. We will look more closely at the discrepancies between English ‘if...then’ and the TFL connective ‘ \rightarrow ’ in §3.3.

■ Exercises 2.5

A. Symbolize each English sentence in TFL.

1. Liz will go to the party if Megan and Ben both go.
2. If Megan doesn’t go to the party, Liz won’t go either.
3. If Megan goes to the party, then Ben will go if Liz does too.
4. Ben will go only if Liz does.

5. Russia will attend the summit only if Japan does not.
6. Sam will get a raise as long as he keeps working hard.
7. Sam's getting a raise is contingent on his not getting fired first.
8. Laura will take biology provided Bea does, but Bea will take it only if Alexis doesn't.
9. If either Alice or Bob is a spy, then the code has been broken.
10. If neither Alice nor Bob is a spy, then the code remains unbroken.

2.6 Biconditional

Consider the following sentence:

- (25) Bucephalus is a horse if and only if he is a mammal

Let's use the following symbolization key:

H : Bucephalus is a horse

M : Bucephalus is a mammal

Sentence 25 can be paraphrased as 'Bucephalus is a horse *if* he is a mammal, and Bucephalus is a horse *only if* Bucephalus is a mammal'. This is just the conjunction of two conditionals we already know how to symbolize. So we can symbolize it as ' $((H \rightarrow M) \wedge (M \rightarrow H))$ '. We call this a BICONDITIONAL, because it amounts to stating both directions of the conditional. That is, (25) says (falsely, as it happens) that Bucephalus' being a mammal is both necessary and sufficient for his being a horse.

We could treat every biconditional this way, as a conjunction of two conditionals. So, just as we do not need a new TFL symbol to deal with *exclusive or*, we do not really need a new TFL symbol to deal with biconditionals. However, since biconditionals occur a lot in logic and philosophy, we'll use the dedicated connective ' \leftrightarrow ' to symbolize them. We can then can symbolize sentence 25 with the TFL sentence ' $(H \leftrightarrow M)$ '.

Since 'if and only if' gets used so much in logic and philosophy, it is often abbreviated with a single, snappy word, 'iff'. So 'if' with only *one* 'f' is the English conditional. But 'iff' with *two* 'f's is the English biconditional. Armed with this we can say:

A sentence can be symbolized as $(\phi \leftrightarrow \psi)$ if it can be paraphrased in English as ' ϕ if and only if ψ ', that is, ' ϕ iff ψ '.

Another way to express a biconditional relationship in English is with the words 'just in case', as in:

- (26) The triangle is equilateral just in case all of its sides have the same length.

This says that the triangle's having sides of the same length is both necessary and sufficient for its being equilateral. If we use ' E ' for 'the triangle is equilateral' and ' S ' for 'the triangle's sides have the same length', it can be symbolized as the biconditional ' $(E \leftrightarrow S)$ '.

2.7 ‘Unless’

We have now seen all of the connectives of TFL. We can use them together to symbolize many kinds of sentences. But some cases are harder than others. And a typically nasty case is the English-language connective ‘unless’:

- (27) Unless you wear a jacket, you will catch a cold.
 (28) You will catch a cold unless you wear a jacket.

These two sentences are clearly equivalent. Let’s use the following symbolization key:

- J : You will wear a jacket.
 D : You will catch a cold.

Both sentences mean that if you do not wear a jacket, then you will catch a cold. So we could symbolize them as ‘ $(\neg J \rightarrow D)$ ’. In general, you can think of ‘unless’ as having the meaning ‘if not’. So sentences of the form ‘Unless ϕ , ψ ’ or ‘ ψ unless ϕ ’ mean ‘If not ϕ , then ψ ’ or again, ‘ ψ if not ϕ ’, and can be symbolized as $(\neg\phi \rightarrow \psi)$.

As we shall see in the next chapter, in TFL ‘ $(\neg J \rightarrow D)$ ’ is equivalent to the disjunction ‘ $(J \vee D)$ ’. So this is another way to symbolize (27) or (28). And that makes some sense: intuitively they do say that either you will wear a jacket, or you will catch a cold. So we can use the following guideline:

If a sentence can be paraphrased as ‘Unless ϕ , ψ ’ or ‘ ψ unless ϕ ’, then it can be symbolized as $(\neg\phi \rightarrow \psi)$, or simply $(\phi \vee \psi)$.

One caveat: although ‘Unless’ can be symbolized as a conditional or an inclusive disjunction, ordinary speakers of English often use ‘unless’ to mean something more like the biconditional, or like exclusive disjunction. Suppose I say: ‘I will go running unless it rains’. I probably mean that I’ll go running if it doesn’t rain, and *also* that I’ll go running *only if* it doesn’t rain, i.e. what we could put by saying ‘I will go running if and only if it does not rain’ (a biconditional), or ‘either I will go running or it will rain, but not both’ (an exclusive disjunction). However, in this book we’ll always use ‘unless’ in the strict sense in which it can be symbolized using a conditional or an inclusive disjunction. (As it happens, our guideline here is also the one used on the LSAT.)

■ Exercises 2.7

A. Symbolize each English sentence in TFL.

1. Unless something terrible happens, the team will win the playoffs.
2. Lua will win the race if and only if both Emily and Bill sit out.
3. Annie will mow the grass just in case her sister does the dishes, provided that there are dishes to be done.
4. Neither Li nor Simon will go the party unless Grace does.
5. Unless those creatures are men in costumes, they are either chimpanzees or gorillas.

2.8 Symbolizing Whole Arguments

So far we've been concerned with symbolizing individual statements. But logic is ultimately about the analysis of *arguments*, so we need to be able to symbolize whole arguments as well. Luckily, as we learned in §1.1, an argument is just a collection of statements, so symbolizing an argument is just a matter of symbolizing each of the statements (the premises and the conclusion) that comprise the argument.

Take again the simple example from §2.1:

(1) It's raining outside.	A
If it's raining outside, Jenny is miserable.	$(A \rightarrow C)$
\therefore Jenny is miserable.	$\therefore C$

The argument consists of two premises and a conclusion, and involves two atomic sentences. Using 'A' for 'It's raining outside' and 'C' for 'Jenny is miserable', the argument as a whole gets symbolized as you see on the right. By symbolizing an argument in TFL like this, we reveal its *logical form*, specifically its *truth-functional* logical form.

Although symbolizing an argument just involves symbolizing the individual statements it consists of, there are some complications to be aware of. Take the following English argument:

The murder either occurred in the attic or in the basement. Furthermore, if Prof. Plum was awake, the murder can't have occurred in the basement or the dining room. So if the murder didn't occur in the attic, Prof. Plum must have been asleep.

There is no \therefore symbol in this argument to indicate the conclusion, so we have to figure out what the conclusion is from context. In this case it's pretty clear: the word 'so' indicates that the person presenting the argument regards the last sentence as following from the others, so it is the conclusion. But this needn't always be the case: sometimes people present an argument by stating the conclusion first, and afterwards telling you what the premises are that demonstrate that conclusion.

Another thing we have to be careful about is correctly identifying the reappearance of the same atomic sentences in different parts of the argument. Take the following symbolization key:

- A: The murder occurred in the attic.
- B: The murder occurred in the basement.
- D: The murder occurred in the dining room.
- W: Prof. Plum was awake.
- S: Prof. Plum was asleep.

Using this key, our argument would be symbolized as: $(A \vee B)$, $(W \rightarrow \neg(B \vee D)) \therefore (\neg A \rightarrow S)$. But as we'll learn how to show in Chapter 3, this TFL argument is not valid: the conclusion $(\neg A \rightarrow S)$ does not follow from these premises (notice that none of the premises contains the letter 'S'). However, we shouldn't conclude that the English argument must therefore be invalid too — perhaps we just didn't symbolize it correctly.

We here used separate atomic TFL sentences to represent the statement ‘Prof. Plumb was awake’ in the second premise and the statement ‘Prof. Plumb was asleep’ in the conclusion. But we could instead symbolize the second of these statements as the negation of the first, that is, use $\neg W$ in place of S . This would now give us the following symbolization: $(A \vee B), (W \rightarrow \neg(B \vee D)) \therefore (\neg A \rightarrow \neg W)$. And as we’ll soon be able to show, this TFL argument now *is* valid. So this is a better way to represent the logical structure of the original English argument, since the person giving it presumably meant to draw a conclusion that follows from the premises. In general, when interpreting an argument — in this class or elsewhere — it is always a good idea to abide by the PRINCIPLE OF CHARITY:

If an argument has more than one *plausible* interpretation or symbolization, but only one of them yields a valid argument, then we should go with the symbolization on which the argument is valid.

■ Exercises 2.8

A. Symbolize each argument in TFL.

1. If Dorothy plays the piano in the morning, then Roger wakes up cranky. Dorothy plays piano in the morning unless she is distracted. So if Roger does not wake up cranky, then Dorothy must be distracted.
2. It will either rain or snow on Tuesday. If it rains, Neville will be sad. If it snows, Neville will be cold. Therefore, Neville will either be sad or cold on Tuesday.
3. If Zoog remembered to do his chores, then things are clean but not neat. If he forgot, then things are neat but not clean. Therefore, things are either neat or clean; but not both.

2.9 The Syntax of TFL

In the course of learning to symbolize English sentences in TFL, we’ve gotten a pretty good intuitive sense of how to build up complex TFL sentences from atomic ones using our five connectives. But as we move ahead, it will be necessary for us to be a bit more precise about the structure of our logical language. Because just as not every string of English words counts as a grammatical English sentence (e.g. ‘shelf on book red lies’ is just gibberish), so not every string of symbols counts as grammatical, or well-formed sentence of TFL.

We have seen that there are three kinds of symbols in TFL:

Atomic sentences: $A, B, C, \dots, Z, \dots, A_1, B_1, Z_1, A_2, A_{25}, J_{375}, \dots$

Connectives: $\neg, \wedge, \vee, \rightarrow, \leftrightarrow$

Brackets: $(,)$

This constitutes the LEXICON of TFL. Now define an EXPRESSION OF TFL to be any string of symbols of TFL. That is: write down any sequence of symbols from the lexicon of TFL, in any order, and you have an expression of TFL.

Strings like ' $(A \leftarrow B)$ ' or ' $(p \vee C)$ ' or ' $\neg(\phi \wedge A)$ ' are not expressions of TFL because they contain symbols like ' \leftarrow ' (leftward arrows) ' p ' (lowercase letters) and ' ϕ ' (Greek letters) that are not even in the Lexicon of TFL. On the other hand, ' $(A \wedge B)$ ' is an expression of TFL, and so are ' $(\neg)A \rightarrow$ ' and ' $\neg(\vee()) \wedge (\neg\neg())((B)$ '. However, whereas the first of these expressions also counts as a *sentence* of TFL, the rest are just *gibberish*. What we want are some rules to tell us precisely which TFL expressions count as sentences.

Obviously, individual atomic sentences like ' A ' and ' G_{13} ' should count as sentences. We can form further sentences out of these by using the various connectives. Using negation, we can get ' $\neg A$ ' and ' $\neg G_{13}$ '. Using conjunction, we can get ' $(A \wedge G_{13})$ ', ' $(G_{13} \wedge A)$ ', ' $(A \wedge A)$ ', and ' $(G_{13} \wedge G_{13})$ '. We could also apply negation repeatedly to get sentences like ' $\neg\neg A$ ' and ' $\neg\neg\neg A$ ', or apply negation to one of our conjunctions to get sentences like ' $\neg(A \wedge G_{13})$ ' and ' $\neg(G_{13} \wedge \neg G_{13})$ '. There are infinitely many possible combinations, even starting with just these two sentence letters. And of course there are infinitely many sentence letters. So there is no point in trying to list all of the sentences of TFL one by one.

Instead, we will describe the process by which sentences can be *constructed*. Consider negation: given any sentence ϕ of TFL, putting a negation in front of it gives us a sentence $\neg\phi$. We can say similar things for each of the other connectives. For instance, if ϕ and ψ are sentences of TFL, then $(\phi \wedge \psi)$ is a sentence of TFL. (What's up with the funny Greek letters, which are *not* in the lexicon of TFL? We'll get to that in §2.9 below.)

Providing clauses like this for all of the connectives, we arrive at the following SYNTAX, or formal definition of what counts as a SENTENCE OF TFL:

THE SYNTAX OF TFL:

1. Every atomic sentence is a sentence.
2. If ϕ is a sentence, then $\neg\phi$ is a sentence.
3. If ϕ and ψ are sentences, then $(\phi \wedge \psi)$ is a sentence.
4. If ϕ and ψ are sentences, then $(\phi \vee \psi)$ is a sentence.
5. If ϕ and ψ are sentences, then $(\phi \rightarrow \psi)$ is a sentence.
6. If ϕ and ψ are sentences, then $(\phi \leftrightarrow \psi)$ is a sentence.
7. Nothing else is a sentence.

Definitions like this are called *recursive*. Recursive definitions begin with some list of base elements (in this case, atomic sentences), and then present ways to generate indefinitely many more elements by compounding together previously generated elements. We can then determine if any given TFL expression counts as a sentence by checking whether it can be generated by applying these recursive rules of syntax.

For example, suppose we want to know whether or not ' $\neg\neg D$ ' is a sentence of TFL. Looking at clause 2 of the definition, we know that ' $\neg\neg D$ ' is a sentence *if* ' $\neg D$ ' is a sentence. So now we need to ask whether or not ' $\neg D$ ' is a sentence. Again looking at clause 2 of the definition, ' $\neg D$ ' is a sentence *if* ' D ' is. And ' D ' is an atomic sentence of TFL, so we know that ' D ' is a sentence by the clause 1 of the definition. So by applying the clauses of our definition repeatedly, we see that our original sentence ' $\neg\neg D$ ' can be generated by applying the rules of our syntax to atomic sentences, and thus counts as a TFL sentence.

Next, consider a more complex example: ' $\neg(P \wedge \neg(\neg Q \vee R))$ '. Looking at clause 2 of the definition, this is a sentence if ' $(P \wedge \neg(\neg Q \vee R))$ ' is. And by clause 3 the latter is a sentence

if both ' P ' and ' $\neg(\neg Q \vee R)$ ' are sentences. The former is an atomic sentence, and the latter is a sentence if ' $(\neg Q \vee R)$ ' is a sentence. Looking at clause 4, we see this is a sentence if both ' $\neg Q$ ' and ' R ' are sentences. And both are! So we've shown that the expression we started with is indeed a sentence.

Notice that negation differs from our other operators. It attaches to a *single* sentence to form a new sentence, and is therefore called a **UNARY OPERATOR**. By contrast, the other operators all operate on *pairs* of sentences to form new sentences, and are therefore called **BINARY OPERATORS**.

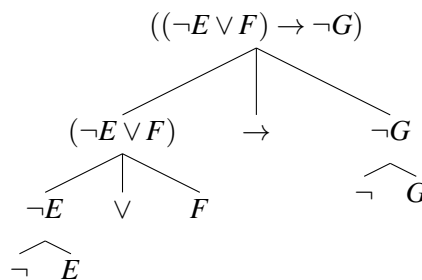
Our syntactic rules tell us that any sentence formed by applying a binary operator to a pair of sentences must be enclosed by parentheses. For example, when putting ' S ' and ' R ' together using a conjunction, the resulting sentence ' $(S \wedge R)$ ' must have brackets around it. So ' $S \wedge R$ ' is not technically a sentence of TFL, but a mere expression. This is *not* the case for negation, however! Putting a negation in front of a sentence never requires adding parentheses. So ' $\neg\neg D$ ' and ' $\neg(S \wedge R)$ ' are well-formed sentences, but ' $(\neg(\neg(D)))$ ' or ' $(\neg(S \wedge R))$ ' are not sentences. (We'll return to the rationale behind these bracketing rules in §2.9 below.)

Some Syntactic Notions

Ultimately, every TFL sentence is constructed nicely out of atomic sentences on the basis of our syntactic rules. When we are dealing with any complex (i.e. non-atomic) sentence, we can see that there must be some connective that was introduced *most recently* when constructing that sentence. We call that the **MAIN OPERATOR** of the sentence:

The **MAIN OPERATOR** of a complex TFL sentence is the one that was introduced *most recently* in the process of constructing that sentence.

In the case of ' $\neg\neg D$ ', the main operator is the very first ' \neg ' sign. In the case of ' $((\neg E \vee F) \rightarrow \neg G)$ ', the main operator is ' \rightarrow ' because the last step in constructing this sentence is to connect ' $(\neg E \vee F)$ ' and ' $\neg G$ ' using ' \rightarrow ' (and putting brackets around the result). One can visually represent the process in which a sentence is constructed from its parts via a **SYNTACTIC TREE** for the sentence. For example, the syntactic tree for ' $((\neg E \vee F) \rightarrow \neg G)$ ' looks like this:



This shows that ' $((\neg E \vee F) \rightarrow \neg G)$ ' was constructed by connecting ' $(\neg E \vee F)$ ' and ' $\neg G$ ' using ' \rightarrow ', and ' $(\neg E \vee F)$ ' was in turn constructed by connecting ' $\neg E$ ' and ' F ' using ' \vee ', and ' $\neg E$ ' was constructed by putting ' \neg ' in front of ' E ', and so on. If we represent the construction process in terms of a tree structure like this, then the main operator of a sentence is whichever operator occurs on a branch of its own at the *first level* below the sentence as a whole (which again, in this case, is ' \rightarrow ').

The syntactic structure of sentences in TFL also allows us to give a formal definition of the *scope* of a negation (mentioned in §2.3). The scope of a ‘ \neg ’ in a given sentence is whatever subsentence of that sentence has ‘ \neg ’ as its main logical operator. For example, consider the complex TFL sentence:

$$(\neg(R \wedge B) \leftrightarrow Q)$$

This was constructed by putting a biconditional between ‘ $\neg(R \wedge B)$ ’ and ‘ Q ’. So ‘ $\neg(R \wedge B)$ ’ is a *subsentence* of the sentence as a whole. And the main logical operator for this subsentence is ‘ \neg ’. So the scope of the negation in ‘ $(\neg(R \wedge B) \leftrightarrow Q)$ ’ is just ‘ $\neg(R \wedge B)$ ’. More generally:

The SCOPE of any connective in a sentence is the subsentence for which that connective is the main logical operator.

So again in the case of ‘ $(\neg(R \wedge B) \leftrightarrow Q)$ ’, the scope of ‘ \leftrightarrow ’ is the sentence as a whole (since it is the main operator of the whole sentence), and the scope of ‘ \wedge ’ is ‘ $(R \wedge B)$ ’, since ‘ $(R \wedge B)$ ’ is the subsentence of which ‘ \wedge ’ is the main operator.

Bracketing

As mentioned in §2.3 and §2.9 above, brackets are an important part of the syntax of TFL. This is because they demarcate the scope of connectives. For example, there is an important difference between ‘ $\neg(P \wedge Q)$ ’ and ‘ $(\neg P \wedge Q)$ ’. In the case of ‘ $\neg(P \wedge Q)$ ’ the scope of the negation operator is the whole sentence, that is, it is the main operator of the sentence and it serves to negate the entire conjunction it is attached to. In the case of ‘ $(\neg P \wedge Q)$ ’, the scope of the negation is just the subsentence ‘ $\neg P$ ’, and the main operator of the sentence as a whole is ‘ \wedge ’.

Strictly speaking, therefore, a string like ‘ $\neg P \wedge Q$ ’ is *not* a sentence of TFL, but a mere expression, because it is missing brackets. As things stand, it is not clear where in ‘ $\neg P \wedge Q$ ’ the brackets are supposed to go, that is, whether it is supposed to be a negated conjunction, i.e. ‘ $\neg(P \wedge Q)$ ’, or rather a conjunction with a negated left conjunct, i.e. ‘ $(\neg P \wedge Q)$ ’. When working with TFL, however, it will make our lives easier if we are sometimes a little less strict. So, here are some convenient conventions.

First, we’ll allow ourselves to omit the *outermost* brackets on a sentence. Thus we allow ourselves to write ‘ $Q \wedge R$ ’ instead of ‘ $(Q \wedge R)$ ’. However, we have to put the brackets back in when we want to embed this sentence into another, larger sentence! So we cannot write ‘ $P \rightarrow Q \wedge R$ ’, but must write ‘ $P \rightarrow (Q \wedge R)$ ’ instead. With this convention in place, something like ‘ $\neg P \wedge Q$ ’ can now be interpreted as missing its outermost parentheses, and thus being a shorthand for ‘ $(\neg P \wedge Q)$ ’ rather than ‘ $\neg(P \wedge Q)$ ’.

Second, it can be a bit painful to stare at long sentences with many nested pairs of brackets. To make things a bit easier on the eyes, we will allow ourselves to use square brackets, ‘[’ and ‘]’, instead of rounded ones. So there is no logical difference between ‘ $(P \vee Q)$ ’ and ‘ $[P \vee Q]$ ’, for example. Combining this convention with the first one, we can rewrite the unwieldy sentence:

$$(((H \rightarrow I) \vee (I \rightarrow H)) \wedge (J \vee K))$$

more simply as follows:

$$[(H \rightarrow I) \vee (I \rightarrow H)] \wedge (J \vee K)$$

The scope of each connective is now much more visually apparent.

Metalanguage and Metavariables

Our recursive definition of TFL sentences included clauses like the following:

3. If ϕ and ψ are sentences, then $(\phi \wedge \psi)$ is a sentence.

But notice that ' $(\phi \wedge \psi)$ ' is *not* a sentence of TFL. In fact, it isn't even an expression of TFL! After all, it includes Greek letters, and these are not among the symbols that constitute the the lexicon of TFL. Atomic TFL sentences are just ordinary uppercase roman letters (possibly subscripted) like ' A ' or ' B ' or ' S_{21} ', not Greek letters. So what's going on in our recursive clauses?

The answer is that in these clauses, we are using Greek letters as *variables* that "range over" arbitrary expressions of TFL. Consider how, in a math class, you might explain to someone that if m and n are any two positive integers, it holds that $m + n \geq m$. In this case we are using ' m ' and ' n ' as variables that range over arbitrary positive integers, and saying that $m + n \geq m$ holds no matter what positive integers m and n are.

In the same way, we are here using Greek letters as variables, except that we are using them as variables that range over arbitrary expressions in the language of TFL (rather than over integers, say). The language of TFL has been the object of our study for the past several sections, and we therefore call it the OBJECT LANGUAGE. But we have been conducting our discussion of TFL in English, so English is our METALANGUAGE: it is the language in which we talk about the object language.⁴ For this reason, variables like ' ϕ ' and ' ψ ' are called METAVARIABLES: they form part of our metalanguage, English, and they range over arbitrary expressions in our object language, TFL. But again, these metavariables are *only* part of our metalanguage, and not themselves included in the object language TFL.

Greek letters ϕ , ψ , χ etc. are METAVARIABLES in our METALANGUAGE, used to talk about arbitrary expressions of TFL. Roman letters $A, B, C \dots$ etc. are atomic sentences of our OBJECT LANGUAGE TFL .

So what clause 3 in our definition says is that if ϕ and ψ are any arbitrary sentences of TFL, then the result of writing whatever ϕ is, followed by ' \wedge ', followed by whatever ψ is, and enclosing the result in parentheses, produces a sentence of TFL. For example, if ϕ is the TFL sentence ' $\neg D$ ' and ψ is the TFL sentence ' $(S \rightarrow T)$ ', then clause 3 tells us that ' $(\neg D \wedge (S \rightarrow T))$ ' is a sentence of TFL. But again, the string of symbols ' $(\phi \wedge \psi)$ ' is not itself a sentence (nor even an expression) of TFL.

■ Exercises 2.9

A. Which of the following are sentences of TFL? If it's not, how could you rewrite it to make it a grammatical sentence?

1. $(A \vee \psi)$

⁴Of course these notions aren't fixed. If we had been discussing the syntax of Korean, say, rather than TFL, then our metalanguage would still have been English, but our object language would have been Korean.

2. $R \wedge \neg S_2$
3. $\neg(P \rightarrow Q)$
4. $((\neg P) \wedge Q)$
5. $(\neg P \rightarrow Q \vee R)$
6. $(R \vee (P \leftarrow Q))$
7. $S \rightarrow R \rightarrow T$
8. $(A \vee B \wedge C \vee D)$
9. $\neg\neg\neg\neg F$
10. $\neg\neg\neg(\neg F)$
11. $\neg \wedge S$
12. $((A \rightarrow (A \wedge \neg F)) \vee (D \leftrightarrow E))$
13. $(A \wedge (B \wedge ((C \wedge D) \wedge E)))$

B. What is the main operator in each of the following?

1. $\neg(P \rightarrow Q)$
2. $(\neg P \rightarrow Q)$
3. $((A \wedge B) \vee \neg C)$
4. $(A \wedge (B \vee \neg C))$
5. $(\neg(P \wedge R) \rightarrow (S \rightarrow (Q \vee V)))$
6. $\neg((P \wedge R) \rightarrow (S \rightarrow (Q \vee V)))$
7. $\neg(A \rightarrow \neg(\neg C \leftrightarrow B))$
8. $\neg A \rightarrow \neg(\neg C \leftrightarrow B)$
9. $(A \wedge (B \vee ((C \vee D) \wedge E)))$
10. $[(H \rightarrow I) \vee (I \rightarrow H)] \wedge (J \vee K)$

C. Are there any TFL sentences TFL that contain no atomic sentences? Does any TFL sentence contain more binary connectives than atomic sentences? Explain your answers.

D. Here are some more complex English sentences. Symbolize each and say what the main operator is:

1. If Lisa got paid, she will go to the mall only if she has enough money for a shirt or a phone case or a pair of shoes.
2. Weiting will win the race if either Lily or Sam drop out; otherwise, she will lose.
3. Neither Sweden nor Ireland will attend the summit if Russia and China don't both attend.
4. Either Sweden or Ireland will not attend the summit if Russia and China both don't attend.
5. Sarah isn't going to the party unless Richard and Pam are both going, and Tim is going iff neither Pam nor Quincy are going.
6. If Canada subsidizes exports, then the US will raise tariffs if Mexico opens new factories.
7. Hanyu will go hiking as long as Liam comes too, unless the weather turns bad — in that case she'll go on a bike ride.
8. If evolutionary biology is correct, higher life forms arose by chance, and if that's so, then there isn't any design or divine intervention in nature.

The Semantics of TFL

3

We ended the last chapter by looking at the SYNTAX, or grammar, for the language of TFL. In this chapter, we'll be concerned with the SEMANTICS, or meaning, of TFL sentences. More specifically, we're going to look at the meanings of our five TFL connectives, and see how the meaning of a complex TFL sentence is determined by the connectives it contains. Once we've done that, we can use our semantics to give a precise definition of various important logical notions, like validity.

3.1 Meanings for TFL Connectives

Negation Let's begin with negation. The meaning of the TFL connective ' \neg ' should roughly resemble that of the English word 'not'. But what does 'not' mean? This might seem like a baffling question. What are meanings, anyhow?

To make the issue more tractable, let's ask a simpler question: if you put 'not' into a sentence, what does that do to the *truth-value* of the sentence? Take a true sentence, like 'Frida Kahlo was a painter'. If you add a 'not' into it, you get the false sentence 'Frida Kahlo was not a painter'. And similarly, if you take a false sentence and negate it, you get a true sentence.

So we can characterize the meaning of the TFL connective ' \neg ' as a mapping between truth values: given something true, it returns something false, and given something false, it returns something true. We'll abbreviate 'True' with 'T' and 'False' with 'F'. We can then represent the meaning of the TFL connective ' \neg ' via the following *characteristic truth table* for negation:

ϕ	$\neg\phi$
T	F
F	T

What this says is that for any TFL sentence ϕ : if ϕ is true, then $\neg\phi$ is false, and if ϕ is false, then $\neg\phi$ is true.

Conjunction A similar line of thought goes for conjunction. If you take two true sentences, and put an 'and' between them, the conjunction you've formed is true. On the other hand, if even one of the two conjoined sentences is false, the entire conjunction is false. So the characteristic truth table for conjunction looks like this:

ϕ	ψ	$(\phi \wedge \psi)$
T	T	T
T	F	F
F	T	F
F	F	F

Note that conjunction is *symmetrical*. The truth value for $(\phi \wedge \psi)$ is always the same as the truth value for $(\psi \wedge \phi)$.

Disjunction Recall that ‘ \vee ’ represents inclusive or. So, for any sentences ϕ and ψ , $(\phi \vee \psi)$ is true iff at least one of ϕ and ψ is true. This gives us the following characteristic truth table for disjunction:

ϕ	ψ	$(\phi \vee \psi)$
T	T	T
T	F	T
F	T	T
F	F	F

Like conjunction, disjunction is symmetrical: ‘ $(\phi \vee \psi)$ ’ always has the same truth value as ‘ $(\psi \vee \phi)$ ’.

As we saw in §2.4, the English construction ‘either ... or’ is sometimes used to express *exclusive* disjunction, which “excludes” the possibility of both disjuncts’ being true. If we liked, we could expand TFL by introducing a new connective \oplus (sometimes called XOR) with the following characteristic truth table, which differs from table for \vee in the first row, the case where ϕ and ψ are both true:

ϕ	ψ	$(\phi \oplus \psi)$
T	T	F
T	F	T
F	T	T
F	F	F

However, as we discussed, we don’t need to (and won’t) go this route, since the effect of the exclusive $(\phi \oplus \psi)$ can be achieved using the TFL connectives we already have, via $(\phi \vee \psi) \wedge \neg(\phi \wedge \psi)$, i.e. by saying that one of ϕ or ψ is true, but not both.

Conditional Conditionals are considerably more contentious. In fact, we might as well be up front about it: they are a mess. We’ll simply stipulate that, in TFL, $(\phi \rightarrow \psi)$ is false if ϕ is true and ψ is false, but is true in *all other circumstances*. This gives us the following characteristic truth table for the conditional:

ϕ	ψ	$(\phi \rightarrow \psi)$
T	T	T
T	F	F
F	T	T
F	F	T

Notice that the conditional is *asymmetrical*: $(\varphi \rightarrow \psi)$ and $(\psi \rightarrow \varphi)$ need not have the same truth value. For example, if φ is true and ψ false, then $(\varphi \rightarrow \psi)$ is false but $(\psi \rightarrow \varphi)$ is true.

You can perhaps already see that this is a controversial way to symbolize English ‘if ... then’ constructions. The above truth table tells us that any TFL conditional with a false antecedent is true (and similarly, any TFL conditional with a true consequent is true). But it’s far from clear that any English conditional whose antecedent turns out to be false is therefore automatically true. This is known as “the paradox of material implication.” In the case of conditionals, there in other words isn’t just a worry about TFL bypassing certain *subtleties* of meaning, but of missing out on the meaning of the corresponding English expression altogether. We’ll look at some of these issues in §3.3 below.

Biconditional Since a biconditional is the same as the conjunction of two conditionals running in either directions, $(\varphi \leftrightarrow \psi)$ is true iff both $(\varphi \rightarrow \psi)$ and $(\psi \rightarrow \varphi)$ are true. This gives us the following table:

φ	ψ	$(\varphi \leftrightarrow \psi)$
T	T	T
T	F	F
F	T	F
F	F	T

An easy way to remember this is that a biconditional is true when both sides have the *same* truth value, and false when the two sides have different truth values. The biconditional is therefore symmetrical. Its truth table is in effect the opposite of exclusive disjunction. As we’ll soon be able to show, this truth table is indeed the same as the one we would get for $(\varphi \rightarrow \psi) \wedge (\psi \rightarrow \varphi)$.

3.2 Truth-Functionality

The fact that we can give characteristic truth tables like these for our TFL connectives means that they are *truth-functional*:

A connective is TRUTH-FUNCTIONAL iff the truth value of a sentence with that connective as its main logical operator is uniquely determined by the truth value(s) of the constituent sentence(s).

Indeed, this is what gives TFL its name: *truth-functional logic*.

Many languages have connectives that are not truth-functional. In English, for example, we can form a new sentence from any simpler sentence by prefixing it with the unary connective ‘It is necessarily the case that...’. The truth value of this new sentence is not fixed solely by the truth value of the original sentence. For consider two true sentences:

1. $2 + 2 = 4$
2. Shostakovich wrote fifteen string quartets

Whereas it is necessarily the case that $2 + 2 = 4$, it is not *necessarily* the case that Shostakovich wrote fifteen string quartets. If he had died earlier or later, he might have

written fewer or more quartets than he in fact did. So the English unary connective ‘It is necessarily the case that...’ is not *truth-functional*. TFL cannot represent non-truth-functional connectives like these; it can only represent truth-functional connectives like e.g. ‘It is not the case that...’ or ‘...and...’.

Since TFL’s connectives are all truth-functional, we have to ignore everything except the truth-functional aspects of English when symbolizing English sentences or arguments into TFL. A lot is inevitably lost in the process. There are subtleties to our ordinary claims that far outstrip their mere truth values: sarcasm, poetry, snide implicature, emphasis. These are all important parts of everyday discourse, but none of it is retained in TFL.

For example, as already remarked in §2.3, TFL cannot capture the subtle differences between the following English sentences:

1. Adam is energetic and Adam is not athletic.
2. Although Adam is energetic, he is not athletic.
3. Despite being energetic, Adam is not athletic.
4. Adam is energetic, albeit not athletic.

They all get symbolized with the same TFL sentence, perhaps ‘ $(E \wedge \neg A)$ ’. Similarly, in symbolizing ‘Adam is energetic’ as ‘ E ’, we are ignoring all aspects of its meaning except its truth value.

This is why we talk of *symbolizing* English sentences. Some logic textbooks talk about *translating* English sentences into TFL. But a good translation should preserve more than mere truth values and truth-functional aspects of meaning. So we can’t really *translate* English into TFL, properly speaking.

3.3 Conditionals in TFL and English

When we introduced the truth table for ‘ \rightarrow ’, we didn’t provide any justification for it. In fact, we noticed that it seems problematic as a symbolization of English ‘if...then’. But there are some things to be said in favor of the truth table we provided.

First, the TFL conditional has some attractive *logical* features given our truth table.¹ The following all seem correct for English ‘if...then’ statements:

- Arguments of the form ‘If φ then ψ ; φ ; therefore ψ ’ are valid. (This form of argument is called *modus ponens*.)
- Arguments of the form ‘If φ then ψ ; ψ ; therefore φ ’ are not valid. (This is called the fallacy of *affirming the consequent*.)
- Statements of the form ‘If φ , then φ ’ are necessarily true.

As we will see once we define validity and other logical notions in TFL, the same holds for the corresponding TFL symbolizations using \rightarrow in place of ‘if...then’: arguments of the form $(\varphi \rightarrow \psi), \varphi \therefore \psi$ will be valid in TFL, ones of the form $(\varphi \rightarrow \psi), \psi \therefore \varphi$ will not be valid (at least if φ and ψ are atomic sentences), and any sentence of the form $(\varphi \rightarrow \varphi)$ is a logical necessity (or “tautology”) in TFL. And importantly, out of the sixteen possible binary truth functions, the one we have assigned to ‘ \rightarrow ’ is the *only* one that has all these

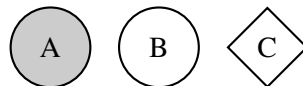
¹I owe this observation to Branden Fitelson.

logical properties! So if we have to pick a truth-functional connective to symbolize English ‘if ... then’, then ‘ \rightarrow ’ is the best choice among the sixteen available options.

Second, there’s an argument to suggest that the truth table we’ve given for ‘ \rightarrow ’ captures at least certain uses of English ‘if ... then’.² Suppose Lara has drawn several shapes on a piece of paper, and colored some of them grey. I have not seen them, but I claim:

If any shape is grey, then it is also circular.

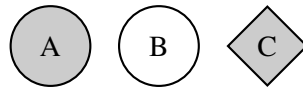
As it happens, Lara has drawn the following:



In this case, my general conditional claim is true. And this in turn means that each of its *instances* must be true:

- If A is grey, then it is circular (true antecedent, true consequent)
- If B is grey, then it is circular (false antecedent, true consequent)
- If C is grey, then it is circular (false antecedent, false consequent)

However, if Lara had drawn the following:



then my claim would have been false, because it would then have had a false instance:

- If C is grey, then it is a circular (true antecedent, false consequent)

Notice that this distribution of truth values exactly matches that in our truth-table for ‘ \rightarrow ’: the only false conditional is the one with a true antecedent but a false consequent. So this suggests that the truth values of at least some English ‘if ... then’ statements match those predicted by our truth table.

At the same time, it’s clear that there are other uses of ‘if ... then’ in English that aren’t adequately symbolized using ‘ \rightarrow ’. Consider the following two sentences:

- (1) If Hillary Clinton had won the 2016 US election, then she would have been the first female president of the US.
- (2) If Hillary Clinton had won the 2016 US election, then she would have turned into a helium balloon and floated away into the sky.

Intuitively, sentence 1 is true and sentence 2 is false. But both have false antecedents and false consequents. (Hillary did not win; she did not become the first female president of the US; and she of course did not turn into a helium balloon.) So our truth table would incorrectly count both sentence true.

²Versions of this argument are given by Dorothy Edgington (2014), ‘Conditionals’, in the *Stanford Encyclopedia of Philosophy* (<http://plato.stanford.edu/entries/conditionals/>) and Warren Goldfarb (2003), in his textbook *Deductive Logic*.

These are examples of *subjunctive conditionals*, because they are in the subjunctive mood (that is, they involve words like ‘had’ and ‘would’). They ask us to imagine something contrary to fact—a world in which Hillary won the 2016 election—and then ask us to evaluate what *would* have happened in that case. What we’ve seen is that subjunctive conditionals are not adequately symbolized using ‘ \rightarrow ’. In fact, we’ve seen that subjunctive conditionals aren’t even *truth functional*! After all, (1) and (2) have antecedents and consequents with the same truth values (all false), but the two conditionals themselves have different truth values. Since subjunctive conditionals aren’t truth-functional, there is no hope of symbolizing them in the truth-functional language of TFL.

Still, for the reasons given earlier, ‘ \rightarrow ’ is the best candidate we have for symbolizing at least certain uses of English ‘if ... then’. We’ll therefore continue to symbolize them this way, while remaining mindful of the simplification involved.

3.4 Complete Truth Tables

We’ve seen what the characteristic truth tables for the five TFL connectives are. Our next step is to use these truth tables to build up truth tables for complex TFL sentences that contain multiple connectives. To construct a truth table for a complex sentence like ‘ $(H \wedge I) \rightarrow H$ ’ we have to start with truth-values for the atomic sentences, and then calculate the truth value of the complex sentence.

So far, we’ve used symbolization keys to assign truth values to TFL sentences. For example, we might say that the TFL sentence ‘ B ’ is to symbolize ‘Big Ben is in London’. Since Big Ben *is* in London, this symbolization would make ‘ B ’ true. But we can also assign truth values *directly*. We could simply stipulate that ‘ B ’ is to be true, or stipulate that it is to be false. Such stipulations are called *valuations*:

A VALUATION is any assignment of truth values to particular atomic sentences of TFL.

To construct the COMPLETE TRUTH TABLE for a complex TFL sentence, we will have to calculate its truth value on every possible valuation of the atomic sentences it contains.

Let’s look at an example. Take the sentence ‘ $(H \wedge I) \rightarrow H$ ’. There are four possible ways to assign True and False to the atomic sentence ‘ H ’ and ‘ I ’—four possible valuations. We can represent these as follows:

H	I	$(H \wedge I) \rightarrow H$
T	T	
T	F	
F	T	
F	F	

To calculate the truth value of the entire sentence ‘ $(H \wedge I) \rightarrow H$ ’, we first copy the truth values for the atomic sentences and write them underneath the letters in the sentence:

H	I	$(H \wedge I) \rightarrow H$		
T	T	T	T	T
T	F	T	F	T
F	T	F	T	F
F	F	F	F	F

Next we have to consider the subsentence ‘ $(H \wedge I)$ ’. This is a conjunction, and the characteristic truth table for conjunction tells us that a conjunction is true iff both conjuncts are true. Since ‘ H ’ and ‘ I ’ are both true on (and only on) the first line of the truth table, the conjunction ‘ $(H \wedge I)$ ’ is true on the first row of the table and false on the rest:

H	I	$(H \wedge I) \rightarrow H$		
T	T	T	T	T
T	F	T	F	T
F	T	F	F	F
F	F	F	F	F

Notice how we’ve recorded the truth value for the subsentence ‘ $(H \wedge I)$ ’ on each row underneath its main operator, ‘ \wedge ’.

Now, our TFL sentence as a whole is a conditional, $\phi \rightarrow \psi$, with ‘ $(H \wedge I)$ ’ as ϕ and ‘ H ’ as ψ . So to determine the truth-value of the whole conditional, we have to look at the truth values of ‘ $(H \wedge I)$ ’ and ‘ H ’ on each row. On the second row, for example, ‘ $(H \wedge I)$ ’ is false and ‘ H ’ is true. Since a conditional is true when the antecedent is false, we write a ‘T’ in the second row underneath the conditional symbol. We continue for the other three rows and get this:

H	I	$(H \wedge I) \rightarrow H$		
T	T	T	T	T
T	F	F	T	T
F	T	F	T	F
F	F	F	T	F

The conditional is the main logical operator of this sentence. The column of ‘T’s underneath ‘ \rightarrow ’ therefore tells us the truth value for the sentence as a whole on each of the four possible valuations of its atomic constituents ‘ H ’ and ‘ I ’. What the table shows is that ‘ $(H \wedge I) \rightarrow H$ ’ is true regardless of the truth values of ‘ H ’ and ‘ I ’. They can be true or false in any combination, and the complex sentence still comes out true. Since we have considered all four possible valuations, this means that ‘ $(H \wedge I) \rightarrow H$ ’ is true on *every* valuation.

In this example, I erased some ‘T’s and ‘F’s as we went along to make things more readable. When actually writing truth tables on paper, however, it is impractical to erase whole columns or rewrite the whole table for every step. Although it is more crowded, the complete truth table with no columns erased looks like this:

H	I	$(H \wedge I) \rightarrow H$		
T	T	T	T	T
T	F	T	F	T
F	T	F	F	T
F	F	F	F	T

Most of the columns underneath the sentence are only there for bookkeeping purposes. The column that matters most is the column underneath the *main logical operator* for the sentence, since this tells you the truth value of the entire sentence. I have emphasized this column by putting it in bold. When you work through truth tables yourself, you should similarly emphasize the column under the main operator (perhaps by highlighting it or circling it).

As you can see from this example, a complete truth table has a row for every possible assignment of True and False to the relevant atomic sentences. Each of these rows represents a *valuation*. The number of rows depends on the number of different atomic sentences involved. A sentence that contains only one atomic sentence requires only two rows, as in the characteristic truth table for negation. This is true even if the same letter is repeated many times, as in the sentence $[(C \leftrightarrow C) \rightarrow C] \wedge \neg(C \rightarrow C)$. The complete truth table requires only two rows because there are only two possibilities: ‘C’ can be true or it can be false. The truth table for this sentence looks like this:

C	$((C \leftrightarrow C) \rightarrow C) \wedge \neg(C \rightarrow C)$									
T	T	T	T	T	T	F	F	T	T	T
F	F	T	F	F	F	F	F	F	T	F

Looking at the column underneath the main logical operator, we see that the sentence is false on both rows, i.e. on every valuation, no matter whether ‘C’ is true or false.

