

## The Central Theorems of Kripke's Theory of Truth

We start with an interpreted language  $\mathcal{L}$ , and we add a new predicate "Tr," which is supposed to represent the set of true sentences of the enlarged language. We need to determine how the new predicate is to behave. In classical semantics, without truth-value gaps, this means assigning a set  $E$  as the extension of "Tr," so that, if the closed term  $\tau$  denotes something in  $E$ ,  $\text{Tr}(\tau)$  will be true, whereas, if  $\tau$  denotes something outside  $E$ ,  $\text{Tr}(\tau)$  will be false and  $\neg\text{Tr}(\tau)$  will be true. Because of the Liar Paradox, it isn't possible to pick the set  $E$  so that  $E =$  the set of sentences that are made true by assigning  $E$  as the extension of "Tr." In other words, we can't arrange things so that  $\text{Tr}([\ulcorner S \urcorner])$  is true or false according as  $S$  is true or false.

If we allow truth-value gaps, we shall need to assign to "Tr" an *extension*  $E$ , consisting of things that definitely satisfy the predicate, and an *antiextension*  $A$ , consisting of things that are definitely fail to satisfy it.  $E$  and  $A$  cannot overlap, but there may be things that aren't in either of them. If  $\text{Den}(\tau)$  is in  $E$ ,  $\text{Tr}(\tau)$  will be true, whereas if  $\text{Den}(\tau)$  is in  $A$ ,  $\text{Tr}(\tau)$  will be false. If  $\text{Den}(\tau)$  is in neither  $E$  nor  $A$ ,  $\text{Tr}(\tau)$  will be neither true nor false. The central theorem underlying Kripke's theory of truth (which was proven also by Robert Martin and Peter Woodruff) tells us that we can find a pair of nonoverlapping sets  $E$  and  $A$  with the following properties:

$E =$  the set of sentences that are made true by setting the extension of "Tr" equal to  $E$  and setting the antiextension equal to  $A$ .

$A =$  the union of the set of nonsentences with the set of sentences that are made false by setting the extension of "Tr" equal to  $E$  and the antiextension equal to  $A$ .

In other words,  $\text{Tr}([\ulcorner S \urcorner])$  will be true, false, or unsettled according as  $S$  is true, false, or unsettled. We have thus arranged things so that  $\text{Tr}([\ulcorner S \urcorner])$  and  $S$  always have the same semantic status, even though we haven't arranged things so that sentences of the form

$$\text{Tr}([\ulcorner S \urcorner]) \leftrightarrow S$$

are all true.

We'll prove the theorem for  $\mathcal{L}$  the language of arithmetic, but we're using the language of arithmetic as a stand-in for a vast range of other languages we might equally well have used. In particular, the language of arithmetic lacks truth-value gaps, so that any truth value gaps we encounter occur because of the indeterminacy of "Tr." It doesn't have to be this way. The construction works just as well if the base language has truth-value gaps.

## I. The Classical Situation

The language  $\mathcal{L}^{\text{Tr}}$  is obtained from the language of arithmetic  $\mathcal{L}$  by adjoining the new predicate "Tr." Thus the *terms* of  $\mathcal{L}^{\text{Tr}}$  are specified by the following stipulation:

"0" is a term.

Each of the variables "x," "x'," "x'," "x''", "x''''", "x''''''", and so on, is a term.

[These are the variables for official purposes. In practices, the notation is cumbersome, so I'll ask other letters from the end of the alphabet, with or without numerical subscripts, to play the role of variables.]

If  $\tau$  is a term,  $S\tau$  is a term.

If  $\tau$  and  $\rho$  are terms,  $(\tau + \rho)$  is a term.

If  $\tau$  and  $\rho$  are terms  $(\tau \times \rho)$  is a term.

Nothing is a term unless it is required to be by the five clauses above.

A *closed term* is a term with no variables. The *denotation* of closed terms is defined inductively, as follows (for any closed terms  $\tau$  and  $\rho$ ):

$\text{Den}("0") = 0.$

$\text{Den}(S\tau) = \text{Den}(\tau) + 1$

$\text{Den}((\tau + \rho)) = \text{Den}(\tau) + \text{Den}(\rho).$

$\text{Den}((\tau \times \rho)) = \text{Den}(\tau) \cdot \text{Den}(\rho).$

The *formulas* of  $\mathcal{L}^{\text{Tr}}$  are specified as follows:

If  $\tau$  and  $\rho$  are terms, then

$\tau = \rho$

$\tau < \rho$

and

$\text{Tr}(\tau)$

are formulas.

If P and Q are formulas, so are

$(P \wedge Q)$

$\neg P$   
and  
 $(\forall v)P,$

for each variable  $v$ .

Nothing is a formula unless it's required to be by the clauses above.

