

TRUTH AND INDUCTIVE DEFINABILITY

RICHARD G. HECK

The purpose of this note is to give an exposition of Kripke's theory of truth [4] within the context of the theory of induction on abstract structures. The results are not presented in full generality: That is not needed for the applications made here. For that, see Moschovakis [6], Kechris and Moschovakis, [2], and McGee [5, ch. 5]. Useful papers here, besides Kripke's original, include Feferman [1]. A background in basic set theory, together with some mathematical sophistication, should be enough for the reader to understand this note.

1. ORDINALS

We shall need some basic facts about ordinal numbers.¹ Intuitively, an *ordinal number* answers the question which position in a sequence is occupied by a given element. The sequences in question must be *non-repeating*—no element may occur more than once—and *well-ordered*. A sequence is well-ordered if it is linearly ordered—transitive, irreflexive, and totally ordered—and the sequence must satisfy an analogue of the least number principle: Given any non-empty set S of elements of the sequence, there must be a least one, that is, a member of S that comes before every other member of S . If we think of the sequence as ordered by the relation $<$, then the conditions just mentioned may be formalized as:

$$\begin{aligned} &\forall x \forall y \forall z (x < y \wedge y < z \rightarrow x < z) \\ &\forall x \forall y (x < y \rightarrow \neg y < x) \\ &\forall x \in S \forall y \in S (x < y \vee x = y \vee y < x) \\ &\forall x (x \subseteq S \rightarrow \exists y \in x \forall z \in x (y \leq z)) \end{aligned}$$

The natural numbers 1, 2, 3, and so forth can be construed as ordinals, pronounced 'first', 'second', 'third', and so forth.

There are also infinite ordinals. Consider the following sequence:

1, 2, 3, . . . , Julius Caesar, Brutus, Claudius

This is clearly a non-repeating, well-ordered sequence. Each natural number n occupies the n^{th} position in the sequence. What position does Caesar occupy? This position is called the ω^{th} position: The first infinite ordinal is ω (omega). The next position, occupied by Brutus, is the $(\omega + 1)^{\text{st}}$ position; the next, occupied by Claudius, the $(\omega + 2)^{\text{nd}}$. And so on.²

Consider now the following sequence:

1, 2, 3, . . . , -1, -2, -3, . . . , Caesar, Brutus, Claudius.

¹The theory of ordinal numbers can be developed in set theory, which treats ordinal numbers as sets of a certain kind. See any decent textbook.

²Note, by the way, that ordinal addition is not commutative: $2 + \omega$ is just ω again, since the thing that occupies the ω^{th} place after the 2^{nd} thing is just the ω^{th} thing.

Again, each natural number n occupies the n^{th} position of the sequence; each negative integer $-n$ occupies the $(\omega + n)^{\text{th}}$ position. What position does Caesar occupy here? This is the $(\omega + \omega)^{\text{th}}$ or $(2\omega)^{\text{th}}$ position; Brutus occupies the $(2\omega + 1)^{\text{st}}$; Claudius, the $(2\omega + 2)^{\text{nd}}$. And so on, through $3\omega, 4\omega, \dots, n\omega$, and onward, until we reach $\omega \times \omega$ or ω^2 , which is the position Caesar occupies in this sequence:

$\langle 1, 1 \rangle, \langle 1, 2 \rangle, \dots, \langle 2, 1 \rangle, \langle 2, 2 \rangle, \dots, \langle n, 1 \rangle, \langle n, 2 \rangle, \dots, \text{Caesar}.$

And so on through $\omega^2 + \omega, \omega^2 + 2\omega, \omega^2 + \omega^2 = 2(\omega^2), 3(\omega^2), \omega(\omega^2) = \omega^3, \omega^4, \omega^\omega$, and so on (and on) for a very long time.

One can see here that there are two kinds of ordinals: There are *successor* ordinals and there are *limit* ordinals. We say that α is a successor ordinal if $\alpha = \beta + 1$, for some ordinal β . Otherwise, α is a limit. Examples of limit ordinals are 0 (which is the only finite limit),³ ω , and ω^2 .

