# 1 Introduction to model-theoretic inferentialism

The purpose of this book is to explore what rules of logic express about the meanings of the logical symbols they govern. Suppose that the only thing you know about the symbol '*' is that the following rules govern its behavior. Given an English sentence of the form A*B, it follows that A, and it also follows that B. Given A and B together, it follows that A*B. Can you tell what the symbol '*' means? Did you think that '*' must mean what we mean by 'and' in English, and that the truth behavior of sentences involving '*' must conform to the standard truth table for conjunction? If you did, can you be certain that the same deductive role for '*' specified by these rules might not also allow some alternative (or unintended) interpretation for '*'?

## 1.1 The broader picture

Questions like this are special cases of a general concern in the philosophy of language. To what extent can the meanings of expressions of a language be defined by the roles they play in our reasoning? Does knowing the meaning of a sentence simply amount to knowledge of which sentences entail it and which ones it entails? Would it be possible at least in principle for alien anthropologists from a planet circling the star Alpha Centauri (who know initially nothing about our language) to learn what our sentences mean by simply investigating the way we reason from one to another in different circumstances?

The predominant research methodology for understanding natural language in the field of artificial intelligence assumes the answer is: Yes. It is presumed that the project of designing computers capable of understanding a natural language such as English amounts to the development of systems capable of drawing the conclusions humans typically draw from a body of available text. So if a computer program can summarize correctly

1

articles from the *New York Times*, or cogently answer questions about a visit to a restaurant (Shank and Abelson, 1977), then the system counts as one that understands, and in fact models how understanding is possible in the natural world. If this research program is right, the problem of what sentences mean is reduced to resolving the issues involved in storing, efficiently retrieving, and drawing conclusions from relevant information found in a massive data-base incorporating knowledge of a typical human (Davis and Lenat, 1982).

Of course there is wide disagreement about the viability of such an answer to the problem of understanding meaning. Everyone from Fodor (1988, Chapter 2) to Searle (1980) objects, and for widely varying reasons. (See Cummins (1989, Chapters 4–6) for a useful discussion.) Among the many challenges that such a view faces is the complaint that the inferential roles of an expression are by themselves too weak to determine the semantical interpretations we intuitively assign to them. Worries about underdetermination of meaning motivate externalist strategies, such as covariance and causal role theories. These hold that part of what determines meaning is connections to the outside world mediated by our perceptual abilities. So to have a full understanding of what 'red' means, you would have to associate this word with the visual sensation you have when you see blood. Others point out the importance of connections to our actions as well as to our perceptions. If they are right, the aliens who wanted to know what we mean would have to study not only the logical relations between our utterances, but also the connections between those utterances and what we sense, and how we act.

Nevertheless, it is fairly widely held that inferential relations between sentences will contribute an important part of the answer about how meaning is determined. The idea is that at least part of what helps fix the meaning of a sentence like 'McKinley was assassinated' is the deductive connections it has with other sentences such as 'McKinley was killed', 'McKinley died', 'McKinley was shot', and even 'McKinley was president'. A person who does not understand how the truth of the latter sentences are inferentially related to the truth of 'McKinley was assassinated' doesn't really know what that sentence means.

One of the difficulties in trying to resolve questions about how inferential relations play a role in fixing meaning is that working out a theory in detail for a given language is such a massive project. If you are truly serious

about supporting such an inferentialist account, the best plan would be to show in detail that it has the resources to determine the meanings of hundreds of thousands of words and billions of expressions. Unfortunately, that project is still far beyond us. In the meantime, a sensible research strategy for controlling some of the complexity starts by defining solvable "toy" problems. Once these are understood, the lessons learned can be applied to more complex cases. The question of how and whether the deductive roles of the logical symbols can determine their meanings is a good example of such a toy problem. Here it is less plausible to think that relations to perception and action are essential to the meanings of the logical symbols. If there is any hope for a theory that presumes that meaning can be defined by logical relations alone, the best place to look would be in a study of the meanings of the logical connectives. One of the purposes of this book is to explore that question in detail and to report on the lessons learned. The realm of logic is especially well chosen since here it is possible to give mathematical demonstrations that answer such questions as whether connective meaning is or is not underdetermined by a set of rules. So the question on the table will be whether the alien anthropologists, knowing what will be revealed in this book, can determine what our logical symbols mean from a study of the rules we use.