In the past, we've included other symbols in what we were calling the language of arithmetic: the disjunction sign " $\vee$ ," the conditional " $\rightarrow$ ," the biconditional " $\leftrightarrow$ ," and the existential quantifier " $\exists$ ." We're leaving them out here, but it's no real loss, since we can treat them as defined:

$$(P \vee Q) =_{\text{Def}} \neg(\neg P \wedge \neg Q)$$

$$(P \rightarrow Q) =_{\text{Def}} \neg(P \wedge \neg Q)$$

$$(P \leftrightarrow Q) =_{\text{Def}} (\neg(P \wedge \neg Q) \wedge \neg(\neg P \wedge Q))$$

$$(\exists v)P =_{\text{Def}} \neg(\forall v)\neg P$$

The reason for leaving them out of the official language is that the proofs we are about to give are very long and tedious, leaving out unneeded symbols makes the proofs not quite as bad, by avoiding repetitions.

An occurrence of a variable  $v$  in a formula is *bound* if it occurs within some subformula that begins with " $(\forall v)$ ." (We count a formula as a subformula of itself.) If not bound, *free*. A formula with no free variables is a *sentence*. It is a sentence that can say something that is either true or false.

Let  $\mathbb{N}$  be the standard model of  $\mathcal{L}$ , the model whose domain is the natural numbers, in which "0" denotes the number 0 and "+" stands for addition and "×" for multiplication. We extend  $\mathbb{N}$  to a *classical model*  $(\mathbb{N}, E)$  of  $\mathcal{L}^{\text{Tr}}$  by designating a set  $E$  of natural numbers to serve as the extension of "Tr." We assume that we have fixed a Gödel numbering, so that we'll talk about sentences being members of  $E$ . We specify the conditions under which a sentence is *true* in  $(\mathbb{N}, E)$ , as follows:

$$\tau = \rho \text{ is true in } (\mathbb{N}, E) \text{ iff } \text{Den}(\tau) = \text{Den}(\rho).$$

$$\tau < \rho \text{ is true in } (\mathbb{N}, E) \text{ iff } \text{Den}(\tau) < \text{Den}(\rho).$$

$$\text{Tr}(\tau) \text{ is true in } (\mathbb{N}, E) \text{ iff } \text{Den}(\tau) \in E.$$

$(P \wedge Q)$  is true in  $(\mathbb{N}, E)$  iff  $P$  and  $Q$  are both true in  $(\mathbb{N}, E)$ .

$\neg P$  is true in  $(\mathbb{N}, E)$  iff  $P$  isn't true in  $(\mathbb{N}, E)$ .

$(\forall v)P(v)$  is true in  $(\mathbb{N}, E)$  iff, for each closed term  $\tau$ ,  $P(\tau)$  is true in  $(\mathbb{N}, E)$ .

Here " $P(\tau)$ " means the result of substituting  $\tau$  for each free occurrence of " $v$ " in " $P(v)$ ."

A sentence is said to be *false* in  $(\mathbb{N}, E)$  iff its negation is true in  $(\mathbb{N}, E)$ . Consequently,

$\tau = \rho$  is false in  $(\mathbb{N}, E)$  iff  $\text{Den}(\tau) \neq \text{Den}(\rho)$ .

$\tau < \rho$  is false in  $(\mathbb{N}, E)$  iff  $\text{Den}(\tau) \not\prec \text{Den}(\rho)$ .

$\text{Tr}(\tau)$  is false in  $(\mathbb{N}, E)$  iff  $\text{Den}(\tau) \notin E$ .

$(P \wedge Q)$  is false in  $(\mathbb{N}, E)$  iff either  $P$  or  $Q$  is false in  $(\mathbb{N}, E)$ .

$\neg P$  is false in  $(\mathbb{N}, E)$  iff  $P$  isn't false in  $(\mathbb{N}, E)$  iff  $P$  is true in  $(\mathbb{N}, E)$ .

$(\forall v)P(v)$  is false in  $(\mathbb{N}, E)$  iff, for some closed term  $\tau$ ,  $P(\tau)$  is false in  $(\mathbb{N}, E)$ .

The defined terms " $\vee$ ," " $\neg$ ," " $\rightarrow$ ," and " $\exists$ " behave just the way you'd expect:

A disjunction  $(P \vee Q)$  is true in  $(\mathbb{N}, E)$  iff one or both disjuncts are true in  $(\mathbb{N}, E)$ .

A disjunction is false in  $(\mathbb{N}, E)$  iff both its disjuncts are false in  $(\mathbb{N}, E)$ .

A conditional  $(P \rightarrow Q)$  is true in  $(\mathbb{N}, E)$  iff its antecedent  $P$  is false in  $(\mathbb{N}, E)$  or its consequent  $Q$  is true in  $(\mathbb{N}, E)$ .

A conditional is false in  $(\mathbb{N}, E)$  iff its antecedent is true in  $(\mathbb{N}, E)$  and its consequent is false in  $(\mathbb{N}, E)$ .

A biconditional  $(P \leftrightarrow Q)$  is true in  $(\mathbb{N}, E)$  iff its components are either both true or both false in  $(\mathbb{N}, E)$ .

A biconditional  $(P \leftrightarrow Q)$  is false in  $(\mathbb{N}, E)$  iff one of its components is true in  $(\mathbb{N}, E)$ , and the other is false in  $(\mathbb{N}, E)$ .

An existential sentence  $(\exists v)P(v)$  is true in  $(\mathbb{N}, E)$  iff, for some closed term  $\tau$ ,  $P(\tau)$  is true in  $(\mathbb{N}, E)$ .