All of the ordinals mentioned so far are *countable*, meaning that there are only as many objects in any such sequence as there are natural numbers: Indeed, any such sequence can be constructed simply from natural numbers; for example, the elements $\langle m, k \rangle$ of the last displayed sequence may be replaced by $2^m 3^k$; Caesar may be replaced by 5.

The countable ordinals themselves, under the obvious ordering, form a non-repeating, well-ordered sequence. If we add Caesar to the end of this sequence, we may ask what position he occupies. The answer is ω_1 , the first *uncountable* ordinal. And of course the ordinals do not stop there. But we shall not need to go further. The one fact we do need is that the set of countable ordinals is not itself a countable set. More precisely, what we need to know is just this.

Fact 1. *Let C be the set of all countable ordinals. Suppose that f is a one-one function from C into (not necessarily onto) some set S . Then S is uncountable.*

For a proof, see any decent textbook on set theory.

The other fact we shall need about the ordinals is that, just as one can show that all natural numbers have some property by means of mathematical induction, it is possible to show that all ordinals have some property by means of *transfinite* induction. In the finite case, one typically argues as follows: One shows that 0 is F —the basis case—and one then shows that, if a given number n is F , so is $n + 1$ —the induction step. An equivalent method of proof replaces that sort of induction step with this one: If all numbers less than n are F , then so is n .⁴ One uses this kind of induction, for example, in proving that every number has a unique prime factorization: One does not simply assume that n has a unique prime factorization and then show that $n + 1$ does; rather, one assumes that *every* number less than n has a unique prime factorization and then shows that n does.

Transfinite induction is simply this same method of proof, but applied to all ordinals (or all ordinals less than some given ordinal, e.g., to all countable ordinals). That is, to show that all ordinals are F , one shows that $F0$ and that, if $\forall \beta < \alpha (F\beta) \rightarrow F\alpha$. The induction step often divides into different cases for successor ordinals and limit ordinals.

³Strictly speaking, there is no ordinal 0, since there is no zeroth position of a sequence. But we allow 0 for convenience.

⁴The latter follows from the former: Suppose we want to show that $\forall x Fx$ and we know that $\forall n [\forall x < n (Fx) \rightarrow Fn]$. Then consider $\forall x < n (Fx)$ and argue by induction in the usual way.

Similarly, in the finite case, we can define functions on the natural numbers by recursion: One says what $\phi(0)$ is supposed to be and then, assuming one knows what $\phi(n)$ is, one says what $\phi(n+1)$ is to be. Again, there is a similar but somewhat different method: Assuming one knows what $\phi(k)$ is, for all $k < n$, one says what $\phi(n)$ is to be. Functions on the ordinals can be defined by this method, as well, in which case it is known as *transfinite recursion*. We'll see examples below.

2. MONOTONE OPERATORS

Let A be a set. A^n is the n^{th} Cartesian product of A with itself, that is, the set of sequences $\langle a_1, \dots, a_n \rangle$ with elements from A . It is easy to show that A^n is countable if A is. The *power set* of A , $\wp(A)$ is the set of all subsets of A . If A is countably infinite, then $\wp(A)$ is *uncountably infinite*, by a simple diagonalization due to Cantor.

Let A be a set of sets. We say that an operator $\Phi : A \rightarrow A$ is *monotonic* if, whenever $S \subseteq T$, $\Phi(S) \subseteq \Phi(T)$. Our attention will be focused upon operators $\Phi : \wp(A^n) \rightarrow \wp(A^n)$ that map sets of n -tuples whose elements are taken from A to sets of n -tuples whose elements are taken from A .

The key theorem is that monotonic operators always have fixed points, in the following sense:

Definition. S is a *fixed point* of Φ if $\Phi(S) = S$.

The idea behind the proof is simple. We start with \emptyset and start applying Φ , forming the sequence $\emptyset \subseteq \Phi(\emptyset) \subseteq \Phi(\Phi(\emptyset)) \subseteq \dots$ and continue in this way. The various inclusions hold because Φ is monotone. Obviously, the first one must hold; but then the second one holds; and so forth. If we don't hit a fixed point before we run out of finite ordinals, we take the union of everything we've had so far and continue applying Φ , taking unions at any limits we reach. Eventually, we will reach a fixed point, for reasons we shall see below.

We now give a real proof of this result.