An important conclusion established here is that the answers depend on decisions about the format of the rules, and decisions about how to define what rules express about connective meaning. As different choices are made, the view that the patterns of inference set up by the rules of deduction are sufficient for determining connective meaning are in some cases vindicated and in others undermined. By noting the strategies that can be employed to resolve problems of meaning underdetermination for the logical connectives, new insights may be obtained about how one might resolve problems of meaning underdetermination for language in general.

## 1.2  Proof-theoretic and model-theoretic inferentialism

A quick way to summarize what this book is about is to say that it explores the viability of a brand of inferentialism in logic. A logical inferentialist, or inferentialist for short, is somebody who believes that the rules of logic alone determine the meanings of the logical connectives. (For an excellent review of inferentialism including comparisons with the work of Brandom

(1994, 2001), see Tennant (forthcoming, Section 2).) However, there are at least two rather different ways of identifying the foundational commitments of the view. One is to think of inference as the rival of reference for serving as the foundation for a theory of meaning. A referentialist theory of meaning takes the idea of denotation of an expression as basic. Names refer to objects, and in general, all words refer to something. In the model-theoretic tradition, the idea is fleshed out by taking the reference (extension) of predicate letters to be sets of n-tuples of objects and the reference of well-formed formulas to be the truth-values t and f. One can even go so far as to claim that the reference of a connective is its truth function.

St. Augustine advocates an idea like this in the famous quotation that opens Wittgenstein's *Investigations*. "Thus as I heard words repeatedly used in their proper places in various sentences, I gradually learnt to understand what objects they signified; and after I had trained my mouth to form these signs, I used them to express my own desires" (Wittgenstein, 1969, p. 2). Certainly this was Wittgenstein's doctrine on what explained the meaning of names in the *Tractatus Logico-Philosophicus* (Wittgenstein, 1961, 3.203). Wittgenstein's point in placing such a referential theory of language on the table was to launch an argument that a theory of meaning based on word reference accounts, at best, for a very limited set of cases. To fully understand how language works, one must appreciate the rich variety in the ways in which words have a *use* in the social activities. Even reference itself makes no sense apart from those wider concerns. For those who take this to heart in the case of logic, it is natural to found a theory of the use of logical connectives on their inferential roles, for it is the activity of assessing inferential relations that appears relevant to the social role of logic.

Inferential theories will be particularly attractive to those who have epistemological worries about reference, or about how reference can be fixed. Those theories will also attract people with anti-realist sympathies, and who advocate pragmatic or coherentist rather than correspondence theories of truth. Therefore it is no surprise that semantical theories for logic based on proof-theoretic role have historical roots in the work of intuitionists such as Brouwer. (See Sundholm (1986 especially Section 2) for a good discussion of the motives in this tradition.) Such proof-theoretic inferentialists reject model-theoretic semantics in the Tarskian tradition as misguided. Its truth conditions are non-constructive, and attempts to characterize the notion of truth without constructivist constraints leads to

paradox. Therefore a tradition has grown up among logical inferentialists to develop proof-theoretical semantics, which strictly avoids any mention of concepts from model theory such as reference, truth, and validity as preservation of truth.

Strictly speaking, however, the demand that model-theoretic notions be laid aside is not part of the central inferentialist doctrine. Proof-theoretic inferentialists subscribe to two doctrines that can and should be distinguished.

(Inferentialism)   The connectives obtain their meanings from the proof-theoretic roles that are established by the rules that govern them.

(Proof-Theoretic Semantics)   The meanings determined for the connectives must be characterized using only concepts found in proof theory. In particular, notations like denotation and truth are not to be employed.