There will be four rows in the complete truth table for a sentence containing two atomic sentences, like $(H \wedge I) \rightarrow H$. And there will be eight rows in the complete truth table for a sentence containing three atomic sentences, e.g.:

M	N	P	$M \wedge (N \vee P)$			
T	T	T	T	T	T	T
T	T	F	T	T	T	F
T	F	T	T	T	F	T
T	F	F	T	F	F	F
F	T	T	F	F	T	T
F	T	F	F	F	T	F
F	F	T	F	F	F	T
F	F	F	F	F	F	F

So truth tables grow quickly! Four atomic sentences would require 16 rows, five atomic sentences require 32 rows, six atomic sentences require 64 rows, and so on.

A complete truth table for a sentence with n atomic sentences must have 2^n rows, representing the 2^n possible valuations.

In order to fill in the columns under the atomic sentences, begin with the right-most atomic sentence (‘P’ in the table above) and alternate between ‘T’ and ‘F’. In the next column to the left (the one under ‘N’ in our table), write two ‘T’s followed by two ‘F’s, and repeat. For the third atomic sentence (‘M’ in our table), write four ‘T’s followed by four ‘F’s. This yields an eight line truth table like the one above. For a 16 line truth table, the next column of atomic sentences should have eight ‘T’s followed by eight ‘F’s. For a 32 line table, the next column would have 16 ‘T’s followed by 16 ‘F’s. And so on. In general, you should construct your truth tables according to the following rules:

1. Write down the complex sentence you are working with, and to its left list the atomic sentences it contains *in alphabetical order*.
2. Determine how many rows your table will require given how many atomic sentences are involved. Again, for n atomic sentences you need 2^n rows.
3. Fill in the truth values for each atomic sentence according to the pattern described above. The column under the right-most atomic sentence will follow the pattern T F T F T F, the next column to the left will have the pattern T T F F T T F F, the column to the left of that the pattern T T T T F F F F, and so on.
4. Then calculate the truth value for the complex sentence as a whole on every row of the truth table (i.e. on every possible valuation). When you're done, remember to highlight or circle the column under the sentence's main logical operator.

These rules give us a canonical format for truth tables, which makes it easier to compare (and grade) truth tables written by different people.

■ Exercises 3.4

A. Construct complete truth tables in canonical format (i.e. by following the rules we gave) for each of the following:

1. $A \rightarrow A$
2. $C \rightarrow \neg C$
3. $A \rightarrow \neg(B \wedge \neg A)$
4. $(P \rightarrow Q) \vee (Q \rightarrow P)$
5. $(B \rightarrow (C \wedge A)) \vee (C \wedge \neg A)$
6. $(A \leftrightarrow B) \leftrightarrow \neg(A \leftrightarrow \neg B)$
7. $(A \rightarrow B) \vee (B \rightarrow A)$
8. $(A \wedge B) \rightarrow (B \vee A)$
9. $\neg(A \vee B) \leftrightarrow (\neg A \wedge \neg B)$
10. $[(A \wedge B) \wedge \neg(A \wedge B)] \wedge C$
11. $[(A \wedge B) \wedge C] \rightarrow B$

B. Show that ' $((A \vee B) \wedge \neg(A \wedge B))$ ' has a truth table that matches that for exclusive disjunction given in §3.1 above.

3.5 Semantic Concepts

Now that we know how to construct complete truth tables for complex sentences, we can define some important semantic concepts and see how to use truth tables to test whether they apply. In §1.2, we looked at the notions of *necessary truth* and *necessary falsity*. Both notions have surrogates in TFL. We'll start with a surrogate for necessary truth.

A sentence ϕ is a TF TAUTOLOGY iff it is true on every valuation.

That is, a TFL sentence is a TF tautology if it is true on every row of its complete truth table, since rows represent valuations. If you look back at the truth table for ' $(H \wedge I) \rightarrow H$ ' from §3.4 you'll see that it's a tautology. Other tautologies include sentences of the form $(\phi \rightarrow \phi)$ and ones of the form $(\phi \vee \neg\phi)$; the latter is called the LAW OF EXCLUDED MIDDLE.

Notice that this is only a *surrogate* for necessary truth. There are some necessary truths that we cannot adequately symbolize in TFL. For example, ' $2 + 2 = 4$ ' and 'Every city either is or is not in France' are both necessary truths, but if we symbolize them in TFL, the best we can offer is an atomic sentence, and no atomic sentence is a tautology. Still, if we can adequately symbolize some English sentence using a TFL sentence which is a tautology, then that English sentence expresses a necessary truth.

We have a similar surrogate for necessary falsity:

A sentence ϕ is a TF CONTRADICTION iff it is false on every valuation.

A sentence is a TF contradiction if it is false on every line of its complete truth table. The truth table for ' $[(C \leftrightarrow C) \rightarrow C] \wedge \neg(C \rightarrow C)$ ' we constructed in §3.4 shows that this sentence is a contradiction. Sentences of the form $(\phi \wedge \neg\phi)$ or $\neg(\phi \rightarrow \phi)$ are other examples of contradictions. Notice that the negation of any tautology is a contradiction. Lastly, we have a surrogate for contingency:

A sentence ϕ is TF CONTINGENT iff it is neither a tautology nor a contradiction, i.e. if it is true on at least one valuation and false on at least one.

These notions all apply to *single* sentences of TFL. Another useful notion—which we've occasionally already made use of—is that of *equivalence*. This is a property that applies to *pairs* of sentences of TFL:

ϕ and ψ are TF EQUIVALENT iff they have the same truth value on every valuation.

Notice that by this definition, any two tautologies, and any two contradictions, are equivalent. Equivalence is more interesting when we find pairs of contingent sentences that are equivalent. For example, we've observed that a biconditional like ' $(A \leftrightarrow B)$ ' is equivalent to a conjunction of two conditionals, ' $((A \rightarrow B) \wedge (B \rightarrow A))$ '. Similarly, since $(\phi \rightarrow \psi)$ is true if ϕ is true or if ψ is false (and false otherwise), TFL conditionals are equivalent to disjunctions of the form $(\neg\phi \vee \psi)$. DEMORGAN'S LAWS, which we looked at in 2.4, are another example of equivalences.

Again, we can test for TF equivalence using truth tables. Consider the sentences ' $\neg(P \vee Q)$ ' and ' $\neg P \wedge \neg Q$ '. To find out whether they're TF equivalent, we construct JOINT TRUTH TABLE that includes both sentences at once:

P	Q	$\neg(P \vee Q)$	$\neg P \wedge \neg Q$
T	T	F	T
T	F	F	F
F	T	F	F
F	F	T	T

Joint truth tables like these are constructed by listing to the left (in alphabetical order) all the atomic sentences that occur in *any* of the sentences being compared. We then include two separate sections in the table, one for each sentence, and calculate the truth value of each sentence in every row, i.e. on every valuation.

Looking at the columns under the main operators of the two sentences (negation for the first sentence, conjunction for the second) we see that both are false on the first three rows and true on the last row. Since they match on every row, they have the same truth value on every valuation, and are therefore TF equivalent. This pair of sentences is an instance of one of DeMorgan's Laws, so the table shows that the law does indeed hold in TFL.

Another notion that applies to pairs of sentences is the following:

ϕ and ψ are TF CONTRADICTORY iff they have opposite truth values on every valuation.

We can again test for this property by drawing a joint truth table for the two sentences to be compared, and checking to see that the truth values listed underneath their main connectives are different in every row of the table. It's important not to confuse the notion of two sentences being *contradictory* with the notion of a sentence's being a *contradiction*: the former is a property of pairs of sentences, the latter a property of a single sentence. But the two notions are connected (as the names suggest): if ϕ and ψ are contradictory, then their conjunction ($\phi \wedge \psi$) is a contradiction.

Next, here is a notion that applies to arbitrarily large *collections* of sentences:

ϕ_1, \dots, ϕ_n are JOINTLY TF CONSISTENT iff there is at least one valuation which makes them all true.

Collections of sentences are also said to be JOINTLY TF INCONSISTENT iff they are not TF consistent, i.e. iff there is no valuation on which they are all true. Again, it is easy to test for joint TF consistency (and inconsistency) using joint truth tables:

P	Q	$P \wedge Q$	$P \vee Q$	$P \rightarrow Q$
T	T	T	T	T
T	F	F	T	F
F	T	F	T	T
F	F	F	F	T

We can see from this that ' $P \wedge Q$ ', ' $P \vee Q$ ', and ' $P \rightarrow Q$ ' are consistent, because there is a row in their joint table on which they are all true, namely the first.

■ Exercises 3.5

A. Determine if each of the following is a tautology, a contradiction, or contingent:

1. $(S \rightarrow R) \wedge (S \wedge \neg R)$
2. $((E \rightarrow F) \rightarrow F) \rightarrow E$
3. $P \rightarrow (P \rightarrow P)$
4. $(P \rightarrow P) \rightarrow P$

$$5. \neg(A \vee B) \leftrightarrow (\neg A \wedge \neg B)$$

B. Determine if the following are equivalent, contradictory, consistent, or inconsistent:

1. ' $\neg(A \wedge B)$ ' and ' $(\neg A \vee \neg B)$ '.
2. ' $(F \wedge M)$ ' and ' $\neg(F \vee M)$ '
3. ' $(R \vee \neg S)$ ' and ' $(S \wedge R)$ '
4. ' $(A \vee B) \wedge \neg(A \wedge B)$ ' and ' $\neg(A \leftrightarrow B)$ '

C. Answer the following:

1. If ϕ is a tautology, must $(\phi \vee \psi)$ be a tautology?
2. If ϕ is a contradiction, must $(\phi \vee \psi)$ be a contradiction?
3. If ϕ is a tautology, must $(\phi \wedge \psi)$ be a tautology?
4. If ϕ is contingent, must $(\phi \wedge \psi)$ be contingent?
5. If ϕ is a tautology, must $(\psi \rightarrow \phi)$ be a tautology?
6. If ϕ and ψ are both contingent, must $(\phi \wedge \psi)$ be contingent?
7. If ϕ and ψ are equivalent, must they be consistent?
8. If ϕ and ψ are inconsistent, must they be contradictory?
9. If ϕ and ψ are contradictory, must they be inconsistent?
10. If ϕ and ψ are equivalent, what property does $(\phi \rightarrow \psi)$ have?
11. If ϕ and ψ are contradictory, what property does $(\phi \leftrightarrow \psi)$ have?

3.6 Validity in TFL

We've defined various important logical concepts using the notion of a valuation. But we have yet to discuss our most important logical concept, that of *validity*. Recall that in §1.1 we said that an argument is valid iff it's *impossible* for the premises to be true and the conclusion false. In TFL, we can spell out this concept of impossibility using our notion of a valuation, and define validity for TFL as follows:

An argument $\phi_1, \dots, \phi_n \therefore \psi$ is TF VALID iff *no valuation* makes all the premises ϕ_1, \dots, ϕ_n true but the conclusion ψ false.

Instead of saying that the argument $\phi_1 \dots \phi_n \therefore \psi$ is *valid*, we can also say that the premises $\phi_1 \dots \phi_n$ *entail* the conclusion ψ :

ϕ_1, \dots, ϕ_n TF ENTAIL ψ iff no valuation makes all of ϕ_1, \dots, ϕ_n true and ψ false.

We can again test for TF entailment, or validity, with joint truth tables. To test whether ' $\neg L \rightarrow (J \vee L)$ ' and ' $\neg L$ ' TF entail ' J ' we construct the following joint truth table:

J	L	$\neg L \rightarrow (J \vee L)$	$\neg L$	J
T	T	F	T	T
T	F	T	F	T
F	T	F	T	F
F	F	T	F	F

The only row on which both ' $\neg L \rightarrow (J \vee L)$ ' and ' $\neg L$ ' are true is the second row. But that is a row on which ' J ' is also true! So there is no valuation that makes both ' $\neg L \rightarrow (J \vee L)$ ' and ' $\neg L$ ' true and ' J ' false, meaning that ' $\neg L \rightarrow (J \vee L)$ ' and ' $\neg L$ ' entail ' J ', i.e. that the argument from premises

Because TF entailment is such an important concept, we'll introduce some new notation in connection with it. Rather than saying that the sentences $\varphi_1, \dots, \varphi_n$ TF entail ψ , we will abbreviate this by writing:

$$\varphi_1, \dots, \varphi_n \models \psi$$

The symbol ' \models ' is known as *the double-turnstile*, since it looks like a turnstile with two horizontal beams.

There is no limit on the number of sentences that can be mentioned before the symbol ' \models '. Indeed, we can even consider the limiting case where none are mentioned:

$$\models \psi$$

This says that there is no valuation which makes all the sentences mentioned on the left side of ' \models ' true while making ψ false. Since *no* sentences are mentioned on the left side of ' \models ' in this case, this just means that there is no valuation which makes ψ false. But that just means that ψ is a tautology! So writing $\models \psi$ gives us a short way to say that ψ is a tautology.

Sometimes we will want to deny that a TF entailment holds. We will write:

$$\varphi_1, \dots, \varphi_n \not\models \psi$$

to say that $\varphi_1, \dots, \varphi_n$ *do not* TF entail ψ , i.e. to say that there *does* exist a valuation that makes all of $\varphi_1, \dots, \varphi_n$ true but ψ false. Similarly $\not\models \psi$ means that ψ is *not* a tautology, i.e. that there exists a valuation that makes ψ false.

Lastly, we can abbreviate the claim that φ and ψ are equivalent as follows:

$$\varphi \models \psi$$

This encapsulates the idea that φ and ψ are equivalent iff each entails the other. That makes sense: if $\varphi \models \psi$ then there is no valuation that makes φ true but ψ false, and if $\psi \models \varphi$ then there is no valuation that makes ψ true but φ false. Putting it together, this means that there's no valuation on which φ and ψ have different truth values, meaning that φ and ψ are equivalent. Similarly reasoning holds the other direction, from equivalence to mutual entailment

It's important to be clear that ' \models ' is not a symbol in the language of TFL. Rather, it is a symbol of our *metalinguage* (recall the difference between object language and metalanguage from §2.9). The following is a claim in our metalanguage, not in the language TFL:

- $P, P \rightarrow Q \models Q$

This is shorthand for the following metalinguistic claim:

- There is no valuation that makes ' P ' and ' $P \rightarrow Q$ ' true but ' Q ' false

For this reason it is also important not to confuse the symbols ‘ \rightarrow ’ and ‘ \models ’. The conditional ‘ \rightarrow ’ is a symbol in our object language, TFL. The TFL sentence ‘ $P \rightarrow Q$ ’ just says that it’s not the case that ‘ P ’ is true and ‘ Q ’ is false. By contrast ‘ $P \models Q$ ’ is a sentence in the metalanguage. It doesn’t just claim that it is not the case that ‘ P ’ is true and ‘ Q ’ is false, but makes the stronger — and as it happens false — claim that *there exists no valuation at all* that makes ‘ P ’ true and ‘ Q ’ false.

Despite this important difference, there are some close connections between conditionals in TFL and claims about entailment in the metalanguage. Observe the following:

- $\phi \models \psi$ iff no valuation makes ϕ true and ψ false.
- $\phi \rightarrow \psi$ is a tautology iff no valuation makes $\phi \rightarrow \psi$ false. Since a conditional is only false if its antecedent is true and its consequent false, this means that $\phi \rightarrow \psi$ is a tautology iff no valuation makes ϕ true and ψ false.

Combining these two observations, we see that:

$$\phi \models \psi \text{ iff } \models \phi \rightarrow \psi$$

This means that the entailment $\phi \models \psi$ holds if and only if the corresponding conditional $\phi \rightarrow \psi$ is a tautology. Note: the mere *truth* of the conditional $\phi \rightarrow \psi$ does not suffice for the entailment $\phi \models \psi$ to hold. For the latter to hold, ‘ $\phi \rightarrow \psi$ ’ doesn’t just have to be *true*, it has to be a *tautology*. More generally, we have:

$$\phi_1, \dots, \phi_n, \psi \models \chi \text{ iff } \phi_1, \dots, \phi_n \models \psi \rightarrow \chi$$

■ Exercises 3.6

A. Use joint truth tables to determine if the following entailments hold:

1. $A \rightarrow A \models A$
2. $A \rightarrow (A \wedge \neg A) \models \neg A$
3. $A \vee (B \rightarrow A) \models \neg A \rightarrow \neg B$
4. $A \vee B, B \vee C, \neg A \models B \wedge C$
5. $(B \wedge A) \rightarrow C, (C \wedge A) \rightarrow B \models (C \wedge B) \rightarrow A$
6. $P \rightarrow (P \vee Q), P \models Q$
7. $(P \rightarrow Q), (\neg R \rightarrow \neg Q), P \models R$
8. $A \vee (B \vee C), \neg A, B \rightarrow C \models C$

B. Symbolize the following argument and use a joint truth table to determine if it is valid:

If Bugs Bunny is a rabbit, then he is a mammal.
 Bugs Bunny is a mammal.
 \therefore Bugs Bunny is a rabbit.

C. Construct joint truth tables to determine whether the following hold:

1. $(A \leftrightarrow B) \models ((A \rightarrow B) \wedge (B \rightarrow A))$

2. $(A \rightarrow B) \models (B \rightarrow A)$
3. $(A \rightarrow B) \models \neg A \vee B$
4. $((A \wedge B) \wedge C) \models (A \wedge (B \wedge C))$ (That is: is \wedge is associative?)
5. $((A \vee B) \vee C) \models (A \vee (B \vee C))$ (That is: is \vee is associative?).
6. $((A \rightarrow B) \rightarrow C) \models (A \rightarrow (B \rightarrow C))$ (That is: is \rightarrow associative?)
7. $((A \leftrightarrow B) \leftrightarrow C) \models (A \leftrightarrow (B \leftrightarrow C))$ (That is: is \leftrightarrow is associative?)

D. Here are some important logical equivalence laws:

Double Negation (DN):

$$\neg\neg\phi \models \phi$$

DeMorgan's Laws (DeM):

$$\neg(\phi \wedge \psi) \models \neg\phi \vee \neg\psi$$

$$\neg(\phi \vee \psi) \models \neg\phi \wedge \neg\psi$$

Laws of Redundancy (R):

$$\phi \wedge \phi \models \phi$$

$$\phi \vee \phi \models \phi$$

Distributive Laws (Dist):

$$\phi \wedge (\psi \vee \chi) \models (\phi \wedge \psi) \vee (\phi \wedge \chi)$$

$$\phi \vee (\psi \wedge \chi) \models (\phi \vee \psi) \wedge (\phi \vee \chi)$$

Material Conditional Laws:

Imp: $\phi \rightarrow \psi \models \neg\phi \vee \psi$

NegImp: $\neg(\phi \rightarrow \psi) \models \phi \wedge \neg\psi$

Cont: $\phi \rightarrow \psi \models \neg\psi \rightarrow \neg\phi$

For each of the following, say which equivalence law it is an instance of:

1. $(P \vee \neg Q) \wedge (P \vee \neg Q) \models (P \vee \neg Q)$
2. $(P \vee \neg Q) \wedge (P \vee \neg Q) \models (P \vee (\neg Q \wedge \neg Q))$
3. $\neg A \wedge \neg\neg(B \wedge \neg C) \models \neg(A \vee \neg(B \wedge \neg C))$
4. $\neg C \vee (B \wedge \neg A) \models (\neg C \vee B) \wedge (\neg C \vee \neg A)$
5. $\neg B \vee \neg(A \rightarrow \neg C) \models B \rightarrow \neg(A \rightarrow \neg C)$
6. $(A \vee C) \rightarrow \neg B \models \neg\neg B \rightarrow \neg(A \vee C)$

E. Consider the following principle, and explain whether it is correct or not:

- Suppose ϕ and ψ are equivalent. Then given any argument that contains ϕ (either as a premise or as its conclusion), replacing ϕ with ψ will not affect that argument's validity.

The Limits of Our Tests

This is an important milestone: a test for the validity of arguments! But we shouldn't get carried away. It is important to understand the *limits* of our achievement. We can illustrate these limits with a few examples. First, consider the argument:

1. Daisy has four legs. \therefore Daisy has more than two legs.

To symbolize this argument in TFL, we would have to use two different atomic sentences — perhaps ' F ' and ' T ' — for the premise and the conclusion respectively. Obviously ' F ' does not TF entail ' T ', and the argument $F \therefore T$ is therefore not TF valid. And yet the English argument is valid, i.e. it's impossible for the premise to be true and the conclusion false!

This case perhaps isn't so bad. As we discussed in §1.4, logic only aims to identify arguments that are *formally* valid, that is, valid in virtue of their logical structure. And the above argument isn't formally valid. But now consider this argument:

2. Some parallelograms are squares. All squares are equilateral. \therefore Some parallelograms are equilateral.

To symbolize this in TFL, we'd again have to use different atomic sentences for the premises and conclusion, something like:

$$P, S \therefore E$$

Again this argument is not TF valid. And yet the English argument is valid, and indeed formally valid this time! So here TFL really does fail us. To capture the structure in virtue of which this argument is valid, we need a stronger system of logic, namely FOL, which we'll look at in the second part of this book.

Similar shortcomings beset our other truth table tests. Consider the sentence:

3. Jan is neither bald nor not-bald.

We could symbolize this in TFL as ' $\neg(J \vee \neg J)$ ', or (given DeMorgan's Laws) as ' $\neg J \wedge \neg \neg J$ '. If we constructed a truth-table for this, we'd see that it's a TF contradiction. But the English sentence (3) does not seem like a contradiction: maybe Jan is on the borderline between being bald and not-bald, and (3) is in fact true! So from the fact that an English sentence receives a contradictory symbolization in TFL, we can't always conclude that the English sentence is necessarily false.

Lastly, consider the following sentence:

4. It's not the case that, if God exists, then God answers malevolent prayers.

Symbolising this in TFL, we would offer something like ' $\neg(G \rightarrow M)$ '. Now, ' $\neg(G \rightarrow M)$ ' is TF equivalent to ' $G \wedge \neg M$ ' and therefore TF entails ' G ' (check this with a truth table). So if we symbolize the English sentence (4) in TFL, it seems to entail that God exists. But that's strange: surely even the atheist can accept (4), and not thereby commit herself to the existence of God!

In different ways, all of these examples highlight some of the limits of working with a language like TFL that can *only* handle truth-functional connectives. These limits give rise to some interesting questions in the field of *philosophical logic*. Our first two arguments raise questions about the distinction between validity and formal validity, and how these are related to our surrogate notion in TFL. The case of Jan's baldness raises the question of what logic we should use when dealing with *vague* language. And the case of the atheist raises the question of how to deal with the *paradoxes of material implication*, which arise from differences between English 'if ... then' constructions and the TFL conditional ' \rightarrow ' (see §3.3). Part of the purpose of this course is to equip you with the tools needed to explore these questions in philosophical logic. But we have to walk before we can run. We have to become proficient in using TFL before we can adequately discuss its limits, and consider alternatives.

3.7 Truth Table Shortcuts

As you become better at constructing truth tables, you will quickly notice that you can use shortcuts to lighten your work. For example, you know for sure that a disjunction is true whenever one of the disjuncts is true. So once you find one true disjunct, there is no need to work out the truth values of the other disjuncts. Thus you might offer:

P	Q	$(\neg P \vee \neg Q) \vee \neg P$		
T	T	F	FF	FF
T	F	F	TT	TF
F	T			TT
F	F			TT

We don't need to know what the truth value of ' $(\neg P \vee \neg Q)$ ' is on the third and fourth row, because we already know that ' $\neg P$ ' is true on these rows, meaning that sentence as a whole must be true too. What we ultimately care about is the column under the main connective, so you only need to do as much work as is needed to determine that column.

Similarly, you know for sure that a conjunction is false whenever one of the conjuncts is false. So if you find one false conjunct, there is no need to work out the truth value of the other conjunct. Thus you might offer:

P	Q	$\neg(P \wedge \neg Q) \wedge \neg P$		
T	T			FF
T	F			FF
F	T	T	F	TT
F	F	T	F	TT

There's no need to look at the truth value of ' $\neg(P \wedge \neg Q)$ ' on the first and second row since we already know that the second conjunct ' $\neg P$ ' is false on these rows.

A similar short cut is available for conditionals. You immediately know that a conditional is true if either its consequent is true, or its antecedent is false. Thus you might present:

P	Q	$((P \rightarrow Q) \rightarrow P) \rightarrow P$		
T	T			T
T	F			T
F	T	T	F	T
F	F	T	F	T

So ' $((P \rightarrow Q) \rightarrow P) \rightarrow P$ ' is a tautology. In fact, it is an instance of *Peirce's Law*, named after Charles Sanders Peirce.

We can apply shortcuts when testing for entailment, or validity, too. To test for entailment, we need to identify "bad" lines in the joint truth table: lines where the premises are all true but the conclusion is false. A line like this is a COUNTEREXAMPLE to the entailment. Now:

- If the conclusion is true on a line, then that line can't be counterexample. (And we don't need to evaluate *anything else* on that line to confirm this.)
- If any premise is false on a line, then that line can't be a counterexample. (And we don't need to evaluate *anything else* on that line to confirm this.)

With this in mind, we can speed up our tests for validity considerably. Consider how we might test the following argument for validity:

$$\neg L \rightarrow (J \vee L), \neg L \therefore J$$

The *first* thing we should do is evaluate the conclusion. If we find that the conclusion is *true* on some row, then that row is not a counterexample, and we can ignore it. In this case that leaves us with only the third and fourth rows to consider:

J	L	$\neg L \rightarrow (J \vee L)$	$\neg L$	J
T	T			T
T	F			T
F	T	?	?	F
F	F	?	?	F

with the question-marks indicating where we need to keep digging. The easiest premise to evaluate is the second, so we do that next. Filling in rows three and four for it gives us:

J	L	$\neg L \rightarrow (J \vee L)$	$\neg L$	J
T	T			T
T	F			T
F	T		F	F
F	F	?	T	F

Since ‘ $\neg L$ ’ is false on row three, that row is certainly not a counterexample and we can ignore it. So this leaves us with only row four to consider. Filling it in gives us:

J	L	$\neg L \rightarrow (J \vee L)$	$\neg L$	J
T	T			T
T	F			T
F	T		F	F
F	F	T F F	T	F

The truth table has no counterexample rows, so the argument is valid. Any valuation which makes the conclusion false also makes at least one premise false.

■ Exercises 3.7

A. Using shortcuts, check whether each of the following is a tautology, a contradiction, or contingent:

1. $\neg D \vee D$
2. $(A \wedge B) \vee (B \wedge A)$
3. $\neg[A \rightarrow (B \rightarrow A)]$
4. $A \leftrightarrow [A \rightarrow (B \wedge \neg B)]$
5. $\neg(A \wedge B) \leftrightarrow A$
6. $(A \wedge \neg A) \rightarrow (B \vee C)$
7. $(B \wedge D) \leftrightarrow [A \leftrightarrow (A \vee C)]$

3.8 Partial Truth Tables

Using shortcuts like these can save a lot of work. But we can get even more efficient. Recall that truth tables grow exponentially: to test an argument involving n atomic sentences, we have to consider a joint truth table with 2^n rows. So if an argument involves 5 atomic sentences, for example, that would mean setting up a 32 row table!

We can be more efficient by using the method of *constructing partial truth tables*. To show that an entailment fails, it suffices to find a single counterexample, i.e. a single valuation that makes all the premises true and the conclusion false. So rather than to set up a complete joint truth table and determine whether any row meets this condition, it is often quicker to try and actively construct a truth table row that does the trick. There are two possible outcomes:

- ▷ We might *succeed* in constructing a counterexample. We can then conclude that the entailment fails and the argument is *invalid*.
- ▷ We might discover that it's *impossible* to construct a counterexample. In this case, we can conclude that the entailment holds and the argument is *valid*.

Example 1 Suppose we have to test whether the following is valid:

$$P \leftrightarrow \neg R, (P \vee Q) \rightarrow \neg S \therefore P \rightarrow (S \vee Q)$$

Rather than set up a sixteen row joint truth table, we'll see if we can "reverse engineer" an assignment of truth values to the atomic sentences ' P ', ' Q ', ' R ', and ' S ' that makes both premises true and the conclusion false. That is, we want to know whether there's a way to fill in truth-values for atomic sentences so as to get a truth-table row that looks as follows:

P	Q	R	S	$P \leftrightarrow \neg R$	$(P \vee Q) \rightarrow \neg S$	$P \rightarrow (S \vee Q)$
?	?	?	?	T	T	F

Let's begin with the conclusion. To make ' $P \rightarrow (S \vee Q)$ ' false, we have to make ' P ' true and ' $(S \vee Q)$ ' false, which means making both ' S ' and ' Q ' false:

P	Q	R	S	$P \leftrightarrow \neg R$	$(P \vee Q) \rightarrow \neg S$	$P \rightarrow (S \vee Q)$
T	F	?	F	T	T	T F F F F

Next, given that ' P ' is true, in order to make ' $P \leftrightarrow \neg R$ ' true we have to make $\neg R$ true as well, meaning ' R ' has to be false:

P	Q	R	S	$P \leftrightarrow \neg R$	$(P \vee Q) \rightarrow \neg S$	$P \rightarrow (S \vee Q)$
T	F	F	F	T T T F	T	T F F F F

At this point we have a valuation that covers all the atomic sentences. And we know that it makes the conclusion false and the first premise true. To make sure that our valuation really constitutes a counterexample to the entailment, we have to be sure that it makes the second premise true as well. And it does:

P	Q	R	S	$P \leftrightarrow \neg R$	$(P \vee Q) \rightarrow \neg S$	$P \rightarrow (S \vee Q)$
T	F	F	F	T	T	F

Since the valuation we've constructed succeeds as a counterexample, we can conclude that the entailment does not hold and that the argument is not TF valid. Constructing this partial truth table was a lot quicker than calculating a complete 16 row table!

Example 2 For another example, consider following TFL argument:

$$A \rightarrow (D \wedge C), B \leftrightarrow \neg D \therefore A \rightarrow (\neg B \wedge C)$$

Again, to test if the premises entail the conclusion, we have to determine whether there is a truth-table row that looks like this:

A	B	C	D	$A \rightarrow (D \wedge C)$	$B \leftrightarrow \neg D$	$A \rightarrow (\neg B \wedge C)$
?	?	?	?	T	T	F

We begin with the conclusion: to make ' $A \rightarrow (\neg B \wedge C)$ ' false, we have to make ' A ' true and ' $(\neg B \wedge C)$ ' false. We could make ' $(\neg B \wedge C)$ ' false by either making ' $\neg B$ ' false or making ' C ' false. We don't know which yet, so all we have at this stage is:

A	B	C	D	$A \rightarrow (D \wedge C)$	$B \leftrightarrow \neg D$	$A \rightarrow (\neg B \wedge C)$
T	?	?	?	T	T	F

However, since we're trying to make the first premise ' $A \rightarrow (D \wedge C)$ ' true, and we've made ' A ' true, we have to make both ' D ' and ' C ' true:

A	B	C	D	$A \rightarrow (D \wedge C)$	$B \leftrightarrow \neg D$	$A \rightarrow (\neg B \wedge C)$
T	?	T	T	T	T	F

Next, since ' D ' is true, ' $\neg D$ ' must be false, so given that we're trying to make ' $B \leftrightarrow \neg D$ ' true, we have to make ' B ' false as well:

A	B	C	D	$A \rightarrow (D \wedge C)$	$B \leftrightarrow \neg D$	$A \rightarrow (\neg B \wedge C)$
T	F	T	T	T	F	F

We now have truth values assigned to all our atomic sentences. But there's a problem: we've had to make ' B ' false and ' C ' true, which means that ' $(\neg B \wedge C)$ ' is true. But that now makes our conclusion ' $A \rightarrow (\neg B \wedge C)$ ' true:

A	B	C	D	$A \rightarrow (D \wedge C)$	$B \leftrightarrow \neg D$	$A \rightarrow (\neg B \wedge C)$
T	F	T	T	T	F	T

whereas we were trying to construct a valuation that makes it false! What we've discovered is that it's impossible to construct such a valuation. So we can conclude that the entailment holds, and that the argument is TF valid.

Notice that the truth table row we've constructed technically does not, by itself, show that the entailment holds. All it shows is that the particular valuation on which ' A ', ' C ', and ' D ' are true and ' B ' is false does not succeed in making premises true and the conclusion false.

To show that the argument is valid, we have to show that *no other* valuation can make the premises true and the conclusion false either.

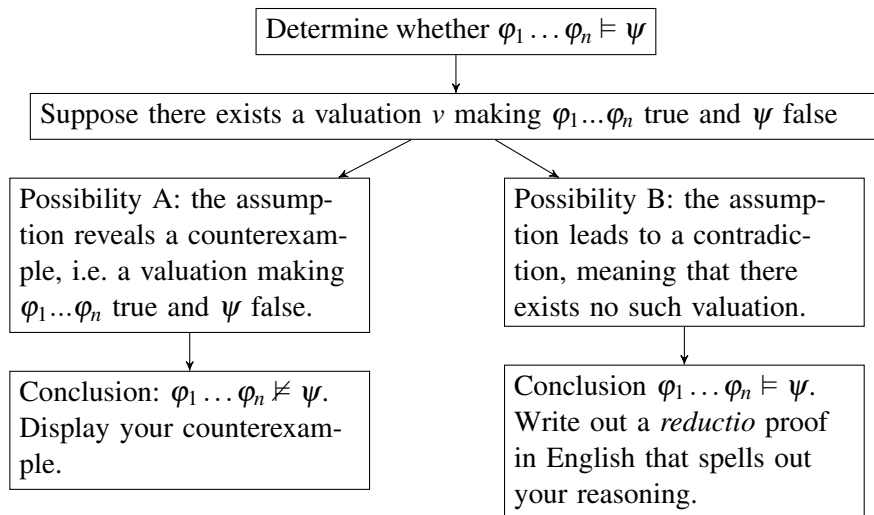
To show that the entailment holds, we can give a verbal proof in English that recapitulates the reasoning we went through in constructing our (failed) truth table row. The proof looks like this:

Claim: $A \rightarrow (D \wedge C), B \leftrightarrow \neg D \models A \rightarrow (\neg B \wedge C)$

Proof: assume (for reduction) that there exists a valuation, let's call it v , that makes ' $A \rightarrow (D \wedge C)$ ' and ' $B \leftrightarrow \neg D$ ' true but ' $A \rightarrow (\neg B \wedge C)$ ' false. Since ' $A \rightarrow (\neg B \wedge C)$ ' is false, ' A ' must be true. And since ' $A \rightarrow (D \wedge C)$ ' is true and ' A ' is true, we know ' $(D \wedge C)$ ' must be true, meaning that both ' D ' and ' C ' are true. Further, since ' D ' is true, ' B ' must be false in order for ' $B \leftrightarrow \neg D$ ' to be true. But now if ' B ' is false and ' C ' is true, ' $\neg B \wedge C$ ' is true, meaning that the conclusion ' $A \rightarrow (\neg B \wedge C)$ ' is true as well. This contradicts our original assumption that ' $A \rightarrow (\neg B \wedge C)$ ' is false on v . Since our assumption lead to a contradiction, we can conclude that the assumption is false, that is, that there *does not* exist a valuation that makes ' $A \rightarrow (D \wedge C)$ ' and ' $B \leftrightarrow \neg D$ ' true but ' $A \rightarrow (\neg B \wedge C)$ ' false. So the entailment holds. QED³

This style of proof is called a proof by *reductio ad absurdum*: we began with an assumption (that there exists a valuation that makes the premises true and the conclusion false), showed that a contradiction (or “absurdity”) results from it, and concluded that the assumption is false (that there exists no such valuation).

Instead of giving a reductio proof like this in English, we could instead construct a full 16 row truth table, and demonstrate validity that way (showing that none of the 16 rows constitute counterexamples to the entailment). But giving the proof in English is often quicker, and gives us more insight, so we'll generally give proofs like this to demonstrate that an entailment does hold. The following diagram summarizes the method of constructing counterexamples via partial truth tables:



³Here 'QED' abbreviates the latin phrase "*quod erat demonstrandum*," meaning "which was to be proven." Writing QED at the end of a proof is a signal that the proof is complete, and establishes the claim we set out to prove.

■ Exercises 3.8

A. Use the partial truth table method to determine whether the following entailments hold. Remember: if you find that the entailment holds, you need to give a *reductio* proof in English. A one-row partial table only suffices to demonstrate that an entailment *fails* to hold.

1. $A \rightarrow (C \vee E), B \rightarrow D \models (A \vee B) \rightarrow (C \rightarrow (D \vee E))$
2. $D \vee \neg A, \neg(B \vee C) \rightarrow \neg D \models A \rightarrow (B \wedge C)$
3. $(D \rightarrow H) \rightarrow P, D \rightarrow \neg(C \vee G), C \vee H \models D \rightarrow P$
4. $A \rightarrow (B \wedge E), D \rightarrow (A \vee C), \neg E \models D \rightarrow B$
5. $\neg A \vee (B \rightarrow C), E \rightarrow (B \wedge A), C \rightarrow E \models C \leftrightarrow A$
6. $\neg C \rightarrow (\neg B \wedge \neg D), C \rightarrow \neg A, B \vee A \models A \leftrightarrow \neg C$
7. $P \vee Q, P \rightarrow R, Q \rightarrow \neg R \models P \leftrightarrow R$
8. $A \vee [A \rightarrow (A \leftrightarrow A)] \models A$
9. $A \leftrightarrow \neg(B \leftrightarrow A) \models A$
10. $A \rightarrow B, B \models A$
11. $A \vee B, B \vee C, \neg B \models A \wedge C$
12. $A \leftrightarrow B, B \leftrightarrow C \models A \leftrightarrow C$

Testing for Other Semantic Notions

We can use partial truth tables to test for other semantic notions, besides entailment.

Tautology To test whether ‘ $(U \wedge T) \rightarrow (S \wedge W)$ ’ is a tautology, we can set up a partial truth table and see whether it’s possible to make the sentence false:

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
				F

Since this is a conditional, and we’re trying to make it false, the antecedent must be true and the consequent must be false:

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
	T	F	F	

In order for the ‘ $(U \wedge T)$ ’ to be true, both ‘ U ’ and ‘ T ’ must be true.

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
	T	T		T T T F F

Now we just need to make ‘ $(S \wedge W)$ ’ false. To do this, we need to make at least one of ‘ S ’ and ‘ W ’ false. We can make both ‘ S ’ and ‘ W ’ false if we want. All that matters is that the whole sentence turns out false on this line. Making an arbitrary decision, we finish the table in this way:

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
F	T	T	F	T T T F F F F

So we now have a partial truth table which shows that there is a valuation which makes ' $(U \wedge T) \rightarrow (S \wedge W)$ ' false, namely, the valuation which makes ' S ' false, ' T ' true, ' U ' true and ' W ' false. So we can conclude that ' $(U \wedge T) \rightarrow (S \wedge W)$ ' is not a tautology.

Our partial truth table suffices to show that this sentence is *not* a tautology. But a partial truth table does not suffice to show that a sentence *is* a tautology, just as it doesn't suffice to show that an argument is valid. To show that a sentence is a tautology, i.e. to show that it's true on *every* valuation, we'd have to either give a full truth table, or a *reductio* argument in English showing that it's *impossible* to construct a valuation that makes the sentence false.

Contradiction. To test whether a sentence is a contradiction, we see whether we can construct a valuation that makes it true:

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
				T

To make the sentence true, it will suffice to ensure that the antecedent is false. Since the antecedent is a conjunction, we can just make one of them false. Making an arbitrary choice, let's make ' U ' false; we can then assign any truth value we like to the other atomic sentences.

S	T	U	W	$(U \wedge T) \rightarrow (S \wedge W)$
F	T	F	F	F F T T F F F

Since there is a valuation that makes the sentence true, our partial table shows that it is *not* a contradiction. Again, to show that something *is* a contradiction (false on *every* valuation), we'd have to give a full truth table or a *reductio* argument showing it is impossible to make the sentence true.

Consistency. To test some sentences for consistency, we would test whether we can construct a partial truth table which makes all of the sentence true. If we succeed, that is sufficient to demonstrate consistency. To demonstrate *inconsistency* we'd have to give a full truth table, or a *reductio* argument showing that it is impossible to make all the sentences in question true.

Equivalence To test two sentences for equivalence, we would test whether we can construct a partial truth table on which the two sentences have different truth values. If we succeed, that is sufficient to show that the sentences are *not* equivalent. To demonstrate equivalence, we'd have to give a full truth table, or a *reductio* argument showing that it's impossible to make the sentences have different truth values (or alternatively, two *reductio* arguments showing that the sentences mutually entail each other).

This table summarises what is required:

	Yes	No
Entailment?	complete table or <i>reductio</i>	partial truth table
Tautology?	complete table or <i>reductio</i>	partial truth table
Contradiction?	complete table or <i>reductio</i>	partial truth table
Consistent?	partial truth table	complete table or <i>reductio</i>
Equivalent?	complete table or <i>reductio</i>	partial truth table

■ Exercises 3.8

A. Use the partial truth table method to determine whether these pairs of sentences are equivalent. And remember, if you find they *are* equivalent you need to give a *reductio* proof (or a full truth table).

1. $A, \neg A$
2. $A, A \vee A$
3. $A \rightarrow A, A \leftrightarrow A$
4. $A \vee \neg B, A \rightarrow B$
5. $A \wedge \neg A, \neg B \leftrightarrow B$
6. $\neg(A \wedge B), \neg A \vee \neg B$
7. $\neg(A \rightarrow B), \neg A \rightarrow \neg B$
8. $(A \rightarrow B), (\neg B \rightarrow \neg A)$

B. Use the partial truth table method to determine whether these sentences are jointly consistent or inconsistent. Again, to show they're *inconsistent* you need to give a *reductio* proof (or a full truth table).

1. $A \wedge B, C \rightarrow \neg B, C$
2. $A \rightarrow B, B \rightarrow C, A, \neg C$
3. $A \vee B, B \vee C, C \rightarrow \neg A$
4. $A, B, C, \neg D, \neg E, F$

Natural Deduction for TFL

4

4.1 The Idea Behind Natural Deduction

We’ve seen how to use truth tables to determine whether a TFL argument is valid. Truth tables are nice because they give us a completely mechanical test for validity: we just crunch through the table and see whether there is any valuation that makes all the premises true and the conclusion false. But truth tables don’t give us much *insight* into why arguments are valid.

A rather different approach to logic is to try to *deduce* the conclusion from the premises via a series of simple inferences. Think of this as the Sherlock Holmes approach to logic: given some evidence (premises), Holmes methodically draws one inference after another, until he arrives at a conclusion about who committed the crime. If all the inferences in the deduction are correct, the conclusion must follow from the premises Holmes started with.

What we will do in this chapter is to introduce a NATURAL DEDUCTION system that formalizes this process of deducing things from premises. We’ll introduce a few very basic rules of inference, which can then be combined into more complicated chains of reasoning. Indeed, with just a small set of rules, we will be able to capture *all* valid arguments. Whereas truth tables are completely mechanical, natural deduction requires insight and ingenuity. This makes it harder, but also more interesting and rewarding.

The move to natural deduction can be motivated by more than the search for insight, however. It might also be motivated by necessity. In TFL, truth tables give us a completely mechanical test for validity. Of course they can get unmanageably big as the number of atomic sentences increase, but you could in principle program a computer to crunch through them for you. When we get to FOL, in the second part of this book, things will look very different. There is nothing like the truth table test for validity available in FOL. The fact that there is no completely mechanical test for validity in FOL is a deep mathematical result, independently proved by Alan Turing and Alonzo Church in 1936. So in FOL, using methods that require ingenuity and insight become indispensable, and we will have to rely on natural deduction to prove arguments valid.