Definition. If Φ is monotonic, then the function I_Φ^x is defined by transfinite recursion thus:

$$\begin{aligned} I_\Phi^0 &= \emptyset \\ I_\Phi^\alpha &= \Phi\left(\bigcup_{\xi < \alpha} I_\Phi^\xi\right), \text{ for } \alpha \neq 0 \end{aligned}$$

Theorem 2 (Fixed-point Theorem). *Let $\Phi : A \rightarrow A$ be monotonic. Then Φ has a fixed point. Moreover, if A is countable, then there is a countable ordinal α such that $\Phi(I_\Phi^\alpha) = I_\Phi^\alpha$, and I_Φ^α is the minimal fixed point of Φ , in the sense that it is a subset of every other fixed point of Φ .⁵*

Lemma 3. *If Φ is monotonic, then $\forall \beta \forall \alpha (\beta < \alpha \rightarrow I_\Phi^\beta \subseteq I_\Phi^\alpha)$. In particular, $I_\Phi^\alpha \subseteq I_\Phi^{\alpha+1}$.*

Proof. By transfinite induction on α . This holds vacuously for $\alpha = 0$. So suppose $\alpha \neq 0$. We want to see that $\forall \beta < \alpha (I_\Phi^\beta \subseteq I_\Phi^\alpha)$. So let $\beta < \alpha$. Then obviously $\bigcup_{\xi < \beta} I_\Phi^\xi \subseteq \bigcup_{\xi < \alpha} I_\Phi^\xi$, since the latter union includes all the sets the former one

⁵More generally, the α in question will be less than $|A|^+$, which is the least cardinal greater than the cardinality of A .

does. So, since Φ is monotone, we have $I_\Phi^\beta = \Phi(\bigcup_{\xi < \beta} I_\Phi^\xi) \subseteq \Phi(\bigcup_{\xi < \alpha} I_\Phi^\xi) = I_\Phi^\alpha$, as claimed. \square

Corollary 4. $I_\Phi^{\alpha+1} = \Phi(I_\Phi^\alpha)$

Proof. By definition, $I_\Phi^{\alpha+1} = \Phi(\bigcup_{\xi < \alpha+1} I_\Phi^\xi)$. I claim that $\bigcup_{\xi < \alpha+1} I_\Phi^\xi = I_\Phi^\alpha$. Now, $\bigcup_{\xi < \alpha+1} I_\Phi^\xi = I_\Phi^\alpha \cup \bigcup_{\xi < \alpha} I_\Phi^\xi$. But if $\xi < \alpha$, then $I_\Phi^\xi \subseteq I_\Phi^\alpha$, so $\bigcup_{\xi < \alpha} I_\Phi^\xi \subseteq I_\Phi^\alpha$. Hence, $I_\Phi^\alpha \cup \bigcup_{\xi < \alpha} I_\Phi^\xi = I_\Phi^\alpha$. \square

Proof of Theorem 2. Suppose that for no countable α do we have $\Phi(I_\Phi^\alpha) = I_\Phi^\alpha$. Then for every α , $I_\Phi^\alpha \subsetneq \Phi(I_\Phi^\alpha) = I_\Phi^{\alpha+1}$. Hence, for each countable α , we can find $\varphi(\alpha) \in I_\Phi^{\alpha+1} \setminus I_\Phi^\alpha$.⁶ Now suppose $\alpha \neq \beta$; either $\alpha < \beta$ or conversely; without loss of generality, suppose $\alpha < \beta$. Then $\varphi(\beta) \in I_\Phi^{\beta+1} \setminus I_\Phi^\beta$, whence $\varphi(\beta) \notin I_\Phi^\alpha$, whereas $\varphi(\alpha) \in I_\Phi^{\alpha+1}$. And since $\alpha < \beta$, $\alpha + 1 \leq \beta$, whence $I_\Phi^{\alpha+1} \subseteq I_\Phi^\beta$, and so $\varphi(\alpha) \in I_\Phi^\beta$. But that means that φ is one-to-one, whence it is a function from all the countable ordinals into A , which is impossible.

So there is a countable ordinal α such that $\Phi(I_\Phi^\alpha) = I_\Phi^\alpha$. Henceforth let α be the least such ordinal. So I_Φ^α is a fixed point of Φ .