(Proof-Theoretic Semantics) is not essential to inferentialism, despite its centrality in the historical tradition going back to intuitionism. This book demonstrates that an inferentialism that investigates how proof-theoretic roles determine model-theoretic readings is of technical and philosophical interest. Sundholm (1986, p. 478) recommends a project of this kind. The idea is to find a way to "read off a [model-theoretic] semantics from the ... rules." He is describing what he takes to be a failed attempt by Hacking to carry this out in the case of classical logic, and goes on to say "the problem still remains open how to find a workable proposal along these lines." This book shows how to solve that problem.

Tennant (forthcoming, Section 2, note 7) notes that Brandom's views on inferentialism have evolved towards a kind of quietism with respect to epistemological and metaphysical commitments. Recognizing that inferentialism is orthogonal to those concerns opens room for brands of model-theoretic inferentialism. The model-theoretic tradition has important intuitions in its favor. The role of semantics in logic is to provide a definition of validity, and arguably validity amounts to the preservation of truth. Given this standard, proof-theoretic semantics, for all its technical interest, looks like an oxymoron. Although proof-theoretic notions of validity are in the offing (Schroeder-Heister, 2006) they are complex, and it is not clear how they meet the concerns for limning correctness of reasoning that motivated the concept of validity in the first place. From the model-theoretic side, however, defining validity is straightforward.

I do not wish to dwell on possible failings of proof-theoretic semantics here. It is sufficient for my purposes to argue for a pluralism in the style of semantics chosen, so that an inferentialism that uses model-theoretic notions counts as a live option. As we will see, such a view need not be a rival of proof-theoretic inferentialism. It is a way of thinking that can actually be of service to the proof-theoretic side. So a *model-theoretic inferentialist* subscribes to these two theses: (Inferentialism), of course, and (Model-Theoretic Semantics).

(Inferentialism)   The connectives obtain their meanings from the proof-theoretic roles that are established by the rules that govern them.

(Model-Theoretic Semantics)   When characterizing the meanings of the connectives, it is of interest (and even helpful to those in the proof-theoretic tradition) to employ concepts from model theory such as truth, reference, and validity understood as preservation of truth.

The purpose of this book will not be to argue for model-theoretic inferentialism. Although, some comfort for the view is found in the result that rules for most connectives can be proven to fix exactly one model-theoretic interpretation for the connectives, there are also reasons for worry. In some cases (notably disjunction (Section 7.3), and the failure of referentialism in predicate logic (Section 14.7)), it can be argued that inferential roles of the connectives fail to determine the expected model-theoretic counterparts. In those negative cases, proof-theoretic inferentialists may find some support for their view.

Results of this book can also be helpful to the proof-theoretic tradition in another way. One of the major challenges to inferentialism in logic is Prior's (1960) famous demonstration that there are sets of rules that do not define an acceptable logical operation. His example of *tonk* showed that not all logical systems determine corresponding meanings for the connectives they regulate. The response of those with inferentialist sympathies has been to invent harmony constraints on the rules designed to guarantee that the rules are successful in defining connective meaning. Many definitions of harmony have been proposed. One of the first was Belnap's (1962) requirement that the rules be conservative and unique. Notions involving inversion and normalization (Prawitz, 1965, p. 33; Dummett, 1978, pp. 220–222; Weir, 1986; Tennant, 1997, pp. 308ff. and forthcoming) have also been introduced as the missing constraint. The difficulty

proof-theoretic inferentialists face in motivating these responses to Prior is to provide independent evidence that the constraints proposed are necessary conditions for defining connective meaning. Results found in this book will help inferentialists motivate such constraints. For example, it will be shown (in Sections 13.1 and 13.4) that any set of rules that determines a connective meaning from the model-theoretic point of view meets Belnap's conservation and uniqueness requirements. The approach taken here will also help us develop a new understanding of the inversion principle and normalization. (See Sections 13.5–13.7.)