$(\exists v)P(v)$  is false in  $(\mathbb{N}, E)$  iff, for every closed term  $\tau$ ,  $P(\tau)$  is false in  $(\mathbb{N}, E)$ .

We want to see that, whatever set  $E$  we pick as the extension of "Tr," we'll never have

$$E = \{\text{sentences true in } (\mathbb{N}, E)\}.$$

We'll use the Gödel Self-Referential Lemma. In its general form, the lemma tells us that, for any formula  $Q(v_1, v_2, \dots, v_n)$ , there is a formula  $P(v_1, v_2, \dots, v_n)$  so that the biconditional

$$(\forall v_1)(\forall v_2)\dots(\forall v_n)(Q([\text{Tr}(P(v_1, v_2, \dots, v_n))] , v_1, v_2, \dots, v_n) \leftrightarrow P(v_1, v_2, \dots, v_n))$$

is a theorem of Peano Arithmetic (PA). In case  $n = 0$ , this tells us that, for any formula  $Q(v)$ , we can find a sentence  $P$  so that

$$(Q([\text{Tr}(P)] ) \leftrightarrow P)$$

is a theorem of PA. Thus we can find a sentence  $L$  so that

$$(\neg \text{Tr}([\text{Tr}(L)] ) \leftrightarrow L)$$

is a theorem of PA.  $L$  is the formal equivalent of Eubulides' "This statement is not true." For any set of numbers  $E$ ,  $(\neg \text{Tr}([\text{Tr}(L)] ) \leftrightarrow L)$  is true in  $(\mathbb{N}, E)$ . Consequently,

$$\begin{aligned} L \in \{\text{sentences true in } (\mathbb{N}, E)\} \\ \text{iff } L \text{ is true in } (\mathbb{N}, E) \\ \text{iff } \neg \text{Tr}([\text{Tr}(L)] ) \text{ is true in } (\mathbb{N}, E) \\ \text{iff } \text{Tr}([\text{Tr}(L)] ) \text{ isn't true in } (\mathbb{N}, E) \\ \text{iff } \text{Den}([\text{Tr}(L)] ) \notin E \\ \text{iff } L \notin E \end{aligned}$$

So we cannot have  $E = \{\text{sentences true in } (\mathbb{N}, E)\}$ .

## II. Truth-value Gaps

We get truth-value gaps when we allow "Tr" to be partially undefined. We assign to "Tr" a pair  $(E, A)$ , with  $E \cap A = \emptyset$ . The extension  $E$  is intended to consist of things that are definitely true, while the antiextension  $A$  contains things that are definitely not true. We stipulate what it is for a sentence to be true or false in  $(\mathbb{N}, (E, A))$ :

$$\begin{aligned} \tau = \rho \text{ is true in } (\mathbb{N}, (E, A)) \text{ iff } \text{Den}(\tau) = \text{Den}(\rho). \\ \tau = \rho \text{ is false in } (\mathbb{N}, (E, A)) \text{ iff } \text{Den}(\tau) \neq \text{Den}(\rho). \end{aligned}$$

- $\tau < \rho$  is true in  $(\mathbb{N},(E,A))$  iff  $\text{Den}(\tau) < \text{Den}(\rho)$ .  
 $\tau < \rho$  is false in  $(\mathbb{N},(E,A))$  iff  $\text{Den}(\tau) \not< \text{Den}(\rho)$ .  
 $\text{Tr}(\tau)$  is true in  $(\mathbb{N},(E,A))$  iff  $\text{Den}(\tau) \in E$ .  
 $\text{Tr}(\tau)$  is false in  $(\mathbb{N},(E,A))$  iff  $\text{Den}(\tau) \in A$ .  
 $(P \wedge Q)$  is true in  $(\mathbb{N},(E,A))$  iff  $P$  and  $Q$  are both true in  $(\mathbb{N},(E,A))$ .  
 $(P \wedge Q)$  is false in  $(\mathbb{N},(E,A))$  iff one or both of  $P$  and  $Q$  are false in  $(\mathbb{N},(E,A))$ .  
 $\neg P$  is true in  $(\mathbb{N},(E,A))$  iff  $P$  is false in  $(\mathbb{N},(E,A))$ .  
 $\neg P$  is false in  $(\mathbb{N},(E,A))$  iff  $P$  is true in  $(\mathbb{N},(E,A))$ .  
 $(\forall v)P(v)$  is true in  $(\mathbb{N},(E,A))$  iff, for each closed term  $\tau$ ,  $P(\tau)$  is true in  $(\mathbb{N},(E,A))$ .  
 $(\forall v)P(v)$  is false in  $(\mathbb{N},(E,A))$  iff, for some closed term  $\tau$ ,  $P(\tau)$  is false in  $(\mathbb{N},(E,A))$ .