Let F be any fixed point of Φ . We show by transfinite induction that $I_\Phi^\beta \subseteq F$, for all β , from which it follows that $I_\Phi^\alpha \subseteq F$. Obviously, if $\beta = 0$, then $I_\Phi^0 = \emptyset \subseteq F$. So suppose $I_\Phi^\xi \subseteq F$, for all $\xi < \beta$. Then $\bigcup_{\xi < \beta} I_\Phi^\xi \subseteq F$,⁷ so $I_\Phi^\beta = \Phi(\bigcup_{\xi < \beta} I_\Phi^\xi) \subseteq \Phi(F) = F$, by monotonicity. \square

Notation. We write: I_Φ , for the minimal fixed point of Φ .

3. INDUCTIVE DEFINITIONS

Definition. A predicate letter R is said to *occur positively* in a formula ϕ if, and only if, it is in the smallest class C of formulae containing:

- (1) All formulae in which R does not occur;
- (2) $R(t_1, \dots, t_n)$, for any terms t_i ;
- (3) the conjunctions, disjunctions, and existential or universal quantifications of any formulae in C .

What this means, in essence, is that R does not occur inside negation or in the *antecedent* of a conditional (since $p \rightarrow q$ can be defined as: $\neg p \vee q$). We call a formula that is logically equivalent to a formula in which R occurs positively *R-positive*.

Let \mathcal{L} be a language; \mathcal{T} , a theory in \mathcal{L} . We will be interested in languages \mathcal{L}^* that result from adding to \mathcal{L} some new relation symbols of various arities. *Par abuse de langue*, we shall also think of \mathcal{T} as a theory in \mathcal{L}^* , with the same axioms. We shall write these new relation constants as lower case Greek letters, σ , τ , and the like, and indicate their arity with superscripts as necessary.

Our main theorem is this:

Theorem 5 (Inductive Definitions). *Let \mathcal{T} be a theory in \mathcal{L} ; let \mathcal{L}^* be \mathcal{L} plus a new n -place relation symbol σ^n ; and let $\phi(x_1, \dots, x_n, \sigma^n)$ be a formula of \mathcal{L}^* in*

⁶The Axiom of Choice is needed for this step.

⁷If $\forall x \in S(x \subseteq T)$, then $\cup S \subseteq T$.

which σ^n occurs positively and in which only the variables shown are free. Then the theory \mathcal{T}^* , which is the result of adding the new axiom

$$(\sigma) \quad \forall x_1 \dots \forall x_n [\sigma^n(x_1, \dots, x_n) \equiv \phi(x_1, \dots, x_n, \sigma^n)]$$

to \mathcal{T} , is a conservative extension of \mathcal{T} . A fortiori, \mathcal{T}^* is consistent if \mathcal{T} is consistent.

Definition. We say that σ^n is *inductively defined* by the formula $\phi(x_1, \dots, x_n, \sigma^n)$.

Before we give the proof, a couple examples. Let \mathcal{L} be the language of arithmetic, and let \mathcal{T} be PA. First, we can inductively define 2^x . Let $\psi(x, y, \text{exp})$ be:

$$(x = 0 \wedge y = 1) \vee (x = n \wedge \exists z (\text{exp}(n, z) \wedge y = 2 \times z))$$

Then exp occurs positively in ψ so, by the theorem, the result of adding the axiom:

$$\text{exp}(x, y) \equiv (x = 0 \wedge y = 1) \vee (x = n \wedge \exists z (\text{exp}(n, z) \wedge y = 2 \times z))$$

to PA is consistent. This, as said, gives an inductive characterization of 2^x .

Second, we can give an inductive definition of the notion of a *term* of the language. Assume that we have already introduced the basics of Gödel numbering into the system. Let $\phi(x, \tau)$ be the formula:

$$\begin{aligned} x &= \ulcorner 0 \urcorner \vee \\ &\exists y (x = \ulcorner Sy \urcorner \wedge \tau(y)) \vee \\ &\exists z \exists y (x = \ulcorner y + z \urcorner \wedge \tau(y) \wedge \tau(z)) \vee \\ &\exists z \exists y (x = \ulcorner y \times z \urcorner \wedge \tau(y) \wedge \tau(z)) \end{aligned}$$