The modern development of natural deduction dates from simultaneous papers from 1934 by Gerhard Gentzen and Stanisław Jaśkowski. Later, in 1952, Frederic Fitch introduced the graphical “Fitch notation” for natural deduction proofs that we will use here.

4.2 Setting up Natural Deduction Proofs

The NATURAL DEDUCTION system we will develop includes a pair of rules for every connective in the language. INTRODUCTION RULES allow us to prove a sentence that has that connective as the main logical operator, and ELIMINATION RULES allow us to prove something *from* a sentence that has that connective as the main logical operator.

Our natural deduction proofs will be *formal proofs*. They will consist of a sequence of lines, with the premises listed at the top and the conclusion at the bottom. All the lines in between have to be justified as following from earlier lines via some rule of inference. As an illustration, consider the following instance of DeMorgan's Law:

$$\neg(A \vee B) \therefore \neg A \wedge \neg B$$

We would start this proof by writing the premise:

$$\begin{array}{l|l} 1 & \neg(A \vee B) & \text{Premise} \\ \hline \end{array}$$

Note that we have numbered the premise, since we will want to refer back to it. Indeed, every line on a proof is numbered so that we can refer back to it. Note also that we have drawn a vertical line to the left and a horizontal line underneath the premise. Everything written above the horizontal line is an *assumption* — so the premise is introduced into the proof as an assumption. Everything written below the horizontal line will either be something which follows from this assumption, or it will be some new assumption.

We are hoping to conclude that ' $\neg A \wedge \neg B$ '. So we are hoping ultimately to end our proof with a line that looks like this:

$$\begin{array}{l|l} n & \neg A \wedge \neg B \end{array}$$

for some line number n . It doesn't matter what line number we end on, but we would obviously prefer a short proof to a long one.

Or to take another example, suppose we wanted to prove that the following is valid:

$$A \vee B, \neg(A \wedge C), \neg(B \wedge \neg D) \therefore \neg C \vee D$$

The argument has three premises, so we start by writing them all down, numbered, and drawing a vertical line the the left and a horizontal line underneath:

$$\begin{array}{l|l} 1 & A \vee B & \text{Premise} \\ 2 & \neg(A \wedge C) & \text{Premise} \\ 3 & \neg(B \wedge \neg D) & \text{Premise} \\ \hline \end{array}$$

We are aiming to conclude with a line that looks like this:

$$\begin{array}{l|l} n & \neg C \vee D \end{array}$$

What we have to learn are rules of inference, and how to chain them together to move in a step-by-step fashion from the premises to the conclusion.

Before we look at the rules, however, we'll introduce some new terminology and notation having to do with proofs. We will use the following expression:

$$\varphi_1, \dots, \varphi_n \vdash \psi$$

to mean that ψ is *provable* from $\varphi_1, \dots, \varphi_n$. That is, that there *exists a proof* which ends with ψ and whose premises include at most $\varphi_1, \dots, \varphi_n$. We'll call a provability claim of this form a SEQUENT. By providing a natural deduction proof, we can demonstrate that a sequent holds, i.e. demonstrate that ψ is indeed provable from $\varphi_1, \dots, \varphi_n$. When we want to say that ψ is *not* provable from $\varphi_1, \dots, \varphi_n$, we write:

$$\varphi_1, \dots, \varphi_n \not\vdash \psi$$

Natural deduction does not give us a way to verify claims like these, about the *non-existence* of a proof. More complicated reasoning would be required to show this kind of thing.

The symbol ' \vdash ' is called the *single turnstile*. This is *not* the same as the double turnstile symbol ' \models ' that we used to symbolize entailment in chapter 3. The single turnstile ' \vdash ' says something about the *existence* of a certain kind of *proof* (one that begins with certain premises and ends in a certain conclusion). The double turnstile ' \models ' says something about the *non-existence* of a certain kind of *valuation* (one that makes the premises true and the conclusion false). Valuations are completely different kinds of things from proofs, so it's important not to confuse ' \vdash ' (the *proof theoretic* notion of *provability*) with ' \models ' (the *semantic* notion of *entailment*).

That said, the system of natural deduction that we will develop is designed to deliver a proof whenever a semantic entailment holds. That is, it is designed to ensure:

COMPLETENESS: If $\varphi_1, \dots, \varphi_n \models \psi$ then $\varphi_1, \dots, \varphi_n \vdash \psi$

meaning that if ψ is semantically entailed by $\varphi_1, \dots, \varphi_n$, then there must also exist a proof of ψ from $\varphi_1, \dots, \varphi_n$. Our natural deduction system is also designed to guarantee the other direction:

SOUNDNESS: If $\varphi_1, \dots, \varphi_n \vdash \psi$ then $\varphi_1, \dots, \varphi_n \models \psi$

meaning that whenever ψ is provable from $\varphi_1, \dots, \varphi_n$, then ψ is also semantically entailed by $\varphi_1, \dots, \varphi_n$.¹ All good proof systems should be both sound and complete, to ensure that the proof-theoretic notion of provability matches up perfectly with the semantic notion of entailment. In a more advanced logic class, you learn how to provide "meta-logical" proofs showing that a proof system is both sound and complete, but here you'll just have to take my word for it that our natural deduction system will have both of these features.

We also have a proof-theoretic analogue of the semantic notion of TF equivalence:

Two sentences φ and ψ are PROVABLY EQUIVALENT iff each is provable from the other; i.e., both $\varphi \vdash \psi$ and $\psi \vdash \varphi$, also written $\varphi \dashv\vdash \psi$.

Given that our natural deduction system is both sound and complete, we will again have it that any two TF equivalent sentences are also provably equivalent, and vice versa. Let's now look at the rules that constitute our system of natural deduction.

¹Note that this is a different notion of soundness from the one we discussed in §1.3

4.3 Conjunction Rules

Suppose I want to show that Alice is both a logician and a tennis player. One obvious way to do this would be as follows: first I show that Alice is a logician; then I show that Alice is a tennis player; then I put these two demonstrations together to obtain the conjunction.

Our natural deduction system will capture this thought via the rule of \wedge -Introduction, or $\wedge I$ for short. Perhaps I am working through a proof, and have obtained ' H ' on line 8 and ' R ' on line 15. Then on any subsequent line I can obtain ' $H \wedge R$ ' thus:

8		H	
		\vdots	
15		R	
		\vdots	
		$H \wedge R$	$\wedge I$ 8, 15

Every line of our proof must either be a premise (or an assumption, as we'll see later), or must be justified by some rule like this. We cite ' $\wedge I$ 8, 15' here to indicate that ' $H \wedge R$ ' is obtained by the rule of conjunction introduction applied to lines 8 and 15. We could equally well have conjoined the conjuncts in the opposite order to infer ' $R \wedge H$ ' rather than ' $H \wedge R$ ', though then we should also adjust our rule citation to read ' $\wedge I$ 15, 8', with the line numbers of the two conjuncts listed in the opposite order.

More generally, our conjunction introduction rule is:

m		ϕ	
n		ψ	
		$\phi \wedge \psi$	$\wedge I$ m, n

Here lines m and n can occur in either order, i.e. ϕ could occur first in the proof followed by ψ later on, or ψ could occur first followed by ϕ later on.

The rule is called "conjunction *introduction*" because it introduces the symbol ' \wedge ' into our proof where it may have been absent. Correspondingly, we have a rule that *eliminates* that symbol. Suppose you have shown that Alice is both a logician and a tennis player. Then you're entitled to infer that Alice is a logician, and you're also entitled to infer that Alice is a tennis player. This gives us our conjunction elimination rule(s):

m		$\phi \wedge \psi$	
		ϕ	$\wedge E$ m

and equally:

m	$\varphi \wedge \psi$	
	ψ	$\wedge E\ m$

The point is simply that, when you have a conjunction on some line in a proof, you can obtain either of its two conjuncts by applying the rule of $\wedge E$.

One point is worth emphasizing: you can only apply this rule (as well as the other rules we'll introduce) to the *main logical operator* of a sentence. So the following would not be a legitimate use of $\wedge E$:

1	$P \wedge (Q \wedge R)$	
2	R	$\wedge E\ 1$

I can only apply $\wedge E$ to the main operator of ' $P \wedge (Q \wedge R)$ ', giving me either ' P ' or ' $(Q \wedge R)$ '. I could then apply $\wedge E$ to the latter to get ' R ' itself; but I can't get R *directly* from ' $P \wedge (Q \wedge R)$ '.

Here's an example that illustrates this. The following argument is valid (showing that \wedge is associative):

$$A \wedge (B \wedge C) \therefore (A \wedge B) \wedge C$$

To provide a proof for this argument, we start by writing the premise:

1	$A \wedge (B \wedge C)$	Premise
---	-------------------------	---------

From the premise, we can get both ' A ' and ' $(B \wedge C)$ ' by applying $\wedge E$ twice. And we can then apply $\wedge E$ twice more to $(B \wedge C)$ to get a proof that looks like this:

1	$A \wedge (B \wedge C)$	Premise
2	A	$\wedge E\ 1$
3	$B \wedge C$	$\wedge E\ 1$
4	B	$\wedge E\ 3$
5	C	$\wedge E\ 3$

Again: we *cannot* get ' B ' or ' C ' by applying $\wedge E$ directly to line 1. We first have to get ' $(B \wedge C)$ ' from 1, and then get ' B ' and ' C ' out of *this* by applying $\wedge E$ again. To get to our desired conclusion, we now just put the various atomic sentences back together using $\wedge I$:

1	$A \wedge (B \wedge C)$	Premise
2	A	$\wedge E$ 1
3	$B \wedge C$	$\wedge E$ 1
4	B	$\wedge E$ 3
5	C	$\wedge E$ 3
6	$A \wedge B$	$\wedge I$ 2, 4
7	$(A \wedge B) \wedge C$	$\wedge I$ 6, 5

Notice that whereas $\wedge E$ gets applied to a single line, $\wedge I$ gets applied to two lines. However, $\wedge I$ doesn't necessarily have to be applied to two *different* lines. If we wanted, for example, we could formally prove ' $A \wedge A$ ' from ' A ' as follows:

1	A	
2	$A \wedge A$	$\wedge I$ 1, 1

And we could now apply $\wedge E$ to line 2 to prove the rather uninteresting fact that ' A ' follows from ' A ':

1	A	
2	$A \wedge A$	$\wedge I$ 1, 1
3	A	$\wedge E$ 2

4.4 Conditional Rules

Consider the following argument:

If Jane is smart then she is fast. Jane is smart. So Jane is fast.

This argument is certainly valid. And it suggests a conditional elimination rule ($\rightarrow E$):

m	$\varphi \rightarrow \psi$	
n	φ	
	ψ	$\rightarrow E$ m, n

This rule implements the *modus ponens* form of inference mentioned earlier: given a conditional, and given its antecedent, we can infer its consequent. Again, this is an elimination rule, because it allows us to obtain a sentence that may not contain ' \rightarrow ', having started with a sentence that did contain ' \rightarrow '. Note that the conditional and its antecedent can appear in

either order in our proof. However, in the citation for $\rightarrow E$, we should cite the conditional first, followed by the antecedent.

The rule for conditional introduction is also quite easy to motivate. The following argument should be valid:

Alice is a chemist. Therefore, if Alice is a biologist, then Alice is both a chemist and a biologist.

If someone doubted that this was valid, we could try to convince them otherwise by explaining ourselves as follows:

We are given that Alice is a chemist. Now, assume additionally, for the sake of argument, that Alice is also a biologist. Then by conjunction introduction, Alice is both a chemist and a biologist. Of course, this only follows given our assumption that Alice is a biologist. So what we've shown is that *if* Alice is a biologist, then she is both a chemist and a biologist.

Transferred into natural deduction format, here is the pattern of reasoning that we just used. We started with one premise, 'Alice is a chemist':

1	C	Premise
---	-----	---------

The next thing we did is to make an *temporary assumption* ('Alice is a biologist'), for the sake of argument. To indicate that we are no longer dealing merely with our original premise 'C', but with an additional assumption, we continue our proof as follows:

1	C	Premise
2	<div style="border-left: 1px solid black; padding-left: 10px; vertical-align: middle;">B</div>	Assumption

Introducing ' B ' as a temporary assumption opens up a *subproof*. Inside this subproof we can now reason under the assumption, or hypothesis, that ' B ' holds. We indicate this by drawing a line under ' B ' (to indicate that it is an assumption) and by indenting it with a further vertical line (to indicate that we have entered a new subproof headed by this assumption).

With this extra assumption in place, we are in a position to use $\wedge I$:

1	C	Premise
2	<div style="border-left: 1px solid black; padding-left: 10px; vertical-align: middle;">B</div>	Assumption
3	<div style="border-left: 1px solid black; padding-left: 10px; vertical-align: middle;">$C \wedge B$</div>	$\wedge I$ 1, 2

So we have now shown that, under the assumption that ' B ' holds, we can infer ' $C \wedge B$ '. We can therefore conclude that, *if* ' B ' obtains, so does ' $C \wedge B$ '. Or, put another way, we can conclude ' $B \rightarrow (C \wedge B)$ ':

1		C	
2			B
3			$C \wedge B$ $\wedge I$ 1, 2
4		$B \rightarrow (C \wedge B)$	$\rightarrow I$ 2–3

Notice that we have popped back out of the subproof opened by our assumption. This indicates that we have now *discharged* the temporary assumption ‘ B ’, and concluded that the conditional itself follows just from our original premise ‘ C ’. Although you can in principle always make an assumption, it is absolutely essential that any assumptions you do make be discharged by the end of the proof. This is why they are *temporary* assumptions.

The general pattern at work here is the following: to prove a conditional $\phi \rightarrow \psi$, we *assume* the antecedent ϕ temporarily or “for the sake of argument” (thereby opening a new subproof), and then try to prove the consequent ψ from that assumption. If we succeed, we can discharge the assumption (popping out of the subproof) and conclude that the conditional $\phi \rightarrow \psi$ holds:

m			ϕ	
			\vdots	
n			ψ	
		$\phi \rightarrow \psi$		$\rightarrow I$ m – n

Here’s another illustration of $\rightarrow I$ in action. Suppose we want to prove:

$$P \rightarrow Q, Q \rightarrow R \vdash P \rightarrow R$$

We start by listing both of our premises. Then, since we are aiming to prove a conditional, namely, ‘ $P \rightarrow R$ ’, we assume the antecedent to that conditional as an additional assumption, which opens a new subproof:

1		$P \rightarrow Q$	Premise
2		$Q \rightarrow R$	Premise
3			P Assumption
			\vdots
			R

Our goal now is to prove R from this temporary assumption, inside of our subproof. Given ‘ P ’, we can use $\rightarrow E$ on the first premise. This will yield ‘ Q ’. And we can then use $\rightarrow E$ on the second premise to get ‘ R ’. So, by assuming ‘ P ’ we were able to prove ‘ R ’! We can now apply the $\rightarrow I$ rule, thereby discharging ‘ P ’ and popping out of the subproof::

1	$P \rightarrow Q$	Premise
2	$Q \rightarrow R$	Premise
3	P	Assumption
4	Q	$\rightarrow E$ 1, 3
5	R	$\rightarrow E$ 2, 4
6	$P \rightarrow R$	$\rightarrow I$ 3–5

Notice that when applying $\rightarrow I$ to obtain ' $P \rightarrow R$ ', we have to cite the *entire* subproof that begins with ' P ' and ends with ' R '. So we use a dash, rather than just a comma, between the two line numbers (writing '3–5' rather than '3,5').

4.5 Additional assumptions and subproofs

The rule $\rightarrow I$ invoked the idea of opening subproofs via additional assumptions. This needs to be handled with some care. Consider this proof:

1	A	Premise
2	B	Assumption
3	$B \wedge B$	$\wedge I$ 2, 2
4	B	$\wedge E$ 3
5	$B \rightarrow B$	$\rightarrow I$ 2–4

This is a perfectly legitimate, if somewhat unusual, proof. What it shows is that the argument $A \therefore B \rightarrow B$ is valid. This is as it should be: ' $B \rightarrow B$ ' is a tautology, and any argument with a tautology as its conclusion is valid. But suppose we now tried to continue the proof as follows:

1	A	Premise
2	B	Assumption
3	$B \wedge B$	$\wedge I$ 2, 2
4	B	$\wedge E$ 3
5	$B \rightarrow B$	$\rightarrow I$ 2–4
6	B	No!! $\rightarrow E$ 5, 4

If we were allowed to do this, it would be a disaster: our proof would now purport to show that the argument $A \therefore B$ is valid. We could in this way prove any conclusion we liked from any premise whatsoever. That would obviously destroy the soundness of our natural deduction system.

What has gone wrong here is that on line 6 we've illegitimately tried to apply $\rightarrow E$ to line 4, which occurs inside a subproof we've already closed off. A subproof can be thought of as showing what *would* follow if the assumption that opens it held. While we are working within the subproof, we can refer to the assumption that we made to open the subproof, and to anything that we obtained from our premises. After all, those premises still hold. But once we close the subproof and return to the main proof, the assumption that opened it has been DISCHARGED, and it becomes illegitimate to draw upon anything that depends upon that assumption, i.e. on anything inside the subproof opened by that assumption. Thus we stipulate:

To cite any individual line when applying a rule, that line must (1) occur before the application of the rule, but (2) not occur within a closed subproof.

The application of $\rightarrow E$ in the faulty proof above involves citing a line (namely line 4) that occurs within a subproof that has (by line 6) been closed. This is illegitimate.

Once we have started thinking about what we can show by making additional assumptions, nothing stops us from posing the question of what we could show if we were to make *even more* assumptions. We can in other words introduce a subproof within a subproof. Here is an example of such nested subproofs:

1	A	Premise
2	B	Assumption
3	C	Assumption
4	A \wedge B	$\wedge I$ 1, 2
5	C \rightarrow (A \wedge B)	$\rightarrow I$ 3–4
6	B \rightarrow (C \rightarrow (A \wedge B))	$\rightarrow I$ 2–5

This proof gets set up as follows: we begin with the premise 'A' and the goal of proving ' $B \rightarrow (C \rightarrow (A \wedge B))$ '. Since the conclusion is a conditional, we assume its antecedent 'B' and set ourselves the new goal of proving its consequent ' $(C \rightarrow (A \wedge B))$ ' using this additional assumption. But our new goal ' $(C \rightarrow (A \wedge B))$ ' is *itself* a conditional, so we repeat the same process: assume its antecedent 'C', and try to prove its consequent ' $A \wedge B$ ' in the subproof we've opened. Proving ' $A \wedge B$ ' is easy: we can just apply $\wedge I$ to our original premise from line 1 and our first assumption from line 2. Referring back to lines 1 and 2 in step 4 of the proof in this manner is legitimate, since neither line occurs in a subproof that has been closed by the time of step 4.

But it would now *not* be legitimate to continue the proof as follows:

1	A	
2	B	
3	C	
4	A ∧ B	∧I 1, 2
5	C → (A ∧ B)	→I 3–4
6	B → (C → (A ∧ B))	→I 2–5
7	C → (A ∧ B)	No!! →I 3–4

This would be awful. This proof would purport to show that ' $C \rightarrow (A \wedge B)$ ' can be deduced from the premise ' A '. But the argument ' $A \therefore C \rightarrow (A \wedge B)$ ' is certainly not valid. Again, if we were allowed to do this kind of thing, our proof system would no longer be sound.

The problem is that the subproof that began with the assumption ' C ' occurs within the scope of (i.e. within the subproof opened by) assumption ' B ' on line 2. By line 6, we have *discharged* assumption ' B '. So it is cheating to try to help ourselves (on line 7) to a subproof that occurs within the scope of an assumption that has already been discharged. Here the problem isn't that we cited an individual *line* that occurs inside a closed subproof, but that we cited an entire *subproof* that occurs inside a closed subproof. So we expand our stipulation to cover rules that cite entire subproofs:

To cite a subproof when applying a rule, the subproof must (1) come before the application of the rule, but (2) not occur within some other closed subproof.

Our proof above violates this stipulation, since the subproof of lines 3–4 occurs within the subproof spanning lines 2–5, which has already been closed by the point we get to line 7.

And to again emphasize a point from earlier: although you can, in principle, always make any temporary assumption you like, you always have to ultimately have a way of discharging that assumption and popping out of the subproof that it opens (otherwise it wouldn't be temporary any more!). For this reason it is very important to *only* make assumptions when you have a discharge strategy in mind. At this point, we have only one rule that allows us to make an assumption, and that is the rule of \rightarrow I: if your goal is to prove a conditional $\phi \rightarrow \psi$, you should assume its antecedent ϕ and try to prove its consequent ψ inside the subproof opened by that assumption. But at this point, this is the only time at which you should be making assumptions — when aiming to prove conditionals.

■ Exercises 4.5

A. Prove the following sequents (these require only conjunction rules and \rightarrow E):

1. $A \wedge (B \wedge C) \vdash C \wedge B$
2. $P \wedge Q, (Q \wedge P) \rightarrow R \vdash R$
3. $A \wedge B, B \rightarrow (A \rightarrow C) \vdash C \wedge B$
4. $(A \rightarrow ((C \wedge A) \rightarrow D)), (A \wedge C) \vdash D$

B. Prove the following sequents (these now also require \rightarrow I, and note that the last one asks you to prove an equivalence, meaning you have to give proofs in both directions):

1. $(A \rightarrow (B \wedge D)), (A \rightarrow C) \vdash (A \rightarrow (C \wedge D))$
2. $A \rightarrow B, B \rightarrow C \vdash A \rightarrow C$
3. $A \rightarrow B, B \rightarrow C \vdash A \rightarrow (B \wedge C)$
4. $(A \wedge B) \rightarrow C, A \rightarrow B \vdash A \rightarrow C$
5. $A \rightarrow B \vdash (A \wedge C) \rightarrow (B \wedge C)$
6. $A \rightarrow (B \rightarrow C) \dashv\vdash (A \wedge B) \rightarrow C$

4.6 Proving Theorems and Reiterating

We said that the sequent $\varphi_1, \dots, \varphi_n \vdash \psi$ means that there exists a proof ending in ψ whose premises include at most $\varphi_1, \dots, \varphi_n$. Similarly, we will write:

$$\vdash \varphi$$

to say that there is a proof of φ with no premises whatsoever. Sentences which are provable with no premises are called **THEOREMS**. Notice the similarity with our notation in semantics, where we used $\models \varphi$ to say that φ is a tautology. This is no accident: given that our system of natural deduction is both sound and complete, every *tautology* should be a *theorem* of our proof system, and vice versa.

For example, since $(A \wedge B) \rightarrow A$ is a tautology, we should be able to prove it with no premises. Here's how that looks like:

1		$A \wedge B$	Assumption
2		A	\wedge E 1
3		$(A \wedge B) \rightarrow A$	\rightarrow I 1-2

Notice that, unlike in any of the other proofs we've looked at, the leftmost vertical line, which appears next to our conclusion, has no premise listed at its top. This graphically indicates that ' $(A \wedge B) \rightarrow A$ ' is a theorem in our proof system, i.e. something that is provable without any premises.

As another example, take the tautology ' $A \rightarrow (B \rightarrow A)$ '. This is provable as a theorem as follows:

1		A	Assumption
2		B	Assumption
3		A ∧ B	∧I 1, 2
4		A	∧E 3
5		B → A	→I 2–4
6		A → (B → A)	→I 1–5

This proof is a bit odd: since we're trying to prove ' $(B \rightarrow A)$ ' on line 5, the subproof that begins with ' B ' on line 2 has to end with ' A '. We already have ' A ' as an assumption on line 1, but the only way to get it to appear at the end of our second subproof is to first conjoin it with ' B ' to get ' $(A \wedge B)$ ' on line 3, and then to use $\wedge E$ to get it back on its own on line 4.

In order to avoid having to use the trick of using $\wedge I$ and $\wedge E$ in this way to repeat earlier lines at later stages in a proof, we'll allow ourselves to use the following shortcut rule:

REITERATION RULE: at any point in a proof, we may write down a sentence occurring on any line that (i) appears before that point in the proof, and (ii) isn't inside a closed subproof.

This lets us shorten the above proof to:

1		A	Assumption
2		B	Assumption
3		A	Reit 1
4		B → A	→I 2–3
5		A → (B → A)	→I 1–4

Reiteration also gives us a quick way to prove that ' $B \rightarrow B$ ' is a theorem:

1		B	Assumption
2		B	Reit 1
3		B → B	→I 1–2

In fact, just as $\wedge I$ can be applied to a single line to go from ' A ' to ' $A \wedge A$ ', so $\rightarrow I$ can in principle be applied to a subproof that consists of just one line. So an even shorter proof of the theorem ' $B \rightarrow B$ ' can be given like this:

1		B	Assumption
2		$B \rightarrow B$	\rightarrow I 1–1

■ Exercises 4.6

A. Prove the following theorems:

1. $\vdash (A \wedge B) \rightarrow (B \wedge A)$
2. $\vdash (A \rightarrow B) \rightarrow ((B \rightarrow C) \rightarrow (A \rightarrow C))$
3. $\vdash (A \rightarrow (B \rightarrow C)) \rightarrow (B \rightarrow (A \rightarrow C))$
4. $\vdash A \rightarrow (B \rightarrow A)$
5. $\vdash (A \wedge B) \rightarrow (B \wedge A)$
6. $\vdash P \rightarrow P$
7. $\vdash Q \rightarrow (P \rightarrow P)$

4.7 Biconditional Rules

The rules for the biconditional will be like double-barrelled versions of the rules for the conditional. In order to prove ' $F \leftrightarrow G$ ', for instance, you must be able to prove ' G ' on the assumption ' F ' and prove ' F ' on the assumption ' G '. The biconditional introduction rule \leftrightarrow I therefore requires two subproofs. Schematically, the rule works like this:

i		ϕ	
		\vdots	
j		ψ	
k		ψ	
		\vdots	
l		ϕ	
		$\phi \leftrightarrow \psi$	\leftrightarrow I $i-j, k-l$

There can be as many lines as you like between i and j , and as many lines as you like between k and l . Moreover, the subproofs can come in any order, and the second subproof does not need to come immediately after the first.

The biconditional elimination rule \leftrightarrow E is a bit like \rightarrow E in both directions. If you have the left-hand subsentence of the biconditional, you can obtain the right-hand subsentence, and if you have the right-hand subsentence, you can obtain the left-hand subsentence:

m	$\varphi \leftrightarrow \psi$	
n	φ	
	ψ	$\leftrightarrow E\ m, n$

and:

m	$\varphi \leftrightarrow \psi$	
n	ψ	
	φ	$\leftrightarrow E\ m, n$

As usual, lines m and n can occur in either order, but in the citation for $\leftrightarrow E$, we always cite the line number of the biconditional first.

Here's an example involving $\leftrightarrow I$:

1	$A \wedge B$	Premise
2	A	Assumption
3	B	$\wedge E\ 1$
4	B	Assumption
5	A	$\wedge E\ 1$
6	$A \leftrightarrow B$	$\leftrightarrow I\ 2-3, 4-5$

My reasoning for setting this up is that since it's my goal to prove ' $A \leftrightarrow B$ ', I will have to create two subproofs, one for each direction of the biconditional. My first subproof (lines 2–3) begins with the assumption ' A ' and ends with ' B ', and my second subproof (lines 4–5) goes the other way, beginning with the assumption ' B ' and ending with ' A '.

■ Exercises 4.7

A. Prove the following:

1. $A \wedge B \vdash A \leftrightarrow B$
2. $A \leftrightarrow B \vdash B \leftrightarrow A$
3. $A \leftrightarrow B \vdash (A \wedge C) \leftrightarrow (B \wedge C)$
4. $(A \wedge B) \leftrightarrow (A \wedge C) \vdash A \rightarrow (B \leftrightarrow C)$
5. $A \leftrightarrow B, B \leftrightarrow C \vdash A \leftrightarrow C$
6. $K \wedge L \vdash K \leftrightarrow L$
7. $A \leftrightarrow B \vdash (A \wedge C) \leftrightarrow (B \wedge C)$
8. $(Z \wedge K) \leftrightarrow (Y \wedge M), D \wedge (D \rightarrow M) \vdash Y \rightarrow Z$
9. $\vdash P \leftrightarrow P$

4.8 Negation Rules

Our next connective is negation. In the context of natural deduction, negation is unusual because the rules governing it involve another notion, that of *contradiction*.

Consider this: an effective way to argue against someone is to show that the assumptions they are making collectively lead to a contradiction. At that point, you have your opponent in a bind: since their assumptions lead to a contradiction, they can't all be true! So your opponent has to give up at least one of their assumptions. This argumentative strategy involves the style of *reductio ad absurdum* reasoning that we encountered in §3.8. We are going to make use of this idea in our proof system by adding the absurdity symbol ' \perp '. You can think of it as officially declaring 'contradiction!' or 'reductio!' or 'but that's absurd!'

We can introduce this symbol into a proof whenever we explicitly contradict ourselves, i.e. whenever we have both a sentence and its negation appearing in the proof:

m	φ	
n	$\neg\varphi$	
	\perp	$\perp I\ m, n$

This is traditionally called the rule of $\neg E$ (e.g. by Gentzen in his original formulation of these rules) since the negation sign in $\neg\varphi$ is eliminated in favor of \perp . But students generally find it more intuitive to think of the rule as introducing the absurdity symbol \perp into the proof, so we'll call this the $\perp I$ rule. Notice that it doesn't matter in what order the sentence and its negation appear in the proof, and they also don't need to appear on adjacent lines—as long as a sentence and its negation appear in your proof (and neither is trapped inside a closed subproof), you can declare a contradiction with \perp .

Next, we turn to our rule of negation introduction. This rule will be a formal implementation of the *reductio ad absurdum* proof strategy described earlier: if making an assumption leads you to a contradiction, then you know that assumption must be wrong, and you can infer its negation. The rule looks like this:

i	φ	
	\vdots	
j	\perp	
	$\neg\varphi$	$\neg I\ i-j$

There can be as many lines between i and j as you like. As with other rules that require subproofs, you need to cite the entire subproof when applying $\neg I$. Notice that since the subproof has to end with the contradiction symbol \perp , we will usually have to use $\perp I$ in the course of using $\neg I$.

Here's an example of how this works. Suppose we want to show:

$$\neg A \vdash \neg(A \wedge B)$$

This should strike you as intuitively valid: if A isn't true, then of course any conjunction containing A can't be true either. To prove this, what we'll do is assume that ' $(A \wedge B)$ ' is true, derive a contradiction from this assumption together with our premise, and then conclude ' $\neg(A \wedge B)$ ' by \neg I. That looks like this (the Reiteration step could be skipped too):

1	$\neg A$	Premise
2	$A \wedge B$	Assumption
3	A	\wedge E 2
4	$\neg A$	Reit 1
5	\perp	\perp I 3, 4
6	$\neg(A \wedge B)$	\neg I 2–5

The strategy of reasoning by *reductio ad absurdum* can take another form too, however. The \neg I rule says that in order to show that some sentence ϕ is false (i.e. that $\neg\phi$ holds), we have to show that assuming ϕ leads to a contradiction. But instead of using *reductio* to show that something is *false*, we can also use it to show that something is *true*: to prove that ϕ is true, we suppose that it is false, i.e. assume $\neg\phi$, and reduce *that* to a contradiction.

Interestingly, the rules we have assembled so far won't yet let us replicate this second style of *reductio* reasoning. So we'll add it to our proof system as another primitive rule, which we'll call *indirect proof*:

i	$\neg\phi$	
	\vdots	
j	\perp	
	ϕ	IP i – j

This is called indirect proof because it lets us prove ϕ “indirectly,” by assuming its negation and deriving a contradiction from that assumption. It's quite similar to \neg I, except that the location of the negation symbol gets reversed: instead of *deriving* a negated sentence (as with \neg I), we instead *assume* a negated sentence and then infer its un-negated counterpart.

It bears emphasis that IP is a very powerful rule, because any proof whatsoever can in principle be done using IP as the overall strategy! Just assume the negation of *whatever conclusion* you're trying to prove, derive a contradiction from that, and then infer your conclusion by IP.² But be careful: indirect proofs tend to be longer and more complicated than direct proofs. So IP should only be used as a *last resort*, when you're sure that there's no other way to complete the proof.

Here is an example of something that can *only* be proven using IP:

$$P \wedge \neg P \vdash Q$$

²In section 3.8, on partial truth tables, we essentially did this kind of proof in English: we assumed that the premises are true and that the conclusion false (i.e. that its negation is true), and then showed that a contradiction resulted. We can now do the same reasoning formally, using our rule of IP.

The corresponding entailment $P \wedge \neg P \models Q$ holds: there is no valuation that makes the premise true and the conclusion false. That's simply because no valuation makes the premise true! So if our system of deduction is to be complete, we had better be able to provide a natural deduction proof of this as well. We can do that with IP, as follows:

1	$P \wedge \neg P$	Premise
2	$\neg Q$	Assumption
3	P	$\wedge E$ 1
4	$\neg P$	$\wedge E$ 1
5	\perp	$\perp I$ 3,4
6	Q	IP 2-5

This proof illustrates the EXPLOSION PRINCIPLE: from a contradiction, like $P \wedge \neg P$, anything whatsoever can be proven! We here proved Q , but exactly the same sequence of steps would have equally well let us prove A , or D , or $(S \rightarrow L)$, or anything else. The Explosion Principle can be stated in general terms as $\perp \vdash \phi$.

Another thing that we can only prove using IP is the LAW OF EXCLUDED MIDDLE, which says that $\vdash \phi \vee \neg\phi$, i.e. that any statement of the form $\phi \vee \neg\phi$ is a theorem. Proponents of *intuitionistic logic* reject the Law of Excluded Middle, because they reject the assumption we've been making throughout this book, that for any statement, it or its negation must be true. So they must reject our rule IP, since it lets us prove Excluded Middle.

There are also logicians who reject the Explosion Principle. For example, proponents of *relevance logic* hold that there must always be some "relevant connection" between the premises and conclusion of a valid argument — something the Explosion Principle violates. And proponents of *paraconsistent logic* hold the view that some contradictions are true, and that accepting a contradiction should therefore not allow you to infer anything whatsoever. So they reject the Explosion Principle, and therefore have to use a different set of rules that doesn't prove this principle.

Intuitionistic logic, relevance logic, and paraconsistent logic are all varieties of NON-CLASSICAL LOGIC. However, in CLASSICAL LOGIC, which is what we are here studying, both the Law of Excluded Middle and the Explosion Principle hold. Since we can't prove these without IP, we have to add this rule into our natural deduction system in order to render it complete with respect to classical logic. You can learn more about negation rules and completeness in Exercise 4.8 below.

Let's look at one more proof with IP. CONTRAPOSITION is the following general equivalence law: $\phi \rightarrow \psi \models \neg\psi \rightarrow \neg\phi$. Let's prove an instance of the right-to-left direction: $\neg B \rightarrow \neg A \vdash A \rightarrow B$. Since we're trying to prove a conditional, we'll use $\rightarrow I$: we'll assume A and try to prove B . But how can we get B from our premise $\neg B \rightarrow \neg A$ together with our assumption A ? The only way to do it is by indirect proof, like this:

1		$(\neg B \rightarrow \neg A)$	Premise
2			
3			
4			
5			
6			
7			

To prove B via IP, I assumed $\neg B$, and showed that this, together with my assumption of A and our premise 1, leads to a contradiction. The other direction of contraposition holds too: $A \rightarrow B \vdash \neg B \rightarrow \neg A$. Try proving this as well; here you'll be able to use $\neg I$ instead of IP.

■ Exercises 4.8

A. Prove the following:

1. $A \rightarrow B, A \rightarrow \neg B \vdash \neg A$
2. $P \rightarrow \neg Q \vdash Q \rightarrow \neg P$
3. $A \wedge \neg B \vdash \neg(A \rightarrow B)$
4. $\neg(P \wedge Q) \vdash P \rightarrow \neg Q$
5. $\vdash \neg(A \wedge \neg A)$
6. $A \rightarrow \neg B \vdash B \rightarrow \neg A$
7. $A \rightarrow B, B \rightarrow \neg A \vdash \neg A$
8. $(A \wedge B) \rightarrow \neg A \vdash A \rightarrow \neg B$
9. $\neg A \vdash A \rightarrow B$
10. $(A \wedge \neg B) \rightarrow B \vdash A \rightarrow B$
11. $(A \leftrightarrow \neg B) \vdash \neg(A \leftrightarrow B)$
12. $\neg(A \wedge \neg B) \vdash A \rightarrow B$
13. $P \rightarrow Q \dashv\vdash \neg Q \rightarrow \neg P$ [Do both directions]
14. $\neg(P \rightarrow Q) \dashv\vdash (P \wedge \neg Q)$ [Do both directions]
15. $P \dashv\vdash \neg\neg P$ [Do both directions]
16. $\vdash \neg P \rightarrow (P \rightarrow Q)$
17. $\vdash (\neg A \rightarrow A) \rightarrow A$
18. $\vdash \neg(A \wedge \neg A)$

B. We've seen that if we only have $\perp I$ (aka $\neg E$) and $\neg I$, our proof system is incomplete: we can't prove every classically valid argument to be valid. So we added IP as an additional rule. In the following questions, we'll look more closely at what's needed to get a complete proof system.

1. It turns out that once we add IP into our proof system, we don't really need $\neg I$ any longer, because anything we can prove using $\neg I$ can be proven using IP instead (though again, not the other way around!). The core of $\neg I$ can be expressed as:

$$A \rightarrow \perp \vdash \neg A$$

Your challenge: prove $A \rightarrow \perp \vdash \neg A$ using IP instead of $\neg I$. This shows that we never really have to use $\neg I$.

2. Instead of adding IP into our system, a more conservative approach would have been to add the rule of *double negation elimination*:

i	$\neg\neg\phi$ ϕ	DNE i
-----	--------------------------	---------

That would have been another way to get a complete proof system, but without rendering $\neg I$ idle. The core of IP can be expressed as:

$$\neg A \rightarrow \perp \vdash A$$

Your challenge: prove $\neg A \rightarrow \perp \vdash A$ using DNE (together with our other rules) instead of IP. This shows we could replace IP with DNE.

3. Another alternative to adding IP into our proof system would have been to add a pair of rules, one for the Explosion Principle:

i	\perp ϕ	EX i
-----	-------------------	--------

and another that lets us write down any instance of the Law of Excluded Middle:

$\phi \vee \neg\phi$	LEM
----------------------	-----

Your challenge: prove $\neg A \rightarrow \perp \vdash A$ using EX and LEM (together with our other rules) instead of IP. This shows we could replace IP with this pair of rules. Note: to do this problem, you'll have to know how to work with disjunction rules, so read the next section on the $\vee E$ rule before attempting it.

4.9 Disjunction Rules

Our last connective to deal with is disjunction. The $\vee I$ rule is pretty straightforward. Suppose Alice is a logician. Then certainly Alice is either a logician or a chemist. After all, to say that Alice is either a logician or a chemist is to say something weaker than to say that Alice is a logician. In fact, we can weaken the claim however we like. Suppose Alice is a logician. It follows that Alice is *either* a logician *or* a kumquat. Equally, it follows that *either* Alice is a logician *or* the earth is flat. Many of these are strange inferences to draw. But there is nothing *logically* wrong with them: in each case the conclusion has to be true if the premise that Alice is a logician is true.

Our disjunction introduction rules implement this idea of arbitrarily weakening a claim:

m	φ	
	$\varphi \vee \psi$	$\vee I m$

and

m	φ	
	$\psi \vee \varphi$	$\vee I m$

Notice that ψ can be *any* sentence. So the following is a perfectly kosher use of $\vee I$:

1	M	
2	$M \vee ([(A \leftrightarrow B) \rightarrow (C \wedge D)] \leftrightarrow [E \wedge F])$	$\vee I 1$

The disjunction elimination rule is slightly trickier. Suppose you know that Alice is either a logician or a chemist. What can you conclude? Not that Alice is a logician; she might be a chemist instead. And equally, not that Alice is a chemist; she might be a logician instead. Disjunctive premises, just by themselves, are hard to work with!

But suppose that we could somehow show both of the following: first, that Alice's being a logician implies that she has a PhD; second, that Alice's being a chemist also implies that she has a PhD. Then if we know that Alice is either a logician or a chemist, we know that, whichever she happens to be, she has a PhD. Our disjunction elimination rule $\vee E$ formalizes this insight:

m	$\varphi \vee \psi$	
i	φ	
	\vdots	
j	χ	
k	ψ	
	\vdots	
l	χ	
	χ	$\vee E m, i-j, k-l$

This is a bit more complicated than our previous rules, but the idea is fairly simple. Suppose we have some disjunction $\varphi \vee \psi$ and our goal is to prove some claim χ . If we can give two subproofs, one showing that our goal χ follows from the assumption that φ holds, and another showing that the *same* conclusion χ *also* follows from the assumption that ψ holds, then we can infer χ itself by $\vee E$. This rule formally implements a proof strategy called *argument by cases*: the disjunction $\varphi \vee \psi$ tells us that one of two cases obtains, either φ

holds or ψ does; if it can now be shown that χ must hold in *either case*, then we can conclude that χ holds on the basis of the original disjunction.

Notice that the citation for $\vee E$ is quite complex. We have to cite *three* things: the line number of the original disjunction, and the two subproofs. As usual, there can be as many lines as you like between i and j , and as many lines as you like between k and l . Moreover, the subproofs and the disjunction can come in any order, and do not have to be adjacent.

Some examples will help illustrate the rule. Consider this problem:

$$(P \wedge Q) \vee (P \wedge R) \vdash P$$

The premise tells us that either ' $(P \wedge Q)$ ' holds or ' $(P \wedge R)$ ' holds. But in either case, ' P ' must hold, so ' P ' follows from our disjunctive premise. Here's how this looks as a proof:

1	$(P \wedge Q) \vee (P \wedge R)$	Premise
2	$P \wedge Q$	Assumption
3	P	$\wedge E$ 2
4	$P \wedge R$	Assumption
5	P	$\wedge E$ 4
6	P	$\vee E$ 1, 2–3, 4–5

$\vee E$ is similar to rules like $\rightarrow I$, $\leftrightarrow I$, and $\neg I$, in that it requires temporary assumptions and subproofs. But it differs in that, unlike the others, it's an *elimination* rule! This has an important consequence for the kind of strategy to use in relation to $\vee E$. In the case of rules like $\rightarrow I$ and $\neg I$, we always reasoned "from the bottom up," that is, we reasoned *backward* from our goal. E.g. if our goal is to prove a conditional $\phi \rightarrow \psi$, we know to use $\rightarrow I$ as the overall strategy, which means assuming ϕ and then proving ψ .

The $\vee E$ rule is the exception to working backwards like this. Here, we have to reason "from the top down." That is, if you have a disjunction $\phi \vee \psi$ as a premise (or as an assumption, or as something you can easily derive from your premises and assumptions), it's usually a good idea to prove your goal, *whatever it may be*, using $\vee E$ as your overall strategy. This means you have to open up two subproofs, and prove your goal χ first from the left-disjunct ϕ and then from the right-disjunct ψ .