In the presence of truth-value gaps, the Gödel Self-Referential Lemma takes a slightly different form. Given a formula  $Q(v_0, v_1, v_2, \dots, v_n)$ , there is a formula  $P(v_1, v_2, \dots, v_n)$  such that, for any closed terms  $\tau_1, \tau_2, \dots, \tau_n$  and any nonoverlapping sets  $E$  and  $A$ ,  $Q([\text{ }^+ P(v_1, v_2, \dots, v_n) \text{ }], \tau_1, \tau_2, \dots, \tau_n)$  is true, false, or unsettled under  $(\mathbb{N},(E,A))$  according as  $P(\tau_1, \tau_2, \dots, \tau_n)$  is true, false, or unsettled under  $(\mathbb{N},(E,A))$ .

### III. The Smallest Fixed Point

A *fixed point* is a pair  $(E,A)$  with

$$\begin{aligned}
 E &= \{\text{sentences true in } (\mathbb{N},(E,A))\} \\
 A &= \{\text{nonsentences}\} \cup \{\text{sentences false in } (\mathbb{N},(E,A))\}
 \end{aligned}$$

We intend to show that there is a fixed point  $(E_\infty, A_\infty)$ .  $(E_\infty, A_\infty)$  is the *smallest* fixed point, in the following sense: If  $(E,A)$  is another fixed point, we have  $E_\infty \subseteq E$  and  $A_\infty \subseteq A$ .

If you look at Kripke's paper, you will see that he constructs  $E_\infty$  and  $A_\infty$  simultaneously. The argument he gives uses infinite ordinal numbers, and we haven't talked about infinite ordinal numbers, so I'd like to give a more direct, but also more laborious, argument here. It will simplify matters if we begin by constructing  $E_\infty$ , without worrying about  $A_\infty$ . We construct  $E_\infty$  by including within it all the sentences that have to be true, without fretting over how to make  $A_\infty$  include all the sentences that have to be false. Once we're done, we can obtain  $A_\infty$  as  $\{\text{nonsentences}\} \cup \{\text{sentences } S: \neg S \in E_\infty\}$ . Constructing  $E_\infty$  by itself simplifies matters, but it doesn't simplify matters very much, since falsity conditions are just truth conditions for negations. Thus, instead of having to deal separately with the falsity conditions for atomic sentences, conjunctions, negations, and universal sentences, we deal with the truth conditions for negated atomic sentences, negated conjunctions, negated negations, and negated universal sentences, in addition to having to deal with the truth conditions for atomic sentences, conjunctions, and universal sentences. So the definition of  $E_\infty$  is still unpleasantly complicated.

$E_\infty$  is the smallest collection of sentences that satisfies the following twelve conditions:

- i) If  $\text{Den}(\tau) = \text{Den}(\rho)$ , then  $\tau = \rho$  is in  $E_\infty$ .
- ii) If  $\text{Den}(\tau) \neq \text{Den}(\rho)$ , then  $\neg \tau = \rho$  is in  $E_\infty$ .
- iii) If  $\text{Den}(\tau) < \text{Den}(\rho)$ , then  $\tau < \rho$  is in  $E_\infty$ .
- iv) If  $\text{Den}(\tau) \not< \text{Den}(\rho)$ , then  $\neg \tau < \rho$  is in  $E_\infty$ .
- v) If  $\text{Den}(\tau) \in E_\infty$ , then  $\text{Tr}(\tau)$  is in  $E_\infty$ .
- vi) If  $\text{Den}(\tau)$  isn't a sentence, then  $\neg \text{Tr}(\tau)$  is in  $E_\infty$ .
- vii) If  $\text{Den}(\tau)$  is a sentence whose negation is in  $E_\infty$ , then  $\neg \text{Tr}(\tau)$  is in  $E_\infty$ .
- viii) If  $P$  and  $Q$  are in  $E_\infty$ , then  $(P \wedge Q)$  is in  $E_\infty$ .
- ix) If either  $\neg P$  or  $\neg Q$  is in  $E_\infty$ ,  $\neg(P \wedge Q)$  is in  $E_\infty$ .
- x) If  $P$  is in  $E_\infty$ , then  $\neg \neg P$  is in  $E_\infty$ .
- xi) If  $F(\tau)$  is in  $E_\infty$ , for each closed term  $\tau$ , then  $(\forall v)F(v)$  is in  $E_\infty$ .
- xii) If  $\neg F(\tau)$  is in  $E_\infty$ , for some closed term  $\tau$ , then  $\neg(\forall v)F(v)$  is in  $E_\infty$ .

When I say that  $E_\infty$  is the "smallest" set that meets that twelve conditions, I mean that  $E_\infty$  is a set that meets the twelve conditions that is included in every other set that meets the twelve conditions. To see that there have to be such a set, notice that there is a set that meets the twelve conditions, namely, the set of all sentences. Moreover, the intersection of all sets of sentences that meet the twelve conditions meets the twelve conditions. That intersection is  $E_\infty$ .

The definition is set us in such a way that any sentence that is in  $E_\infty$  has been put into  $E_\infty$  by one of the twelve conditions. Moreover, for any sentence in  $E_\infty$ , there is only one of the twelve conditions that could have put it there, and you can tell which condition that is by the syntactic form of the sentence.