Then τ occurs positively in ϕ and so, by the theorem, the result of adding the axiom:

$$\begin{aligned} (\text{Def } \tau) \quad \forall x [\tau(x) \equiv & \quad x = \ulcorner 0 \urcorner \vee \\ & \exists y (x = \ulcorner Sy \urcorner \wedge \tau(y)) \vee \\ & \exists z \exists y (x = \ulcorner y + z \urcorner \wedge \tau(y) \wedge \tau(z)) \vee \\ & \exists z \exists y (x = \ulcorner y \times z \urcorner \wedge \tau(y) \wedge \tau(z))] \end{aligned}$$

to PA is a conservative extension of PA. This axiom characterizes the set of terms inductively.

Before we prove Theorem 5, we shall first need to establish some properties of positive formulae. Let an interpretation \mathcal{M} of \mathcal{L} be given, with domain \mathcal{D} . Let S be a subset of \mathcal{D}^n —i.e., a possible extension for σ^n —and let \mathcal{M}_S be the expansion of \mathcal{M} to an interpretation of \mathcal{L}^* that we get by assigning S to σ^n . Now suppose we begin with an interpretation \mathcal{M}_S of \mathcal{L}^* . In this interpretation, some sequences will satisfy $\phi(x_1, \dots, x_n, \sigma^n)$ and some will not. The set of sequences that do satisfy it is *another* possible extension for σ^n . So there is a natural operator $\Sigma_\phi^{\mathcal{M}}$ on \mathcal{D}^n that we can define as follows:

$$\Sigma_\phi^{\mathcal{M}}(S) = \{ \langle y_1, \dots, y_n \rangle : \langle y_1, \dots, y_n \rangle \text{ satisfies } \phi(x_1, \dots, x_n, \sigma^n) \text{ in } \mathcal{M}_S \}$$

Thus, $\Sigma_\phi^{\mathcal{M}}(S)$ is the extension of ϕ in \mathcal{M}_S .

Now for the key fact:

Lemma 6. *If $\phi(x_1, \dots, x_n, \sigma^n)$ is σ^n -positive, then $\Sigma_\phi^{\mathcal{M}}$ is monotonic and so has a minimal fixed point $I_\phi^{\mathcal{M}}$.*

Proof. The monotonicity of $\Sigma_\phi^{\mathcal{M}}$ will follow immediately from the next result. The existence of a fixed point then follows from Theorem 2. \square

Proposition 7. *Let $\phi(x_1, \dots, x_n, \sigma^n)$ be a σ^n -positive formula in \mathcal{L}^* , and let \mathcal{M} be an interpretation of \mathcal{L} with domain \mathcal{D} . Suppose that $S \subseteq T \subseteq \mathcal{D}^n$. Then if a sequence s satisfies $\phi(x_1, \dots, x_n, \sigma^n)$ in \mathcal{M}_S , s also satisfies $\phi(x_1, \dots, x_n, \sigma^n)$ in \mathcal{M}_T .*

Proof. The proof is by induction on the complexity of σ^n -positive formulas.

The basis case concerns formulas in which σ^n does not occur and atomic formulae of the form $\sigma^n(x_{k_1}, \dots, x_{k_n})$, for some k_1, \dots, k_n . In the former case, the change from \mathcal{M}_S to \mathcal{M}_T is irrelevant. In the latter, if s satisfies $\sigma^n(x_{k_1}, \dots, x_{k_n})$ in \mathcal{M}_S , then $\langle s(k_1), \dots, s(k_n) \rangle \in S \subseteq T$, whence s satisfies $\sigma^n(x_{k_1}, \dots, x_{k_n})$ in \mathcal{M}_T .

So suppose the proposition holds for A and B ; we want to show it holds for their conjunction and disjunction. But if s satisfies $A \wedge B$ in \mathcal{M}_S , then it satisfies both A and B in \mathcal{M}_S ; hence, by the induction hypothesis, it satisfies both A and B in \mathcal{M}_T ; so it satisfies $A \wedge B$ in \mathcal{M}_T . The argument for disjunction is similar.