Here's a more complex example, demonstrating one of the Distributive Laws:

$$P \vee (Q \wedge R) \vdash (P \vee Q) \wedge (P \vee R)$$

1	$P \vee (Q \wedge R)$	Premise
2	P	Assumption
3	$(P \vee Q)$	$\vee I$ 2
4	$(P \vee R)$	$\vee I$ 2
5	$(P \vee Q) \wedge (P \vee R)$	$\wedge I$ 3,4
6	$(Q \wedge R)$	Assumption
7	Q	$\wedge E$ 6
8	R	$\wedge E$ 6
9	$(P \vee Q)$	$\vee I$ 7
10	$(P \vee R)$	$\vee I$ 8
11	$(P \vee Q) \wedge (P \vee R)$	$\wedge I$ 9,10
12	$(P \vee Q) \wedge (P \vee R)$	$\vee E$ 1,2-5,6-11

To identify my overall strategy in this case, I did *not* look at my goal $(P \vee Q) \wedge (P \vee R)$ and think about what introduction rule I might use to prove it (e.g. $\wedge I$). Rather, I noticed that my premise $P \vee (Q \wedge R)$ is a disjunction. And whenever I have a disjunction as a premise like this (or as an assumption, or as something easily derivable from my premises and assumptions), I will prove my current goal, whatever it might be, using $\vee E$ as the overall strategy. So in this case, that means first assuming the left disjunction P of my premise, and proving my conclusion from that (lines 2–5), and then assuming the right disjunct $(Q \wedge R)$ of my premise, and proving my conclusion from that (lines 6–11).

■ Exercises 4.9

A. Prove the following:

1. $A \rightarrow B \vdash A \rightarrow (C \vee B)$
2. $(B \vee A) \rightarrow C \vdash A \rightarrow C$
3. $(A \wedge B) \vee (A \wedge C) \vdash A$
4. $A \vee B \vdash B \vee A$
5. $A \vee B \vdash (A \rightarrow B) \rightarrow B$
6. $A \vee (B \wedge C) \vdash (A \vee B) \wedge (A \vee C)$
7. $A \vee (B \vee C) \vdash (A \vee B) \vee C$
8. $S \leftrightarrow T \vdash S \leftrightarrow (T \vee S)$
9. $(C \wedge D) \vee E \vdash E \vee D$
10. $A \rightarrow C \vdash (A \vee (B \wedge C)) \rightarrow C$
11. $A \wedge (B \vee C) \vdash (A \wedge B) \vee (A \wedge C)$
12. $(A \wedge B) \vee (A \wedge C) \vdash A \wedge (B \vee C)$
13. $A \vee B \vdash (A \rightarrow B) \rightarrow B$
14. $(Z \wedge K) \vee (K \wedge M), K \rightarrow D \vdash D$

4.10 Proof strategies

There is no simple recipe for proofs, and there is no substitute for practice. But here are some strategies to keep in mind.

Work backwards from your goal. Your ultimate goal is arrive at the conclusion. Look at the conclusion and ask what the introduction rule is for its main logical operator. This gives you an overall strategy, and tells you what assumptions (if any) to make, and what your new goal is. Now ask what the main operator of that new goal is, thereby identifying a strategy to prove it, and so on.

For example: if your conclusion is a conditional $\phi \rightarrow \psi$, you should plan to use the \rightarrow I as your strategy. This requires opening a subproof in which you assume ϕ and then set ψ as your new goal. Or if your goal is to prove $\phi \leftrightarrow \psi$, use \leftrightarrow I as your overall strategy (which involves giving two subproof, one in each direction of the arrow). Or if your conclusion is a negated sentence $\neg\phi$, plan to use \neg I as your strategy (which means assuming ϕ and proving a contradiction).

The Exception: \vee E The one important exception to the strategy of working backwards involves the \vee E rule: if you see a disjunction $\phi \vee \psi$ among your premises or assumptions (or as something that you can easily derive from them), it's almost always a good idea to prove your current goal (whatever it may be) by setting up an \vee E proof. That means opening up two subproofs, one that begins with an assumption of ϕ and another that begins with ψ . Inside each of them, you now have to prove whatever your current goal is.

Try an indirect proof. If you can't find any way to prove your goal ϕ directly, try an indirect proof: assume $\neg\phi$ and try to derive a contradiction. If you succeed, you can infer your goal ϕ by IP. This strategy should only be used as a last resort, however! Indirect proofs are often longer and more complicated than direct proofs.

Persist. Try different things. If one approach fails, try something else. I'll never ask you to prove something that cannot be proven.

Don't make random assumptions. Finally, never make an assumption unless you have a strategy in mind for discharging it (i.e. for ultimately closing the subproof that the assumptions opens). That means you should only make an assumption if your strategy is to use one of the following discharge rules: \rightarrow I, \leftrightarrow I, \neg I, IP, or \vee E.

Let's look at one more example. Take the following English argument:

If Guatemala is in Canada, then it is in North America. So if Guatemala is not in North America, it also isn't in Canada.

This is intuitively valid, so we should be able to give a proof of it. We can symbolize the argument as: $C \rightarrow A \therefore \neg A \rightarrow \neg C$. Our goal here is to prove a conditional, ' $\neg A \rightarrow \neg C$ '. So we use \rightarrow I as our strategy, meaning we assume ' $\neg A$ ' and set ourselves the new goal of proving ' $\neg C$ ' from that assumption. And now, since ' $\neg C$ ' has a negation as its main operator, we

use $\neg I$ as our strategy, meaning we assume ‘ C ’ and prove a contradiction from that. Notice how I reasoned “backwards” from the conclusion in order to discover this strategy. Written out, the proof looks like this:

1	$C \rightarrow A$	Premise
2	$\neg A$	Assumption (for $\rightarrow I$)
3	C	Assumption (for $\neg I$)
4	A	$\rightarrow E$ 1, 3
5	\perp	$\perp I$ 2, 4
6	$\neg C$	$\neg I$ 3–5
7	$\neg A \rightarrow \neg C$	$\rightarrow I$ 2–6

■ Exercises 4.10

A. The following three proofs are missing their rule citations. Write them in. Additionally, write down the sequent (i.e. single-turnstile \vdash statement) that each proof demonstrates.

1	$P \wedge S$
2	$S \rightarrow R$
3	P
4	S
5	R
6	$R \vee E$

1	$A \rightarrow D$
2	$A \wedge B$
3	A
4	D
5	$D \vee E$
6	$(A \wedge B) \rightarrow (D \vee E)$

1	$\neg L \rightarrow (J \vee L)$
2	$\neg L$
3	$J \vee L$
4	J
5	$J \wedge J$
6	J
7	L
8	\perp
9	J
10	J

B. Prove each of the following:

1. $J \rightarrow \neg J \vdash \neg J$
2. $Q \rightarrow (Q \wedge \neg Q) \vdash \neg Q$
3. $P \vee Q, \neg P \vdash Q$

4. $\neg R \vee (P \rightarrow Q) \vdash (R \wedge P) \rightarrow Q$
5. $\neg F \rightarrow G, F \rightarrow H \vdash G \vee H$
6. $D \vdash \neg\neg D$
7. $P \wedge (Q \vee R), P \rightarrow \neg R \vdash Q \vee E$
8. $\neg C \vee (A \rightarrow B) \vdash (C \wedge A) \rightarrow B$
9. $C \rightarrow (E \wedge G), \neg C \rightarrow G \vdash G$
10. $M \wedge (\neg N \rightarrow \neg M) \vdash (N \wedge M) \vee \neg M$
11. $(W \vee X) \vee (Y \vee Z), X \rightarrow Y, \neg Z \vdash W \vee Y$

C. Show that the following are provably equivalent:

1. $\neg(P \wedge Q) \dashv\vdash \neg P \vee \neg Q$
2. $\neg(P \vee Q) \dashv\vdash \neg P \wedge \neg Q$
3. $P \vee Q \dashv\vdash \neg(\neg P \wedge \neg Q)$
4. $P \rightarrow Q \dashv\vdash \neg Q \rightarrow \neg P$
5. $P \rightarrow Q \dashv\vdash \neg P \vee Q$
6. $\neg(P \rightarrow Q) \dashv\vdash P \wedge \neg Q$
7. $P \leftrightarrow \neg Q \dashv\vdash \neg(P \leftrightarrow Q)$

D. Prove the following theorems:

1. $\vdash (P \rightarrow Q) \vee (Q \rightarrow P)$
2. $\vdash A \vee \neg A$
3. $\vdash ((P \rightarrow Q) \rightarrow P) \rightarrow P$
4. $\vdash \neg A \rightarrow (A \rightarrow B)$
5. $\vdash J \leftrightarrow [J \vee (L \wedge \neg L)]$

4.11 Derived Rules

We have provided introduction and elimination rules for each of our five connectives. Together with IP, this gives us a complete proof system: every valid argument can be proven using just these few basic rules! In this section, we're going to introduce some additional rules to shorten our proofs and make our proof system easier to work with. It's important to note at the outset that these additional rules are not necessary. They represent a *conservative* extension of our proof system: anything proven using these new rules can also be proven using just our basic set of rules.

To illustrate the motivation for additional rules, consider the following argument:

Alice is either a logician or a chemist. She is not a chemist. So she is a logician.

This involves a very natural form of inference called *Disjunctive Syllogism*. We could symbolize the argument as $L \vee C, \neg C \therefore L$, and we then give a natural deduction proof using $\vee E$ to show that it is valid.

But now consider this: by giving a proof of ' L ' from ' $L \vee C$ ' and ' $\neg C$ ', we have implicitly shown that given *any* sentences of the form $\phi \vee \psi$ and $\neg\psi$, it is possible to prove ϕ . If we substitute the metavariables ϕ and ψ for the sentences ' L ' and ' C ' in our proof, we get a

PROOF TEMPLATE for the disjunctive syllogism form of inference:

m	$(\phi \vee \psi)$	
n	$\neg\phi$	
k_0	ϕ	
k_1	$\neg\psi$	
k_2	\perp	$\perp I\ n, k_0$
k_3	ψ	IP k_0-k_3
k_4	ψ	
k_5	$\psi \wedge \psi$	$\wedge I\ k_4, k_4$
k_6	ψ	$\wedge E\ k_5$
k_7	ψ	$\vee E\ m, k-k_3, k_4-k_6$

Now, if at any time, in the context of any proof whatsoever, we need to prove some sentence ϕ from two sentences of the form $\phi \vee \psi$ and $\neg\psi$, we can simply “slot in” an instance of the above proof template. In other words, once we’ve proven one instance of disjunctive syllogism using our basic rules, we can use that as a template to prove disjunctive syllogism again in the context of any other proof.

Given this, we might as well just introduce a DERIVED RULE into our proof system that lets us skip the actual proof and make the disjunctive syllogism inference *directly*:

m	$\phi \vee \psi$	
n	$\neg\phi$	
	ψ	DS m, n

DS is a DERIVED RULE in the sense that it can be shown to hold using only the *primitive* rules of our system. You can think of derived rules like promissory notes: “I am here justifying my inference by writing ‘DS’, but I promise that, if you asked for it, I could slot in a series of steps using only the primitive rules of our natural deduction system.” Derived rules shorten our proofs, but add no power into our proof system: any proof that appeals to a derived rule could be expanded into one that only appeals to primitive rules.

In fact, we implicitly already added a derived rule to our system in §4.6 when we introduced Reiteration. Reiteration is just a shortcut to let us skip the $\wedge I$ -plus- $\wedge E$ trick that I used in steps k_5 and k_6 in the above proof template.³ There are many further useful derived rules we can add to our proof system. For example, consider the following argument:

If Alice is a chemist, then she has a PhD. Alice doesn’t have a PhD. So she isn’t a chemist.

³Indeed, in this particular case we could have avoided using the $\wedge I$ + $\wedge E$ trick, and shortened our proof template, by treating line $k+4$ as a whole subproof that begins with ψ and ends with ψ (see the end of §4.6).

This inference pattern is called *modus tollens*, and we can introduce a derived rule for it:

m	$\varphi \rightarrow \psi$	
n	$\neg\psi$	
	$\neg\varphi$	MT m, n

Again, this adds no power to our system because it is simply a shortcut for a series of steps involving only primitive rules, as illustrated by the following proof template:

m	$\varphi \rightarrow \psi$	
n	$\neg\psi$	
k_0	φ	
k_1	ψ	$\rightarrow E\ m, k_0$
k_2	\perp	$\perp I\ k_1, n$
k_3	$\neg\varphi$	$\neg I\ k_0-k_2$

In §4.10 we gave a seven step proof showing that $C \rightarrow A \therefore \neg A \rightarrow \neg C$ is valid. Using our derived rule MT, we can now shorten this to just four steps:

1	$C \rightarrow A$	Premise
2	$\neg A$	Assumption
3	$\neg C$	MT 1, 2
4	$\neg A \rightarrow \neg C$	$\rightarrow I\ 2-3$

Here is a complete list of the derived rules that you'll be able to use:

Sequent	Derived Rule
$\phi \rightarrow \psi, \neg\psi \vdash \neg\phi$	MT
$\phi \vee \psi, \neg\psi \vdash \phi$	DS
$\phi \vee \psi, \neg\phi \vdash \psi$	DS
$\phi \vdash \psi \rightarrow \phi$	PMI
$\neg\phi \vdash \phi \rightarrow \psi$	PMI
$\phi \rightarrow \psi \dashv\vdash \neg\phi \vee \psi$	Imp
$\neg(\phi \rightarrow \psi) \dashv\vdash \phi \wedge \neg\psi$	NegImp
$\neg(\phi \wedge \psi) \dashv\vdash \neg\phi \vee \neg\psi$	DeM
$\neg(\phi \vee \psi) \dashv\vdash \neg\phi \wedge \neg\psi$	DeM
$\phi \dashv\vdash \neg\neg\phi$	DN
$(\phi \# \psi) \dashv\vdash (\neg\neg\phi \# \neg\neg\psi) \dashv\vdash (\neg\neg\phi \# \psi) \dashv\vdash (\phi \# \neg\neg\psi)$	SDN
$\neg(\phi \# \psi) \dashv\vdash \neg(\neg\neg\phi \# \neg\neg\psi) \dashv\vdash \neg(\neg\neg\phi \# \psi) \dashv\vdash \neg(\phi \# \neg\neg\psi)$	SDN
$\phi @ \psi \vdash \psi @ \phi$	Com
$\perp \vdash \phi$	EX
$\vdash \phi \vee \neg\phi$	LEM

(where $\#$ in SDN can be any binary connective, and $@$ Com can be any of the three commutative connectives $\vee, \wedge, \leftrightarrow$). The way understand this list of derived rules is as follows:

- ▷ For any sequent $\phi_1 \dots \phi_n \vdash \psi$ matching one from the above list, if $\phi_1 \dots \phi_n$ occur on some earlier lines $j_1 \dots j_n$ in your proof (none of them inside a closed subproof), then you may directly infer ψ and justify it by citing the name of the relevant derived rule followed by the numbers of lines $j_1 \dots j_n$.

For example, if you have the sentence ' P ' on line m in your proof, then you are allowed to directly infer ' $Q \rightarrow P$ ' with the justification 'PMI m ' (for "Paradox of Material Implication"). Or if you have a sentence ' $\neg K \vee \neg L$ ' on line m in your proof, you may directly infer ' $\neg(K \wedge L)$ ' with the justification 'DeM m ' (for "DeMorgan's Law"), or ' $K \rightarrow \neg L$ ' with the justification 'Imp m ', or ' $\neg L \vee \neg K$ ' with the justification 'Com m '. Notice that these latter three are examples of derived rules that hold in *both* directions, so you could also start with ' $\neg(K \wedge L)$ ', or ' $K \rightarrow \neg L$ ', or ' $\neg L \vee \neg K$ ' and derive ' $\neg K \vee \neg L$ '.

The second-to-last derived rule, EX, is the Explosion Principle: it lets you infer any sentence whatsoever from a contradiction. The last derived rule, LEM, lets you write down any instance of the Law of Excluded Middle $\phi \vee \neg\phi$ at any point in your proof with the justification 'LEM'. Here no line number needs to be cited because you're introducing a theorem.

For another example of these rules in action, consider the following theorem:

$$\vdash (P \rightarrow Q) \vee (Q \rightarrow P)$$

This was one of the exercises in §4.10. Proving this using only basic rules is quite difficult, as you will have noticed if you tried that exercise. With derived rules, we can give a much quicker and more intuitive proof of this theorem, by starting out with $P \vee \neg P$ as an instance of the Law of Excluded Middle, and then pursuing an $\vee E$ strategy from there:

1	$P \vee \neg P$	LEM
2	P	Assumption (for $\vee E$)
3	$Q \rightarrow P$	PMI 2
4	$(P \rightarrow Q) \vee (Q \rightarrow P)$	$\vee I$ 3
5	$\neg P$	Assumption (for $\vee E$)
6	$P \rightarrow Q$	PMI 5
7	$(P \rightarrow Q) \vee (Q \rightarrow P)$	$\vee I$ 6
8	$(P \rightarrow Q) \vee (Q \rightarrow P)$	$\vee E$ 1,2-4,5-7

As an exercise, you might try to re-write this proof by “slotting in” a subproof involving only primitive rules wherever the above proof appeals to a derived rule. This involves showing how LEM, and the two versions of PMI, can be proven using only primitive rules of our system.

■ Exercises 4.11

A. The following proofs are missing their citations (rule and line numbers). Add them wherever they are required:

1	$W \rightarrow \neg B$
2	$A \wedge W$
3	$B \vee (J \wedge K)$
4	W
5	$\neg B$
6	$J \wedge K$
7	K

1	$L \leftrightarrow \neg O$
2	$L \vee \neg O$
3	$\neg L$
4	$\neg O$
5	L
6	\perp
7	L

1	$Z \rightarrow (C \wedge \neg N)$
2	$\neg Z \rightarrow (N \wedge \neg C)$
3	$\neg(N \vee C)$
4	$\neg N \wedge \neg C$
5	$\neg N$
6	$\neg N \vee \neg \neg C$
7	$\neg(N \wedge \neg C)$
8	$\neg \neg Z$
9	Z
10	$\neg C$
11	$\neg C \vee \neg \neg N$
12	$\neg(C \wedge \neg N)$
13	$\neg Z$
14	\perp
15	$N \vee C$

B. Prove the following using derived rules:

1. $(A \vee B) \rightarrow C, \neg C \vdash \neg A$
2. $E \vee F, F \vee G, \neg F \vdash E \wedge G$
3. $\neg(A \rightarrow (B \vee \neg C)) \vdash (B \vee C) \rightarrow A$
4. $(Q \rightarrow P) \rightarrow R, \neg Q \vee \neg S \vdash S \rightarrow \neg(\neg R \vee \neg S)$
5. $M \vee (N \rightarrow M) \vdash \neg M \rightarrow \neg N$
6. $A \rightarrow (B \vee C) \vdash (A \rightarrow B) \vee (A \rightarrow C)$
7. $(M \vee N) \wedge (O \vee P), N \rightarrow P, \neg P \vdash M \wedge O$
8. $(B \wedge C) \vee \neg(A \rightarrow \neg D), \neg C \vdash A \wedge D$
9. $(X \wedge Y) \vee (X \wedge Z), \neg(X \wedge D), D \vee M \vdash M$
10. $\vdash (A \rightarrow B) \vee (B \rightarrow A)$

If you want more practice, you can of course also re-do any of the earlier proofs in this chapter using derived rules.

C. Provide proof templates (like those I provided for DS and MT) that justify the addition of the De Morgan rules, the Imp and NegImp rules, and LEM as derived rules. If you don't want to bother with metavariables, you can just prove instances of the corresponding sequents, but in any case, be sure to only use primitive rules in your proofs.

II

First-order logic

Symbolization in FOL

5

Consider the following argument, which is obviously valid:

Willard is a logician.
All logicians wear funny hats.
 \therefore Willard wears a funny hat.

But how would we symbolize it in TFL? We could offer the following symbolization key:

L : Willard is a logician.
 A : All logicians wear funny hats.
 F : Willard wears a funny hat.

and then symbolize the argument as:

$$L, A \therefore F$$

But the truth-table test will indicate that this is *invalid*. What has gone wrong?

The problem is not that we have made a mistake while symbolizing the argument. This is the best symbolization we can give *in TFL*. The problem lies with TFL itself! This argument is not valid in virtue of its *truth-functional* structure, but rather in virtue of its *subsential* structure. For example, ‘All logicians wear funny hats’ establishes a certain relationship between being a logician and hat-wearing. But in TFL, the best we can do is symbolize it as an atomic sentence A . Because TFL doesn’t let us represent any subsential structure, we lose the connection between Willard’s being a logician and Willard’s wearing a hat in the TFL symbolization..

To symbolize arguments like the preceding one, we will have to develop a new logical language which will allow us to *split the atom*. That is to say: it will let us represent logically significant structure *inside* atomic sentences. This will be the language of *first-order logic*, or *FOL*.

The details of FOL will be explained throughout this chapter, but here is the basic idea for splitting the atom. A sentence like ‘Willard is a logician’ is internally composed of a name, ‘Willard’, and a predicate, ‘___ is a logician’. In FOL, we’ll use lowercase letters to symbolize names, and uppercase letters to symbolize predicates. So we might use ‘ a ’ to symbolize the name ‘Willard’, and ‘ L ’ to symbolize the predicate ‘___ is a logician’. The whole sentence can then be symbolized as ‘ La ’, thereby representing the fact that this atomic sentence is internally composed of a name and a predicate.

A sentence like ‘All logicians wear funny hats’ involves two predicates: ‘___ is a logician’ and ‘___ wears funny hats’. It also involves the word ‘all’, which relates the two

predicates. This is called a *quantifier*, because it tells us something about “quantities” — in this case that *all* individuals who are logicians wear funny hats, rather than just some of them. FOL will have two quantifiers, ‘ \forall ’ and ‘ \exists ’. ‘ \exists ’ will roughly convey ‘There is at least one thing such that ...’ and ‘ \forall ’ will convey ‘Every thing is such that ...’.

So FOL has three new ingredients: names, predicates, and quantifiers. And we’ll be able to use these ingredients to represent the internal structure of atomic sentences. That is the general idea. But FOL is significantly more complex than TFL, so we’ll build up slowly.

5.1 Names and Predicates

In English, a *singular term* is a word or phrase that refers to a *specific* person, place, or thing. The word ‘dog’ is not a singular term, because there are many dogs. But ‘Fido’ is a singular term, because it refers to a specific dog. Likewise, the phrase ‘Philip’s dog Fido’ is a singular term, because it also refers to that specific terrier.

Proper names are a particularly important kind of singular term. These are expressions that, unlike e.g. ‘Philip’s dog Fido’, pick out individuals without describing them. The name ‘Emerson’ is a proper name, and it alone does not tell you anything about Emerson. Of course, some names are traditionally given to boys or girls. If ‘Hilary’ is used as a singular term, you might guess that it refers to a woman. But then again you might be wrong. Indeed, the name does not necessarily mean that the person referred to is even a person: Hilary might be a giraffe, for all you could tell just from the name.

In FOL, our NAMES are lowercase letters ‘*a*’ through to ‘*t*’. We can also add subscripts if we want to use some letter more than once. So the following are all names in FOL:

$$a, b, c, \dots, s, t, a_1, f_{32}, j_{390}, m_{12}$$

These should be thought of along the lines of proper names in English. But with one difference. ‘Syracuse’ is a proper name, but it is the name of both a city in New York State and of a city in Italy. And there are over thirty towns in the US that have the name ‘Springfield’. We live with this kind of ambiguity in English, allowing context to determine that ‘Syracuse’ is being used to refer to a city in the US rather than to one in Italy. In FOL, we do not tolerate any such ambiguity. Each name must refer to *exactly* one thing. (However, two different names may refer to the same thing, like ‘Mark Twain’ and ‘Sam Clemens’ refer to the same person.)

As with TFL, we’ll provide symbolization keys. These indicate, temporarily, what object a name will refer to. So we might offer the following:

e: Elsa
g: Gregor
m: Marybeth

The second ingredient in FOL are PREDICATES. The simplest predicates express properties of individuals. Here are some examples of English predicates:

____ is a dog
 ____ is a member of Monty Python
 A piano fell on ____

In general, you can think about predicates as things which combine with names and other singular terms to make sentences. Conversely, you can start with sentences and make predicates out of them by removing terms. Consider the sentence, ‘Vinnie borrowed the family car from Nunzio.’ By removing a singular term, we can obtain three different predicates:

____ borrowed the family car from Nunzio
 Vinnie borrowed ____ from Nunzio
 Vinnie borrowed the family car from ____

FOL predicates are capital letters *A* through *Z*, with or without subscripts. We might write a symbolization key for predicates like this:

A: ____ is angry
H: ____ is happy

If we combine our two symbolization keys, we can start to symbolize some English sentences that use these names and predicates in combination. For example:

- (1) Elsa is angry.
- (2) Gregor and Marybeth are angry.
- (3) If Elsa is angry, then so are Gregor and Marybeth.

Sentence (1) is just symbolized as ‘*Ae*’. Sentence (2) is a conjunction of two simpler sentences. The simple sentences can be symbolized as ‘*Ag*’ and ‘*Am*’. Then we help ourselves to our resources from TFL, and symbolize the entire sentence as ‘(*Ag* \wedge *Am*)’.

This illustrates an important point: FOL has all of the truth-functional connectives of TFL! Lastly, sentence (3) is a conditional, whose antecedent is sentence (1) and whose consequent is sentence (2). So we can symbolize it as ‘*Ae* \rightarrow (*Ag* \wedge *Am*)’.

We can also use TFL connectives to symbolize sentences that involve COMPOUND PREDICATES, that is, predicates formed out of simpler ones. Consider the following sentence:

- (4) Herbie is a white car

It involves the compound predicate ‘____ is a white car’. But we can paraphrase the sentence as a conjunction involving simpler predicates: ‘Herbie is white and Herbie is a car’. Using the following symbolization key:

W: ____ is white
C: ____ is a car
h: Herbie

we can thus symbolize (4) as ‘(*Wh* \wedge *Ch*)’.

In this case, the compound predicate was formed out of an adjective and a noun. But there are other ways to do this too:

- (5) Herbie is a car from Germany.
- (6) Herbie is a car from Germany that is fast.

Sentence (5) involves a compound predicate formed from a noun and the prepositional phrase ‘from Germany’, and (6) involves a compound predicate formed from a noun, a prepositional phrase, and the relative clause ‘that is fast’. These can also be symbolized as conjunctions of simple predications. Using ‘ G ’ for ‘____ is from Germany’ and ‘ F ’ for ‘____ is fast’, (5) can be symbolized as ‘ $Ch \wedge Gh$ ’ and (6) as ‘ $(Ch \wedge Gh) \wedge Fh$ ’.

One does, occasionally have to be careful when symbolizing compound predicates. Suppose I have a violin that I’ve named ‘Lucy’, and now consider the sentence ‘Lucy is a fake Stradivarius’. You might think that using ‘ c ’ for ‘Lucy’, ‘ F ’ for ‘____ is fake’ and ‘ S ’ for ‘____ is a Stradivarius’, we can symbolize this as ‘ $Fc \wedge Sc$ ’. But that wouldn’t be right: ‘ $Fc \wedge Sc$ ’ entails ‘ Sc ’, but our English sentence does not entail ‘Lucy is a Stradivarius’! The word ‘fake’ is what’s called a *non-intersective* adjective; in a case like this, we’d have to use a single FOL predicate, say ‘ S ’, for the whole English predicate ‘____ is a fake Stradivarius’.

■ Exercises 5.1

A. Symbolize the following in FOL:

1. Ada was both a mathematician and a computer scientist.
2. Ada was a mathematician from either England or Wales.
3. Susan will attend the party only if neither Tom nor Ella does.
4. Jen is a talented composer, and Fritz is also a composer, but not talented.
5. Susan and Tom will not both attend the party if Ella does.
6. If Ella attends the party, Susan and Tom both will not attend.

5.2 Quantifiers and Quantifier Scope

Next up are quantifiers. Consider these sentences:

- (7) Everyone is happy.
- (8) Someone is angry.

Sentence (7) superficially looks like it has the same kind of structure as something like ‘Elsa is happy’. So you might be tempted to symbolize (7) as ‘ He ’, with the explanation that ‘ e ’ is to symbolize ‘everyone’. But that would be a serious mistake.

‘Everyone’ is not a proper name — it doesn’t pick out any particular individual — and so it should not be symbolized using a name like ‘ e ’ in FOL. The word ‘everyone’ is rather a quantifier. Logically, quantifiers behave very differently from names. For example, whereas ‘Either Elsa is happy or Elsa is not happy’ is a necessary truth, ‘Either everyone is happy or everyone is not happy’ is not a necessary truth. In fact, it’s presumably false: some people are happy and others are not.

To express claims about every individual in a set, we’ll use the FOL symbol ‘ \forall ’. This is called the UNIVERSAL QUANTIFIER. A quantifier in FOL always has to be followed by a variable. FOL variables are lowercase letters ‘ u ’ through ‘ z ’, with or without subscripts. So we could symbolize sentence (7) as ‘ $\forall x Hx$ ’. The variable ‘ x ’ is serving as a kind of placeholder. The expression ‘ $\forall x$ ’ intuitively means that you can pick anyone and put them

in as ‘ x ’. The subsequent ‘ Hx ’ indicates, of that thing you picked out, that it is happy. So ‘ $\forall xHx$ ’ can be read as saying “every individual x is such that: x is happy.”

There is no special reason to use ‘ x ’ rather than some other variable here. The sentences ‘ $\forall xHx$ ’, ‘ $\forall yHy$ ’, ‘ $\forall zHz$ ’, and ‘ $\forall x_5Hx_5$ ’ use different variables, but they are all logically equivalent, and any of them could be used to symbolize (7).

To symbolize sentence (8), we introduce another new symbol: the EXISTENTIAL QUANTIFIER, ‘ \exists ’. Like the universal quantifier, the existential quantifier requires a variable. Sentence (8) can be symbolized by ‘ $\exists xAx$ ’. You can read ‘ $\exists xAx$ ’ as saying “there is some individual x such that: x is angry”. Again, the variable is just a kind of placeholder; we could just as well have symbolized sentence (8) as ‘ $\exists zAz$ ’, ‘ $\exists w_{256}Aw_{256}$ ’, or whatever.

In FOL, quantifiers always range over a DOMAIN of objects. So really, ‘ $\forall xHx$ ’ says that every object x *in the domain* is such that x is happy, and ‘ $\exists xAx$ ’ says that some object x *in the domain* is such that x is angry. So whenever we symbolize something with a quantifier, we should also specify a domain of objects in our symbolization key. English quantifiers like ‘everyone’ and ‘someone’ are used to talk about people, so in this case we can specify our domain as consisting of people. We’ll look more closely at picking domains later on.

Next, consider the following two sentences:

(9) Someone is a logician and someone is an architect.

(10) Someone is a logician and an architect.

These sentences clearly mean different things, and must therefore receive different symbolizations. Sentence (9) says that there is someone who is a logician, and also that there is someone (maybe a different person) who is an architect. So it is a conjunction of two existential sentences. Using the following symbolization key:

domain: people

L : ____ is a logician

A : ____ is an architect

sentence (9) can be symbolized as ‘ $\exists xLx \wedge \exists xAx$ ’.

Notice that we could equally well have symbolized it as ‘ $\exists xLx \wedge \exists yAy$ ’ or ‘ $\exists zLz \wedge \exists vAv$ ’. As we observed earlier, sentences like ‘ $\exists xAx$ ’ and ‘ $\exists yAy$ ’ that differ only in which variable gets used are equivalent to each other, so ‘ $\exists xLx \wedge \exists xAx$ ’ is therefore equivalent to ‘ $\exists xLx \wedge \exists yAy$ ’. In particular, using the same variable in both conjuncts, as in ‘ $\exists xLx \wedge \exists xAx$ ’, does *not* mean that the person who is a logician is the same as the one who is an architect; and using different variables, as in ‘ $\exists xLx \wedge \exists yAy$ ’, does *not* mean it’s a different person. Both of these FOL sentences just say that there is someone who is a logician, and that there is someone (maybe the same, maybe different) who is an architect.

Sentence (10), on the other hand, says that there exists some *one* individual who is *both* a logician and an architect. It gets symbolized as ‘ $\exists x(Lx \wedge Ax)$ ’. You can read this as: there is some object x in the domain of people such that x is a logician and x is also an architect. The difference between ‘ $\exists x(Lx \wedge Ax)$ ’ and the earlier ‘ $\exists xLx \wedge \exists xAx$ ’ has to do with the SCOPE of the quantifiers.

The way scope works with quantifiers is very similar to it works with negation, which we discussed in §2.9. In the sentence ‘ $\neg P \wedge \neg Q$ ’, the first \neg has scope over just the first

conjunct, and the second \neg has scope over just the second conjunct. It is the conjunction \wedge that is the main logical operator. Similarly, in $\exists x Lx \wedge \exists x Ax$, the first quantifier has scope over just the first conjunct, and the second quantifier over just the second conjunct, with the conjunction \wedge functioning as the main logical operator.

On the other hand, in $\exists x(Lx \wedge Vx)$ the quantifier is the main logical operator, and has scope over the entire sentence, with the conjunction \wedge occurring inside the scope of the quantifier. This is similar to how in $\neg(P \wedge Q)$, the negation is the main logical operator, with the conjunction occurring inside its scope. We'll give a more precise definition of the notion of scope in relation to quantifiers in §5.10, but the analogy with negation should give you a working handle for now.

Closely connected to the notion of scope are the notions of FREE and BOUND variables. A quantifier is said to *bind* the variables that occur in its scope. So in $\exists x(Lx \wedge Vx)$, the quantifier ' $\exists x$ ' binds the variable in ' Lx ' and also the one in ' Vx ', since they both occur inside its scope. On the other hand, in $\exists x Lx \wedge \exists x Ax$, the first quantifier binds the variable in ' Lx ', since it occurs in its scope, but does not bind the variable in ' Ax ' — that variable is bound by the second quantifier, in whose scope it occurs.

Or compare the following two FOL sentence:

$$(11) \forall x(Hx \vee Wx)$$

$$(12) \forall x Hx \vee Wx$$

In (11), the quantifier ' $\forall x$ ' is the main operator, and it has scope over the entire sentence; it therefore binds both occurrences of the variable ' x '. On the other hand, in (12), the main operator is the disjunction, with the quantifier only taking scope over the left disjunct. So here, the quantifier binds the variable in ' Hx ', but not the one in ' Wx '. Variables like this, that are not bound by any quantifier, are said to be FREE VARIABLES. Formulas like (12) that contain free variables are called OPEN FORMULAS. Formulas like (11) where all the variables are bound are called CLOSED FORMULAS.

When symbolizing English statements, you should not have any free variables in your symbolization, i.e. your symbolizations should always be closed formulas. The closed formula (11) could be used to symbolize the English statement 'Everyone is either happy or wise'. The open formula (12), on the other hand, says something like: either everyone is happy, or x is wise. But who is x ? Since ' x ' isn't a name, it doesn't refer to any particular thing, and we can't determine whether it's true that x is wise. English statements always have truth values, whereas open formulas like (12) do not, so open formulas shouldn't be used to symbolize English statements.

We could turn (12) into a close formula either by adding parentheses as in (11), or by adding another quantifier to bind the variable in ' Wx ', as in $\forall x Hx \vee \forall x Wx$. The latter is a closed formula that could be used to symbolize the English statement 'Either everyone is happy or everyone is wise'. The moral is that parentheses are important! They indicate the scope of quantifiers (just like they indicate the scope of negation), and therefore indicate what variables a quantifier binds.

■ Exercises 5.2

A. Using a domain of people, H for '____ is happy', and W for '____ is wise', symbolize the following statements:

1. Someone is happy, and someone is not happy.
2. Everyone is both friendly and honest.
3. Everyone is either friendly or honest.
4. Everyone is friendly and someone is honest.
5. If everyone is friendly then Liz is happy.
6. Liz is not happy, but someone is.
7. Someone is neither honest nor friendly.
8. Someone is honest but not friendly.
9. If Liz is friendly but not happy, then someone is friendly but not happy.
10. Someone is happy, but not everyone is wise.

B. Which of the following are open formulas, and which are closed formulas? What variables does each quantifier bind? (Remember: letters *a*-*s* are names in FOL, *t*-*z* are variables.)

1. Wb
2. $\exists x Wx$
3. $\forall y Wx$
4. $\exists x Hx \wedge Wx$
5. $\exists x (Hx \wedge Wx)$
6. $\exists x Hx \wedge Wb$
7. $\forall x (Hx \rightarrow Wy)$
8. $\forall x (Hx \rightarrow Wy)$
9. $\forall x (Hx \wedge \exists y (Wy \wedge Hx))$
10. $\forall x Hx \wedge \exists y (Wy \wedge Hx)$

5.3 Common Quantifier Phrases and Domains

Consider these sentences:

- (13) Some dogs are poodles.
- (14) Every dog is a canine.

Let's use the following symbolization key:

Domain: animals

D : ____ is a dog

P : ____ is a poodle

C : ____ is a canine

Sentence (13) gets symbolized using an existential quantifier as ' $\exists x(Dx \wedge Px)$ '. You can read this as "there is some object x (in the domain of animals) such that x is a dog and x is a poodle", which does capture the intent of (13).

Sentence (14) gets symbolized using a universal quantifier. You might be tempted to symbolize it as ' $\forall x(Dx \wedge Cx)$ ', using a universal quantifier together with a conjunction, just as we used an existential quantifier together with a conjunction for (13). But that would be a mistake: ' $\forall x(Dx \wedge Cx)$ ' means "every object x (in the domain) is such that x is a dog and

x is a canine”, or more simply, “every object (in the domain) is both a dog and a canine.” That’s not at all what (14) says!

To see the route towards the correct symbolization, notice that (14) can be paraphrased as “for any object x in the domain, *if* x is a dog, *then* x is a canine.” So (14) gets symbolized using a universal quantifier together with a conditional, as ‘ $\forall x(Dx \rightarrow Cx)$ ’. As general guidelines, we have the following:

An English sentence that can be paraphrased as ‘Some F is G ’ can be symbolized as $\exists x(\mathcal{F}x \wedge \mathcal{G}x)$.

An English sentence that can be paraphrased as ‘Every F is G ’ can be symbolized as $\forall x(\mathcal{F}x \rightarrow \mathcal{G}x)$.

The same patterns apply to quantified sentences that involve compound predicates:

(15) Some tame dogs are poodles.

(16) Every wild dog is a canine.

Let’s use ‘ T ’ for ‘____ is tame’ and ‘ W ’ for ‘____ is wild’. Recall from §5.1 that in general, compound predicates get symbolized as conjunctions of the component predicates. So we can symbolize (15) as $\exists x((Tx \wedge Dx) \wedge Px)$, with the conjunction ‘ $(Tx \wedge Dx)$ ’ symbolizing the compound predicate ‘is a tame dog’. Similarly, (16) gets symbolized as ‘ $\forall x((Wx \wedge Dx) \rightarrow Cx)$ ’, with ‘ $(Wx \wedge Dx)$ ’ symbolizing the compound predicate ‘is a wild dog’. This can be read as: every object x (in the domain) is such that *if* x is wild and x is a dog, *then* x is a canine.

When symbolizing more complex sentences like these, it is useful to distinguish the RESTRICTOR PREDICATE from the MAIN PREDICATE. Intuitively, the restrictor predicate tells you what class of things the sentence says something about — e.g. pet dogs, or wild dogs, or wild dogs from Africa — and the main predicate then says what is true of some or all of them (e.g. that they’re poodles, or canines). The general pattern is that, in the case of an existential quantifier, the restrictor predicate and main predicate always get connected by a conjunction. And in the case of a universal quantifier, the restrictor predicate and main predicate always get connected by a conditional:

Some <u>tame dogs</u> are poodles Restrictor Main	Every <u>wild dog</u> is a canine Restrictor Main
$\exists x(\underbrace{(Tx \wedge Dx)}_{\text{Restrictor}} \wedge \underbrace{Px}_{\text{Main}})$	$\forall x(\underbrace{(Wx \wedge Dx)}_{\text{Restrictor}} \rightarrow \underbrace{Cx}_{\text{Main}})$

It doesn’t matter how complex the restrictor and main predicate get, the pattern is always the same. Using A for ‘____ is from Africa’ and M for ‘____ is a mammal’ we get:

Every <u>wild dog from Africa</u> is both a mammal and a canine Restrictor Main
$\forall x(\underbrace{(Wx \wedge (Dx \wedge Ax))}_{\text{Restrictor}} \rightarrow \underbrace{(Mx \wedge Cx)}_{\text{Main}})$

Again: since it's a universal sentence the restrictor predicate and main predicate get connected by a \rightarrow . If the quantifier had been 'Some' instead of 'Every', they would have been connected by a \wedge .

So far we've been pretty loose about picking domains. But in practice, picking a domain can be a delicate matter, and can affect what the correct symbolization of a sentence is. Take (14) from earlier, 'Every dog is a canine.' In a domain of animals, this gets symbolized as ' $\forall x(Dx \rightarrow Cx)$ '. But suppose we instead use a domain that consists of just dogs. In this domain, 'Every dog is a canine' can just be symbolized as ' $\forall xCx$ ': every object x in the domain (of dogs) is a canine. In other words, by restricting our *domain* to dogs, we no longer need the explicit restrictor *predicate* ' Dx ' in our symbolization.

You could in principle always avoid complex symbolizations by just restricting your domain appropriately. A sentence that standardly gets symbolized as $\forall x(\mathcal{F}x \rightarrow \mathcal{G}x)$ could just be symbolized $\forall x\mathcal{G}x$ by making the domain consist of whatever things the restrictor predicate $\mathcal{F}x$ is true of. And similarly, a sentence that standardly gets symbolized as $\exists x(\mathcal{F}x \wedge \mathcal{G}x)$ could just be symbolized as $\exists x\mathcal{G}x$ by picking a domain that consists of whatever things $\mathcal{F}x$ is true of.

However, this gain in convenience comes at a cost. If we pick a domain of dogs in our symbolization key for (14), then we can no longer use the same symbolization key to symbolize sentences that talk about things other than dogs. That could be a problem: 'Every dog is a canine' might appear as the first premise in an argument, whose second premise is 'No pelican is a canine', and whose conclusion is 'No pelican is a dog'. Since the second premise and the conclusion talk about pelicans, we can't symbolize them using a domain of just dogs. We'll instead have to go back to using a larger domain, like all animals, and revert to the more complex symbolization $\forall x(Dx \rightarrow Cx)$ for (14).

In practice, we're going to be fairly restrictive about what domains to use. In order to standardize symbolizations, you should always pick one of two domains in your symbolization key: either a domain that consists of things in general (dogs, people, other animals, plants, stars, numbers etc.), or a domain that consists of people. You can use the following guidelines to determine which domain to pick:

- ▷ If a sentence or argument only contains quantifiers restricted to people (like 'everyone', 'someone', 'anyone', 'no one'), and no standalone names of things other than people, you may (but aren't required to) use a domain consisting of just people.
- ▷ For any other kind of sentence or argument, use the domain of things.

So in the case of something like 'Every dog is a canine', you should use a domain of things and symbolize it as $\forall x(Dx \rightarrow Cx)$. For 'Every logician is wise', you'd also use a domain of things and symbolize it as ' $\forall x(Lx \rightarrow Wx)$ '. On the other hand, for something like 'everyone is happy', you may use a domain of just people and symbolize it as $\forall xHx$.

It's important to notice, though, that you don't *have* to use a domain of people to symbolize 'Everyone is happy.' This can be symbolized in a domain of things too, by making the restriction to people (carried by the 'one' in 'everyone') explicit with an additional predicate ' P ' for '____ is a person'. This would give us the symbolization ' $\forall x(Px \rightarrow Hx)$ '. In fact, using an extra predicate like this is what you would have to do if the sentence occurred as part of a larger sentence or argument that involves reference to things other than people.