Let  $A_\infty = \{\text{nonsentences}\} \cup \{\text{sentences } S: \neg S \in E_\infty\}$ . Our first task is to show that  $E_\infty$  and  $A_\infty$  are disjoint. We'll do this by showing that  $(E_\infty \sim A_\infty)$  satisfies the twelve conditions. Since  $E_\infty$  is the smallest set that meets the twelve conditions, if  $(E_\infty \sim A_\infty)$  is also a set that meets the twelve conditions, then  $E_\infty$  must be included in  $(E_\infty \sim A_\infty)$ , which means that no member of  $E_\infty$  is in  $A_\infty$ . So we have to verify the twelve conditions.

**i)** If  $\text{Den}(\tau) = \text{Den}(\rho)$ , then clause i) puts  $\tau = \rho$  in  $E_\infty$ , and none of the clauses puts  $\neg \tau = \rho$  into  $E_\infty$ . So  $\tau = \rho$  isn't in  $A_\infty$ .

**ii)** If  $\text{Den}(\tau) \neq \text{Den}(\rho)$ , then clause ii) puts  $\neg \tau = \rho$  into  $E_\infty$ . None of the clauses puts  $\tau = \rho$  into  $E_\infty$ . The only way  $\neg \neg \tau = \rho$  could be in  $E_\infty$  would be for it to be put there by clause x), and clause x) would only put  $\neg \neg \tau = \rho$  into  $E_\infty$  if  $\tau = \rho$  were in  $E_\infty$ . So  $\neg \neg \tau = \rho$  isn't in  $E_\infty$ , and  $\neg \tau = \rho$  isn't in  $A_\infty$ .

**iii)** Similar to i).

**iv)** Similar to ii).

**v)** Suppose that  $\text{Den}(\tau)$  is in  $E_\infty \sim A_\infty$ . Then, because  $E_\infty$  satisfies clause v),  $\text{Tr}(\tau)$  is in  $E_\infty$ . Also, because  $\text{Den}(\tau)$  is in  $E_\infty$ ,  $\text{Den}(\tau)$  must be a sentence, and so  $\neg\text{Tr}(\tau)$  isn't put into  $E_\infty$  by clause vi). The negation of  $\text{Den}(\tau)$  isn't in  $E_\infty$ , and so  $\neg\text{Tr}(\tau)$  wasn't put into  $E_\infty$  by clause vii). So  $\neg\text{Tr}(\tau)$  isn't in  $E_\infty$  at all; hence  $\text{Tr}(\tau)$  isn't in  $A_\infty$ .

**vi)** If  $\text{Den}(\tau)$  isn't a sentence, then  $\neg\text{Tr}(\tau)$  is in  $E_\infty$ , because  $E_\infty$  satisfies clause vi).  $\text{Tr}(\tau)$  isn't in  $E_\infty$ , because it isn't put into  $E_\infty$  by clause v). So  $\neg\neg\text{Tr}(\tau)$  isn't put into  $E_\infty$  by clause viii). So  $\neg\neg\text{Tr}(\tau)$  isn't in  $E_\infty$ , and  $\neg\text{Tr}(\tau)$  isn't in  $A_\infty$ .

**vii)** If  $\text{Den}(\tau)$  is a sentence whose negation is in  $E_\infty \sim A_\infty$ , then  $\neg\text{Tr}(\tau)$  is put into  $E_\infty$  by clause vii). Because the negation of  $\text{Den}(\tau)$  isn't in  $A_\infty$ , the negation of the negation of  $\text{Den}(\tau)$  isn't in  $E_\infty$ . So  $\text{Den}(\tau)$  isn't in  $E_\infty$ . So clause v) doesn't put  $\text{Tr}(\tau)$  into  $E_\infty$ . So  $\text{Tr}(\tau)$  isn't in  $E_\infty$ , which means that clause x) doesn't put  $\neg\neg\text{Tr}(\tau)$  into  $E_\infty$ . So  $\neg\text{Tr}(\tau)$  isn't in  $A_\infty$ .

**viii)** If  $P$  and  $Q$  are both in  $E_\infty \sim A_\infty$ , then  $(P \wedge Q)$  is in  $E_\infty$ . Neither  $\neg P$  nor  $\neg Q$  is in  $A_\infty$ , so clause ix) doesn't put  $\neg(P \wedge Q)$  in  $E_\infty$ . So  $(P \wedge Q)$  isn't in  $A_\infty$ .

**ix)** If  $\neg P$  is in  $E_\infty \sim A_\infty$ , then  $\neg(P \wedge Q)$  is in  $E_\infty$ . Because  $\neg P$  isn't in  $A_\infty$ ,  $\neg\neg P$  isn't in  $E_\infty$ . So  $P$  isn't in  $E_\infty$ . So clause viii) doesn't put  $(P \wedge Q)$  into  $E_\infty$ . So clause x) doesn't put  $\neg\neg(P \wedge Q)$  into  $E_\infty$ . So  $\neg(P \wedge Q)$  isn't in  $A_\infty$ . The argument that, if  $\neg Q$  is in  $E_\infty \sim A_\infty$ , then  $\neg(P \wedge Q)$  is in  $E_\infty \sim A_\infty$  is similar.