So suppose the proposition holds for $A(x_i)$, where x_i occurs free in $A(x_i)$, as may additional variables. We want to show that the proposition holds for $\exists x_i A(x_i)$ and $\forall x_i A(x_i)$. If s satisfies $\exists x_i A(x_i)$ in \mathcal{M}_S , then there is a sequence t such that $\forall j \neq i (t(j) = s(j))$ and t satisfies $A(x_i)$ in \mathcal{M}_S . By the induction hypothesis, t satisfies $A(x_i)$ in \mathcal{M}_T ; but then s satisfies $\exists x_i A(x_i)$ in \mathcal{M}_T . The argument for the universal quantifier is similar. \square

Our goal, recall, is to show that the result \mathcal{T}^* of adding the new axiom (σ) to \mathcal{T} is a conservative extension of \mathcal{T} . To do so, we show that any model \mathcal{M} of \mathcal{T} can be expanded to a model of \mathcal{T}^* , that is, that we can expand \mathcal{M} in such a way as to make (σ) true.⁸ The model in question will be the one that assigns $I_\phi^{\mathcal{M}}$ to σ^n . That this does the trick is the content of the following.

Lemma 8. *Let \mathcal{M} be an interpretation of \mathcal{L} ; let $\phi(x_1, \dots, x_n, \sigma^n)$ be a σ^n -positive formula in \mathcal{L}^* ; and let I be any fixed point of $\Sigma_\phi^{\mathcal{M}}$. Then*

$$(\sigma) \quad \forall x_1 \dots \forall x_n [\sigma^n(x_1, \dots, x_n) \equiv \phi(x_1, \dots, x_n, \sigma^n)]$$

is true in \mathcal{M}_I .

Proof. The formula mentioned is true in \mathcal{M}_I iff

$$\{ \langle y_1, \dots, y_n \rangle : \langle y_1, \dots, y_n \rangle \text{ satisfies } \sigma^n(x_1, \dots, x_n) \text{ in } \mathcal{M}_I \} =$$

$$\{ \langle y_1, \dots, y_n \rangle : \langle y_1, \dots, y_n \rangle \text{ satisfies } \phi(x_1, \dots, x_n, \sigma^n) \text{ in } \mathcal{M}_I \}$$

But

$$\{ \langle y_1, \dots, y_n \rangle : \langle y_1, \dots, y_n \rangle \text{ satisfies } \sigma^n(x_1, \dots, x_n) \text{ in } \mathcal{M}_I \} = I,$$

by the definition of \mathcal{M}_I ;

$$I = \Sigma_\phi^{\mathcal{M}}(I),$$

since I is a fixed point of $\Sigma_\phi^{\mathcal{M}}$; and

$$\Sigma_\phi^{\mathcal{M}}(I) = \{ \langle y_1, \dots, y_n \rangle : \langle y_1, \dots, y_n \rangle \text{ satisfies } \phi(x_1, \dots, x_n, \sigma^n) \text{ in } \mathcal{M}_I \},$$

by the definition of $\Sigma_\phi^{\mathcal{M}}$. \square

We can now complete the proof of theorem 5.

⁸Thus, \mathcal{T}^* is, as it is said, *semantically* conservative over \mathcal{T} .

Proof of Theorem 5. Let \mathcal{M} be any model of \mathcal{T} . We shall show that $\mathcal{M}_{I_\phi^\mathcal{M}}$, which is an expansion of \mathcal{M} , is a model of \mathcal{T}^* . All axioms of \mathcal{T} continue to be true in $\mathcal{M}_{I_\phi^\mathcal{M}}$, since we have not changed the interpretation of \mathcal{L} . And by lemma 8, (σ) is also true in $\mathcal{M}_{I_\phi^\mathcal{M}}$. So we are done. \square

A simple generalization of the above arguments yields the following.

Theorem 9 (Simultaneous Inductive Definitions). *Let \mathcal{T} be a theory in \mathcal{L} ; let $\psi(y_1, \dots, y_m, \sigma^m, \tau^n)$ and $\phi(x_1, \dots, x_n, \sigma^m, \tau^n)$ be formulae of \mathcal{L}^* in which both σ^m and τ^n occur positively. Then the theory \mathcal{T}^* which is the result of adding the new axioms*

$$(\sigma) \quad \forall y_1 \dots \forall y_m [\sigma^m(y_1, \dots, y_m) \equiv \psi(y_1, \dots, y_m, \sigma^m, \tau^n)]$$

$$(\tau) \quad \forall x_1 \dots \forall x_n [\tau^n(x_1, \dots, x_n) \equiv \phi(x_1, \dots, x_n, \sigma^m, \tau^n)]$$

to \mathcal{T} , is a conservative extension of \mathcal{T} . A fortiori, \mathcal{T}^* is consistent if \mathcal{T} is consistent.