Similarly, if you wanted to symbolize:

(17) Everyone in Chicago is happy.

in a domain of people, you would symbolize it as $\forall x(Cx \rightarrow Hx)$, using C for ‘____ is in Chicago’. But if we use a domain of things in general, the implicit restriction to people carried by ‘everyone’ has to be made explicit by adding ‘ Px ’, for ‘ x is a person’, into the restrictor predicate, giving us ‘ $\forall x((Px \wedge Cx) \rightarrow Hx)$ ’.

One last point about picking domains: sentence (17) contains the name ‘Chicago’, which refers to a city, not a person, so you might think that our guidelines require a domain of things here. In fact, though, we can use the domain of people, because we aren’t symbolizing the name ‘Chicago’ in (17) *as a name*, but rather treating it as part of the larger predicate ‘____ is in Chicago’, which we symbolize as ‘ C ’. What the guidelines mean by a “standalone name” is a name that will get *symbolized as a name*. It is only if a sentence contains a standalone name (in this sense) of something other than a person that you shouldn’t use a domain of people. The reason is that every object referred to by a name *in our symbolization* must be part of the domain.

■ Exercises 5.3

A. Symbolize the following (using a domain of things):

1. Some dogs are cute but not friendly.
2. Some dogs are neither friendly nor cute.
3. Some ancient manuscripts are priceless.
4. Every wealthy artist is happy.
5. Every artist is both wealthy and happy.
6. Every politician from London will get re-elected.
7. All red mushrooms are deadly if eaten.

B. The following contain quantifiers restricted to people, and can be symbolized in either a domain of people or a domain of things. Symbolize each first in a domain of people, and then in a domain of things:

1. Someone is wise but not happy.
2. Everyone is both wise and happy.
3. Everyone from Sweden is either wise or happy.
4. Someone from Norway is wise, but not everyone in Scandinavia is wise.

5.4 Quantifiers and Negation

So far we’ve only looked at English quantifier phrases involving ‘every’ and ‘some’. But our two FOL quantifiers \forall and \exists can also be used to symbolize other quantifier phrases in English. Consider the following sentence:

(18) No one is angry.

This could be paraphrased as ‘It is not the case that someone is angry’. We can therefore symbolize it using negation together with an existential quantifier: $\neg\exists xAx$. Interestingly, though, this is not the only option. If you think about it, (18) could also be paraphrased as ‘Everyone is not-angry’. So we can also symbolize our sentence using negation and a universal quantifier: $\forall x\neg Ax$ (“every individual x is such that x is not angry”). Indeed, as we will see, it holds in general that $\forall v\neg\phi$ is logically equivalent to $\neg\exists v\phi$. (Notice that I have here returned to the practice of using ‘ ϕ ’ as a metavariable, now over FOL sentences; and v as a metavariable over FOL variables.) So whenever you have a negation in front of an existential quantifier in FOL, you can move the negation over the quantifier, and flip the quantifier into a universal to obtain an equivalent sentence.

A similar pattern emerges if we consider the following.

(19) Not everyone is happy.

This can be symbolized as $\neg\forall xHx$. But again, if you think about it, (19) could be paraphrased as ‘Someone is not-happy.’ So another way to symbolize this sentence is as $\exists x\neg Hx$. This illustrates that in general $\neg\forall v\phi$ is equivalent to $\exists v\neg\phi$, meaning that if a negation occurs in front of a universal quantifier, we can move it over the quantifier and flip the quantifier into an existential. So we have the following logical laws:¹

QUANTIFIER EQUIVALENCE LAWS:

$\forall v\neg\phi$ is equivalent to $\neg\exists v\phi$

$\neg\forall v\phi$ is equivalent to $\exists v\neg\phi$

Next consider the following examples involving more complex quantifier phrases:

(20) No dog is a poodle.

(21) Not all dogs are poodles.

Sentence (20) says that there does not exist a dog that is a poodle. So it can be symbolized as a negated existential sentence, $\neg\exists x(Dx \wedge Px)$. Whereas (20) says something which is in fact false, (21) says something true: that it’s not the case that every dog is a poodle. It can therefore be symbolized as a negated universal sentence, $\neg\forall x(Dx \rightarrow Px)$.

Again, though, each of these can be symbolized another way too. Sentence (21) could also be paraphrased as saying ‘Some dog (i.e. at least one) is a non-poodle’. So instead of symbolizing it as a negated universal ‘ $\neg\forall x(Dx \rightarrow Px)$ ’, it could also be symbolized as an existential: $\exists x(Dx \wedge \neg Px)$. Similarly, sentence (20) could also be paraphrased as ‘every dog is a non-poodle’. So instead of symbolizing it as a negated existential ‘ $\neg\exists x(Dx \wedge Px)$ ’, it could also be symbolized as a universal: $\forall x(Dx \rightarrow \neg Px)$. In both of these cases, the two possible symbolizations are equivalent to each other.

We’ll learn how to use natural deduction to prove these equivalences later. But you can already see the reason for it given the Quantifier Equivalence Laws from above, together with some of the equivalences we know from TFL. For example, our symbolization ‘ $\neg\forall x(Dx \rightarrow Px)$ ’ of (21) is, by the Quantifier Equivalence Laws, equivalent to

¹Notice that this in turns means that $\forall v\phi$ is equivalent to $\neg\exists v\neg\phi$. So we don’t really need \forall in our language in addition to \exists , we could just define one in terms of the other and always write $\neg\exists v\neg\phi$ when we mean $\forall v\phi$.

$\exists x \neg(Dx \rightarrow Px)$. Then, applying the NegImp law inside the scope of the existential, the latter is equivalent to $\exists x(Dx \wedge \neg Px)$, which was our second possible symbolization!

In general, we've seen the following patterns:

An English sentence that can be paraphrased as 'No F is G ' can be symbolized as $\neg \exists x(\mathcal{F}x \wedge \mathcal{G}x)$, or as $\forall x(\mathcal{F}x \rightarrow \neg \mathcal{G}x)$

An English sentence that can be paraphrased as 'Not every F is G ' can be symbolized as $\neg \forall x(\mathcal{F}x \rightarrow \mathcal{G}x)$ or as $\exists x(\mathcal{F}x \wedge \neg \mathcal{G}x)$.

These patterns also apply if we want to symbolize our earlier examples (18) and (19) in a domain of things rather than people. In a domain of things, 'No one (i.e. no person) is angry' gets symbolized as $\neg \exists x(Px \wedge Ax)$ (or alternatively, as $\forall x(Px \rightarrow \neg Ax)$, i.e. "every person is non-angry"). And 'Not everyone is angry' would become $\neg \forall x(Px \rightarrow Ax)$ (or alternatively, $\exists x(Px \wedge \neg Ax)$, i.e. "at least one person is non-angry").

■ Exercises 5.4

A. Symbolize the following:

1. No honest politician is rich.
2. Some honest politicians aren't rich.
3. No logician is both rich and famous.
4. Not every logician is both rich and famous.
5. No one from Sweden is famous.

5.5 The Utility of Paraphrase

As we've seen, it is important to get the structure of the sentences you want to symbolize right. Sometimes you will be able to move from English directly to a sentence of FOL. Other times, it helps to paraphrase the sentence one or more times. Each successive paraphrase should move from the original sentence closer to something that you can finally symbolize directly in FOL.

For the next several examples, we will use this symbolization key:

domain: people

B : ____ is a bassist

R : ____ is a rock star

k : Kim Deal

Now consider these sentences:

- (22) If Kim Deal is a bassist, then she is a rock star.
- (23) If someone is a bassist, then she is a rock star.

These sentences look similar, and even have the same words in the consequent ('... she is a rock star'), but they mean very different things and will require different symbolizations.

Sentence (22) can be paraphrased as, ‘If Kim Deal is a bassist, then *Kim Deal* is a rockstar’. This can obviously be symbolized as the conditional ‘ $Bk \rightarrow Rk$ ’.

Sentence (23) is more tricky. You might think that it’s a conditional, with an existential quantifier in its antecedent, and should be symbolized along the lines of ‘ $\exists x Bx \rightarrow Rx$ ’. This isn’t correct, however, since the quantifier in this symbolization isn’t binding the variable in ‘ Rx ’, that is, it’s an open sentence. In the English (23), the pronoun ‘she’ is referring back to the bassist. You might try to fix it by adding parentheses, giving us ‘ $\exists x (Bx \rightarrow Rx)$ ’. This is now at least a closed sentence, with the quantifier binding both variables. Unfortunately it doesn’t capture the meaning of (23). Recall that in TFL, $(P \rightarrow Q)$ is equivalent (by the Imp Law) to $(\neg P \vee Q)$. Similarly, in FOL ‘ $\exists x (Bx \rightarrow Rx)$ ’ is equivalent to ‘ $\exists x (\neg Bx \vee Rx)$ ’. This will be true as long as there is some non-bassist in the world, but that clearly doesn’t suffice for (23) to be true!

Paraphrase will help us reach our goal. (23) can be paraphrased as ‘If some person is a bassist, then that person is a rock star’. This sentence is not about any particular person, but rather says something about every person who is a bassist. It can be paraphrased as ‘For any person x , if x is a bassist, then x is a rockstar’, which can be symbolized as ‘ $\forall x (Bx \rightarrow Rx)$ ’. This is now the correct symbolization of (23). There is a surprising but important upshot to this: sometimes, English sentences that involve the quantifier ‘someone’ get symbolized using the universal quantifier \forall in FOL!²

Next, consider these sentences:

- (24) If anyone is a bassist, then Kim Deal is a rock star.
 (25) If anyone is a bassist, then she is a rock star.

The same words appear as the antecedent in sentences (24) and (25) (‘If anyone is a bassist...’). But again, they mean very different things, and will have to be symbolized differently. Paraphrase will help us.

Sentence (24) can be paraphrased, ‘If there exists at least one bassist, then Kim Deal is a rock star’. This is a conditional whose antecedent is an existentially quantified sentence. We can symbolize the entire sentence with a conditional as the main logical operator: ‘ $\exists x Bx \rightarrow Rk$ ’. (Notice that this is not an open sentence, because k is a name, not a variable.)

Sentence (25) again has a pronoun ‘she’ referring back to the bassist. It can be paraphrased as ‘For all people x , if x is a bassist, then x is a rock star’, or just ‘All bassists are rock stars’. So it gets symbolized as a universally quantified sentence, ‘ $\forall x (Bx \rightarrow Rx)$ ’, just like sentence (23) from earlier. What these examples illustrate is that the English quantifier ‘anyone’ sometimes gets symbolized as an existential quantifier in FOL, and at other times as a universal quantifier. To determine which, try paraphrasing the sentence using words *besides* ‘any’ or ‘anyone’.

■ Exercises 5.5

A. Symbolize the following:

1. If everyone is wealthy, economists are happy.

²Sentence (23) involves what linguists call a “donkey anaphor.” You can read more about this in the Stanford Encyclopedia of Philosophy: <https://plato.stanford.edu/entries/anaphora/>.

2. If anyone is wealthy, economists are happy.
3. If someone is wealthy, they are an economist.
4. If a politician is corrupt, they are also dishonest.
5. All violins and cellos are stringed instruments.

B. Below are the syllogistic figures identified by Aristotle and his successors, along with their medieval names. These formed the foundation of formal logic for over two millennia, until the end of the 19th century. Formalize each syllogistic figure in FOL.

- **Barbara.** All G are F. All H are G. So: All H are F
- **Celarent.** No G are F. All H are G. So: No H are F
- **Ferio.** No G are F. Some H is G. So: Some H is not F
- **Darii.** All G are F. Some H is G. So: Some H is F.
- **Camestres.** All F are G. No H are G. So: No H are F.
- **Cesare.** No F are G. All H are G. So: No H are F.
- **Baroko.** All F are G. Some H is not G. So: Some H is not F.
- **Festino.** No F are G. Some H are G. So: Some H is not F.
- **Datisi.** All G are F. Some G is H. So: Some H is F.
- **Disamis.** Some G is F. All G are H. So: Some H is F.
- **Ferison.** No G are F. Some G is H. So: Some H is not F.
- **Bokardo.** Some G is not F. All G are H. So: Some H is not F.
- **Camenes.** All F are G. No G are H So: No H is F.
- **Dimaris.** Some F is G. All G are H. So: Some H is F.
- **Fresison.** No F are G. Some G is H. So: Some H is not F.

C. In §5.3 we noted that English sentences of the form ‘No \mathcal{F} is \mathcal{G} ’ can be symbolized either as $\neg\exists x(\mathcal{F}x \wedge \mathcal{G}x)$ or as $\forall x(\mathcal{F}x \rightarrow \neg\mathcal{G}x)$, and ones of the form ‘Not all \mathcal{F} are \mathcal{G} ’ can be symbolized as either $\neg\forall x(\mathcal{F}x \rightarrow \mathcal{G}x)$ or as $\exists x(\mathcal{F}x \wedge \neg\mathcal{G}x)$. Following these templates, give two different symbolizations for each of the following:

1. No spy is famous.
2. Not all spies are famous.
3. Not every famous villain is a spy.
4. Not all famous spies are villains.
5. No spy is both famous and a villain.
6. No villain is both famous and a spy.
7. Not every spy is both famous and a villain.
8. Some spies are not villains.

D. For each argument, write a symbolization key and symbolize the argument in FOL.

1. Willard is a logician. All logicians wear funny hats. So Willard wears a funny hat
2. Nothing on my desk escapes my attention. There is a computer on my desk. As such, there is a computer that does not escape my attention.
3. All my dreams are black and white. Old TV shows are in black and white. Therefore, some of my dreams are old TV shows.
4. Neither Holmes nor Watson has been to Australia. A person could see a kangaroo only if they had been to Australia or to a zoo. Although Watson has not seen a kangaroo, Holmes has. Therefore, Holmes has been to a zoo.

5. No one expects the Spanish Inquisition. No one knows the troubles I've seen. Therefore, anyone who expects the Spanish Inquisition knows the troubles I've seen.
6. All babies are illogical. Nobody who is illogical can manage a crocodile. Berthold is a baby. Therefore, Berthold is unable to manage a crocodile.

5.6 Many-Place Predicates

So far, we have only considered sentences with one-place predicates and one quantifier. The full power of FOL really comes out when we start to use many-place predicates and multiple quantifiers, however. Whereas the logic of singly-quantified sentences has been well known for over two millennia since Aristotle, it took until the work of Gottlob Frege in the late 19th for a logic capable of handling sentences with multiple quantifiers to be developed. The system of FOL we are here studying is a fragment of the logic Frege developed in his book *Begriffsschrift* (1879).

ONE-PLACE PREDICATES concern *properties* that objects might have. They have one argument place, or gap, in them. To make a sentence, we simply slot a name into that gap. Other predicates concern *relations* between things. Here are some examples of relational predicates in English:

____ loves ____
 ____ is to the left of ____
 ____ is in debt to ____

These are TWO-PLACE PREDICATES: they need to be filled in with two terms in order to make a sentence. Conversely, if we start with an English sentence containing many singular terms, we can remove two singular terms, to obtain different two-place predicates. Consider the sentence 'Vinnie borrowed the family car from Nunzio'. By deleting two singular terms, we can obtain any of three different two-place predicates

Vinnie borrowed ____ from ____
 ____ borrowed the family car from ____
 ____ borrowed ____ from Nunzio

And by removing all three singular terms, we obtain a THREE-PLACE PREDICATE:

____ borrowed ____ from ____

Indeed, there is in principle no upper limit on the number of argument places that our predicates may contain.

It's important to realize that the multiple argument places in a predicate can be filled either with the same term, or with different terms, and in various different orders. For example, if we begin with the two-place predicate '____ loves ____', we can fill the gaps with the names 'Karl' and 'Imre' in various different ways, to obtain different English sentences:

- (26) Karl loves Imre.
- (27) Imre loves Karl.
- (28) Karl loves Karl.

In FOL, many-place predicates are symbolized via uppercase letters, just like one-place predicates. To symbolize the above sentences, we can use the following symbolization key:

domain: people
i: Imre
k: Karl
L: ____ loves ____

As in the case of one-place predicates, FOL names appear *after* the predicate letter. So sentence (26) will be symbolized as '*Lki*', sentence (27) as '*Lik*', and sentence (28) as '*Lkk*'. You can think of the FOL predicate letter '*L*' as having two invisible argument places after it, into which we can slot the names '*i*' and '*k*'. The convention is that the first gap after the predicate letter represents the first gap in the corresponding English predicate, and the second gap represents the second gap in the English predicate. So since sentence (28) results from putting 'Karl' into the first gap in '____ loves ____' and 'Imre' into the second, its symbolization in FOL has '*k*' in the first gap after '*L*' and '*i*' in the second.

Another way to put it is that the first gap in the English predicate '____ loves ____' is for the *agent* of the relation — the lover, the person doing the loving — and the second gap is for the *patient*, or direct object, of the relation — the beloved, the person who is loved. So given our symbolization key, the first name in the FOL sentence '*Lki*' represents the agent, the lover, and the second name represents the patient, the beloved.

Here are some more sentences that we can symbolize using this key:

- (29) Imre loves himself.
- (30) Karl loves Imre, but not vice versa.
- (31) Karl is loved by Imre.

Sentence (29) can be paraphrased as 'Imre loves Imre', and is symbolized by '*Lii*'. Sentence (30) is a conjunction. We can paraphrase it as 'Karl loves Imre, and Imre does not love Karl', and so symbolize it as '*Lki* \wedge \neg *Lik*'. Sentence (31) is in the passive voice, but it can be paraphrased in the active voice as 'Imre loves Karl', and so symbolized as '*Lik*'. Of course, there are differences of *tone* between the active and passive voice; but we have preserved the truth conditions.

The difference between active and passive voice illustrates something important. Suppose we had instead used the following symbolization key of our predicate '*L*':

L: ____ is loved by ____

Now '*L*' symbolizes an English predicate in the passive voice, meaning that the first gap now represents the patient, the person who is loved, and the second gap the agent, the person doing the loving. Using this symbolization key, the FOL sentence '*Lki*' now means that Karl is loved by Imre, that is to say, that Imre loves Karl. So Sentence (26) — which says that Karl loves Imre, i.e. that Imre is loved by Karl — can no longer be symbolized as '*Lki*', but must be symbolized as '*Lik*'. The overall moral is simple: *differences in the order of names matter*.

■ Exercises 5.6

A. symbolize the following using two-place predicates.

1. Tom chased Jerry, and Tom was also chased by Jerry.
2. Tom chased Jerry, but Jerry didn't chase Tom.
3. If Tai admires Meg, then Meg admires herself.
4. Tom didn't chase Jerry unless Jerry was chased by Tom.

5.7 Multiple Generality

Now that we have two-place predicates to work with, we can also symbolize sentences that combine such predicates with quantifiers. Suppose we're again working with the following symbolization key:

domain: people
i: Imre
k: Karl
L: ____ loves ____

Now consider the following sentences:

- (32) Everyone loves Imre.
- (33) Imre loves everyone.
- (34) Imre is loved by everyone

Starting with (32), we could paraphrase this as: every person x is such that x loves Imre. So here the variable x will now go into the first slot after L , since x is the lover, and the name i into the second, giving us $\forall x Lxi$. There's a common mistake students make with examples like this. In the English sentence (32), it looks like the quantifier phrase 'everyone' occupies the first argument place of 'loves', so students are sometimes tempted to put a quantifier into the corresponding argument place in the FOL symbolization, producing something like ' $L\forall xi$ '. However, this is not a grammatical sentence of FOL! In FOL, only *names* and *variables* may occur in the argument places of predicate letters. So again, the way to symbolize this is as ' $\forall x Lxi$ ', with the universal quantifier out front binding a variable in the first argument place of ' L '.

Sentence (33) could be paraphrased as: every person x is such that Imre loves x . So here the variable now goes into the second argument place of ' L ', giving us the symbolization ' $\forall x Lix$ '. As for (34), notice that this is equivalent to (32), so it can also be symbolized as ' $\forall x Lxi$ '.

As already mentioned, the real power of FOL lies in its ability to treat sentences with many-place predicates and *multiple* quantifiers. Take the following:

- (35) Everyone loves someone.

On its most natural interpretation, this sentence says that for every person x , there is some person y whom x loves. It can therefore be symbolized as ' $\forall x \exists y Lxy$ '. This would be true, for example, in a love triangle situation where Karl loves Imre, Imre loves Juan, and Juan loves Karl (and no one loves anyone else): no matter which person x we consider, we can find some person y whom x loves.

However, English sentences like (35) that contain multiple quantifiers are *ambiguous*. There is another interpretation of (35) on which it says that every x loves some one particular person y , a claim that can also be expressed via the following English sentence:

(36) There is someone who is loved by everyone.

Since this now claims that there is some particular person y who is loved by every x , we can symbolize it as ' $\exists y\forall xLxy$ '. This would be false in our earlier love triangle situation, since there is no lucky individual y who is loved by *everyone* in that scenario. A situation where ' $\exists y\forall xLxy$ ' is true would be one where, for example, Karl loves Imre, Juan loves Imre, and Imre also loves himself. Now Imre is the lucky individual y who's loved by everyone.

The two examples we've looked at — ' $\forall x \exists y Lxy$ ' and ' $\exists y \forall x Lxy$ ' — differ in the order, or *scope*, of the two quantifiers: in the first, the universal $\forall x$ has the larger scope (i.e. comes first, and is the main operator), whereas in the second the existential $\exists y$ has the larger scope.³ As we've seen, this difference in scope results in a difference in meaning between these two FOL sentences. Accidentally switching the scope of quantifiers gives rise to the so-called *quantifier shift fallacy*. For example, the following argument is not valid:

Everything is caused by something. (E \forall)
 \therefore There is some one thing that caused everything. (E \exists)

Using ‘C’ for ‘____ caused ____’ (and a domain of things), the premise can be symbolized as ‘ $\forall x \exists y Cyx$ ’: for every x there exists some y such that y caused x . The conclusion, on the other hand, can be symbolized as ‘ $\exists y \forall x Cyx$ ’: there exists some y such that for every x , y caused x . The latter is not implied by ‘ $\forall x \exists y Cyx$ ’. We’ll leave it as an exercise for you to describe a situation that would make the premise true, but the conclusion false.

Such fallacies, though, arise only when we swap around universal with existential quantifiers. With strings of the same quantifier, the order doesn't matter. For example, ' $\exists x \exists y Lxy$ ' and ' $\exists y \exists x Lxy$ ' would naturally be used to symbolize English sentences 'there is someone who loves someone' and 'there is someone who is loved by someone', respectively. But, though these differ in nuance, they are true in exactly the same situations; ' $\exists x \exists y Lxy$ ' and ' $\exists y \exists x Lxy$ ' are therefore equivalent. Also, to return to a point from §5.2, notice that ' $\exists x \exists y Lxy$ ' does not require that x and y be different individuals. This sentence, as well as ' $\exists y \exists x Lxy$ ', would be true in a situation where Imre loves himself (and no one loves anyone else). After all, that would be a situation where someone loves someone, and also one where someone is loved by someone.

Similar comments apply to pairs like ‘ $\forall x\forall yLxy$ ’ and ‘ $\forall y\forall xLxy$ ’: if everyone loves everyone (as per the first sentence) then it follows that everyone is loved by everyone (as per the second), and vice versa. So there can be no situation that makes one true but not the other. And notice that for either to be true, everyone has to, among other things, love themselves. So both of these imply $\forall xLxx$, as we’ll be able to show soon using natural deductions.

Lastly, multiply quantified sentences can of course also involve “negative quantifiers.” For example:

³Of course we could use different variables: (36) could also be paraphrased as saying that there is some person x who is loved by every y , and thus symbolized as $\exists x\forall yLyx$ instead of $\exists y\forall xLxy$. These FOL sentences look different, but they are equivalent: in both cases the existential quantifier binds the variable in the *second* argument position of L , and the universal quantifier binds the variable in the *first* argument position. This is what's crucial to capturing the meaning of (36).

- (37) No one loves everyone.
 (38) There's someone who loves no one.

Sentence (37) is the denial of 'there is someone who loves everyone'; since the latter gets symbolized as ' $\exists x\forall yLxy$ ', sentence (37) gets symbolized as its negation, ' $\neg\exists x\forall yLxy$ '. Sentence (38) says that there is some person x such that no matter what y we pick, x does not love y . It therefore gets symbolized as ' $\exists x\forall y\neg Lxy$ '.

Another way to think about (38) is as saying: there is some x such that there does not exist any y whom x loves. So we can also symbolize it as ' $\exists x\neg\exists yLxy$ '. This illustrates that the Quantifier Equivalence Laws from §5.2, which govern the movement of negation across quantifiers, continue to hold in multiply quantified sentences. Similarly, if we start with ' $\neg\exists x\forall yLxy$ ', which was our symbolization of (37), and move the negation across both quantifiers, we end up with ' $\forall x\exists y\neg Lxy$ '. This says that for every x there is at least one y whom x does not love, which is indeed another way to capture the truth conditions of (37).

■ Exercises 5.7

A. Symbolize the following (in a domain of people):

1. Socrates admires someone.
2. Everyone admires Socrates.
3. If everyone admires Socrates, then Socrates also admires himself.
4. Everyone loves someone.
5. Someone loves everyone.
6. Everyone is loved by someone.
7. Someone is loved by everyone.
8. Everyone loves everyone.
9. Someone loves no one
10. No one loves everyone.
11. No one is loved by everyone.
12. Someone is loved by no one.

5.8 Intermediate Steps to Symbolization

As we are starting to see, symbolization in FOL can become tricky. When symbolizing a complex sentence, it is best to proceed by way of several intermediate steps. Let's look at some examples. Consider the following sentences:

- (39) Geraldo owns a dog.
 (40) Someone owns a dog.
 (41) All of Geraldo's friends are dog owners.
 (42) Every dog owner is a friend of a dog owner.

We'll use the following symbolization key:

domain: things

D : ____ is a dog
 F : ____ is a friend of ____
 O : ____ owns ____
 g : Geraldo

Sentence (39) can be paraphrased as, ‘There is a dog that Geraldo owns’. This can be symbolized by ‘ $\exists x(Dx \wedge Ogx)$ ’.

Sentence (40) can be paraphrased as, ‘There is some y such that y owns a dog’. We can begin by just focusing on the initial quantifier, which gives us ‘ $\exists y(y \text{ owns a dog})$ ’. Now the fragment ‘ y owns a dog’ is exactly like sentence (39), except it contains a variable instead of a name. So we could symbolize this fragment as $\exists x(x \text{ is a dog} \wedge y \text{ owns } x)$. Putting it all together we get the following:

$$\exists y \exists x (Dx \wedge Oyx)$$

In working out how to symbolize this sentence, we wrote down things like ‘ $\exists y(y \text{ owns a dog})$ ’ and $\exists x(x \text{ is a dog} \wedge y \text{ owns } x)$. To be very clear: these are *neither* FOL sentence *nor* English sentences. They use bits of FOL (‘ \exists ’, ‘ y ’) and bits of English (‘owns a dog’). These really are just *intermediate steps* on the way to symbolizing the English sentence, a bit of “scratch work” we do on the side as we work through the problem.

Sentence (41) can be paraphrased as, ‘Everyone who is a friend of Geraldo is a dog owner’. Since being a dog owner is the same as owning a dog, we can in turn paraphrase this as ‘Everyone who is a friend of Geraldo owns a dog’. So we can write:

$$\forall x [Fgx \rightarrow x \text{ owns a dog}]$$

as our first intermediate step. Now the consequent of the conditional, ‘ x is a dog owner’, is structurally just like sentence (39). Using our symbolization of (39) as a guide, we get:

$$\forall x [Fgx \rightarrow \exists y (Dy \wedge Oxy)]$$

Notice that it was essential that we used a variable other than ‘ x ’ for the existential quantifier in the consequent. If we had instead written:

$$\forall x [Fgx \rightarrow \exists x (Dx \wedge Oxx)]$$

we would have had a *clash of variables*. The first variable ‘ x ’ after the predicate ‘ O ’ represents the agent, the person doing the owning, who is a friend of Geraldo. Accordingly, this variable should get bound by the initial quantifier ‘ $\forall x$ ’ that also binds the ‘ x ’ in the antecedent ‘ Fgx ’. But if we now use ‘ $\exists x$ ’ in the consequent, then it should bind every ‘ x ’ in its scope, including the first one in ‘ Oxx ’. To avoid this clash of variables, we have to use a different variable for the quantifier in the consequent, as in ‘ $\forall x [Fgx \rightarrow \exists y (Dy \wedge Oxy)]$ ’. The broad moral is that a single variable cannot serve two masters simultaneously.

Moving to sentence (42), it can be paraphrased as ‘For any x , if x is a dog owner, then x is a friend of some dog owner’. As our first intermediate step, we might have:

$$\forall x [x \text{ is a dog owner} \rightarrow \exists y (y \text{ is a dog owner} \wedge Fxy)]$$

Again, being a dog owner is the same as owning some dog, and we know how to symbolize that. To avoid a variable clash, we’ll have to use an existential quantifier that won’t threaten

to bind either the ‘ x ’ or the ‘ y ’ we already have in our intermediate step. So let’s use ‘ $\exists z$ ’ in both cases, giving us:

$$\forall x [\exists z (Dz \wedge Oxz) \rightarrow \exists y (\exists z (Dz \wedge Oyz) \wedge Fxy)]$$

Here ' $\exists z(Dz \wedge Oxz)$ ' just says that x owns a dog, and ' $\exists z(Dz \wedge Oyz)$ ' that y owns a dog.

We here decided to use the same variable, ‘ z ’, in both the antecedent *and* the consequent of the conditional. This is ok, because there is no scope overlap between the two. We might graphically represent the scope of the various quantifiers thus:

$$\overbrace{\forall x \overbrace{\exists z (Dz \wedge Oxz) \rightarrow \exists y \overbrace{(\exists z (Dz \wedge Oyz) \wedge Fyx)}^{\text{scope of '}\exists y\text{'}}}}^{\text{scope of '}\forall x\text{'}}}$$

5.9 Adding Identity

Consider this sentence:

(43) Pavel owes money to everyone

Let's use a domain of people, ' p ' for 'Pavel', an ' O ' for '____ owes money to ____'. We can then symbolize sentence (43) by ' $\forall x Opx$ '. But this has a (perhaps) odd consequence. It requires that Pavel owes money to *every* member of the domain. Since Pavel himself must be a member of the domain, this entails that Pavel owes money to himself. And maybe we did not want to say that. Maybe what we meant to say was:

(44) Pavel owes money to everyone *else*

(45) Pavel owes money to everyone *other than* Pavel

(46) Pavel owes money to everyone *except* himself

But we do not have any way of dealing with the italicized words yet. The solution is to add a new symbol to FOL.

The symbol '=' will be a two-place predicate, denoting the relation of *identity*. Since identity is such a basic logical concept — similar to how e.g. conjunction, negation, or existential quantification are basic logical concepts — '=' functions as a LOGICAL CONSTANT in FOL. This means that the symbol '=' *has* to be interpreted as '____ is identical to ____'; you can't assign it a different meaning in a symbolization key.

To highlight the fact that identity is special in being the only logical constant among two-place predicates, we adopt a different notational convention for it, and write it *between* two terms rather than in front of them. This notation will be familiar to you from math, where you write things like $\frac{1}{2} = \frac{4}{8}$. Note that in saying that some objects x and y are identical, we don't merely mean that they are very similar, or indistinguishable in the way that e.g. two cans of Coca Cola are. We mean that they are *one and the same* object.

To put this to use, suppose we want to symbolize this sentence:

(47) Pavel is Mister Checkov.

Using ' c ' for the name 'Mister Checkov', sentence (47) can be symbolized as ' $p = c$ '. This tells us that Pavel and Mister Checkov are one and the same person, and that the names ' p ' and ' c ' refer to the same individual.

We can also now deal with sentences (44)–(46). All of these sentences can be paraphrased as 'Everyone who is not Pavel is owed money by Pavel'. Paraphrasing some more, we get: 'For all x , if x is not Pavel, then x is owed money by Pavel'. Now that we are armed with our new identity symbol, we can symbolize this as ' $\forall x (\neg x = p \rightarrow Opx)$ '.

This last sentence contains the formula ' $\neg x = p$ '. This might look a bit strange, but it just means that we are negating the entire formula, ' $x = p$ '. From math, you're probably familiar with the notation ' \neq ' for negated identity, so we'll also adopt this notational convention here, though only as a convenient shorthand:

An FOL sentence of the form $\neg t_1 = t_2$ can be abbreviated as $t_1 \neq t_2$

Using this notational shorthand, we can rewrite our symbolization as $\forall x(x \neq p \rightarrow Opx)$.

In addition to sentences that use the words ‘else’, ‘other than’, and ‘except’, identity will be helpful when symbolizing some sentences that contain the words ‘besides’ and ‘only.’ Consider these examples:

(48) No one besides Pavel owes money to Hikaru.

(49) Only Pavel owes Hikaru money.

Letting ‘*h*’ name Hikaru, sentence (48) can be paraphrased as, ‘No one who is not Pavel owes money to Hikaru’. This can be symbolized by $\neg \exists x(x \neq p \wedge Oxh)$. Sentence (48) can be paraphrased as ‘for all *x*, if *x* owes money to Hikaru, then *x* is Pavel’. This can be symbolized as $\forall x(Oxh \rightarrow x = p)$. In fact, these two symbolizations are equivalent to each other; and (48) and (49) do seem to express the same claim.

But there is one subtlety here. Our symbolizations imply that anyone who is not Pavel does not owe money to Hikaru. But (48) and (49) also seem to imply that Pavel does owe money to Hikaru. To capture this, we can add ‘*Oph*’ as a conjunct to either symbolization, giving us e.g. $Oph \wedge \forall x(Oxh \rightarrow x = p)$ as a final symbolization. This, in turn, can be shortened to $\forall x(Oxh \leftrightarrow x = p)$.

Identity can also be used to symbolize claims about *how many* things there are of a particular kind. We’ll go look at three kinds of claims of this sort.

There are at least...

Consider the following ‘at least’ claims:

(50) There is at least one apple

(51) There are at least two apples

(52) There are at least three apples

We’ll use ‘*A*’ for ‘___ is an apple’, and a domain of things. Sentence (50) does not require identity. It can be adequately symbolized by $\exists xAx$: there is some apple; perhaps many, but at least one.

It might be tempting to also translate sentence (51) without identity. But consider the sentence $\exists x \exists y(Ax \wedge Ay)$. This says that there is some apple *x* in the domain and also some apple *y* in the domain. Since nothing precludes these from being one and the same apple, this would be true even if there were only one apple. (Recall here the point from §5.2 that a difference in variables need not indicate a difference in the objects the variables pick out.) To make sure that we are dealing with *different* apples, we need to use identity, and symbolize (51) as $\exists x \exists y(Ax \wedge Ay \wedge x \neq y)$.

Sentence (52) requires talking about three different apples. Now we need three existential quantifiers, and we need to make sure that each will pick out something different: $\exists x \exists y \exists z(Ax \wedge Ay \wedge Az \wedge x \neq y \wedge y \neq z \wedge x \neq z)$. As you can see, by following this pattern, we can symbolize claims of the sort ‘there are at least *n* apples’ for any (finite) number *n*.

There are at most...

Now consider these sentences:

- (53) There is at most one apple
 (54) There are at most two apples

Sentence (53) can be paraphrased as, ‘It is not the case that there are at least *two* apples’. This is just the negation of sentence (51):

$$\neg \exists x \exists y (Ax \wedge Ay \wedge \neg x = y)$$

But sentence (53) can also be approached in another way. It means that if you pick out an object and it’s an apple, and then you pick out an object and it’s also an apple, you must have picked out the same object both times. With this in mind, it can be symbolized by:

$$\forall x \forall y [(Ax \wedge Ay) \rightarrow x = y]$$

The two sentences will turn out to be logically equivalent.

Similarly, sentence (54) can be approached in two equivalent ways. It can be paraphrased as ‘It is not the case that there are at least *three* apples’, which is just the negation of (52) above. Alternatively, we can understand it as saying that if you pick out an apple, and an apple, and an apple, then you will have picked out the same apple at least once. Thus:

$$\forall x \forall y \forall z [(Ax \wedge Ay \wedge Az) \rightarrow (x = y \vee x = z \vee y = z)]$$

Again, by following this pattern we can symbolize claims of the sort ‘there are at most *n* apples’ for any *n*.

There are exactly...

Lastly, there are statements that specify a precise numerical quantity:

- (55) There is exactly one apple.
 (56) There are exactly two apples.
 (57) There are exactly three apples.

Sentence (55) can be paraphrased as ‘There is *at least* one apple and there is *at most* one apple’. This is just the conjunction of (50) and (53) from above:

$$\exists x Ax \wedge \forall x \forall y [(Ax \wedge Ay) \rightarrow x = y]$$

But it is perhaps more straightforward to paraphrase sentence (55) as, ‘There is a thing *x* which is an apple, and everything which is an apple is just *x* itself’. Thought of in this way, we’d symbolize it as:

$$\exists x [Ax \wedge \forall y (Ay \rightarrow x = y)]$$

Similarly, sentence (56) may be paraphrased as, ‘There are *at least* two apples, and there are *at most* two apples’, and thus symbolized as the conjunction of (51) and (54). More efficiently, though, we can paraphrase it as ‘There are at least two different apples, and every apple is one of those two apples’. Then we offer:

$$\exists x \exists y [Ax \wedge Ay \wedge x \neq y \wedge \forall z (Az \rightarrow (x = z \vee y = z))]$$

Continuing with this patten, we could symbolize the claim that there are exactly three apples as follows:

$$\exists x \exists y \exists z [Ax \wedge Ay \wedge Az \wedge x \neq y \wedge x \neq z \wedge y \neq z \wedge \forall v (Av \rightarrow (x = v \vee y = v \vee z = v))]$$

and so on, for any number n of apples.

Finally, consider these sentence:

(58) There are exactly two things

(59) There are exactly two objects

It might be tempting to add a predicate to our symbolization key, to symbolize the English predicate ‘_____ is a thing’ or ‘_____ is an object’. But this is unnecessary. Words like ‘thing’ and ‘object’ apply trivially to everything. So we can symbolize either sentence as:

$$\exists x \exists y [x \neq y \wedge \forall z (x = z \vee y = z)]$$

Logical Truths Involving Identity

We introduced the symbol ‘=’ as an additional logical constant, i.e. as a logical symbol whose meaning remains fixed, just like ‘ \exists ’ or ‘ \neg ’. Part of the motivation for treating identity as a logical constant is that it seems like a very primitive logical concept, much like negation or existence.

Another, related motivation is that there seem to be certain basic logical truths involving identity, just like e.g. the LAW OF EXCLUDED MIDDLE from TFL is a basic logical truth involving negation. One such primitive truth is that everything is identical to itself, which we can express as:

$$\forall x x = x$$

This is sometimes called the LAW OF IDENTITY, and it will be a theorem in the system of natural deduction for FOL that we will develop later.

Another logical truth involving identity sometimes goes by the name of LEIBNIZ’S LAW.⁴ It says that if x and y are one and the same thing, then x and y must share all their properties. So if we use ‘ D ’ for ‘_____ is a dog’, the following would be an instance of Leibniz’s Law:

$$\forall x \forall y [x = y \rightarrow (Dx \leftrightarrow Dy)]$$

This says that if x and y are identical, then x is a dog iff y is. Similarly, if we use ‘ O ’ for ‘_____ owns _____’, then another instance of Leibniz’s Law says that if x and y are identical, then x owns a dog iff y does:

$$\forall x \forall y [x = y \rightarrow (\exists z (Dz \wedge Oxz) \leftrightarrow \exists z (Dx \wedge Oyz))]$$

⁴This law is named after Gottfried Willhelm Leibniz (1646–1716). Leibniz actually endorsed a stronger claim, which says not only that x and y must share all their properties if they are identical, but also (and more controversially!) that if x and y share all their properties, then they are identical. This second claim is called the “Identity of Indiscernibles.” Can you think of a potential example of objects x and y that share all their properties but are still distinct?

But Leibniz's Law itself cannot be captured in FOL. Since it makes a claim about *all properties*, we'd need to have quantifiers that bind variables in *predicate* position to really express Leibniz's Law in full generality, writing something like:

$$\forall x \forall y [x = y \rightarrow \forall P (Px \leftrightarrow Py)]$$

A logic that contains quantifiers like ' $\forall P$ ' is said to be a SECOND-ORDER LOGIC. As its name indicates, FOL is a *first-order* logic. It can express *instances* of Leibniz's Law, concerning particular properties like owning a dog, but it cannot express it in full generality.

One last logical truth involving identity is the following:

$$\exists x x = x$$

This says that there exists at least one thing. This is obviously a controversial case: you might think the claim there exists something (rather than nothing) shouldn't just be a truth of logic. But it is a logical truth in FOL, since FOL requires that quantifier domains have at least one member (see the stipulation in §5.2 above), and it will be a theorem in our system of natural deduction. Systems of logic that avoid having this as a logical truth, and allow for empty domains as well as non-referring names, are called FREE LOGICS.

■ Exercises 5.9

A. Using the following symbolization key:

domain: things

- A: ____ is a card
- B: ____ is black
- C: ____ is a club
- D: ____ is a deuce
- J: ____ is a jack
- O: ____ is one-eyed
- W: ____ is wild

symbolize each sentence in FOL:

1. All clubs are black cards.
2. There are no wild cards.
3. There are at least two clubs.
4. There is more than one one-eyed jack.
5. There are at most two one-eyed jacks.
6. There are two black jacks.
7. There are four deuces.
8. One-eyed jacks and deuces are wild.
9. If one-eyed jacks are wild, then there are exactly two wild cards.

B. Using the following symbolization key:

domain: things

B : ____ is in Farmer Brown's field.
 H : ____ is a horse.
 W : ____ has wings.

symbolize the following sentences in FOL:

1. There are at least three horses.
2. There are at least three things.
3. There is more than one horse in Farmer Brown's field.
4. There are exactly two horses in Farmer Brown's field.
5. There is a single winged horse in Farmer Brown's field, and all other things in the field are wingless.

5.10 The Syntax of FOL

We've been learning to symbolize English sentences in the language of FOL, but it's time to be more precise about the grammar, or syntax of FOL. As in the case of TFL (see §2.9), we will in this section precisely define the notion of a SENTENCE OF FOL.

There are six kinds of symbols that constitute the LEXICON of FOL:

Predicates	A, B, C, \dots, Z
with subscripts, as needed	$A_1, B_1, Z_1, A_2, A_{25}, J_{375}, \dots$
Names	a, b, c, \dots, t
with subscripts, as needed	$a_1, b_{224}, h_7, m_{32}, \dots$
Variables	u, v, w, x, y, z
with subscripts, as needed	$u_1, y_1, z_1, x_2, \dots$
Connectives	$\neg, \wedge, \vee, \rightarrow, \leftrightarrow$
Brackets	$(,)$
Quantifiers	\forall, \exists

We define an EXPRESSION OF FOL as any string of symbols of FOL. Take any of the symbols in the lexicon of FOL and write them down, in any order, and you have an expression. But not all expressions are well-formed *sentences* of FOL, so we'll need rules to tell us which expressions count as sentences.

In the case of TFL, in §2.9, we went straight from the statement of the lexicon to the definition of a sentence of TFL. In FOL, we will have to go via a more indirect route. We will first define the notion of being a *formula* of FOL, and then single out the sentences from this larger class of formulas.