**x)** If  $P$  is in  $E_\infty \sim A_\infty$ , then  $\neg\neg P$  is in  $E_\infty$ , because  $E_\infty$  satisfies clause x). Because  $\neg P$  isn't in  $E_\infty$ ,  $\neg\neg\neg P$  isn't in  $E_\infty$ . So  $\neg\neg P$  isn't in  $A_\infty$ .

**xi)** Suppose that, for each closed term  $\tau$ ,  $F(\tau)$  is in  $E_\infty \sim A_\infty$ . Then  $(\forall v)F(v)$  is in  $E_\infty$ . There is no closed term  $\tau$  such that  $\neg F(\tau)$  is in  $E_\infty$ . So clause xii) didn't put  $\neg(\forall v)F(v)$  into  $E_\infty$ . So  $\neg(\forall v)F(v)$  isn't in  $E_\infty$ , and  $(\forall v)F(v)$  isn't in  $A_\infty$ .

**xii)** Suppose that, for some closed term  $\tau$ ,  $\neg F(\tau)$  is in  $E_\infty \sim A_\infty$ . Because  $E_\infty$  satisfies clause xii),  $\neg(\forall v)F(v)$  must be in  $E_\infty$ . Since  $\neg F(\tau)$  isn't in  $A_\infty$ ,  $\neg\neg F(\tau)$  must not be in  $E_\infty$ . So  $F(\tau)$  must not be in  $E_\infty$ . This means that clause xi) can't put  $(\forall v)F(v)$  into  $E_\infty$ , and so clause x) can't put  $\neg\neg(\forall v)F(v)$  into  $E_\infty$ . So  $\neg(\forall v)F(v)$  isn't in  $A_\infty$ .

This shows that  $E_\infty$  and  $A_\infty$  are disjoint. We now need to show that  $(E_\infty, A_\infty)$  is a fixed point. We do this by proving, by induction on the complexity of  $P$ , that, for any sentence  $P$ , we have:

$$\begin{aligned} P \in E_\infty &\text{ iff } P \text{ is true in } (\mathbb{N}, (E_\infty, A_\infty)) \\ P \in A_\infty &\text{ iff } P \text{ is false in } (\mathbb{N}, (E_\infty, A_\infty)) \end{aligned}$$

That is, we assume, as inductive hypothesis, that these conditions hold for all sentences simpler than  $P$ , and show that they hold for  $P$ . There are six cases:

**Case 1.**  $P$  has the form  $\tau = \rho$ . Then

$P \in E_\infty$   
 iff  $P$  is put into  $E_\infty$  by clause i)  
 iff  $\text{Den}(\tau) = \text{Den}(\rho)$   
 iff  $\tau = \rho$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$ .

This verifies the first clause. For the second, we have:

$P \in A_\infty$   
 iff  $\neg P \in E_\infty$   
 iff  $\neg P$  is put into  $E_\infty$  by clause ii)  
 iff  $\text{Den}(\tau) \neq \text{Den}(\rho)$   
 iff  $\tau = \rho$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$

**Case 2.**  $P$  has the form  $\tau < \rho$ . Similar.

**Case 3.**  $P$  has the form  $\text{Tr}(\tau)$ . We have

$P \in E_\infty$   
 iff  $P$  is put into  $E_\infty$  by clause ii)  
 iff  $\text{Den}(\tau)$  is in  $E_\infty$   
 iff  $\text{Tr}(\tau)$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$

We get the second clause as follows:

$P \in A_\infty$   
 iff  $\neg P \in E_\infty$   
 iff  $\neg P$  is put into  $E_\infty$  either by clause vi) or clause vii)  
 iff either  $\text{Den}(\tau)$  isn't a sentence or  $\text{Den}(\tau)$  is a sentence whose negation is in  $E_\infty$   
 iff  $\text{Den}(\tau) \in \{\text{nonsentences}\} \cup \{\text{sentences } S: \neg S \in E_\infty\}$   
 iff  $\text{Den}(\tau) \in A_\infty$   
 iff  $\text{Tr}(\tau)$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$

**Case 4.**  $P$  has the form  $(Q \wedge R)$ . We have:

$P \in E_\infty$   
 iff  $P$  is put into  $E_\infty$  by clause viii)  
 iff  $Q$  and  $R$  are both in  $E_\infty$   
 iff  $Q$  and  $R$  are both true in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff  $(Q \wedge R)$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$

iff  $P$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$

For falsity, the argument goes like this:

$P \in A_\infty$   
 iff  $\neg P \in E_\infty$   
 iff  $\neg P$  is put into  $E_\infty$  by clause ix)  
 iff either  $\neg Q$  or  $\neg R$  is in  $E_\infty$   
 iff either  $Q$  or  $R$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff  $(Q \wedge R)$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$

**Case 5.**  $P$  has the form  $\neg Q$ . For truth, the argument is as follows:

$P \in E_\infty$   
 iff  $Q \in A_\infty$   
 iff  $Q$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff  $\neg Q$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$

Here's the story for falsity:

$P \in A_\infty$   
 iff  $\neg P \in E_\infty$   
 iff  $\neg \neg Q \in E_\infty$   
 iff  $\neg \neg Q$  is placed into  $E_\infty$  by clause x)  
 iff  $Q \in E_\infty$   
 iff  $Q$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff  $\neg Q$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$