A proof-theoretic inferentialist may still be uncomfortable with the model-theoretic project. Wasn't the whole point of inferentialism to avoid the realism and anti-verificationism that is implicit in the use of model-theoretic notions? There is a simple answer to that worry. It is to argue for metaphysical and epistemological quietism for model-theoretic notions of reference and truth. I submit that those notions, in themselves, commit us to nothing. This book shows that on very modest assumptions, the roles for the connectives set up by logical rules fix exactly certain corresponding truth conditions. These use the notation: 'v(A)=t' (for valuation v assigns to wff A the value t). What should an inferentialist antecedently committed to a pragmatic or coherentist theory of truth make of 'v(A)=t'? The answer is anything he or she likes, but in fairness, why not read 't' as coherence, or assertibility, rather than truth as correspondence? Model-theoretic inferentialism need not saddle one with any particular reading of the set theoretical machinery of model theory. As a matter of fact, as we will see in Section 13.9, 'v(A)=t' has a provability reading, so that what initially looks like model-theoretic semantics is transformed into a proof-theoretic semantics that is new to the literature. The upshot is that the mere use of the notation of model theory is compatible with a very broad range of epistemological and metaphysical views, including those of the founders and followers of the inferentialist tradition.

## 1.3  Three rule formats

This book shows that the answer to the question: 'do rules fix the meanings of the connectives?' is that it depends. One source of variability in the answers is the format in which we frame the logical rules. A lot depends on the details concerning the way the rules are defined. Three main approaches will be

explored here. The first is axiomatic. An *axiomatic system* lays down a set of axioms (or axiom forms) that serve as examples of logical truths, along with a collection of rules taking a formula or formulas into a new formula. For example, here is a economical axiomatic formulation for propositional logic in a language where → (if then) and ~ (not) are the only connectives.

$\vdash A{\to}(B{\to}A)$

$\vdash \ A{\to}(B{\to}C) \to ((A{\to}B) \to (A{\to}C))$

$\vdash \ ({\sim}A{\to}B) \to (({\sim}A{\to}{\sim}B) \to A)$

(Modus Ponens)

$\vdash A$

$\underline{\vdash A{\to}B}$

$\vdash B$

We use 'A', 'B', 'C', and 'D' as metavariables over well-formed formulas (wffs) of propositional logic. The notation '$\vdash ({\sim}A{\to}B) \to (({\sim}A{\to}{\sim}B) \to A)$' indicates that any wff with the displayed form is provable. In presenting formal systems, we use '$\vdash$' for 'is provable'. (We treat the logical symbols '~', '→', '&', etc. in the metalanguage as *used* to refer to symbols with similar shapes. It is also understood that '~A', for example, refers to the result of concatenating ~ with the wff A. This convention avoids the need to use corner quotes.)

As anyone who tries it knows, finding proofs of wffs in axiomatic systems can be difficult. For example, the shortest proof of A→A in the above system comes to six lines and requires the use of complex and non-obvious instances of the first two axioms.

The second tactic for defining a logic is to use *natural deduction* (ND) format. Here a pair of rules is provided for each connective showing how it is introduced into, and eliminated from inference. Proof finding in ND systems is greatly simplified because of an important innovation; one is allowed to make additional assumptions in the course of a deduction, which are then discharged by the application of the appropriate rules. For example, the rule (→ Introduction) asserts that when one is able to derive sentence B from having made an additional assumption A, then that derivation is a warrant for introducing the conditional A→B and eliminating A from the set of active assumptions. This innovation means that ND rules are defined over more complex structures than are axiomatic systems. An axiomatic proof is a sequence of wffs, each of which is an instance of an

axiom or one that follows from previous steps by a rule. However, it is useful to take it that ND rules are defined over arguments, not wffs. For example, (→ Introduction) takes the argument H, A / C (which asserts that C follows from the ancillary hypothesis A along with other hypotheses H) to the new argument H / A→C (which asserts that the conditional A→C follows from hypotheses H leaving aside the assumption A).