Formulas of FOL

To begin, we define the notion of a *term*:

A TERM is any name or any variable.

Using this, we can next define the notion of an *atomic formula*:

1. If \mathcal{R} is an n -place predicate and t_1, \dots, t_n are terms, then $\mathcal{R}t_1 \dots t_n$ is an ATOMIC FORMULA.
2. If t_1 and t_2 are terms, then $t_1 = t_2$ is an atomic formula.
3. Nothing else is an atomic formula.

Notice that we are again using *metavariables* in this definition (see the discussion in §2.9), albeit cursive ones rather than Greek ones. So here the cursive letter ' \mathcal{R} ' is not itself a predicate of FOL. Rather, it is a symbol of our metalanguage which we are using to talk about arbitrary predicates of FOL. Similarly, the cursive ' t_1 ' is not a term of FOL, but a symbol of the metalanguage that we are using to talk about arbitrary terms (i.e. variables or names) of FOL.

If we let ' F ' be a one-place predicate, ' G ' a three-place predicate, and ' S ' a six-place predicate, then the following all count as atomic formulas of FOL by our definition:

$$\begin{aligned} x &= a \\ Fx \\ Fa \\ Gxay_2 \\ Gaaa \\ Sx_1x_2abyx_6 \end{aligned}$$

Given the notion of an *atomic formula*, we can now recursively define the broader class of FOL *formulas* as follows:

1. Every atomic formula is a formula.
2. If ϕ is a formula, then $\neg\phi$ is a formula.
3. If ϕ and ψ are formulas, then $(\phi \wedge \psi)$ is a formula.
4. If ϕ and ψ are formulas, then $(\phi \vee \psi)$ is a formula.
5. If ϕ and ψ are formulas, then $(\phi \rightarrow \psi)$ is a formula.
6. If ϕ and ψ are formulas, then $(\phi \leftrightarrow \psi)$ is a formula.
7. If ϕ is a formula that contains at least one occurrence of the variable v and does *not* contain either $\forall v$ or $\exists v$, then $\forall v \phi$ is a formula.
8. If ϕ is a formula that contains at least one occurrence of the variable v and does *not* contain either $\forall v$ or $\exists v$, then $\exists v \phi$ is a formula.
9. Nothing else is a formula.

The first few clauses are similar to those from TFL. What's new are clauses (7) and (8), which tell us how to construct quantified formulas. Letting ' F ' be a one-place predicate and ' R ' a two-place predicate, the following all count as formulas of FOL:

Fc	By clause 1
Fx	By clause 1
Rxz	By clause 1
$(Fx \rightarrow Rxz)$	By clause 5
$\forall x(Fx \rightarrow Rxz)$	By clause 7
$(Fy \leftrightarrow \forall x(Fx \rightarrow Rxz))$	By clause 6
$\exists y(Fy \leftrightarrow \forall x(Fx \rightarrow Rxz))$	By clause 8
$\forall z\exists y(Fy \leftrightarrow \forall x(Fx \rightarrow Rxz))$	By clause 7

By contrast, the following are *not* formulas:

$$\begin{aligned} &\forall yRxx \\ &\forall x\exists xRxx \end{aligned}$$

Looking at the first of these, we can say that ' Rxx ' is a formula by clause (1), but ' $\forall yRxx$ ' is not a formula, because it results from attaching the quantifier ' $\forall y$ ' to a formula that does not contain at least one occurrence of the variable ' y ', thus contravening clause (7). Similarly, in the case of the second example, ' Rxx ' is again a formula, and by clause (8) ' $\exists xRxx$ ' is also a formula. But ' $\forall x\exists xRxx$ ' is not a formula because we've attached the quantifier ' $\forall x$ ' to a formula ' $\exists xRxx$ ' that *already* contains a quantifier ' $\exists x$ ' involving the same variable, which contravenes clause (7).

These constraints have the effect of preventing the kind of *variable clashes* we discussed in §5.8. And in fact, we can now give a precise definition of the notion of the *scope* of a quantifier (or other logical operator) in terms of the notion of *main logical operator*:

The MAIN LOGICAL OPERATOR in a formula is the operator that was introduced most recently when constructing that formula according to the syntactic rules of FOL.

The SCOPE of an operator in a formula is the subformula for which it is the main logical operator.

Since quantifiers are a kind of logical operator, this definition covers the scope of quantifiers alongside the scope of truth functional connectives. The scope of the various quantifiers in one of our earlier examples can be illustrated as follows:

$$\begin{array}{c} \text{scope of } \forall z \\ \overbrace{\hspace{10em}} \\ \text{scope of } \exists y \\ \overbrace{\hspace{10em}} \\ \text{scope of } \forall x \\ \overbrace{\hspace{10em}} \\ \forall z\exists y(Fy \leftrightarrow \forall x(Fx \rightarrow Rxz)) \end{array}$$

Returning to the problematic example ' $\forall x\exists xRxx$ ' from above, we can now see that the problem with it is that the quantifier ' $\forall x$ ' has, inside its scope, another quantifier ' $\exists x$ ' involving the same variable, leading to a variable clash. This is ruled out by our clause (7).

Sentences of FOL

With these pieces in place, we're now ready to define the notion of a *sentence* of FOL. To see why we need to distinguish sentences from mere formulas of FOL, recall that logic is concerned with *statements*: sentences that can be either true or false. And many formulas are not true or false. For example, consider the formulas Fc and Fx , and suppose we have the following symbolization key:

domain: people
 F : ____ is a philosopher
 c : Confucius

The formula ' Fc ' can be assigned a truth value: we just ask ourselves whether the person ' c ' refers to is a philosopher. Since Confucius is a philosopher, ' Fc ' is true. By contrast, ' Fx ' has no truth value. After all, ' x ' is just a variable, and doesn't name any specific object in the domain.

Of course, if we put an existential quantifier out front to obtain ' $\exists xFx$ ', we now have something that's capable of being true or false, since this now says that at least one person is a philosopher. The point is that we need to *bind* the variable in ' Fx ' with a quantifier to obtain something true or false.

Since we want all sentences of FOL to be either true or false, we need to exclude formulas like ' Fx ' from the class of sentences. We can do this by giving a precise definition of the notions of *bound* and *free* variables that we already informally discussed in §5.2:

A quantifier $\exists v$ or $\forall v$ BINDS every occurrence of the variable v that appears in its scope. Variables that aren't bound by any quantifier are FREE.

An OPEN FORMULA is one that contains at least one free variable.

For example, consider the following formula, which has a conditional as its main operator:

$$(\forall x(Ex \vee Dy) \rightarrow \exists z(Ex \rightarrow Lzx))$$

The universal quantifier ' $\forall x$ ' has scope over the antecedent ' $\forall x(Ex \vee Dy)$ ', so the first ' x ' in the formula is bound. However, the ' y ' is free. The scope of the existential quantifier ' $\exists z$ ' is the consequent ' $\exists z(Ex \rightarrow Lzx)$ ', so ' z ' is bound. But the ' x 's in this subformula are free in the formula as a whole, since they occur outside the scope of the universal quantifier ' $\forall x$ '. Since the formula as a whole contains free variables, it is an *open formula*.

We could transform it into a *closed formula*, one containing no free variables, by giving ' $\forall x$ ' scope over the entire formula, and adding a quantifier to bind the ' y ' variable, as in for example:

$$\forall x((Ex \vee \exists yDy) \rightarrow \exists z(Ex \rightarrow Lzx))$$

This is now a *sentence* of FOL. In general, we give the following definitions:

A SENTENCE of FOL is any formula of FOL that contains no free variables. Sentences are also called CLOSED FORMULAS.

By requiring that every variable in a sentence be bound, we ensure that all sentences of FOL are capable of being true or false.

■ Exercises 5.10

A. Determine whether each string below is a grammatical formula. If it is, say what the main operator is and whether it has any free (unbound) variables:

1. $\forall x(Fx)$
2. $\forall x(Fx \wedge Gx)$
3. $\forall y(Fx \wedge Gx)$
4. $\forall xFx \wedge Gx \rightarrow Hy$
5. $(\forall x(Fx \wedge Gx) \rightarrow Hy)$
6. $\forall x((Fx \wedge Gx) \rightarrow Hy)$
7. $\exists y\forall x((Fx \wedge Gx) \rightarrow Hy)$
8. $\forall xRxy$
9. $\exists y\forall xRxy$
10. $\exists x\forall xRxx$

B. Identify which variables are bound and which are free.

1. $\exists xLxy \wedge \forall yLyx$
2. $\exists x(Lxy \wedge \forall yLyx)$
3. $\forall xAx \wedge Bx$
4. $\forall x(Ax \wedge Bx) \wedge \forall y(Cx \wedge Dy)$
5. $\forall x\exists y[Rxy \rightarrow (Jz \wedge Kx)] \vee Ryx$
6. $\forall x_1(Mx_2 \leftrightarrow Lx_2x_1) \wedge \exists x_2Lx_3x_2$

C. For each of the following formulas, rewrite the formula so that all variables are bound, or leave it alone if all variables are already bound. You may *only* change these formulas by adding and removing parentheses (i.e. you can't add more quantifiers).

1. $\exists xCx \rightarrow \forall yRxy$
2. $\forall x(\exists yRxy \rightarrow Rxx)$
3. $\forall x(\exists yRxy \rightarrow Ryx)$
4. $\forall x((\exists yRxy \wedge Fy) \rightarrow Rxx)$
5. $\exists x\forall y(\forall zLzx \rightarrow \forall u(Jy \wedge Lzu))$
6. $\exists x(\forall y\forall zLzy \rightarrow \forall u(Ju \wedge Lxu))$

Natural deduction for FOL

6

Since sentences in FOL can contain any of the TFL connectives, proofs in FOL will take over all the natural deduction rules from TFL that we studied in chapter 4, as well as the all derived TFL rules introduced in §4.11. We will also continue to use the same proof theoretic notions, in particular, the symbol ‘ \vdash ’. So:

$$\varphi_1, \dots, \varphi_n \vdash \psi$$

will continue to mean that ψ is *provable* from $\varphi_1, \dots, \varphi_n$, i.e. that there exists a natural deduction proof which ends with ψ and whose premises, or more generally, whose undischarged assumptions, include at most $\varphi_1, \dots, \varphi_n$.

What we will need to add to our natural deduction system are rules to govern the new logical symbols that are specific to FOL: the quantifiers ‘ \forall ’ and ‘ \exists ’, and the identity predicate ‘ $=$ ’. As in the case of TFL, there will be an *introduction* and an *elimination* rule associated with each of these logical symbols. Again, we will of course want to make sure that the rules we add end up producing a proof system that is both sound and complete:

SOUNDNESS: If $\varphi_1, \dots, \varphi_n \vdash \psi$ then $\varphi_1, \dots, \varphi_n \models \psi$

COMPLETENESS: If $\varphi_1, \dots, \varphi_n \models \psi$ then $\varphi_1, \dots, \varphi_n \vdash \psi$

Soundness ensures that if ψ is provable from $\varphi_1, \dots, \varphi_n$, then $\varphi_1, \dots, \varphi_n$ logically entail the conclusion ψ . And completeness ensures us that if $\varphi_1, \dots, \varphi_n$ logically entail the conclusion ψ , then ψ is also provable from premises $\varphi_1, \dots, \varphi_n$. We will formally define the concept of logical entailment for FOL in Chapter 7.

As in the case of TFL, we won’t demonstrate that our proof system is in fact sound and complete — this would require a *meta-logical* proof, a proof *about* our proof system. You’ll just have to take my word for it that soundness and completeness hold for our system. The important point for our purposes is that, due to the soundness of our proof system, we can use proofs to show that arguments are valid: if we can give a proof of a conclusion ψ from some premises $\varphi_1, \dots, \varphi_n$, then we can be sure that the premises entail that conclusion, and that the argument $\varphi_1, \dots, \varphi_n \therefore \psi$ is therefore valid.

Similarly, we will continue to write:

$$\vdash \varphi$$

to mean that φ is a THEOREM of our proof system, i.e. a sentence that is provable using no premises, or undischarged assumptions. Given soundness, theorems of our proof system are guaranteed to be *logical truths*, so we can also use natural deduction proofs to show that something is a logical truth.

6.1 Universal elimination

The elimination rule for the universal quantifier has a very simple idea behind it: from the claim that everything is F , you can infer that any particular thing is F . You name it; it's F . So from $\forall xFx$ we can infer Fa and Fb and Fc , and so on. Similarly with two-place predicates:

1	$\forall xRxd$	
2	Rad	$\forall E$ 1

We obtained line 2 by dropping the universal quantifier ' $\forall x$ ' and replacing every occurrence of the variable ' x ' that this quantifier bound with ' a '. Equally, the following is fine:

1	$\forall xRxd$	
2	Rdd	$\forall E$ 1

Here we obtained line 2 by dropping the universal quantifier and replacing every occurrence of ' x ' with ' d '. We could have done the same with any other name we wanted: ' b ', ' c ', ' e ', you name it. The idea is simple: if, as per line 1, *everything* bears relation R to d , then a does, and so does d itself, and so does b or c , or any other thing.

Both ' Rad ' and ' Rdd ' are said to be *instances* of the quantified sentence ' $\forall xRxd$ ', as are ' Rbd ', ' Rcd ', ' Red ' and so on. In general, this notion is defined as follows:

Given a universally quantified FOL sentence $\forall v\phi(\dots v \dots)$ its INSTANCES are all those FOL sentences $\phi(\dots c \dots)$ that are obtained by dropping the quantifier $\forall v$ and replacing *every* occurrence of the variable v bound by that quantifier with some name c .

The universal elimination rule $\forall E$ then just says that from a universal sentence you may infer any of its instances:

m	$\forall v\phi(\dots v \dots)$	
	$\phi(\dots c \dots)$	$\forall E$ m

Here, as well as in the definition of the notion of an instance, we are using the notation $\phi(\dots v \dots)$ to represent any FOL formula ϕ that contains at least one occurrence of the variable v , and $\phi(\dots c \dots)$ to represent the result of replacing every occurrence of the variable v in that formula with some name c .

Notice that $\phi(\dots v \dots)$ may of course itself contain connectives and even other quantifiers. So if we begin with the complex universally quantified sentence:

$$\forall x(Rax \wedge \exists yRxy)$$

the following would be instance of it, and could be inferred via $\forall E$:

$$\begin{array}{l} Raa \wedge \exists yRay \\ Rab \wedge \exists yRby \end{array}$$

The first results from replacing the variable ‘ x ’ with the name ‘ a ’, and the second from replacing ‘ x ’ with ‘ b ’. By contrast, the following are *not* instances:

$$\begin{array}{l} Rab \wedge \exists yRxy \\ Rac \wedge \exists yRdy \end{array}$$

The first isn’t because we forgot to replace the second occurrence of ‘ x ’ with the name ‘ b ’ (thereby leaving this ‘ x ’ as a free variable), and the second isn’t because we replaced different occurrences of ‘ x ’ with different names. Again, the important point is that when using $\forall E$, the sentence you infer must be an *instance* of the universally quantified sentence you applied the rule to.

It’s also important to emphasize that (as with every elimination rule) you can only apply the $\forall E$ rule when the universal quantifier is the *main logical operator*. So the following is *not* a legitimate use of $\forall E$:

$$\begin{array}{l|l} 1 & \forall xBx \rightarrow Bc \\ \hline 2 & Bb \rightarrow Bc \end{array} \quad \text{No! Illegitimate use of } \forall E \text{ 1}$$

This is illegitimate because ‘ $\forall x$ ’ is not the main logical operator on line 1. If line 1 had instead been ‘ $\forall x(Bx \rightarrow Bc)$ ’, then the quantifier would have been the main operator, and it would have been legitimate to infer ‘ $Bb \rightarrow Bc$ ’.

Another way to put it is that this use of $\forall E$ is illegitimate because ‘ $Bb \rightarrow Bc$ ’ is not an *instance* of ‘ $\forall xBx \rightarrow Bc$ ’. This is because the notion of an instance only applies to sentences that have a quantifier as their main operator. It just doesn’t make sense to talk about the “instances” of ‘ $\forall xBx \rightarrow Bc$ ’, because this sentence has a conditional as its main operator, and conditionals do not have instances. By contrast, ‘ $\forall x(Bx \rightarrow Bc)$ ’ does have a quantifier as its main operator, and ‘ $Bb \rightarrow Bc$ ’ is one of its instances, so we can infer it via $\forall E$.

Using this rule, we can, for example, show that the following argument that we looked at way back in §1.1 is valid:

All rabbits are mammals.
Bugs Bunny is a rabbit.
 \therefore Bugs Bunny is a mammal.

First, we provide the following FOL symbolization (using the obvious symbolization key):

$\forall x(Rx \rightarrow Mx)$
 Rb
 $\therefore Mb$

And then we can give a simple natural deduction to show that the conclusion follows:

1	$\forall x(Rx \rightarrow Mx)$	Premise
2	Rb	Premise
3	$Rb \rightarrow Mb$	$\forall E$ 1
4	Mb	$\rightarrow E$ 3, 2

6.2 Existential introduction

The idea behind the existential introduction rules is again very simple: from the claim that some particular thing is F, you can infer that there exists at least one thing that is F. So we ought to allow:

1	Raa	
2	$\exists xRax$	$\exists I$ 1

Here, we have replaced just one occurrence of the name ‘*a*’ with a variable ‘*x*’, and then existentially quantified over it. Equally, we would have done this:

1	Raa	
2	$\exists xRxx$	$\exists I$ 1

Here we have replaced both occurrences of the name ‘*a*’ with the variable ‘*x*’, and then existentially generalized. Both kinds of inferences are fine: if Narcissus loves himself, ‘*Lnn*’, then we can infer that there exists someone who loves themselves, ‘ $\exists xLxx$ ’, but we can also infer that there exists someone who Narcissus loves, ‘ $\exists xLnx$ ’, and that there is someone who loves Narcissus, ‘ $\exists xLxn$ ’. All can be inferred from ‘*Lnn*’ by $\exists I$.

Another way to put this is to note that ‘*Lnn*’ is an *instance* of each of the quantified sentences ‘ $\exists xLxx$ ’, ‘ $\exists xLnx$ ’, ‘ $\exists xLxn$ ’. The notion of an instances as applied to existential sentences is the same as it is for universal ones:

Given any existentially quantified FOL sentence $\exists v\phi(\dots v \dots)$ its INSTANCES are all those FOL sentences $\phi(\dots c \dots)$ that are obtained by dropping the quantifier $\exists v$ and replacing *every* occurrence of the variable *v* bound by that quantifier with some name *c*.

The existential introduction rule $\exists I$ then just says that from a sentence containing one or more occurrences of a name *c*, we may infer *any* existential sentence of which that original sentence is an instance:

m	$\phi(\dots c \dots)$	
	$\exists v\phi(\dots v \dots)$	$\exists I$ m

To take a slightly more complex example, from:

$$Fa \wedge \forall yRay$$

we could infer any of the following using $\exists I$:

$$\begin{aligned} &\exists x(Fx \wedge \forall yRay) \\ &\exists x(Fa \wedge \forall yRxy) \\ &\exists x(Fx \wedge \forall yRxy) \end{aligned}$$

because ' $Fa \wedge \forall yRay$ ' is an instance of any of these. So again, we can existentially generalize on just one occurrence of ' a ', or on both.

On the other hand ' $\exists x(Fx \wedge \forall yRxy)$ ' could *not*, for example, be inferred from:

$$Fa \wedge \forall yRby$$

because this is not one of its instances. What's gone wrong here is that we tried to generalize on two different names, ' a ' and ' b ', at once, which isn't allowed.

What we could have done instead is to generalize on each of the two names separately, via successive uses of $\exists I$:

1	$Fa \wedge \forall yRby$	Premise
2	$\exists x(Fx \wedge \forall yRby)$	$\exists I$ 1
3	$\exists z\exists x(Fx \wedge \forall yRzy)$	$\exists I$ 2

This is fine because line 1 is an instance of line 2, and line 2 is in turn an instance of line 3.

It's important to notice that in moving from line 2 to line 3 here, it was essential that we introduced an existential quantifier involving a *new* variable ' z ', which did not yet appear in line 2. The following would *not* have been legitimate:

1	$Fa \wedge \forall yRby$	Premise
2	$\exists x(Fx \wedge \forall yRby)$	$\exists I$ 1
3	$\exists y\exists x(Fx \wedge \forall yRyy)$	illegitimate use of $\exists I$ 2

The expression on line 3 involves a *variable clash* (see §5.8) between the newly introduced quantifier ' $\exists y$ ' and the quantifier ' $\forall y$ ' that occurs in its scope. It therefore does not even count as a *formula*, let alone a *sentence* of FOL by the definition we gave in §5.10. For the same reason, we could not have inferred ' $\exists x\exists x(Fx \wedge \forall yRxy)$ ' — this again involves a variable clash, and therefore isn't an FOL formula.

Here's a simple proof that combines our two new quantifier rules to show that $\forall xFx \therefore \forall y(Fy \rightarrow Gy)$ is valid:

1	$\forall xFx$	Premise
2	$\forall y(Fy \rightarrow Gy)$	Premise
3	Fa	$\forall E$ 1
4	$Fa \rightarrow Ga$	$\forall E$ 2
5	Ga	$\rightarrow E$ 4, 3
6	$Ga \wedge Fa$	$\wedge I$ 5, 3
7	$\exists z(Gz \wedge Fz)$	$\exists I$ 6

Notice that from line 1 I could have inferred some other instance by $\forall E$ instead, like ' Fb ', and similarly, from line 2 I could have inferred any other instance, like ' $Fc \rightarrow Gc$ '. But if I had used different names in the two instances like this, I could then not have applied $\rightarrow E$ to them. So although $\forall E$ lets you infer any instance, in the context of a proof you'll usually have to infer some specific instance.

■ Exercises 6.2

A. For each of the following FOL sentences, determine what its a -instance and its b -instance are:

1. $\forall x(Fx \wedge Gx)$
2. $\exists x(Fx \wedge Gx)$
3. $\forall x(Fx \rightarrow \exists yRyx)$
4. $\exists x\forall yLxy$
5. $\forall x(Rxa \rightarrow \exists y(Rxy \wedge Ryx))$
6. $\exists x(Lbx \leftrightarrow Lxa)$
7. $\forall xRxa \rightarrow \exists yRby$

B. Given natural deduction proofs for the following (for these you'll only have to use $\forall E$ and $\exists I$, in addition to TFL rules of course):

1. $\forall x\forall yRxy \vdash Raa \wedge Rab$
2. $\forall x(Fx \rightarrow \exists yGy) \vdash \forall xFx \rightarrow \exists yGy$
3. $\forall xRax, \forall x\forall y(Rxy \rightarrow Ryx) \vdash Rba$
4. $\forall x(Fx \wedge \neg Gx), (Gc \vee Hd) \vdash (Hd \wedge Fd)$
5. $\forall x(Fx \wedge Gx), \forall yHy \vdash \exists z(Hz \wedge Gz)$
6. $\forall xFx, \forall y(Fy \rightarrow Gy) \vdash \exists x(Gx \wedge \exists yFy)$
7. $\forall x(Fx \rightarrow \forall yGy), Fc \vdash \exists x(Fx \wedge Gx)$
8. $\forall x\forall yRxy \vdash \exists xRxx$
9. $\forall x(Fx \rightarrow Gx) \vdash \forall xFx \rightarrow \exists yGy$
10. $\forall xRxx \vdash \exists x\exists yRxy$
11. $\exists xFx \rightarrow \forall yGy \vdash \forall zFz \rightarrow \exists zGz$
12. $\forall x(Fx \rightarrow Gx), \neg Gc \vdash \exists x\exists y(\neg Fx \vee Gy)$
13. $\vdash \exists x(Fx \vee \neg Fx)$

6.3 Universal introduction

Suppose you had shown of each particular thing that it is F (and that there are no other things to consider). Then you would be justified in claiming that *everything* is F . This could be used to motivate the following proof rule: if you had established each and every instance of ' $\forall xFx$ ' holds, then you can infer ' $\forall xFx$ '.

Unfortunately, that rule would be utterly unusable. To establish every single instance of ' $\forall xFx$ ' would require proving ' Fa ', ' Fb ', ..., ' Fj_2 ', ..., ' Fr_{791} ', ... and so on. Since there are infinitely many names in FOL, this process would never end! So we need to be more cunning in coming up with our rule for introducing universal quantifiers.

We can motivate our rule by considering the following:

$$\forall x(Fx \wedge Gx) \therefore \forall xFx$$

This argument is obviously valid: if everything is both F and G , then everything is F . But how could we prove this? Suppose we begin a proof like this:

1	$\forall x(Fx \wedge Gx)$	Premise
2	$Fa \wedge Ga$	$\forall E$ 1
3	Fa	$\wedge E$ 2

We have proven ' Fa ', an instance of the conclusion ' $\forall xFx$ ' that we're aiming for. But of course, nothing stops us from using $\forall E$ in combination with $\wedge E$ in the very same way to prove ' Fb ', ' Fc ', ..., ' Fj_2 ', ..., ' Fr_{791} ', ..., and so on until we run out of space, time, or patience. So it's clear that from our premise, we could in principle prove Fc for *any* name c , that is, we could in principle prove *every* instance of our goal ' $\forall xFx$ '. So we should be entitled to infer ' $\forall xFx$ ' by $\forall I$. It's just that we can't *actually* prove every instance, since our proof would never end.

This leads to the following idea: we should be allowed to infer the universal sentence ' $\forall xFx$ ' by the rule of $\forall I$ if we are able prove an *arbitrary* instance Fc , one that involves some arbitrary name c . For if the name c is truly arbitrary, then it doesn't matter that we specifically proved this particular instance Fc — we could have picked any other name instead, and thereby proven any other instance of the universal sentence we're aiming for.

Our universal introduction rule $\forall I$ implements this idea via a “flagged subproof”:

m	c	Flag
	\vdots	
n	$\phi(\dots c \dots)$	
	$\forall v \phi(\dots v \dots)$	$\forall I$ $m-n$

The Flag-ed name c may not occur outside the subproof (including in the inferred sentence $\forall v \phi(\dots v \dots)$ itself!)

The Flag step on line m is just a way of officially signaling that the name c is being introduced as an arbitrary name in the proof, one that we'll use to prove an arbitrary instance of the universal sentence $\forall v \phi(\dots v \dots)$ that we are aiming for. Our proof that $\forall x(Fx \wedge Gx) \vdash \forall xFx$ can now be presented as follows:

1	$\forall x(Fx \wedge Gx)$	Premise
2	a	Flag
3	$Fa \wedge Ga$	$\forall E$ 1
4	Fa	$\wedge E$ 3
5	$\forall xFx$	$\forall I$ 2–4

Again, the idea is that although we only proved the one instance, ' Fa ', we are allowed to infer the universal sentence ' $\forall xFx$ ' because that instance was arbitrary — we could have just as easily proven any other instance we pleased.

The Flag-ing constraint listed at the bottom of the rule — that the Flag-ed name may not occur outside the subproof — is crucial, because it is what insures that the name we've picked is truly arbitrary.¹ To see the constraint in action, consider this terrible argument:

Everyone loves Beyonce. Therefore everyone loves themselves.

This argument is obviously not valid. We might symbolize it as:

$$\forall xLxb \therefore \forall xLxx$$

Now, suppose we tried to offer the following “proof” to vindicate this argument:

1	$\forall xLxb$	Premise
2	b	Flag
3	Lbb	$\forall E$ 1
4	$\forall xLxx$	Illegitimate use of $\forall I$!

It would be bad if this proof were legitimate, since the conclusion doesn't follow. What makes it illegitimate is that the Flag-ed name ' b ' occurs outside the subproof, namely in our premise on line 1. Since ' b ' occurs in the premise, it doesn't have the status of an arbitrary name, and the sentence ' Lbb ' we proved on line 3 doesn't qualify as an arbitrary instance of our goal ' $\forall xLxx$ ': we could *not* have proven any *other* instance of ' $\forall xLxx$ ', like e.g. ' Laa ' or ' Ljj ', from our premise.

Notice that the flagged constant also cannot occur in the universal sentence that's being inferred via $\forall I$, since it occurs outside the subproof. Consider the following, equally terrible argument:

¹This constraint is actually a little more restrictive than strictly necessary. It would be alright if the Flag-ed constant occurred outside the subproof, as long as it doesn't occur in any earlier *premise or undischarged assumption*. But the constraint as we've here formulated it has the advantage of being concise and easier to remember.

Everyone loves themselves. Therefore everyone loves Beyonce.

which we could symbolize as: $\forall xLxx \therefore \forall xLxb$. Now suppose we tried to prove it as follows:

1	$\forall xLxx$	Premise
2	b	Flag
3	Lbb	$\forall E$ 1
4	$\forall xLxb$	illegitimate use of $\forall I$

Again, this proof had better not be legitimate, since the conclusion does not follow. And it isn't legitimate: the Flag-ed name ' b ' occurs outside of the subproof, in the inferred universal sentence ' $\forall xLxb$ '. Again, although ' Lbb ' is an *instance* of $\forall xLxb$, it doesn't qualify as an *arbitrary* instance because we could not have proven any other instance of ' $\forall xLxb$ ' — like e.g. ' Lab ' or ' Ljb ' — from our premise.

For an example of a correct use of $\forall I$, consider how we might prove that $\forall z(Gz \rightarrow Gz)$ is a theorem. To prove this, we have to open up a flagged subproof inside of which we prove some arbitrary instance of this sentence, such as ' $Gd \rightarrow Gd$ ', as follows:

1	d	Flag
2	Gd	Assumption (for $\rightarrow I$)
3	Gd	Reit 2
4	$Gd \rightarrow Gd$	$\rightarrow I$ 2–3
5	$\forall z(Gz \rightarrow Gz)$	$\forall I$ 4

The constraints on the legitimate application of $\forall I$ are met, since the name ' d ' does not occur outside the subproof. Here ' $Gd \rightarrow Gd$ ' qualifies as an arbitrary instance of ' $\forall z(Gz \rightarrow Gz)$ ': we could just as well have flagged some other name, say ' a ', and instead proved the instance ' $Ga \rightarrow Ga$ ' using that name. There was nothing special about ' d '.

■ Exercises 6.3

A. Give proofs for the following:

1. $\forall x\forall y(Gy \rightarrow Fx) \vdash \forall x(\forall yGy \rightarrow Fx)$
2. $\forall x(Fx \wedge Gx) \vdash \forall x(Fx \wedge Ga)$
3. $\forall xFx \vee \forall xGx \vdash \forall x(Fx \vee Gx)$
4. $\forall xLxx \vdash \forall x\exists yLxy$
5. $\forall x\forall yLxy \vdash \forall xLxx$
6. $\forall x\forall y(Rxy \rightarrow \neg Ryx) \vdash \forall x\neg Rxx$
7. $\neg\exists x(Fx \wedge Gx) \vdash \forall x(Fx \rightarrow \neg Gx)$
8. $\forall x(Fx \rightarrow \forall yGy) \vdash \forall x\forall y(Fx \rightarrow Gy)$
9. $\forall x\forall y(Rxy \rightarrow Ryx), \forall x\forall y\forall z((Rxy \wedge Ryz) \rightarrow Rxz) \vdash \forall x\forall y\forall z((Rxy \wedge Rxz) \rightarrow Ryz)$
10. $\neg\forall xFx \vdash \exists x\neg Fx$

6.4 Existential elimination

Suppose we know that *something* is F. The problem is that simply knowing this does not tell us which particular thing is F. So from ' $\exists xFx$ ' we cannot immediately infer ' Fa ', or ' Fd ', or any other instance of the sentence. What can we do? How can we deduce anything from existential premises?

Well, suppose we know that something is F, and that everything which is F is G. In English, we might pursue the following line of reasoning:

Since something is F, there is some particular thing which is F. We do not know anything about it, other than that it's F, but for convenience, let's call it 'Obbie'. So: Obbie is F. Since everything which is F is G, it follows that Obbie is G. But since Obbie is G, it follows that *something* is G. And nothing depended on which object, exactly, Obbie was. Therefore, something is G.

We can capture this reasoning pattern in a proof as follows:

1	$\exists xFx$	Premise
2	$\forall x(Fx \rightarrow Gx)$	Premise
3	Fo	Assumption (flag <i>o</i>)
4	$Fo \rightarrow Go$	$\forall E$ 2
5	Go	$\rightarrow E$ 4, 3
6	$\exists xGx$	$\exists I$ 5
7	$\exists xGx$	$\exists E$ 1, 3–6

Breaking this down: we started by writing down our premises. At line 3, we then made an additional assumption: ' Fo '. The idea here is that premise 1 tell us that *something* is an F . So on line 3 we introduce some arbitrary name ' o ' for that thing, Flag it as arbitrary to the right, and write down the corresponding instance of the existential premise 1. The name we picked is arbitrary, since we've assumed nothing about the object named by ' o ' other than that the predicate ' F ' is true of it. On the basis of the assumption Fo , we can then establish ' $\exists xGx$ '. Since nothing depended on which specific object ' o ' names, our reasoning pattern is perfectly general: we could equally well have proven ' $\exists xGx$ ' by using any other name on line 3. We can therefore discharge the assumption ' Fo ' on line 3, and simply infer ' $\exists xGx$ ' on its own.

Putting this together, we obtain the existential elimination rule ($\exists E$):²

²As in the case of $\forall I$, our formulation of the flagging constraint at the bottom is more restrictive than strictly necessary. It would be alright if the flag-ed constant c occurred outside the subproof, as long as it doesn't occur in any earlier *premise or undischarged assumption*.

m	$\exists x \phi(\dots x \dots)$	
i	$\phi(\dots c \dots)$	Assumption (flag c)
	\vdots	
j	ψ	
	ψ	$\exists E\ m, i-j$

The Flag-ed name c may not occur outside the subproof (including in the original existential $\exists x \phi(\dots x \dots)$ and the inferred sentence ψ !)

So in general, to prove some sentence ψ from an existential sentence $\exists x \phi(\dots x \dots)$, what we do is flag some arbitrary name c , assume an instance of the existential sentence using this name c , and then prove our goal ψ from that instance. Finally, we discharge our assumption and infer ψ on its own via $\exists E$.

As with universal introduction, the Flag-ing constraint on the name c that's listed at the bottom is very important. To see why, consider the obviously bad argument:

Borges is a librarian. Someone is not a librarian. So Borges is both a librarian and not a librarian.

We might symbolize this as follows:

$$Lb, \exists x \neg Lx \therefore Lb \wedge \neg Lb$$

This is clearly a terrible argument: it presumes that the “someone” who is not a librarian according to the second premise is the individual Borges mentioned in the first premise (which can't be, since Borges is a librarian and the “someone” from premise 2 isn't). Now, suppose we tried to offer the following “proof” to vindicate this argument:

1	Lb	Premise
2	$\exists x \neg Lx$	Premise
3	$\neg Lb$	Assumption (flag b)
4	$Lb \wedge \neg Lb$	$\wedge E\ 1, 3$
5	$Lb \wedge \neg Lb$	No! Illegitimate attempt to use $\exists E\ 2, 3-4$

It would be a bad thing if we could prove the conclusion like this, since it doesn't follow from the premises! And the Flag-ing constraint is what prevents us from doing so: the use of $\exists E$ on the last line is not legitimate, because the Flag-ed name on line 3, namely ‘ b ’, appears outside the subproof, on lines 1 and 5.

We could avoid part of the problem by existentially generalizing line 4 in the subproof to obtain $\exists x (Lx \wedge \neg Lx)$, before discharging our assumption:

1	Lb	Premise
2	$\exists x \neg Lx$	Premise
3	$\neg Lb$	Flag b
4	$Lb \wedge \neg Lb$	$\wedge E$ 1, 3
5	$\exists x (Lx \wedge \neg Lx)$	$\exists I$ 4
6	$\exists x (Lx \wedge \neg Lx)$	No! Illegitimate attempt to use $\exists E$ 2, 3–5

Now, the name ‘ b ’ no longer occurs below the subproof. But this is no better. If it were legitimate, this proof would vindicate an argument of the following sort:

Borges is a librarian. Someone is not a librarian. Therefore someone both is and is not a librarian.

This is clearly a bad argument: we can’t assume that the “someone” who is not a librarian according to the second premise is, specifically, the individual Borges mentioned in the first premise. And again, the Flag-ing constraint rules out our supposed “proof”: the use of $\exists E$ on the the last line is not legitimate because although the Flag-ed name ‘ b ’ doesn’t occur on line 6 any longer, it does still occurs outside the subproof, namely in the premise on line 1.

The overarching problem with both proofs is that because the name ‘ b ’ already occurs in one of our premises, it does not have the status of an arbitrary name in our proof, and therefore can’t be used as an arbitrary name for whatever object premise 2 tell us is not a librarian. The moral is: *if you want to squeeze information out of an existential quantifier, choose a new name for your substitution instance.* That way, you will meet the constraints on the rule for $\exists E$.

Let’s work through a more complicated proof that requires both $\exists E$ and $\forall I$ at the same time. We’ll show that the following is valid:

$$\forall x \exists y Lxy, \forall x \forall y (Lxy \rightarrow Lyx) \therefore \forall x \exists y Lyx$$

If we read ‘ L ’ as ‘loves’, this argument says that if everyone loves someone, and loves is always reciprocated — in the sense that if x loves y , then y loves x back — it follows that everyone is loved by someone. We can show that this is valid with the following proof:

1	$\forall x \exists y Lxy$	Premise
2	$\forall x \forall y (Lxy \rightarrow Lyx)$	Premise
3	a	Flag
4	$\exists y Lay$	$\forall E$ 1
5	Lab	Assumption (flag b)
6	$\forall y (Lay \rightarrow Lya)$	$\forall E$ 2
7	$Lab \rightarrow Lba$	$\forall E$ 6
8	Lba	$\rightarrow E$ 7, 5
9	$\exists y Lya$	$\exists I$ 8
10	$\exists y Lya$	$\exists E$ 4, 5–9
11	$\forall x \exists y Lyx$	$\forall I$ 3–10

This is a relatively complex proof, so let's think through it systematically. As usual, we work backward from the conclusion we're aiming for: we are trying to prove ' $\forall x \exists y Lyx$ ', i.e. that everyone is loved by someone. Since this is a universal sentence, we use $\forall I$ as our overall strategy: we pick an arbitrary name, say ' a ', and open up a subproof where we Flag ' a ', and make it our goal to prove the ' a '-instance of our conclusion, $\exists y Lya$, which says that someone loves a . If we are able to complete the subproof, we're allowed to infer since a is loved by someone, and a was arbitrary, everyone is loved by someone.

So what do our premises imply about our object a ? Well, premise 1 says that everyone loves someone, so we can infer by $\forall E$ that a in particular loves someone, as we did on line 4. Since a loves *someone*, we can give that someone a name, say ' b ', in order to reason about them. So we can say a loves b . The way this works in our proof is that given the existential sentence ' $\exists y Lay$ ' on line 4, we assume Lab as an arbitrary instance of it on line 5.

Next, premise 2 tells us that love is reciprocal. So given that a loves b , we can conclude that b loves a . In our proof, we did this by obtaining line 7 via two steps of $\forall E$ on premise 2, and then doing $\rightarrow E$ on that. Alright: so given that b loves a , we can conclude that a is loved by *someone*, as on line 9. And at this point, having gotten rid of the name ' b ', we can pop out of our subproof by $\exists E$. And finally, since a was arbitrary, we can pop out of our Flag-ed subproof and conclude by $\forall I$ that *everyone* is loved by someone, as on line 11.

■ Exercises 6.4

A. Explain why these two 'proofs' are incorrect.

1	$\forall x Rxx$	Premise
2	a	Flag
3	Raa	$\forall E$ 1
4	$\forall y Ray$	$\forall I$ 2-3

1	$\forall x \exists y Rxy$	Premise
2	$\exists y Ray$	$\forall E$ 1
3	Raa	Assumption
4	$\exists x Rxx$	$\exists I$ 3
5	$\exists x Rxx$	$\exists E$ 2, 3–4

B. The following three proofs are missing their citations (rule and line numbers). Add them, to turn them into full proofs.

1	$\forall x \exists y (Rxy \vee Ryx)$
2	$\forall x \neg Rmx$
3	$\exists y (Rmy \vee Rym)$
4	$Rma \vee Ram$
5	$\neg Rma$
6	Ram
7	$\exists x Rxm$
8	$\exists x Rxm$

1	$\forall x (Jx \rightarrow Kx)$
2	$\exists x \forall y Lxy$
3	$\forall x Jx$
4	$\forall y Lay$
5	Laa
6	Ja
7	$Ja \rightarrow Ka$
8	Ka
9	$Ka \wedge Laa$
10	$\exists x (Kx \wedge Lxx)$
11	$\exists x (Kx \wedge Lxx)$

C. Provide a proof of each claim.

- $\forall x (Ax \rightarrow Bx), \exists x Ax \vdash \exists x Bx$
- $\exists x (Fx \wedge \exists y \neg Gy) \vdash \exists x (\neg Gx \wedge \exists y Fy)$
- $\exists x (Fx \rightarrow Ga) \vdash \forall x Fx \rightarrow Ga$
- $\exists x \neg Fx \vdash \neg \forall x Fx$
- $\forall x \forall y (Rxy \rightarrow Fx) \vdash \forall x (\exists y Rxy \rightarrow Fx)$
- $\exists x (Fx \rightarrow \forall y Rxy) \vdash \exists x \forall y (Fx \rightarrow Rxy)$
- $\forall x (Fx \rightarrow \forall y \neg Fy) \vdash \neg \exists x Fx$
- $\exists x \exists y Rxy \vdash \exists y \exists x Rxy$
- $\exists y (\forall x (Gx \rightarrow Gy) \wedge \forall z (Gy \rightarrow Gz)), \exists x Gx \vdash \forall x Gx$
- $\forall x \forall y Lxy \vdash \forall x (\exists y Lxy \wedge \exists y Lyx)$
- $\forall x \exists y Lxy, \forall x \forall y (Lxy \rightarrow Lyx) \vdash \forall x \exists y Lyx$
- $\exists x \forall y (Fx \leftrightarrow Fy) \vdash \neg \forall x Fx \rightarrow \forall x \neg Fx$
- $\forall x \exists y (Fx \rightarrow Gy), \forall x \exists y (\neg Fx \rightarrow Gy) \vdash \exists z Gz$
- $\exists y \forall x (Fx \wedge Gy) \vdash \forall x \exists y (Fx \wedge Gy)$
- $\forall x \exists y (Fx \wedge Gy) \vdash \exists y \forall x (Fx \wedge Gy)$

16. $\forall x \exists y (Fx \wedge Gy) \vdash \exists y \forall x (Fx \wedge Gy)$
17. $\forall x (Mx \leftrightarrow Nx), \exists y (My \wedge \exists x Rxy) \vdash \exists x Nx$
18. $\vdash \forall z (Pz \vee \neg Pz)$
19. $\vdash \forall x \forall y Rxy \rightarrow \forall x Rxx$
20. $\vdash \forall y \exists x (Qy \rightarrow Qx)$
21. $\forall x \forall y (Gxy \rightarrow Gyx) \vdash \forall x \forall y (Gxy \leftrightarrow Gyx)$
22. $\forall x (\neg Mx \vee Ljx), \forall x (Bx \rightarrow Ljx), \forall x (Mx \vee Bx) \vdash \forall x Ljx$

D. In §B problem part A, we considered fifteen syllogistic figures of Aristotelian logic. Provide proofs for each of the argument forms. Note: you will find it *much* easier if you symbolize (for example) ‘No F is G’ as ‘ $\forall x (Fx \rightarrow \neg Gx)$ ’ rather than ‘ $\neg \exists x (Fx \wedge Gx)$ ’.