**Case 6.**  $P$  has the form  $(\forall v)F(v)$ .

$P \in E_\infty$   
 iff clause xi) puts  $P$  into  $E_\infty$   
 iff, for every closed term  $\tau$ ,  $F(\tau)$  is in  $E_\infty$   
 iff, for every closed term  $\tau$ ,  $F(\tau)$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff  $(\forall v)F(v)$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is true in  $(\mathbb{N}, (E_\infty, A_\infty))$

For falsity,

$P \in A_\infty$   
 iff  $\neg P \in E_\infty$   
 iff  $\neg P$  is put into  $E_\infty$  by clause xii)  
 iff, for some closed term  $\tau$ ,  $\neg F(\tau)$  is in  $E_\infty$   
 iff, for some closed term  $\tau$ ,  $F(\tau)$  is in  $A_\infty$   
 iff, for some closed term  $\tau$ ,  $F(\tau)$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$  (by inductive hypothesis)  
 iff,  $(\forall v)F(v)$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$   
 iff  $P$  is false in  $(\mathbb{N}, (E_\infty, A_\infty))$

This shows that  $(E_\infty, A_\infty)$  is a fixed point. It is easy to check that, if  $(E, A)$  is another fixed point,  $E$  satisfies conditions i) - xii). Consequently,  $E_\infty \subseteq E$ , and  $(E_\infty, A_\infty)$  is the *smallest* fixed point.

#### IV. Other Fixed Points

There are many fixed points other than the smallest. To get other fixed points, we use the following:

**Theorem.** Given nonoverlapping sets  $D_0$  and  $B_0$ , with  $D_0 \subseteq \{\text{sentences true in } (\mathbb{N}, (D_0, B_0))\}$  and  $B_0 \subseteq \{\text{nonsentences}\} \cup \{\text{sentences false in } (\mathbb{N}, (D_0, B_0))\}$ . There exists a fixed point  $(D_\infty, B_\infty)$  with  $D_0 \subseteq D_\infty$  and  $B_0 \subseteq B_\infty$ .

The pair  $(D_\infty, B_\infty)$  we construct will, in fact, be the smallest fixed point that extends  $(D_0, B_0)$ . That is, it is a fixed point extending  $(D_0, B_0)$  that is included in every other such fixed point.

The construction of  $(D_\infty, B_\infty)$  will resemble the construction of  $(E_\infty, A_\infty)$ , with the small changes required to ensure that  $D_0$  is in the extension and  $B_0$  in the antiextension.  $D_\infty$  will be the smallest set of sentences that meets the following eleven conditions:

- a) If  $\text{Den}(\tau) = \text{Den}(\rho)$ , then  $\tau = \rho \in D_\infty$ .
  - b) If  $\text{Den}(\tau) \neq \text{Den}(\rho)$ , then  $\neg \tau = \rho \in D_\infty$ .
  - c) If  $\text{Den}(\tau) < \text{Den}(\rho)$ , then  $\tau < \rho \in D_\infty$ .
  - d) If  $\text{Den}(\tau) \not\leq \text{Den}(\rho)$ , then  $\neg \tau < \rho \in D_\infty$ .
  - e) If  $\text{Den}(\tau) \in D_\infty$  or  $\text{Den}(\tau) \in D_0$ , then  $\text{Tr}(\tau) \in D_\infty$ .
  - f) If  $\text{Den}(\tau)$  isn't a sentence or  $\text{Den}(\tau)$  is a sentence whose negation is in  $D_\infty$  or  $\text{Den}(\tau) \in B_0$ , then  $\neg \text{Tr}(\tau) \in D_\infty$ .
  - g) If  $P$  and  $Q$  are in  $D_\infty$ ,  $(P \wedge Q)$  is in  $D_\infty$ .
  - h) If either  $\neg P$  or  $\neg Q$  is in  $D_\infty$ ,  $\neg(P \wedge Q)$  is in  $D_\infty$ .
  - i) If  $P$  is in  $D_\infty$ , so is  $\neg \neg P$ .
  - j) If, for every closed term  $\tau$ ,  $F(\tau)$  is in  $D_\infty$ ,  $(\forall v)F(v)$  is in  $D_\infty$ ;
  - k) If, for some closed term  $\tau$ ,  $\neg F(\tau)$  is in  $D_\infty$ ,  $\neg(\forall v)F(v)$  is in  $D_\infty$ .
- $B_\infty$  is defined to be  $\{\text{nonsentences}\} \cup \{\text{sentences } S: \neg S \in D_\infty\}$ .

The proof that  $(D_\infty, B_\infty)$  is a fixed point is long and tedious, and it doesn't involve any new ideas not already found in the proof that  $(E_\infty, A_\infty)$  is a fixed point, so I won't give it here. Here I'll give only a sketch. The hard part will be to show that  $D_\infty$  and  $B_\infty$  are disjoint. The proof is tricky, just because, if you don't prove things in just the right order, the proof gets bogged down.