Here is an example of a ND system for a propositional logic with → and ~ as its only connectives, using "horizontal" notation that makes apparent the idea that ND rules are defined over arguments.

| (→ Introduction) | (→ Elimination) |
|---|---|
| $\underline{H, A \vdash B}$ | $H \vdash A$ |
| $H \vdash A{\to}B$ | $\underline{H \vdash A{\to}B}$ |
| | $H \vdash B$ |

| (~ Introduction) | (~ Elimination) |
|---|---|
| $H, A \vdash B$ | $H, {\sim}A \vdash B$ |
| $\underline{H, A \vdash {\sim}B}$ | $\underline{H, {\sim}A \vdash {\sim}B}$ |
| $H \vdash {\sim}A$ | $H \vdash A$ |

For ease of comparison with multiple conclusion sequent systems to be presented shortly, it will be assumed that H is a possibly infinite set of wffs. The notation 'H, A' is used as shorthand for $H \cup \{A\}$, and we sometimes omit set braces so that 'A, B' abbreviates '{A, B}'. In the case of ND systems, the symbol '/' is assumed to be in the object language, and a rule takes one from an argument or arguments to a new argument. The symbol '⊢' is used in the metalanguage to indicate the provability of an argument in a system being discussed. Therefore 'H ⊢ C' abbreviates the claim that the object language argument H / C has a proof in that system. (See Hacking (1979, p. 292), who adopts this convention.)

Natural deduction systems will play an important role in this book because of their interesting expressive powers. The results developed for them will help vindicate inferentialist intuitions that natural deduction rules have a special role to play in defining connective meaning.

The third format for presenting rules of logic is multiple conclusion *sequent* notation. A (multiple conclusion) sequent H / G is a generalization of the notion of an argument H / C, where the "conclusion" G is now taken to be a set of wffs. In this book, we will always use 'sequent' to refer to such a multiple conclusion sequent. The sequent H / G is understood to express the

idea that if all the wffs in the hypothesis set H are true, then at least one of the wffs in the conclusion set G is true. Sample rules G→~ for a sequent formulation for a propositional logic with → and ~ as its only connectives follow.

G→~:          (→Left)                     (→Right)
              G ⊢ A, H                     G, A ⊢ B, H
              ‾‾‾‾‾‾‾‾‾                     ‾‾‾‾‾‾‾‾‾‾‾
              G, B ⊢ H                     G ⊢ A→B, H
              ‾‾‾‾‾‾‾‾‾
              G, A→B ⊢ H

              (~Left)                      (~Right)
              G, ⊢ A, H                    G, A ⊢ H
              ‾‾‾‾‾‾‾‾‾                     ‾‾‾‾‾‾‾‾‾
              G, ~A ⊢ H                     G, ⊢ ~A, H

It is useful for comparing systems in different formats to treat (multiple conclusion) sequent format as the most general case, and to define axiomatic and ND systems as special cases of sequent systems. So let an *argument* be defined as a sequent whose conclusion has a single member, and let an *assertion* be an argument with an empty set of hypotheses. Call the items to which a rule is applied the *inputs* to the rule, and let the result of applying the rule be called its *output*. Then a ND system is simply a set of sequent rules whose inputs and outputs are all arguments. Similarly, axiomatic systems are systems whose sequent rules have assertions as their inputs and outputs.

## 1.4  Expressive power and models of rules

This book is about the expressive power of the rules of logic. To what degree does acceptance of the principles of logic force a particular interpretation of its connectives? As we said, the answer depends on how the rules are formulated, but it will also depend on how we define expressive power. Let us explore some of the options.

It helps to start with an analogy from model theory. The idea of a sentence (or group of sentences) expressing a condition on a model should be familiar. For example, the sentence ∃x∃y~x=y expresses that the domain of a model contains at least two objects. The reason is that ∃x∃y~x=y is true on a model exactly when the domain of the model domain contains at least two objects. In general, wff A expresses a property P of models iff A is true on any model M exactly when M has property P. When A is true on a model, we say M is a model of A. So when we say that A expresses P, we mean that for every model