E. Aristotle and his successors identified other syllogistic forms which depended upon ‘existential import’. Symbolize each of these argument forms in FOL and offer proofs.

- **Barbari.** Something is H. All G are F. All H are G. So: Some H is F
- **Celaront.** Something is H. No G are F. All H are G. So: Some H is not F
- **Cesaro.** Something is H. No F are G. All H are G. So: Some H is not F.
- **Camestros.** Something is H. All F are G. No H are G. So: Some H is not F.
- **Felapton.** Something is G. No G are F. All G are H. So: Some H is not F.
- **Darapti.** Something is G. All G are F. All G are H. So: Some H is F.
- **Calemos.** Something is H. All F are G. No G are H. So: Some H is not F.
- **Fesapo.** Something is G. No F is G. All G are H. So: Some H is not F.
- **Bamalip.** Something is F. All F are G. All G are H. So: Some H are F.

F. The following pairs of sentences are all equivalent, showing that we can move quantifiers “across” logical operators under certain circumstances. Give proofs to that they are equivalent:

1. $\forall x (Fx \wedge Ga) \dashv\vdash \forall x Fx \wedge Ga$
2. $\exists x (Fx \vee Ga) \dashv\vdash \exists x Fx \vee Ga$
3. $\forall x (Ga \rightarrow Fx) \dashv\vdash Ga \rightarrow \forall x Fx$
4. $\forall x (Fx \rightarrow Ga) \dashv\vdash \exists x Fx \rightarrow Ga$
5. $\exists x (Ga \rightarrow Fx) \dashv\vdash Ga \rightarrow \exists x Fx$
6. $\exists x (Fx \rightarrow Ga) \dashv\vdash \forall x Fx \rightarrow Ga$

When all the quantifiers occur at the beginning of a sentence, that sentence is said to be in *prenex normal form*. These equivalences are sometimes called *prenexing rules*, since they give us a means for putting any sentence into prenex normal form. For example, ‘ $\exists x Fx \wedge \forall y Gy$ ’ can be put into prenex normal form as ‘ $\exists x \forall y (Fx \wedge Gy)$ ’, or also as ‘ $\forall y \exists x (Fx \wedge Gy)$ ’.

G. Give proofs for the following quantifier equivalence laws involving negation:

1. $\forall x \neg Fx \dashv\vdash \neg \exists x Fx$
2. $\exists x \neg Fx \dashv\vdash \neg \forall x Fx$

6.5 Rules for identity

In §5.9, I mentioned that in saying of some objects a and b that they are *identical*, we don't merely mean that they are very similar to each other, or indistinguishable in the way that e.g. two cans of soda or two pennies might be. Rather, they have to be one and the same object. It follows that no matter how much you tell me about what a and b are like, qualitatively, this won't suffice to conclude that $a = b$. Indeed, no sentences which do not *already* involve an identity claim could justify an inference to ' $a = b$ '

However, we can be sure that every object is identical to *itself*. No premises are required to conclude that much. So this will be the identity introduction rule:

$$\begin{array}{c|c} & c = c \quad =I \end{array}$$

Notice that this rule does not require referring to any prior lines of the proof. For any name c , you can just write $c = c$ at any point, with only the $=I$ rule as justification.

Our elimination rule is more fun. If you have established ' $a = b$ ', then anything that is true of the object named by ' a ' must also be true of the object named by ' b ', since they are one and the same. This means that given any sentence with ' a ' in it, you can replace some or all of the occurrences of ' a ' with ' b '. For example, from ' Raa ' and ' $a = b$ ', you are justified in inferring ' Rab ', ' Rba ' or ' Rbb '. More generally:

$$\begin{array}{c|c} m & a = b \\ n & \varphi(\dots a \dots a \dots) \\ & \varphi(\dots b \dots a \dots) \quad =E\ m, n \end{array}$$

The notation here should be understood as follows: $\varphi(\dots a \dots a \dots)$ is a sentence containing the name a , and $\varphi(\dots b \dots a \dots)$ is a sentence obtained by replacing one or more occurrences of the name a with the name b . Lines m and n can occur in either order, and do not need to be adjacent, but we always cite the statement of identity first.

Symmetrically, we allow:

$$\begin{array}{c|c} m & a = b \\ n & \varphi(\dots b \dots b \dots) \\ & \varphi(\dots a \dots b \dots) \quad =E\ m, n \end{array}$$

That is, if we have established $a = b$, and we have a sentence that contains the name b , then we are allowed to infer any sentence that results from the first by replacing one or more occurrence of b with a .

This rule is closely related to LEIBNIZ'S LAW, which we briefly discussed in §5.9. Leibniz's Law says that if x and y are identical, then x has any given property iff y does too. The following, for example, is an instance of Leibniz's Law:

$$\forall x \forall y (x = y \rightarrow (Dx \leftrightarrow Dy))$$

For example, if ‘ D ’ represents the property of being a dog, Leibniz’s Law tell us that if $x = y$, then x is a dog iff y is also a dog. We can prove this using $=E$ as follows:

1		a	Flag
2		b	Flag
3		$a = b$	Assumption (for $\rightarrow I$)
4		Da	Assumption (for $\leftrightarrow I$)
5		Db	$=E$ 3, 4
6		Db	Assumption (for $\leftrightarrow I$)
7		Da	$=E$ 3, 6
8		$Da \leftrightarrow Db$	$\leftrightarrow I$ 4–5, 6–7
9		$a = b \rightarrow (Da \leftrightarrow Db)$	$\rightarrow I$ 3–8
10		$\forall y(a = y \rightarrow (Da \leftrightarrow Dy))$	$\forall I$ 2–9
11		$\forall x \forall y(x = y \rightarrow (Dx \leftrightarrow Dy))$	$\forall I$ 1–10

The $=I$ rule is also related to a logical law discussed in §5.9, the LAW OF IDENTITY, which says that everything is identical to itself, i.e. that $\forall x x = x$. We can prove this as a theorem using our $=I$ rule:

1		a	Flag
2		$a = a$	$=I$
3		$\forall x x = x$	$\forall I$ 1–2

This shows that identity is *reflexive*, i.e. it is a relation that everything bears to itself. The relation of being at-least-as-tall-as would be another example of a reflexive relation, since everyone is at least as tall as themselves. Using our rules, we can prove other properties of the identity relation as well. For example, we can prove that it is *symmetric*:

1			a	Flag
2				
2			b	Flag
3				
3			$a = b$	Assumption (for \rightarrow I)
4				
4			$a = a$	=I
5				
5			$b = a$	=E 3, 4
6			$a = b \rightarrow b = a$	\rightarrow I 3–5
7			$\forall y(a = y \rightarrow y = a)$	\forall I 2–6
8			$\forall x \forall y(x = y \rightarrow y = x)$	\forall I 1–7

Here we obtain line 5 by replacing one instance of ‘ a ’ in line 4 with an instance of ‘ b ’, which is justified given ‘ $a = b$ ’.

■ Exercises 6.5

A. Here are some important logical properties that a two-place relation R could have:

R is **reflexive** iff $\forall x Rxx$

R is **serial** iff $\forall x \exists y Rxy$

R is **symmetric** iff $\forall x \forall y (Rxy \rightarrow Ryx)$

R is **transitive** iff $\forall x \forall y \forall z ((Rxy \wedge Ryz) \rightarrow Rxz)$

R is **euclidean** iff $\forall x \forall y \forall z ((Rxy \wedge Rxz) \rightarrow Ryz)$

Above we showed that identity is both reflexive and symmetric. Show that identity also has the other three properties: it is transitive, euclidean, and serial. That is, prove each of the following theorems:

1. $\vdash \forall x \forall y \forall z ((x = y \wedge y = z) \rightarrow x = z)$
2. $\vdash \forall x \forall y \forall z ((x = y \wedge x = z) \rightarrow y = z)$
3. $\vdash \forall x \exists y x = y$

B. Provide a proof of each claim. (Remember that $t_1 \neq t_2$ is shorthand for the negated identity sentence $\neg t_1 = t_2$).

1. $Pa \vee Qb, Qb \rightarrow b = c, \neg Pa \vdash Qc$
2. $m = n \vee n = o, An \vdash Am \vee Ao$
3. $\forall x x = m, Rma \vdash \exists x Rxx$
4. $\forall x \forall y (Rxy \rightarrow x = y) \vdash Rab \rightarrow Rba$
5. $\neg \exists x x \neq m \vdash \forall x \forall y (Px \rightarrow Py)$
6. $\exists x Jx, \exists x \neg Jx \vdash \exists x \exists y x \neq y$
7. $\forall x (x = n \leftrightarrow Mx), \forall x (Ox \vee \neg Mx) \vdash On$

8. $\exists x Dx, \forall x (x = p \leftrightarrow Dx) \vdash Dp$
9. $\exists x [(Kx \wedge \forall y (Ky \rightarrow x = y)) \wedge Bx], Kd \vdash Bd$
10. $\vdash Pa \rightarrow \forall x (Px \vee x \neq a)$

C. Show that the following are provably equivalent:

- Fa
- $\exists x (x = a \wedge Fx)$

D. The following are all acceptable ways to symbolize the English sentence ‘there is exactly one F’:

- $\exists x Fx \wedge \forall x \forall y [(Fx \wedge Fy) \rightarrow x = y]$
- $\exists x [Fx \wedge \forall y (Fy \rightarrow x = y)]$
- $\exists x \forall y (Fy \leftrightarrow x = y)$

Show that they are all provably equivalent. (*Hint:* to show that three claims are provably equivalent, it suffices to show that the first proves the second, the second proves the third and the third proves the first; think about why.)

E. Symbolize the following argument, and then give a proof of it:

There is exactly one F. There is exactly one G. Nothing is both F and G. So:
there are exactly two things.

The Semantics of FOL

7

Recall that in TFL, we had the notion of a *valuation*: an assignment of truth values to atomic sentences. We then gave a semantics that allowed us to compute the truth value of an arbitrarily complex TFL sentence on any given valuation. Finally, we used this notion of truth-on-a-valuation to define various logical concepts, like equivalence and entailment.

In this chapter we will do something very similar for FOL, except that the notion of a valuation now gets replaced by the more complex notion of an *interpretation*. We'll first have to see what an FOL interpretation is, and then learn how to compute the truth values of arbitrarily complex FOL sentences in a given interpretation. With that in hand, we can then again define logical concepts — like logical entailment, or validity — for FOL.

Ultimately, we'll be able to use this machinery to show that a given FOL argument is *not* valid, by giving an interpretation that makes its premises true and its conclusion false. For reasons that will emerge in this chapter, we won't use interpretations to show that arguments *are* valid. We will use the natural deduction proofs we learned in chapter 6 to do this.

7.1 Predicates and their Extensions

The connectives of TFL are all truth-functional, and consequently, TFL only cares about what truth-values sentences have. We can assign a truth value to a sentence *directly*, via a valuation that just stipulates that the sentence '*P*', for example, is to have the value *true*. Alternatively, we can do it *indirectly*, by offering a symbolization key, e.g.:

P: Big Ben is in London

This stipulates that the TFL sentence '*P*' is to have the same truth value as the English sentence 'Big Ben is in London' (which, as it happens, is true). But no further aspect of the meaning of the English sentence carries over to TFL.

FOL is similarly impoverished as regards meaning. It goes beyond mere truth values, because it allows us to split atomic sentences into their parts, consisting of terms and predicates. A term is a word that refers to a particular object, and a predicate is a word that is *true of* objects. But FOL doesn't care about any other aspect of predicate's meaning besides what objects it's true of, or any other aspect of a term's meaning besides what it refers to. For example, when we provide a symbolization key for some FOL predicate, such as:

S: ____ is a US state.

this isn't intended to suggest that our FOL predicate '*S*' carries the same *meaning* as the English predicate. It simply tells us that the FOL predicate is to be *true of* exactly those things that the English predicate ' is a US state' is true of.

Alternatively, we can stipulate what objects a predicate is true of *directly*, by just listing those objects. So we might stipulate that '*S*' is to be true of: Alabama, Alaska, Arizona, ... and so on, listing all 50 states. This is a perfectly legitimate interpretation of an FOL predicate, because, again, all we care about is what objects it's true of, and our stipulation settles this. The things a predicate is true of comprise the *EXTENSION* of the predicate. FOL is said to be an *EXTENSIONAL LANGUAGE* because it doesn't care about any aspect of a predicate's meaning besides its extension.

Our stipulations about predicate extensions can be as arbitrary as we like. For example, we could stipulate that '*H*' should have an extension consisting of the following objects:

H: Barack Obama, the number π , the play *Hamlet*

It doesn't matter that these objects have nothing in common, they still form a perfectly good predicate extension. Suppose we add the following names to our symbolization key:

b: Usain Bolt
o: Barack Obama
p: the number π

Together, these stipulations then settle the truth value of any atomic FOL sentence formed from the predicate '*H*' and the names '*d*', '*o*', and '*p*': the sentences '*Ho*' and '*Hp*' will both be true on this interpretation, because Obama and π are in the predicate's extension, but '*Hb*' will be false, because Usain Bolt is not one of the objects '*H*' is true of.

Many-Place Predicates Things get slightly more complicated when we move from one-place predicates to two-place predicates. Consider a symbolization key like:

L: loves

This key should be read as saying something like:

'*L*' is true of *x* and *y* (in that order) iff *x* loves *y*.

The qualifier "in that order" is very important here. Since *x* might love *y* without *y* also loving *x*, a two-place predicate like this can apply to a pair of objects in one order but not another.

How should a direct stipulation of the extension of a two-place predicate like this look? This is a bit tricky. If we just *list* objects that '*L*' applies to, we won't know which of the objects are the lovers and which are the objects they love. A simple list would in other words give us no way to indicate the order in which the predicate holds of objects.

To deal with this, we instead let two-place predicates be true of *pairs* of objects. We could for example stipulate that '*R*' is to be true of, and only of, the following pairs of objects, indicated with angle brackets:

R : $\langle \text{Lenin, Marx} \rangle$
 $\langle \text{Heidegger, Sartre} \rangle$
 $\langle \text{Sartre, Heidegger} \rangle$
 $\langle \text{Marx, Marx} \rangle$

The angle-brackets tell us in what order the predicate R applies to the objects: the first object between the brackets always corresponds to the first argument slot in the predicate, and the second object corresponds to the second argument slot. Suppose, for example, that we add the following stipulations for names:

l : Lenin
 m : Marx
 h : Heidegger
 s : Sartre

Then ' Rlm ' will be true, since the pair $\langle \text{Lenin, Marx} \rangle$ is the extension of ' R '. But ' Rml ' will be false, since $\langle \text{Marx, Lenin} \rangle$ is not in the extension of ' R '. However, both ' Rhs ' and ' Rsh ' will be true, since both $\langle \text{Heidegger, Sartre} \rangle$ and $\langle \text{Sartre, Heidegger} \rangle$ are on our list of pairs, and Rmm will also be true, since the pair $\langle \text{Marx, Marx} \rangle$ is in R 's extension.

If we were dealing with a three-place predicate, its extension would consist not of *pairs* of objects, but of ordered *triples* of objects, like $\langle \text{Heidegger, Marx, Sartre} \rangle$. And a four-place predicate would have ordered *quadruples* of objects in its extension, a five-place predicate would have ordered *quintuples*, and so on. In general, we call ordered things like these TUPLES. So the extension of a many-place predicate can be specified by giving a list of tuples: either of pairs, or of triples, or of quadruples etc. depending on whether we're dealing with a two-, three-, or four-place predicate.

7.2 FOL Interpretations

We defined a *valuation* of a sentence ϕ (or collection of sentences) of TFL to be an assignment of truth values to all the atomic sentences contained in ϕ (or the collection). In FOL, the role of a valuation will be played by an INTERPRETATION. FOL interpretations are more complex than TFL valuations, because they have *three* components:

An FOL INTERPRETATION of a sentence ϕ (or of a collection of sentences C , e.g. an argument) consist of:

1. A specification of a DOMAIN containing at least one object
2. For each name in ϕ (or in C), an assignment of exactly one object in the domain. This object is the name's REFERENT.
3. For each predicate in ϕ (or in C), a specification of what objects in the domain (if any), and in what order, that predicate is true of. This constitutes the predicate's EXTENSION.^a

^aSo notice that a predicate can have an empty extension, in which case it isn't true of any objects in the domain. By contrast, we don't allow "empty names" that lack a referent.

Symbolization keys like those we used in chapter 5 consequently give us one convenient way to present an interpretation. For example, the following counts as one possible FOL interpretation of $(Lw \wedge Tw)$:

Domain: people
 w : Wittgenstein
 L : ____ is a logician
 T : ____ is a school teacher

This has all three components: (i) a specification of a domain, (ii) a specification of a referent for every name in $(Lw \wedge Tw)$, and (iii) a specification of an extension for every predicate in $(Lw \wedge Tw)$. The interpretation then determines a truth value for the sentence. In this case, since Wittgenstein was both a logician and a school teacher, both ' L ' and ' T ' are true of the referent of the name ' w ' on this interpretation, and so $(Lw \wedge Tw)$ as a whole is true.

Alternatively, we can specify interpretations by just directly listing the objects that predicates are true of, as discussed in the previous section. In fact, as we move on, it will often be convenient to consider fairly abstract interpretations where the domain consists of natural numbers, i.e. positive integers, rather than people, or plants, or other objects. One possible interpretation of $\exists x(Fx \wedge Gx)$, for example, would be the following:¹

Domain: 1, 2, 3
 F :
 G : 1, 3

Here the domain contains the numbers 1, 2, and 3, the predicate ' F ' is true of none of those objects, and ' G ' is true of 1 and 3. As you can probably guess, $\exists x(Fx \wedge Gx)$ comes out false on this interpretation, since there's no object in the domain of which both ' F ' and ' G ' are true. We'll look at how to determine truth values more closely in the next two sections.

It will often be useful to represent directly-specified interpretations *diagrammatically*. Interpretations like the one above that involve only one-place predicates can be represented using a MATRIX DIAGRAM:

	F	G
1	—	+
2	—	—
3	—	+

Here we list the objects in the domain on the left, and then put +'s and —'s under each predicate to indicate which objects that predicate is true of. If we were also considering some FOL names, we could include those in our matrix diagram by listing each name to the left of whichever object it refers to.

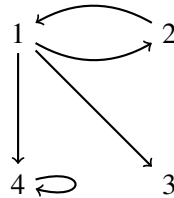
As we discussed in the last section, to directly specify the extension of a *many-place* predicate, we have to give a list of *tuples* of objects, rather than of individual objects, to indicate in what order the predicate holds of the objects. So for example, if we take the sentence $\forall xRxx$ which contains the two-place predicate ' R ', one possible interpretation would be the following:

¹Although this kind of domain officially consists of just numbers, we could in principle regard it as consisting of any objects we like — we're just calling these objects 'object 1', 'object 2', 'object 3' etc. for simplicity.

Domain: 1, 2, 3, 4

R : $\langle 1, 4 \rangle, \langle 4, 4 \rangle, \langle 1, 2 \rangle, \langle 2, 1 \rangle, \langle 1, 3 \rangle$,

To represent interpretations with two-place predicates diagrammatically, we can use ARROW DIAGRAMS like this:



The arrows represent the order in which ' R ' holds of the objects. To indicate that 1 bears the relation represented by ' R ' to 4, we draw an arrow from 1 to 4, and to indicate that 4 bears the relation to itself, we draw an arrow that loops from 4 back to 4, and to indicate that 1 bears the relation to 2, and 2 also bears it to 1, we draw two arrows between them, and so on. So there is one arrow for each pair in the extension of the predicate.

If we wanted, we could make such arrow diagrams more complex. For example, we could label objects in our diagram with FOL names to indicate which object each name refers to. To represent the extension of a one-place predicate in an arrow diagram, we could draw a circle around some objects and label the circle with the predicate. And to represent the extensions of multiple two-place predicates in a single arrow diagram, we might use arrows with dashed as opposed to solid lines, or arrows of different colors.

7.3 Truth in FOL

We have introduced interpretations. Our next task is to give a precise account of what it is for an arbitrarily complex FOL sentence to be true or false in a given interpretation. There are three kinds of sentence in FOL: atomic sentences, sentences whose main logical operator is a sentential connective, and lastly, sentences whose main logical operator is a quantifier. We'll go through each kind in turn.

Our explanation will be completely general, but to make things comprehensible, it will be useful to have a particular interpretation to hand in order to give examples. Let's use the following as our *go-to interpretation*:

Domain: all positive integers

a : 2

b : 3

E : ____ is even

R : ____ is less than ____

I've here specified the extensions of ' E ' and ' R ' indirectly, using English predicates, because I can't list all the numbers, or pairs of numbers, these are true of. But the extension of ' E ' contains even numbers 2, 4, 6, 8, ..., and the extension of ' R ' contains the pair $\langle 2, 3 \rangle$, as well as $\langle 2, 4 \rangle$ and $\langle 2, 5 \rangle$ and $\langle 3, 5 \rangle$, and indeed every pair $\langle x, y \rangle$ where x is less than y .

Truth-Rule for Atomic Sentences Determining the truth value of an atomic sentence on a given interpretation is fairly straightforward. The atomic sentence ' Ea ', for example, is true just in case ' E ' is true of the referent of ' a '. Given our go-to interpretation, this sentence is true since ' a ' refers to 2, and 2 is indeed an even number. By contrast, ' Eb ' would be false on our interpretation, because 3, which is the referent of ' b ', is not in the extension of ' E ', since it is not an even number.

Similarly, ' Rab ' is true on our interpretation just in case the referent of ' a ' is less than the referent of ' b '. Since 2 is indeed less than 3, ' Rab ' is true. By contrast, ' Rba ' is not true, because 3 is not less than 2, or to put it another way, the pair $\langle 3, 2 \rangle$ is not in the extension of ' R ' on our interpretation. Similarly, neither ' Raa ' nor ' Rbb ' are true on our go-to interpretation, since neither 2 nor 3 is less than itself (i.e. the pairs $\langle 2, 2 \rangle$ and $\langle 3, 3 \rangle$ are not in the extension of ' R '). In general, the truth-rule for atomic sentences is:

When \mathcal{R} is an n -place predicate and c_1, c_2, \dots, c_n are names, the sentence $\mathcal{R}c_1c_2\dots c_n$ is true in a given interpretation **iff** \mathcal{R} is true of the objects referred to by c_1, c_2, \dots, c_n (in that order) in that interpretation

Truth-Rules for TFL Connectives The truth rules governing our truth-functional operators are exactly the same as they were in TFL:

$\neg\phi$ is true in a given interpretation **iff** ϕ false in that interpretation

$(\phi \wedge \psi)$ is true in a given interpretation **iff** both ϕ and ψ are true in that interpretation

$(\phi \vee \psi)$ is true in a given interpretation **iff** either ϕ or ψ is true in that interpretation

$(\phi \rightarrow \psi)$ is true in a given interpretation **iff** either ϕ is false or ψ is true in that interpretation.^a

$(\phi \leftrightarrow \psi)$ is true in a given interpretation **iff** ϕ has the same truth value as ψ in that interpretation

^aThis means that $(\phi \rightarrow \psi)$ is false if ϕ is true and ψ is false, and true otherwise.

This is equivalent to the information conveyed by the characteristic truth tables for the connectives; it's just presented in words here rather than truth tables. Some examples will help to illustrate the idea (make sure you understand them!). On our go-to interpretation:

- ' $\neg Raa$ ' is true because 2 is not less than itself
- ' $Rab \wedge Ea$ ' is true because both conjuncts are true
- ' $Rab \wedge Eb$ ' is false because, although ' Rab ' is true, ' Eb ' is false
- ' $Eb \rightarrow \neg Rba$ ' is true because ' Eb ' is false.

7.4 Truth-Rules for Quantified Sentences

The big innovation in FOL is the use of *quantifiers*. And specifying the truth conditions for quantified sentences turns out to be a little tricky. If we only look at simple cases things are pretty straightforward. We can say that ' $\exists xEx$ ' is true iff ' E ' is true of at least one object in the domain, and ' $\forall xEx$ ' is true iff ' E ' is true of every object in the domain. So on our go-to interpretation, ' $\exists xEx$ ' comes out true, and ' $\forall xEx$ ' comes out false.

But what about a more complex sentence like ' $\forall x(Ex \rightarrow Rxb)$ '? This has a universal quantifier as its main operator, so we might again try to say that it is true on an interpretation iff ' $(Ex \rightarrow Rxb)$ ' is true of every object in the domain. The trouble is that ' $(Ex \rightarrow Rxb)$ ' is not a predicate, and our interpretation therefore does not directly specify what objects this complex formula is true of. So while our simple-minded approach worked for ' $\forall xEx$ ', it breaks down when we consider more complex sentences. What we need is some *uniform* and *general* way of specifying the truth conditions of *any* universally (or existentially) quantified sentence, irrespective of how complex it is.

The way we will do this is by temporarily treating variables like ' x ' *as if* they referred to objects in the domain. For example, we will say that the universal sentence ' $\forall x(Ex \rightarrow Rxb)$ ' is true iff the formula ' $(Ex \rightarrow Rxb)$ ' that the quantifier operates on is true *no matter what object in the domain the variable ' x ' is treated as referring to*. Of course, a variable like ' x ' doesn't actually refer to any particular thing, since it's not a name, and ' $(Ex \rightarrow Rxb)$ ' doesn't actually have a truth value. But if we temporarily treated ' x ' *as if* it referred to 1, for example, then the conditional ' $(Ex \rightarrow Rxb)$ ' would be true, since its antecedent ' Ex ' would be false (given that 1 is not even).

We will use the notion of a VARIABLE ASSIGNMENT to implement this idea of temporarily treating variables as if they refer to objects. We can, for example, write $[x : 1]$ for an assignment that treats ' x ' as referring to 1, and $[x : 3]$ for an assignment that treats ' x ' as referring to 3. So although the variable in ' Ex ' is not a name for any object in the domain, if we consider it *relative to an assignment* like $[x : 3]$, the variable works just like a name that refers to the number 3. Such assignments can also cover multiple variables at once. For example, $[x : 1, y : 5, z : 2]$ would be an assignment relative to which ' x ' refers 1, ' y ' refers to 5, and ' z ' refers to 2. On our go-to interpretation, $(Ez \wedge Rxy)$ would then be true relative to this assignment (since 2 is even, and 1 is less than 5), but not relative to e.g. $[x : 5, y : 2, z : 1]$.

Returning to ' $\forall x(Ex \rightarrow Rxb)$ ', our idea is that to determine whether this is true, we go through each object in the domain, and ask whether ' $(Ex \rightarrow Rxb)$ ' would come out true if ' x ' were treated as referring to that object.² So we have to ask ourselves:

is ' $(Ex \rightarrow Rxb)$ ' true on $[x : 1]$?

is ' $(Ex \rightarrow Rxb)$ ' true on $[x : 2]$?

is ' $(Ex \rightarrow Rxb)$ ' true on $[x : 3]$?

²Having brought assignments onto the scene, we should now technically go back and re-do our explanation of the truth conditions for atomic sentences too. That's because we now need to determine the truth values of *formulas* like ' $(Ex \rightarrow Rxb)$ ' relative to an assignment of an object to the free variable ' x ', but our earlier explanation only covered *sentences* (closed formulas), and said nothing about variable assignments. See the Appendix for the technical details.

is ' $(Ex \rightarrow Rxb)$ ' true on $[x : 4]$?
 \vdots

and so on for every positive integer we could assign to ' x '. Now ' $(Ex \rightarrow Rxb)$ ' is true on $[x : 1]$ as well as $[x : 3]$, since in both cases the antecedent ' Ex ' comes out false (because neither 1 nor 3 are even). And ' $(Ex \rightarrow Rxb)$ ' is also true on $[x : 2]$, since both ' Ex ' and ' Rxb ' are true on $[x : 2]$ (because 2 is even, and also less than 3, which is what ' b ' refers to on our go-to interpretation). An object that makes a formula like ' $(Ex \rightarrow Rxb)$ ' true is said to SATISFY that formula. So the numbers 1, 2, and 3 all *satisfy* ' $(Ex \rightarrow Rxb)$ '.

But the number 4 does *not* satisfy it: ' $(Ex \rightarrow Rxb)$ ' is *false* on the assignment $[x : 4]$. That's because although ' Ex ' is true on $[x : 4]$ (since 4 is even), ' Rxb ' is not true on $[x : 4]$ (since 4 is *not* less than 3). So since we found a number that makes ' $(Ex \rightarrow Rxb)$ ' false, we can conclude that the universal sentence ' $\forall x(Ex \rightarrow Rxb)$ ' we started with is itself false — it isn't true of *every* number. And that's of course the result we want: relative to our go-to interpretation, ' $\forall x(Ex \rightarrow Rxb)$ ' says that *every* even number is less than three, which is clearly false, since there are lots of even numbers, including 4, that are not less than 3.

In practice, keeping track of variable assignments while calculating the truth value of a formula can get messy, especially when we're dealing with more than one variable. We'll therefore make use a notational shorthand that will make things a little easier. Instead of explicitly mentioning the variable assignment, we will use superscripts on the variables themselves to indicate what objects we are temporarily treating them as referring to. So instead of saying:

' $(Ex \rightarrow Rxb)$ ' is true on $[x : 2]$

as we did above, we will just add a superscript of 2 to the variable ' x ' itself and write:

$(Ex^2 \rightarrow Rx^2b)$ is true

to indicate that we're treating x as referring to 2. Similarly, instead of saying that ' $(Ez \wedge Rxy)$ ' is true on $[x : 1, y : 5, z : 2]$, we can just say that $(Ez^2 \wedge Rx^1y^5)$ is true.

You can think of $(Ex^2 \rightarrow Rx^2b)$ as a SEMANTIC INSTANCE of $\forall x(Ex \rightarrow Rxb)$, much like $(Ea \rightarrow Rab)$ is a *syntactic instance* of $\forall x(Ex \rightarrow Rxb)$, something we could infer by $\forall E$. In both cases, we delete the quantifier, and then do something with the variable the quantifier bound: in the case of syntactic instances, we replace the variable with a name, and in the case of semantic instances, we assign an object from the domain to the variable.

Given a quantified FOL sentence $\forall v \varphi(\dots v \dots)$ or $\exists v \varphi(\dots v \dots)$, its SEMANTIC INSTANCES in an interpretation \mathcal{I} are obtained by removing the quantifier and assigning some object o from \mathcal{I} 's domain to every occurrence of the variable v bound by that quantifier $\varphi(\dots v^o \dots)$.

It's important to remember, though, that this really is just a notational shorthand. Whereas the syntactic instance $(Ea \rightarrow Rab)$ is a genuine sentence of FOL, the semantic instance $(Ex^2 \rightarrow Rx^2b)$ is not a sentence of FOL (nor even a formula), since the language of FOL doesn't include superscripts on variables. Talk of $(Ex^2 \rightarrow Rx^2b)$ is, again, just shorthand for talking about the FOL formula ' $(Ex \rightarrow Rxb)$ ' relative to an assignment of 2 to ' x '.

With this background in place, we can now state the truth-rules for quantified sentences:

$\forall v \phi(\dots v \dots)$ is true in an interpretation \mathcal{I} **iff** $\phi(\dots v^o \dots)$ is true in \mathcal{I} for *every object* o in \mathcal{I} 's domain (i.e. if *all* of its semantic instances in \mathcal{I} are true).

$\exists v \phi(\dots v \dots)$ is true in an interpretation \mathcal{I} **iff** $\phi(\dots v^o \dots)$ is true in \mathcal{I} for *at least one* object o in \mathcal{I} 's domain (i.e. if *at least one* of its semantic instances in \mathcal{I} are true).

The idea, again, is that we go through each object o in the domain, and check whether $\phi(\dots v \dots)$ is true if v is treated as referring to o ; if so, the universal sentence $\forall v \phi(\dots v \dots)$ is true, if not, it's false. Notice that it's much easier for an existential sentence to be true: for $\exists v \phi(\dots v \dots)$ to be true, it suffices if *one* object in the domain satisfies $\phi(\dots v \dots)$.

This in turn means that that it's very easy for a universal sentence $\forall v \phi(\dots v \dots)$ to be *false*: we just have to find a single object that makes $\phi(\dots v \dots)$ false. For an existential sentence $\exists v \phi(\dots v \dots)$ to be false, on the other hand, we have to make sure that $\phi(\dots v \dots)$ is false for *every* object in the domain. Let's state these "falsity conditions" too:

$\forall v \phi(\dots v \dots)$ is false in a given interpretation \mathcal{I} **iff** $\phi(\dots v^o \dots)$ is false in \mathcal{I} for *at least one object* o in \mathcal{I} 's domain.

$\exists v \phi(\dots v \dots)$ is false in a given interpretation \mathcal{I} **iff** $\phi(\dots v^o \dots)$ is false in \mathcal{I} for *every* object o in \mathcal{I} 's domain.

Again, we're here suppressing explicit mention of variable assignments via our notational shorthand, but you can look at the Appendix for the official version of the semantics, with variable assignments made explicit. Let's now go through a few examples to get a better feel for how to determine the truth values of quantified sentences.

7.5 Truth in an Interpretation: Examples

Since the number of objects we have to consider increases the larger the domain is, we'll here use interpretations with relatively small domains. Let's start with the following interpretation with just three objects in its domain:

Interpretation A

Domain: 1, 2, 3

F : 1

G : 2, 3

H :

	F	G	H
1	+	−	−
2	−	+	−
3	−	+	−

Example 1 $\exists x(Fx \wedge Gx)$

For this to be true, it suffices if a single object in the domain satisfies ' $Fx \wedge Gx$ '. Unfortunately, no matter which object we treat ' x ' as referring to, it comes out false. So our original existential sentence is false. Our explanation in other words goes like this:

- ▷ $\exists x(Fx \wedge Gx)$ is false because ' $(Fx \wedge Gx)$ ' is false for every object in the domain:
 - ▷ $(Fx^1 \wedge Gx^1)$ is false (since Gx^1 is false), and
 - ▷ $(Fx^2 \wedge Gx^2)$ is also false (since Fx^2 is false), and
 - ▷ $(Fx^3 \wedge Gx^3)$ is also false (since Fx^3 is false)

Example 2 $\exists xFx \wedge \exists xGx$

Notice that the main operator in this is \wedge , i.e. it's a conjunction. So here we *first* have to apply the truth-rule for \wedge , and evaluate its two conjuncts $\exists xFx$ and $\exists xGx$ *separately*. And to determine the truth value of each conjunct, we then use the truth-rule for existential sentences. As it turns out, both $\exists xFx$ and $\exists xGx$ are true, meaning that our conjunction as a whole also comes out true:

- ▷ $\exists xFx \wedge \exists xGx$ is true, because
 - ▷ $\exists xFx$ is true
 - ▷ that's because Fx^1 is true
 - ▷ and $\exists xGx$ is also true
 - ▷ that's because Gx^2 is true

Notice that Gx^3 is of course also true, so I could instead have given this as my reason for why $\exists xGx$ is true. It doesn't matter whether I pick 2 or 3 — as long as I can point to at least one object that satisfies ' Gx ', that's enough for $\exists xGx$ to be true.

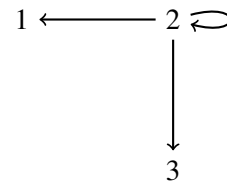
What these two examples illustrate is that we always have to identify the main operator, and then apply the truth-rule appropriate to that operator. In ' $\exists x(Fx \wedge Gx)$ ', the main operator is the existential quantifier, so I apply the truth-rule for the quantifier, and see whether I can find an object to make the whole conjunction ' $(Fx \wedge Gx)$ ' true. In ' $\exists xFx \wedge \exists xGx$ ', by contrast, the main operator is \wedge , so here I do not have to find a single object to make both ' Fx ' and ' Gx ' true. Rather I consider the two conjuncts ' $\exists xFx$ ' and ' $\exists xGx$ ' *separately*. Then I check, first, whether any object satisfies ' Fx ', and second whether any object (perhaps a different one) satisfies ' Gx '. Since I can find an object in each case, both ' $\exists xFx$ ' and ' $\exists xGx$ ' are true; and this means that ' $\exists xFx \wedge \exists xGx$ ' is true by the truth-rule for conjunction.

Things get more complicated with sentences that contain nested quantifiers, like $\exists x\forall yLxy$. Let's determine the truth value of this on the following interpretation:

Interpretation B

Domain: 1, 2, 3

L : $\langle 2,1 \rangle, \langle 2,2 \rangle, \langle 2,3 \rangle$

**Example 3** $\exists x\forall yLxy$

If we paraphrase this in English, it says that there is some object x which bears the relation L to *every* object y . And indeed there is such an object in Interpretation B, namely the number 2: this number bears L to 1, and to itself, and also to 3, that is, to every object in the domain.

To put the same point formally: in order for ' $\exists x\forall yLxy$ ' to be true, there has to be at least one object o that makes $\forall yLx^oy$ true. And there is such an object, namely 2. So ' $\exists x\forall yLxy$ ' is true because $\forall yLx^2y$ is true. Next we ask why $\forall yLx^2y$ is true. That's because no matter which object o we pick for y to refer to, ' Lx^2y^o ' comes out true. That is, Lx^2y^1 , and Lx^2y^2 and Lx^2y^3 are all true. Our official explanation in other words looks like this:

- ▷ $\exists x\forall yLxy$ is true because:
 - ▷ $\forall yLx^2y$ is true. And this is because:
 - ▷ Lx^2y^1 is true, and
 - ▷ Lx^2y^2 is true, and
 - ▷ Lx^2y^3 is true

Although sentences with nested quantifiers occur most commonly with many-place predicates, nested quantifiers can occur in combination with just one-place predicates too. Let's return to Interpretation A, and consider a case like this for our final example.

Example 4 $\exists x(Fx \rightarrow \forall y(Hy \leftrightarrow Gx))$

The main operator is the existential quantifier, so its truth-rule tells us that we have to see whether there is any object o that would make $(Fx^o \rightarrow \forall y(Hy \leftrightarrow Gx^o))$ true on Interpretation A. In fact, an assignment of 1 to ' x ' will do it: the conditional $(Fx^1 \rightarrow \forall y(Hy \leftrightarrow Gx^1))$ comes out true because its antecedent Fx^1 is true and its consequent $\forall y(Hy \leftrightarrow Gx^1)$ is also true.

To explain why ' $\forall y(Hy \leftrightarrow Gx^1)$ ' is true, we need to apply the truth-rule for the universal quantifier: this is true is because no matter what object o we pick, $Hy^o \leftrightarrow Gx^1$ is true. That is: $Hy^1 \leftrightarrow Gx^1$ is true, and $Hy^2 \leftrightarrow Gx^1$ is true, and $Hy^3 \leftrightarrow Gx^1$ is also true. The complete explanation for why Example 4 is true on Interpretation A then looks like this:

- ▷ $\exists x(Fx \rightarrow \forall y(Hy \leftrightarrow Gx))$ is true because:
 - ▷ $Fx^1 \rightarrow \forall y(Hy \leftrightarrow Gx^1)$ is true. And this is because:
 - ▷ Fx^1 is true, and
 - ▷ $\forall y(Hy \leftrightarrow Gx^1)$ is also true. This in turn is because:
 - ▷ $(Hy^1 \leftrightarrow Gx^1)$ is true (since Hy^1 and Gx^1 are both false) and
 - ▷ $(Hy^2 \leftrightarrow Gx^1)$ is true (since Hy^2 and Gx^1 are both false) and
 - ▷ $(Hy^3 \leftrightarrow Gx^1)$ is true (since Hy^3 and Gx^1 are both false)

An explanation of this sort is called a SEMANTIC DEMONSTRATION of the truth or falsity of a given sentence (in an interpretation). In the exercises below, you should give semantic demonstrations like this to justify your claims about the truth values of sentences.

■ Exercises 7.5

A. Take the following interpretation, and determine the truth value of each sentence below on this interpretation (remember to give a semantic demonstration in each case):

Domain: 1, 2

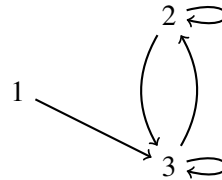
 F : 1, 2 G : 2 H :

	F	G	H
1	+	–	–
2	+	+	–

1. $\exists x(Fx \wedge Gx)$
2. $\forall x(Hx \rightarrow Gx)$
3. $\forall x(Fx \rightarrow \neg Hx)$
4. $\forall x(Fx \leftrightarrow Gx)$
5. $\forall xGx \rightarrow \exists yHy$
6. $\forall xFx \wedge \neg \exists xGx$
7. $\exists x(Gx \wedge \forall y(Gy \rightarrow Hx))$
8. $\forall x(Gx \rightarrow \exists y(Fx \wedge \neg Gy))$

B. Determine the truth values of the sentences below on the provided interpretation

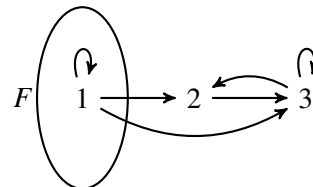
Domain: 1, 2, 3

 R : $\langle 1,3 \rangle, \langle 2,2 \rangle, \langle 2,3 \rangle, \langle 3,2 \rangle, \langle 3,3 \rangle$ 

1. $\exists xRxx$
2. $\forall xRxx$
3. $\forall x\forall yRxy$
4. $\forall x\exists yRxy$
5. $\exists x\forall yRxy$
6. $\exists x\forall y\neg Ryx$
7. $\forall y\exists x\neg Ryx$
8. $\forall x\forall y(Rxy \rightarrow Ryx)$
9. $\exists x\forall y(Rxy \rightarrow Ryx)$
10. $\forall x(\exists yRxy \rightarrow \exists yRyx)$

C. Determine the truth values of the sentences below on the provided interpretation (note that in the diagram, I am using an oval to indicate the extension of the 1-place predicate ' F ', and arrows for the extension of the 2-place predicate ' R ')

Domain: 1, 2, 3

 F : 1 R : $\langle 1,1 \rangle, \langle 1,2 \rangle, \langle 1,3 \rangle, \langle 2,3 \rangle, \langle 3,2 \rangle, \langle 3,3 \rangle$ 

1. $\exists x(Fx \wedge Rxx)$
2. $\exists x(\neg Fx \wedge \neg Rxx)$
3. $\forall x(Rxx \rightarrow Fx)$
4. $\forall x(\exists yRxy \rightarrow Fx)$

5. $\forall x \exists y Rxy$
6. $\forall x \exists y Rxy \rightarrow \exists x Fx$
7. $\exists x (\forall y Rxy \wedge \exists y Ryx)$

7.6 Semantic Concepts

Defining truth in FOL was a bit tricky, due to the presence of quantifiers. But now that we know what determines the truth value of an FOL sentence in an interpretation, we can use that to define various other central logical notions. As you can see, these definitions are basically the same as those from TFL, except that we here use the notion of truth in an *interpretation* rather than truth on a *valuation*. In all the definitions below, the metavariables $\varphi_1 \dots \varphi_n$ and ψ range over arbitrary *sentences* of FOL.³

$\varphi_1, \dots, \varphi_n$ LOGICALLY ENTAIL ψ , written $\varphi_1, \dots, \varphi_n \models \psi$, iff no interpretation makes all of $\varphi_1, \dots, \varphi_n$ true but ψ false.