**Lemma 1.** All the sentences true in  $(\mathbb{N}, (D_0, B_0))$  are in  $D_\infty$ , and all the sentences false in  $(\mathbb{N}, (D_0, B_0))$  are in  $B_\infty$ .

The proof is by induction on the complexity of sentences. Once we have Lemma 1, we'll know that  $D_0 \subseteq D_\infty$  and  $B_0 \subseteq B_\infty$ .

**Lemma 2.** No member of  $D_\infty$  is false in  $(\mathbb{N}, (D_0, B_0))$ .

We show this by proving that  $D_\infty \sim \{\text{sentences false in } (\mathbb{N}, (D_0, B_0))\}$  satisfies conditions a) through k). We use both Lemma 1 and Lemma 2 to prove:

**Lemma 3.**  $D_\infty$  and  $B_\infty$  are disjoint.

The proof of Lemma 3 consists in showing that  $D_\infty \sim B_\infty$  satisfies conditions a) through k). Once we have Lemma 3, an induction on the complexity of sentences proves that, for any sentence  $P$ , we have:

$P \in D_\infty$  iff  $P$  is true in  $(\mathbb{N}, (D_\infty, B_\infty))$   
 $P \in B_\infty$  iff  $P$  is false in  $(\mathbb{N}, (D_\infty, B_\infty))$

Our smallest fixed point,  $(E_\infty, A_\infty)$ , is what we get in the special case in which  $D_0 = \emptyset$  and  $B_0 = \emptyset$ . There are other fixed points too. For example, take  $T$  to be a *Truth teller*, so that  $T$  with  $T$  semantically equivalent to  $\text{Tr}([\text{ }^+T \text{ }])$ . In the smallest fixed point,  $T$  is neither true nor false. However, in the fixed point obtained by setting  $D_0 = \{T\}$  and  $B_0 = \emptyset$ ,  $T$  will be true. In the fixed point obtained by setting  $D_0 = \emptyset$  and  $B_0 = \{T\}$ ,  $T$  is false.

Let  $T_1, T_2$ , and  $T_3$  be three different Truth tellers. Then setting  $D_0 = \{T_1\}$  and  $B_0 = \{T_2\}$  will give us a fixed point in which  $T_1$  is true,  $T_2$  false, and  $T_3$  unsettled. On the other hand, if  $L$  is a Liar sentence –  $L$  is semantically equivalent to  $\neg\text{Tr}([\text{ }^+L \text{ }])$  –  $L$  is unsettled in every fixed point.

How do we get our three Truth tellers? Let  $V(x,y)$  abbreviate the following formula:

$(\exists z)(z \text{ is the result of substituting } [y] \text{ for free occurrences of "y" in the formula whose Gödel number is } x \wedge \text{Tr}(z))$

Use the Self-Referential Lemma to find a formula  $U(y)$  that is semantically equivalent to  $V([\text{ }^+U(y) \text{ }], y)$ , and let our three Truth tellers be  $U([1])$ ,  $U([2])$ , and  $U([3])$ . This procedure gives us infinitely many Truth tellers. If  $R$  and  $S$  are subsets of our infinite set of Truth tellers with  $R \cap S$

$= \emptyset$ , we can find a fixed point in which all the members of R are true, all the members of S are false, and the rest of our Truth-tellers are all undecided.

## V. The Classical Logic Version

The version of the Kripke construction we've been looking at so far does a good job at untangling complicated self-reference, but it has a hard time with simple generalizations. For example, the apparently harmless statement " $(\forall x)(\text{Tr}(x) \rightarrow \text{Tr}(x))$ ," which says, "Every true sentence is a true sentence," is shoved into the gap between truth and falsity.

Kripke gave an alternative construction that avoids this kind of problem. The new construction (which is based upon an idea of Bas van Fraassen) Treats a sentence P as true in  $(\mathbb{N},(E,A))$  iff P is true in every classical model  $(\mathbb{N},S)$  with  $E \subseteq S \subseteq \mathbb{N} \sim A$ . The proof that there are fixed points continues to work with this alternative way of handling truth-value gaps.. The alternative construction treats all the sentences regarded as valid by classical logic as true, and it treats every classical consequence of true sentences as true. It has the odd feature that a disjunction can be true without either disjunct being true An example is  $(L \vee \neg L)$ , where L is the Liar sentence. So we give up the very simple and natural picture of how the truth and falsity of a complex sentence is grounded in the truth or falsity of its simple components, but in return we get to count as true harmless generalizations like "All true sentences are true sentences," that aren't implicated in the Liar paradox.

## VI. References

Kripke, Saul A. : "Outline of a Theory of Truth." *Journal of Philosophy* 72 (1975): 690-716.  
Reprinted in Martin, pp. 53-81.

McGee, Vann. *Truth, Vagueness, and Paradox: An Essay on the Logic of Truth*. Indianapolis: Hackett, 1991.

Martin, Robert L., ed. *Recent Essays on Truth and the Liar Paradox*. Oxford: Oxford University Press, 1984.

Martin, Robert L., and Peter W. Woodruff. "On Representing 'True-in-L' in L." *Philosophia* 5 (1975): 213-17. Reprinted in Martin, pp. 47-51.