φ is a LOGICAL TRUTH, written $\models \varphi$, iff it is true in every interpretation.

φ is a CONTRADICTION iff it is false in every interpretation.

φ and ψ are LOGICALLY EQUIVALENT, written $\varphi \models \psi$ iff they have the same truth value in every interpretation.

$\varphi_1, \dots, \varphi_n$ are JOINTLY CONSISTENT iff there is at least one interpretation which makes them all true. They are JOINTLY INCONSISTENT iff there is no such interpretation.

As in TFL, perhaps the most important of these concepts is entailment, since it is closely related to the notion of a valid argument. We will say that an FOL argument:

$$\varphi_1, \dots, \varphi_n \therefore \psi$$

is VALID just in case its premises $\varphi_1, \dots, \varphi_n$ logically entail its conclusion ψ . Since entailment is defined in terms of the concept of truth-in-an-interpretation, we can use interpretations to investigate the validity of FOL arguments.

In particular, we can use interpretations to show that an argument is *not* valid, or that the corresponding entailment *fails*. To show that a given argument is not valid, it suffices to construct a single interpretation that simultaneously makes all of the premises true but still makes the conclusion false. Such an interpretation is again called a *counterexample* to the validity of the argument, or a COUNTERMODEL.

In the next two sections, we'll look at how to construct countermodels. We'll first look at arguments that involve only one-place predicates, and then at ones that involve many-place predicates. Of course, we can use interpretations to investigate the other logical concepts

³Our logical concepts are in other words only defined for FOL sentences. By contrast, the official explanation of truth in FOL, given in the Appendix, covers formulas in general, not just sentences.

we've defined above too. We'll return to that once we've looked at our central concept of entailment, or validity.

7.7 Countermodels with One-Place Predicates

Example 1 Let's begin with one of the examples we looked at way back in §1.1.

- (1) All rabbits are mammals
 Bugs Bunny is a mammal.
 \therefore Bugs Bunny is a rabbit.

As we noted at the time, this argument is intuitively not valid. We are now in a position to demonstrate this formally. First, we can offer the following FOL symbolization:

$$\begin{aligned} &\forall x(Rx \rightarrow Mx) \\ &Mb \\ &\therefore Rb \end{aligned}$$

Next, we can construct an interpretation which makes the premises true and the conclusion false. One such interpretation would be the following:

Domain: all people
 b: Lady Gaga
 R: ____ is an opera singer
 M: ____ is a vocalist

In this interpretation, ' $\forall x(Rx \rightarrow Mx)$ ' is true, since every opera singer is a vocalist, and ' Mb ' is also true, since Lady Gaga is a vocalist. But ' Rb ' is false, since she is not an opera singer, but a pop singer. So the argument isn't valid. However, using this kind of countermodel has the drawback that it requires us to appeal to real-world knowledge. After all, who knows, maybe, Lady Gaga secretly performs in operas too, in which case ' Rb ' is true in this interpretation, and it doesn't constitute a countermodel after all.

To avoid these kinds of issues, and to give some uniformity to our countermodels, we will again use interpretations with domains that contain just positive integers, and give *direct* specifications of the extensions of predicates, by listing the objects they are true of. Furthermore, since we may have to explain why universal sentences like ' $\forall x(Rx \rightarrow Mx)$ ' are true on our countermodels, and since the number of objects we have to consider to do this increases the larger the domain is, we will try to construct countermodels with the smallest possible domains.

Let's begin by seeing if we can construct a countermodel that has just one object in its domain, say the number 1. First off, let's make sure our conclusion ' Rb ' is false. To that end, we can let ' b ' refer to 1, and let the extension of the predicate ' R ' remain empty. And to make the second premise, ' Mb ' true, we have to put 1 (the referent of ' b ') into the extension of ' M '. So we have an interpretation like this:

Domain: 1
 b : 1
 R :
 M : 1

		R	M
b	1	—	+

And luckily this also makes the first premise ' $\forall x(Rx \rightarrow Mx)$ ' true! After all, $(Rx^1 \rightarrow Mx^1)$ is true (since Rx^1 is false). And since 1 is the only object in our domain, that suffices for ' $\forall x(Rx \rightarrow Mx)$ ' to be true! So our simple interpretation makes both premises true but the conclusion false, showing that the argument isn't valid.

Example 2 Let's look at a slightly more complex example:

$$\exists x(Fx \rightarrow Gx) \therefore \exists xFx \rightarrow \exists xGx$$

We will again begin with the smallest possible domain, containing just 1, and think about what's needed to make the conclusion, ' $\exists xFx \rightarrow \exists xGx$ ' false. Since this is a conditional, we have to make its antecedent ' $\exists xFx$ ' true and its consequent ' $\exists xGx$ ' false. Using just the number 1, we can do that as follows:

	F	G
1	+	—

' $\exists xFx$ ' is now true, since Fx^1 is true. But ' $\exists xGx$ ' is false, because Gx^1 is false. The trouble is that this interpretation will also make our premise ' $\exists x(Fx \rightarrow Gx)$ ' false. After all, $(Fx^1 \rightarrow Gx^1)$ is false, and there's no other object in the domain we could assign to ' x ' to make ' $(Fx \rightarrow Gx)$ ' true.

So let's expand our domain by adding a second object:

	F	G
1	+	—
2	?	?

We don't yet know whether our two predicates should be true or false of this new object, so I've left the cells next to it blank. Let's begin with our conclusion ' $\exists xFx \rightarrow \exists xGx$ ' again: we want it to be false, so ' $\exists xFx$ ' has to be true and ' $\exists xGx$ ' has to be false. Of course, ' $\exists xFx$ ' remains true because Fx^1 is still true, so no change is needed here. But to keep ' $\exists xGx$ ' false, we need to make sure that our new object 2 is not in the extension of ' G ' either:

	F	G
1	+	—
2	?	—

Now, can we make the premise ' $\exists x(Fx \rightarrow Gx)$ ' true? In fact we can, by making sure that ' F ' is not true of 2. For in that case, $(Fx^2 \rightarrow Gx^2)$ will be true, because the antecedent Fx^2 will be false. And that suffices for the truth of ' $\exists x(Fx \rightarrow Gx)$ '. Our final countermodel, and the accompanying semantic demonstrations showing that our premise is true and our conclusion false, then look like this:

Domain: 1, 2

 F : 1 G :

	F	G
1	+	—
2	—	—

- $\exists x(Fx \rightarrow Gx)$ is true because:
 - ▷ $(Fx^2 \rightarrow Gx^2)$ is true (since Fx^2 is false)
- $\exists xFx \rightarrow \exists xGx$ is false because:
 - ▷ $\exists xFx$ is true
 - ▷ since Fx^1 is true
 - ▷ but $\exists xGx$ is false. This is because
 - ▷ Gx^1 is false, and
 - ▷ Gx^2 is also false

We only needed a single object in the domain of our first countermodel, but our second countermodel ended up requiring two objects to simultaneously make the premises true and the conclusion false. Other examples might require you to use three objects, or even four, or more. Is there any upper bound on the number of objects that might be needed to produce a countermodel?

It turns out that for arguments that only involve one-place predicates, the answer is ‘yes’. The logician Leopold Löwenheim showed that if an FOL argument contains just n one-place predicates, then if the argument is invalid, a countermodel with a domain of at most 2^n objects exists. So for something like Example 2, which involves two one-place predicates, we can be sure that we won’t need more than four objects. As we’ll see, there is no such upper bound for arguments with many-place predicates.

Before we turn to many-place predicates, though, one more general observation is in order. Consider the following English argument:

All foxes are mortal.
 \therefore Every vixen is mortal.

This argument is valid in the sense that it’s impossible for its premise to be true but its conclusion to be false: a vixen is just a female fox, so any possible world where every fox is mortal has to be one where every vixen is mortal. But if we symbolize it, the resulting FOL argument is *not* valid:

$\forall x(Fx \rightarrow Mx)$
 $\therefore \forall x(Vx \rightarrow Mx)$

It’s easy to construct an interpretation that makes the premise true and the conclusion false (just put an object in the extension of V but not M or F).

So from the fact that the symbolization of an English argument in FOL is not valid, we can’t straightaway conclude that the original English argument is not valid. What we can conclude is just that the original English argument is not *formally* valid, or more specifically, that it isn’t valid in virtue of the kind of logical form captured in its FOL symbolization. But it might still be valid for other reasons — in this case, because of the connection between

the meanings of ‘fox’ and ‘vixen’. As we discussed in §1.4, logic doesn’t aim to capture the validity of arguments like this; it only cares about formally valid arguments.⁴

■ Exercises 7.7

A. Construct countermodels to demonstrate the following:

1. $\forall xFx \leftrightarrow \forall xGx \not\models \forall x(Fx \leftrightarrow Gx)$
2. $\forall x((Fx \wedge Gx) \rightarrow Hx) \not\models \forall x(Fx \vee Gx) \vee \forall x(Fx \vee Hx)$
3. $\forall xFx \rightarrow \exists xGx \not\models \exists xFx \rightarrow \exists xGx$
4. $(\forall xFx \wedge \forall xGx) \rightarrow \forall xHx \not\models \forall x((Fx \wedge Gx) \rightarrow Hx)$
5. $\exists x\forall y(Fx \rightarrow Gy) \not\models \exists y\forall x(Fx \rightarrow Gy)$
6. $\forall x\exists y(Fy \rightarrow Gx) \not\models \forall x(\exists yFy \rightarrow Gx)$
7. $\forall x(\exists yFy \rightarrow Gx) \not\models \forall x\exists y(Fy \rightarrow Gx)$

7.8 Countermodels with Many-Place Predicates

Let’s next look at how to construct countermodels for arguments with many-place predicates. First, though, there’s some new terminology that will come in handy. Two-place predicates like ‘loves’, ‘respects’, ‘admires’ etc. express RELATIONS between objects (the loving relation, the respecting relation, and so on). And there are some important characteristics that relations like this can have:

A relation R is SERIAL iff $\forall x\exists yRxy$
 A relation R is REFLEXIVE iff $\forall xRxx$
 A relation R is SYMMETRIC iff $\forall x\forall y(Rxy \rightarrow Ryx)$
 A relation R is TRANSITIVE iff $\forall x\forall y\forall z((Rxy \wedge Ryz) \rightarrow Rxz)$

One thing we can do, then, is to show that a relation’s having certain of these characteristics does, or does not, imply its having some other characteristic. For example, being reflexive implies being serial. After all: if every object bears R to itself, then every object bears R to *at least* one thing, namely itself! You could do a natural deduction proof to show that $\forall xRxx \vdash \forall x\exists yRxy$. However, the implication does not hold in the other direction, i.e. being serial does *not* imply being reflexive. Showing this will be our first example.

Example 1 $\forall x\exists yRxy \not\models \forall xRxx$

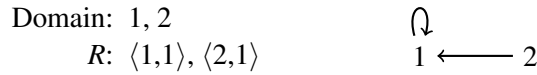
You can perhaps already think of a countermodel that would make ‘ $\forall x\exists yRxy$ ’ true but ‘ $\forall xRxx$ ’ false, but again, we’re going to try to find the smallest countermodel to do the job. Let’s start with a domain that just contains the number 1. For ‘ $\forall x\exists yRxy$ ’ to be true, every object in the domain has to bear the relation R to at least one object. So since we only have one object, that means 1 has to bear R to itself:

⁴Of course we can make the above argument formally valid by adding a second premise: ‘Every vixen is a fox’, which we’d symbolize as ‘ $\forall x(Vx \rightarrow Fx)$ ’. But this is now a *different FOL* argument, one with two premises.



However, this now also makes ' $\forall x Rxx$ ' true, whereas our goal is to make this false.

So let's expand our domain to include two objects, 1 and 2. In order for ' $\forall x Rxx$ ' to be false, at least one of these objects must not bear R to itself, i.e. must not have an arrow looping back to itself. And for ' $\forall x \exists y Rxy$ ' to be true, every object has to bear R to something, i.e. must have at least one outgoing arrow. One easy way to achieve both is to expand our earlier model to look like this:



The accompanying semantic demonstration runs like this:

- $\forall x \exists y Rxy$ is true because:
 - ▷ $\exists y Rx^1y$ is true
 - ▷ since Rx^1y^1 is true
 - ▷ and $\exists y Rx^2y$ is also true
 - ▷ since Rx^2y^1
- $\forall x Rxx$ is false because:
 - ▷ Rx^2x^2 is false

There are of course many other countermodels that would do the job just as well. But what we have in any case discovered is that *at least two* objects are necessary to show that the entailment from seriality to reflexivity fails.

Example 2 Next, let's show the following:

$$\forall x Lxx, \forall x \forall y (Lxy \rightarrow Lyx) \not\models \exists x \forall y Lxy$$

That is: a relation L 's being both reflexive and symmetric does not imply that there is some object x that bears L to everything (there's no official name for this latter characteristic). If we begin with a domain containing just 1, then to make ' $\forall x Lxx$ ' true, we would have to do the following:

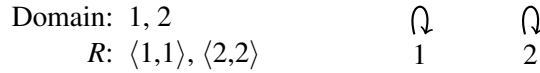


However, since we want to make ' $\exists x \forall y Lxy$ ' *false*, we need to make sure that for every x , there is at least one object it *doesn't* bear L to, i.e. we want the following to hold: $\forall x \exists y \neg Lxy$. And the trouble is that as things stand, every object *does* bear L to something.

So let's try it with two objects. Again, to make ' $\forall x Lxx$ ' true, they both need to bear L to themselves:

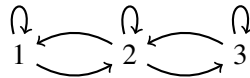


Notice that at this point, ‘ $\forall x \exists y \neg Lxy$ ’ holds. For every object we can find at least one thing it doesn’t bear L to: 1 doesn’t bear L to 2, and 2 doesn’t bear L to 1. We still have to make sure that symmetry holds, i.e. that ‘ $\forall x \forall y (Lxy \rightarrow Lyx)$ ’ is true. But if you think about it, symmetry does hold in our diagram: everything that 1 bears L to — which is just itself — also bears L back to 1. And similarly, everything that 2 bears L to — which is just itself — also bears L back to 2. So we’re done! Our countermodel and accompanying semantic demonstration (which gets pretty involved in the case of symmetry!) looks like this:



- $\forall x Lxx$ is true because:
 - ▷ Lx^1x^1 is true, and
 - ▷ Lx^2x^2 is also true
- $\forall x \forall y (Lxy \rightarrow Lyx)$ is true because:
 - ▷ $\forall y (Lx^1y \rightarrow Lyx^1)$ is true, since:
 - ▷ $(Lx^1y^1 \rightarrow Ly^1x^1)$ is true (because Lx^1y^1 and Ly^1x^1 are both true), and
 - ▷ $(Lx^1y^2 \rightarrow Ly^2x^1)$ is also true (because Lx^1y^2 is false)
 - ▷ and $\forall y (Lx^2y \rightarrow Lyx^2)$ is true, since:
 - ▷ $(Lx^2y^1 \rightarrow Ly^1x^2)$ is true (because Lx^2y^1 is false), and
 - ▷ $(Lx^2y^2 \rightarrow Ly^2x^2)$ is also true (because Lx^2y^2 and Ly^2x^2 are both true)
- $\exists x \forall y Lxy$ is false because:
 - ▷ $\forall y Lx^1y$ is false, since:
 - ▷ Lx^1y^2 is false
 - ▷ and $\forall y Lx^2y$ is also false, since
 - ▷ Lx^2y^1 is false

Again, the countermodel we arrived at is not the only one possible. The following would work too, for example, and might be more intuitive:



But with three objects in the domain, giving a semantic demonstration would require more work. To show that ‘ $\forall x \forall y (Lxy \rightarrow Lyx)$ ’ holds, for example, we’d have to show that ‘ $\forall y (Lxy \rightarrow Lyx)$ ’ is true no matter which of our three objects we assign to ‘ x ’; and then relative to each choice for ‘ x ’, we’d have to show that ‘ $(Lxy \rightarrow Lyx)$ ’ is true no matter what we assign to ‘ y ’. So we’d have to consider nine different assignments of objects to ‘ x ’ and ‘ y ’ in total, whereas the demonstration above only required us to look at four assignments.

■ Exercises 7.8

A. Construct countermodels to demonstrate the following:

1. $\forall x \exists y Lxy \not\models \exists x \forall y Lxy$
2. $\forall x \exists y Lyx \not\models \forall x \exists y Lxy$
3. $\exists x (Fx \wedge \forall y Rxy) \not\models \forall x (Fx \rightarrow \exists y Rxy)$
4. $\exists x (\exists y Axy \wedge \neg \exists y Ayx) \not\models \exists x \forall y (Axy \rightarrow Ayx)$
5. $\forall x (\forall y Rxy \rightarrow \exists z \forall w R wz) \not\models \forall x \exists y \forall z (Rxy \leftrightarrow Rzy)$

7.9 Validity and Decidability

Let's investigate one more argument involving a two-place predicate:

$$\forall x \exists y Rxy, \forall x \forall y \forall z ((Rxy \wedge Ryz) \rightarrow Rxz) \therefore \exists x Rxx$$

For this to be valid, the seriality and transitivity of R would together have to imply that some object bears R to itself (again, there's no official name for this latter characteristic). Let's see if we can construct a countermodel.

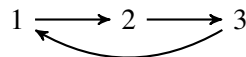
To make the conclusion ' $\exists x Rxx$ ' false, we have to make sure that *no* object in our interpretation bears R to itself, i.e. there can be no arrows that loop from an object onto itself. And as we already noticed at the hand of Example 1 in the previous section, for seriality to hold in a domain of just one object, that object would have to bear the relation to itself. So we cannot make seriality true and ' $\exists x Rxx$ ' false with just one object.

So let's move to a domain with two objects. We can make seriality hold while avoiding any self-looping arrows like this:



But now consider how to make transitivity, i.e. ' $\forall x \forall y \forall z ((Rxy \wedge Ryz) \rightarrow Rxz)$ ', hold. For this to be true, the conditional ' $(Rxy \wedge Ryz) \rightarrow Rxz$ ' has to be true no matter which objects we assign to each of the variables ' x ', ' y ', and ' z '. So suppose we assign 2 to ' y ', and 1 to both ' x ' and ' z ': $(Rx^1y^2 \wedge Ry^2z^1) \rightarrow Rx^1z^1$. Both Rx^1y^2 and Ry^2z^1 are true, so for the latter conditional to be true, the consequent Rx^1z^1 has to be true. But this would now require a self-looping arrow from 1 back to 1! So we can't get what we want with just two objects.

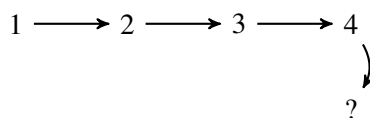
Let's try three objects. What we've just seen is that we can't have any symmetric arrows, since that would introduce a self-looping arrow by transitivity. With three objects, we can secure seriality (making sure every object has an arrow going to some object) while avoiding both symmetric and self-looping arrows like this:



But again, transitivity causes a problem. Let's assign 2 to ' x ', 3 to ' y ', and 1 to ' z ', giving us: $(Rx^2y^3 \wedge Ry^3z^1) \rightarrow Rx^2z^1$. Since both Rx^2y^3 and Ry^3z^1 are true, for the conditional to be true, the consequent Rx^2z^1 has to be true. That would mean adding an arrow going from

2 to 1. But since we already have an arrow going from 1 to 2, that would reintroduce a double-headed arrow!

What this shows is that we can't have any arrows that "go backward" in our interpretation, to an earlier number. We could add a fourth object so that 3's arrow can move forward:



But 4 *also* has to have an arrow going to some object (for seriality). And that arrow can't go backward to 1, 2, or 3, and can't go from 4 to itself, so what's to be done? Does this mean it's impossible to construct a countermodel, and that our argument is therefore valid?

No, it does not: the argument is invalid, and a countermodel exists, but it requires *infinitely many objects* in its domain. The following is a countermodel, for example:

Domain: all positive integers

R : ____ is less than ____ (or: all pairs of integers $\langle m, n \rangle$ where $m < n$)

We can't specify this interpretation "directly" by listing all the objects in the domain and all the pairs in the extension of ' R ', since there are too many of them. But we can specify the extension of ' R ' indirectly, via the English predicate. This interpretation does what we want:

- ' $\forall x \exists y Rxy$ ' is true, since for every positive integer x we can find a y such that $x < y$
- ' $\forall x \forall y \forall z ((Rxy \wedge Ryz) \rightarrow Rxz)$ ' is true, since if $x < y$ and $y < z$, it follows that $x < z$.
- ' $\exists x Rxx$ ' is false, since there exist no integer x such that $x < x$.

What we've seen, then, is that showing an FOL argument to be invalid may require an interpretation with infinitely many objects. As mentioned in §7.7, this is not the case for FOL arguments that contain only *one-place* predicates. If such an argument is invalid, then there exists a countermodel with a domain of at most 2^n objects (where n is the number of one-place predicates the argument contains). Since there are finitely many interpretations for n predicates using at most 2^n objects, this means that a computer could in principle crunch through all of those interpretations and determine if the argument is valid: if it finds a countermodel among them, the argument is invalid, if it doesn't, the argument is valid.

The same holds for TFL. Any TFL argument can be tested for validity by constructing its joint truth table and checking whether any line (i.e. valuation) makes all the premises true and the conclusion false. So a computer can be programmed to mechanically test any TFL argument for validity.

What we've seen is that this does not hold for FOL arguments that contain two-place predicates. Here there is no finite number of interpretations we could program a computer to check such that, if it fails to find a countermodel among them, the argument is guaranteed to be valid. And in fact, as the logicians Alan Turing and Alonzo Church independently proved in 1936, there exists *no* mechanical test for validity in FOL. Validity in FOL is therefore said to be UNDECIDABLE, in contrast to TFL, where validity is DECIDABLE via truth-tables.

Demonstrating FOL arguments to be valid therefore invariably requires some ingenuity and insight — it's not the kind of thing a computer can do. This is why we use natural

deduction proofs to demonstrate validity in FOL: if methods that require ingenuity are going to be required anyhow, we might as well use the method of natural deduction.⁵ To be sure, one *can* also give semantic proofs to demonstrate the validity of FOL arguments. For example to show that the following entailment holds:

$$\exists x(Fx \wedge Gx) \models \exists xFx \wedge \exists xGx$$

we could give the following semantic proof in English:

Semantic Proof: consider some arbitrary interpretation \mathcal{I} , and suppose $\exists x(Fx \wedge Gx)$ is true in \mathcal{I} . This means there is some object in the domain of \mathcal{I} , let's call it a , such that $(Fx^a \wedge Gx^a)$ is true in \mathcal{I} . But then since Fx^a is true in \mathcal{I} , this means $\exists xFx$ is true in \mathcal{I} . And similarly, since Gx^a is true in \mathcal{I} , it follows that $\exists xGx$ must be true in \mathcal{I} . Therefore $\exists xFx \wedge \exists xGx$ is true in \mathcal{I} . Since \mathcal{I} was an arbitrary interpretation, we can conclude that *every* interpretation that makes $\exists x(Fx \wedge Gx)$ true must also make $\exists xFx \wedge \exists xGx$ true, meaning that the entailment above holds.

However, we can also give a natural deduction proof, which goes through a very similar line of reasoning:

1	$\exists x(Fx \wedge Gx)$	Premise
2	$Fa \wedge Ga$	Assumption (flag a)
3	Fa	$\wedge E$ 2
4	$\exists xFx$	$\exists I$ 3
5	Ga	$\wedge E$ 2
6	$\exists xGx$	$\exists I$ 5
7	$\exists xFx \wedge \exists xGx$	$\wedge I$ 4, 6
8	$\exists xFx \wedge \exists xGx$	$\exists E$ 1, 2–8

Natural deduction proofs do not, however, allow use to demonstrate that arguments are *invalid*. To do this, we have to rely on the method of constructing countermodels. And that's what we have been doing in the last several sections.

7.10 Working with Other Semantic Concepts

So far, we've just been focusing on validity, or entailment. But we can use interpretations, as well as natural deduction proofs, in connection with other semantic concepts too.

Take the concept of logical truth first. To show that some FOL sentence φ is *not* a logical truth, it suffices to construct an interpretation that makes it false. And to show that it *is* a logical truth, we can give a natural deduction that proves φ as a theorem. This gives us a

⁵Of course, by using natural deductions to demonstrate validity, we are relying on the fact that our proof system is sound, i.e. only lets us prove valid arguments. See the introduction to Chapter 6 on this concept.

way to approach contradictions as well: since φ is a contradiction iff $\neg\varphi$ is a logical truth, we can show that φ is a contradiction by proving $\neg\varphi$ as a theorem. And to show that φ is *not* a contradiction, it suffices to construct an interpretation in which it is true.

Next consider logical equivalence. To show that φ and ψ are *not* logically equivalent, all we need to do is construct an interpretation on which one of them is true and the other is false. And to show that φ and ψ *are* equivalent, we can give a natural deduction that proves $\varphi \leftrightarrow \psi$ as a theorem, since φ and ψ are equivalent just in case $\varphi \leftrightarrow \psi$ is a logical truth.

Lastly, take the concept of consistency. To show that some sentences are jointly consistent, it suffices to give an interpretation on which they are all true. As for inconsistency, notice that it connects to entailment in the following way:

If $\varphi_1, \dots, \varphi_n \models \perp$, then $\varphi_1, \dots, \varphi_n$ are jointly inconsistent

For suppose $\varphi_1, \dots, \varphi_n \models \perp$. This means there's no interpretation that makes all of $\varphi_1, \dots, \varphi_n$ true but \perp false. However, since every interpretation makes \perp false, this just means that there's no interpretation that makes all of $\varphi_1, \dots, \varphi_n$ true. That is, $\varphi_1, \dots, \varphi_n$ are inconsistent. So if we want to show that $\varphi_1, \dots, \varphi_n$ are inconsistent, we can do that by giving a natural deduction proof showing that $\varphi_1, \dots, \varphi_n \vdash \perp$.

The following table summarizes what is needed to demonstrate that a given concept does or does not apply:

	Yes	No
logical truth?	give a proof	give an interpretation
contradiction?	give a proof	give an interpretation
equivalent?	give a proof	give an interpretation
consistent?	give an interpretation	give a proof
valid?	give a proof	give an interpretation
entailment?	give a proof	give an interpretation

■ Exercises 7.10

A. Show that the following pairs of sentences are not logically equivalent (by constructing an interpretation for each pair that makes one true but the other false):

1. $\exists x Jx, \exists x \neg Jx$
2. $\exists x Jx \wedge \exists x Hx, \exists x (Jx \wedge Hx)$
3. $\forall x Rxx, \exists x Rxx$
4. $\exists x Px \rightarrow \exists y Qy, \exists x (Px \rightarrow \exists y Qy)$
5. $\forall x (Px \rightarrow \neg Qx), \exists x (Px \wedge \neg Qx)$
6. $\exists x (Px \wedge Qx), \exists x (Px \rightarrow Qx)$
7. $\forall x (Px \rightarrow Qx), \forall x (Px \wedge Qx)$
8. $\forall x \exists y Rxy, \exists x \forall y Rxy$
9. $\forall x \exists y Rxy, \forall x \exists y Ryx$

B. Show that the following sentences are jointly consistent (i.e. there's an interpretation that makes them all true):

1. $\forall y Gy, \forall x (Gx \rightarrow Hx), \exists y \neg Iy$
2. $\exists x (Bx \vee Ax), \forall x \neg Cx, \forall x [(Ax \wedge Bx) \rightarrow Cx]$
3. $\exists x Xx, \exists x Yx, \forall x (Xx \leftrightarrow \neg Yx)$
4. $\forall x (Px \vee Qx), \exists x \neg (Qx \wedge Px)$
5. $\exists z (Nz \wedge Ozz), \forall x \forall y (Oxy \rightarrow Oyx)$
6. $\neg \exists x \forall y Rxy, \forall x \exists y Rxy$

7.11 Semantics for Identity

As mentioned in §5.9, FOL is standardly supplemented with a primitive logical symbol $=$ for *identity*. In §6.5 we looked at the deduction rules that govern identity, with which we were then able to prove various logical truths involving identity, like that everything is identical to itself, $\forall x x = x$.

Similarly, we now need to say something about the semantics of identity, specifically how to determine the truth values of identity statements in an interpretation. In this case, the semantics is very simple:

For any names a and b , $a = b$ is true in a given interpretation **iff** a and b refer to very same object in that interpretation.

So if, for example, our interpretation specifies that the name ‘ m ’ refers to 1, that ‘ n ’ refers to 2, and that ‘ o ’ also refers to 2, then ‘ $m = n$ ’ would be false (since 1 and 2 aren’t identical) whereas ‘ $n = o$ ’ would be true (since 2 is identical to 2). And of course, things like ‘ $m = m$ ’ or ‘ $n = n$ ’ will always be true, on any interpretation, since *whatever* ‘ n ’ (or ‘ m ’) refers to, it will be identical to itself.

This clause only covers identity statements involving names. For identity statements involving variables, we again have to bring in variable assignments. E.g. ‘ $x = y$ ’ comes out true if the variables ‘ x ’ and ‘ y ’ are both assigned the value 1, say, but it would come out false if ‘ x ’ were assigned the value 1 and ‘ y ’ the value 2. On the other hand, ‘ $x = x$ ’ comes out true no matter what value is assigned to ‘ x ’. See the Appendix in §7.12 below for the official semantics that makes the role of variable assignments explicit.

We can also consider quantified sentences involving identity. In §5.9 we saw how to symbolize numerical claims using identity. The claim that exactly one object exists can, for example, be symbolized as ‘ $\exists x \forall y x = y$ ’, and the claim that exactly two objects exist can be symbolized as ‘ $\exists x \exists y (x \neq y \wedge \forall z (z = x \vee z = y))$ ’.⁶ Using our semantics, we can now show that ‘ $\exists x \forall y x = y$ ’ is indeed false if the domain contains more than one object:

Domain: 1, 2

$\exists x \forall y x = y$ is false because:

- ▷ $\forall y x^1 = y$ is false
 - ▷ since $x^1 = y^2$ is false.
- ▷ and $\forall y x^2 = y$ is also false
 - ▷ since $x^2 = y^1$ is false.

⁶Again, $x \neq y$ is just shorthand for $\neg(x = y)$

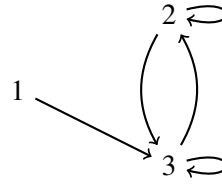
Similarly, we could show that ' $\exists x \exists y (x \neq y \wedge \forall z (z = x \vee z = y))$ ', that is, the claim that exactly two objects exist, is false on a domain that contains 1, 2, and 3. (Though for this, the complete semantic demonstration would be quite long, since there are three possible values to consider for ' x ', and then three possible values for ' y ' in each case!)

■ Exercises 7.11

A. Consider again the following interpretation:

Domain: 1, 2, 3

R : $\langle 1,3 \rangle, \langle 2,2 \rangle, \langle 2,3 \rangle, \langle 3,2 \rangle, \langle 3,3 \rangle$



Determine whether each of the following sentences is true or false in this interpretation, and give a semantic demonstration to justify your answer.

1. $\forall x (\exists y x \neq y \wedge \exists y x = y)$
2. $\forall x \exists y (x \neq y \wedge x = y)$
3. $\exists x \forall y x = y$
4. $\exists x \exists y (x \neq y \wedge (Rxy \wedge Ryx))$
5. $\exists x \forall y (Rxy \leftrightarrow x = y)$
6. $\exists x \forall y (Ryx \leftrightarrow x = y)$
7. $\exists x \exists y ((x \neq y \wedge Rxy) \wedge \forall z (Rzx \leftrightarrow y = z))$

B. Show that each of the following is neither a logical truth nor a contradiction:

1. $\forall x \exists y x \neq y$
2. $\exists x (x = a \wedge x = b)$
3. $\forall x \forall y \forall z ((x = y \vee x = z) \vee y = z)$
4. $\exists x \exists y x \neq y$

Show that the following sentences are jointly consistent (i.e. there's an interpretation that makes them all true):

5. $\neg Raa, \forall x (x = a \vee Rxa)$
6. $\forall x \forall y \forall z (x = y \vee y = z \vee x = z), \exists x \exists y \neg x = y$
7. $\exists x \exists y (Zx \wedge Zy \wedge x = y), \neg Za, a = b$

7.12 Appendix: Semantics with Variable Assignments

In §7.4 we explained the semantics for quantifiers using the notion of a variable assignment: an assignment of an object to a variable, representing our decision to temporarily treat the variable as if it referred to that object. However, once we introduce variable assignments

to deal with quantifiers, we should really go back and re-do our semantics with variable assignments in mind from the beginning, even when we're just looking at atomic sentences.

To see why, consider a simple example like ' $\forall xFx$ ', and suppose our domain contains just the numbers 1 and 2. Then in order for ' $\forall xFx$ ' to be true, the formula ' Fx ' needs to be true on the assignment $[x : 1]$ and also on $[x : 2]$. But what is required for ' Fx ' to be true on $[x : 1]$, say? The answer, of course, is that ' F ' has to be true of 1, the object we're treating ' x ' as referring to. But the semantics we gave for atomic sentences doesn't *technically* tell us this! What we said in §7.3 was just the following:

When \mathcal{R} is an n -place predicate and c_1, c_2, \dots, c_n are names, the sentence $\mathcal{R}c_1c_2\dots c_n$ is true in a given interpretation **iff** \mathcal{R} is true of the objects referred to by c_1, c_2, \dots, c_n (in that order) in that interpretation

This doesn't mention variable assignments anywhere. And furthermore, it only deals with atomic *sentences*, things that only contain *names*. ' Fx ' contains a variable instead of a name. It is merely a *formula*, so the above explanation simply does not apply to it.

To remedy this situation, we need to return to the official explanation of the syntax, or grammar, of FOL that we gave way back in §5.10. Recall that we there defined the general class of *formulas* of FOL, among which we then singled out the smaller class of FOL *sentences*. We did this in four steps. First, we stipulated that a TERM is to be any name *or any variable*. Second, we said that an ATOMIC FORMULA is anything that results from combining any predicate (including the identity symbol) with an appropriate number of *terms*. Third, we said that a FORMULA, more generally, is anything that can be built up from atomic formulas using truth functional operators and quantifiers. And lastly, we then said that a SENTENCE is any formula that contains no free, or unbound, variables.

We here have to do something similar. We will first have to give a general explanation of truth that covers all *formulas*, such as ' Fx ', and then extract from that a notion of truth for just *sentences*. In particular, what we'll do is to explain what is required for any formula whatsoever to be true in an interpretation *relative to a variable assignment*. In what follows, let's use \mathcal{I} for an arbitrary interpretation, and σ for an arbitrary variable assignment. So we will give a general explanation of what is required for any *formula* to be *true on* σ *in* \mathcal{I} . Later we'll then use that to say what's required for a *sentence* to be true in an interpretation.

To begin, we first give an expanded notion of reference that covers both names (like $a, b, c \dots$) and variables (like $x, y, z \dots$). Again: a variable doesn't actually refer to anything, since it isn't a name. But *relative to a variable assignment*, variables can be construed as referring to things. So:

For any term t (name or variable), the REFERENT OF t ON σ IN \mathcal{I} is:

- ▷ whatever object the interpretation \mathcal{I} assigns to t , if t is a name, and
- ▷ whatever object the assignment σ assigns to t , if t is a variable

We can now use this expanded notion of reference to explain truth for all atomic formulas, whether they contain names or variables:

When \mathcal{R} is an n -place predicate and t_1, \dots, t_n are any terms (names *or* variables), the formula $\mathcal{R}t_1\dots t_n$ is true on σ in \mathcal{I} **iff** \mathcal{R} is true of the

objects referred to by t_1, \dots, t_n (in that order) on $[]$ in \mathcal{I} .

And for any terms t_1 and t_2 , the formula $t_1 = t_2$ is true on $[]$ in \mathcal{I} iff t_1 and t_2 refer to the same thing on $[]$ in \mathcal{I} .

This now *does* specify what's required for e.g. ' Fx ' to be true on $[x : 1]$. What's required is that ' F ' be true of the referent of ' x ' on $[x : 1]$. And the referent of ' x ' on this assignment is of course just 1, give our expanded notion of reference. So ' F ' has to be true of 1. With atomic formulas covered, we can go on to explain truth for all complex formulas:

$\neg\phi$ is true on $[]$ in \mathcal{I} **iff** ϕ false on $[]$ in \mathcal{I}

$\phi \wedge \psi$ is true on $[]$ in \mathcal{I} **iff** both ϕ and ψ are true on $[]$ in \mathcal{I}

$\phi \vee \psi$ is true on $[]$ in \mathcal{I} **iff** either ϕ or ψ are true on $[]$ in \mathcal{I}

$\phi \rightarrow \psi$ is true on $[]$ in \mathcal{I} **iff** either ϕ is false or ψ is true on $[]$ in \mathcal{I}

$\phi \leftrightarrow \psi$ is true on $[]$ in \mathcal{I} **iff** ϕ and ψ have the same truth value on $[]$ in \mathcal{I}

$\forall v \phi(\dots v \dots)$ is true on $[]$ in \mathcal{I} **iff** for *every object* o in the domain, $\phi(\dots v \dots)$ is true on $[v:o]$ in \mathcal{I}

$\exists v \phi(\dots v \dots)$ is true on $[]$ in \mathcal{I} **iff** for *at least one object* o in the domain, $\phi(\dots v \dots)$ is true on $[v:o]$ in \mathcal{I} .

In the clauses governing the quantifiers, $[v:o]$ is the assignment that's just like our arbitrary assignment $[]$ in every respect, except that it is stipulated to assign the object o to the variable v . The idea, again, being that we go through each object o in the domain, and check whether the formula $\phi(\dots v \dots)$ is true when that object o is assigned to the variable v . If so, $\forall v \phi(\dots v \dots)$ is true, if not, it is false. On the other hand, for an existential sentence $\exists v \phi(\dots v \dots)$ to be true, it's enough if *at least one* object makes $\phi(\dots v \dots)$ true.

Given all of this, we can now say what is required for any *sentence*, i.e. any closed formula, or formula with no free variables, to be true in an interpretation \mathcal{I} :

For any FOL sentence ϕ and any interpretation \mathcal{I} , ϕ is true in \mathcal{I} iff ϕ is true in \mathcal{I} on *any assignment* $[]$ whatsoever.

With this explanation of truth for sentences in place, the definition of the various logical concepts (like entailment, logical truth, equivalence etc.) now proceeds as it did in §7.6.

Quick Reference

8

Truth Functional Operators (Ch. 3)

ϕ	ψ	$\neg\phi$	$\phi \wedge \psi$	$\phi \vee \psi$	$\phi \rightarrow \psi$	$\phi \leftrightarrow \psi$
T	T	F	T	T	T	T
T	F	F	F	T	F	F
F	T	T	F	T	T	F
F	F	T	F	F	T	T

Deduction Rules for TFL (Ch. 4)

Conjunction Introduction

m	ϕ	
n	ψ	
	$\phi \wedge \psi$	$\wedge I m, n$

Conjunction Elimination

m	$\phi \wedge \psi$	
	ϕ	$\wedge E m$
m	$\phi \wedge \psi$	
	ψ	$\wedge E m$

Conditional Introduction

i	ϕ	Assumption
j	ψ	
	$\phi \rightarrow \psi$	$\rightarrow I i-j$

Conditional Elimination

m	$\phi \rightarrow \psi$	
n	ϕ	
	ψ	$\rightarrow E m, n$

Biconditional Introduction

i	ϕ	Assumption
j	ψ	
k	ψ	Assumption
l	ϕ	
	$\phi \leftrightarrow \psi$	$\leftrightarrow I i-j, k-l$

Biconditional Elimination

m	$\phi \leftrightarrow \psi$	
n	ϕ	
	ψ	$\leftrightarrow E m, n$

m	$\varphi \leftrightarrow \psi$	
n	ψ	
	φ	$\leftrightarrow E\ m, n$

Negation Introduction

m	φ	Assumption
n	\perp	
	$\neg\varphi$	$\neg I\ m-n$

Absurdity Introduction

m	φ	
n	$\neg\varphi$	
	\perp	$\perp I\ m, n$

Indirect Proof

m	$\neg\varphi$	Assumption
n	\perp	
	φ	IP $m-n$

Disjunction Introduction

m	φ	
	$\varphi \vee \psi$	$\vee I\ m$
m	ψ	
	$\psi \vee \varphi$	$\vee I\ m$

Disjunction Elimination

m	$\varphi \vee \psi$	
i	φ	Assumption
j	χ	
k	ψ	Assumption
l	χ	
	χ	$\vee E\ m, i-j, k-l$

Derived Rules for TFL (§4.11)

Sequent	Derived Rule
$\varphi \rightarrow \psi, \neg\psi \vdash \neg\varphi$	MT
$\varphi \vee \psi, \neg\psi \vdash \varphi$	DS
$\varphi \vee \psi, \neg\varphi \vdash \psi$	DS
$\varphi \vdash \psi \rightarrow \varphi$	PMI
$\neg\varphi \vdash \varphi \rightarrow \psi$	PMI
$\varphi \rightarrow \psi \dashv\vdash \neg\varphi \vee \psi$	Imp
$\neg(\varphi \rightarrow \psi) \dashv\vdash \varphi \wedge \neg\psi$	NegImp
$\neg(\varphi \wedge \psi) \dashv\vdash \neg\varphi \vee \neg\psi$	DeM
$\neg(\varphi \vee \psi) \dashv\vdash \neg\varphi \wedge \neg\psi$	DeM
$\varphi \dashv\vdash \neg\neg\varphi$	DN
$(\varphi \# \psi) \dashv\vdash (\neg\neg\varphi \# \neg\neg\psi) \dashv\vdash (\neg\neg\varphi \# \psi) \dashv\vdash (\varphi \# \neg\neg\psi)$	SDN
$\neg(\varphi \# \psi) \dashv\vdash \neg(\neg\neg\varphi \# \neg\neg\psi) \dashv\vdash \neg(\neg\neg\varphi \# \psi) \dashv\vdash \neg(\varphi \# \neg\neg\psi)$	SDN
$\varphi @ \psi \vdash \psi @ \varphi$	Com
$\perp \vdash \varphi$	EX
$\vdash \varphi \vee \neg\varphi$	LEM

Deduction Rules for FOL (Ch. 6)

Universal Elimination

m	$\forall v \varphi(\dots v \dots)$	
	$\varphi(\dots c \dots)$	$\forall E\ m$

Universal Introduction

m	c	Flag
n	$\varphi(\dots c \dots)$	
	$\forall v \varphi(\dots v \dots)$	$\forall I\ m-n$

The Flag-ed name c may not occur outside the subproof.

Existential Introduction

m	$\varphi(\dots c \dots)$	
	$\exists v \varphi(\dots v \dots)$	$\exists I\ m$

Existential Elimination

m	$\exists v \varphi(\dots v \dots)$	
i	$\varphi(\dots c \dots)$	Assumption (flag c)
j	ψ	
	ψ	$\exists E\ m, i-j$

The Flag-ed name c may not occur outside the subproof.

Identity Elimination

m	$a = b$	
n	$\varphi(\dots a \dots a \dots)$	
	$\varphi(\dots b \dots a \dots)$	$=E\ m, n$
m	$a = b$	
n	$\varphi(\dots b \dots b \dots)$	
	$\varphi(\dots a \dots b \dots)$	$=E\ m, n$

Identity Introduction

	$c = c$	$=I$
--	---------	------