Proof. Given an interpretation \mathcal{M} of \mathcal{L} , let $\mathcal{M}_{S,T}$ be the expansion of \mathcal{M} to an interpretation of \mathcal{L}^* in which S is the extension of σ and T is the extension of τ . Define an operator $\Sigma_\psi^\mathcal{M}(S, T)$ as follows: $\Sigma_\psi^\mathcal{M}(S, T)$ is the extension of σ^m in $\mathcal{M}_{S,T}$; that is:

$$\Sigma_\psi^\mathcal{M} = \{ \langle y_1, \dots, y_m \rangle : y_1, \dots, y_m \text{ satisfies } \psi(y_1, \dots, y_m, \sigma^m, \tau^n) \text{ in } \mathcal{M}_{S,T} \}$$

Define $\Sigma_\phi^\mathcal{M}(S, T)$ similarly.

By simultaneous induction, define:

$$\Psi_0 = \Sigma_\psi^\mathcal{M}(\emptyset, \emptyset)$$

$$\Phi_0 = \Sigma_\phi^\mathcal{M}(\emptyset, \emptyset)$$

$$\Psi_\alpha = \Sigma_\psi^\mathcal{M} \left(\bigcup_{\xi < \alpha} \Psi_\xi, \bigcup_{\xi < \alpha} \Phi_\xi \right)$$

$$\Phi_\alpha = \Sigma_\phi^\mathcal{M} \left(\bigcup_{\xi < \alpha} \Psi_\xi, \bigcup_{\xi < \alpha} \Phi_\xi \right)$$

We can then show, just as in the proof of Theorem 5, that $\Sigma_\psi^\mathcal{M}$ and $\Sigma_\phi^\mathcal{M}$ have *simultaneous* fixed points $I_\psi^\mathcal{M}$ and $I_\phi^\mathcal{M}$: That is, $I_\psi^\mathcal{M} = \Sigma_\psi^\mathcal{M}(I_\psi^\mathcal{M}, I_\phi^\mathcal{M})$ and $I_\phi^\mathcal{M} = \Sigma_\phi^\mathcal{M}(I_\psi^\mathcal{M}, I_\phi^\mathcal{M})$. And we can then show that $\mathcal{M}_{I_\psi^\mathcal{M}, I_\phi^\mathcal{M}}$ is a model of \mathcal{T}^* , again as before. \square

Obviously, this can be extended to any finite number of formulas.

Now let $\mathcal{L}_\mathcal{M}$ be an *interpreted* language, consisting of a language \mathcal{L} and an interpretation \mathcal{M} for \mathcal{L} with domain \mathcal{D} ; the interpretation is thus fixed. Suppose that $\phi(x_1, \dots, x_n, \sigma^n)$ is a σ^n -positive formula in $\mathcal{L}_\mathcal{M}^*$. It follows from Theorem 2 that the operator $\Sigma_\phi^\mathcal{M}$ has a fixed point $I_\phi^\mathcal{M}$, which we may just call ' I_ϕ '. And it then follows from Lemma 8 that

$$(\sigma) \quad \forall x_1 \dots \forall x_n [\sigma^n(x_1, \dots, x_n) \equiv \phi(x_1, \dots, x_n, \sigma^n)]$$

is true when σ^n is assigned I_ϕ as its extension.

Definition. If a set S is I_ϕ for some appropriate formula ϕ of $\mathcal{L}_\mathcal{M}^*$, then S is said to be a *fixed point over $\mathcal{L}_\mathcal{M}$* .

Definition. Let $R \subseteq D^m$. We say that R is *inductive* over $\mathcal{L}_{\mathcal{M}}$ if there is a fixed point I_ϕ over $\mathcal{L}_{\mathcal{M}}$ and there are objects a_1, \dots, a_{n-m} such that:

$$\langle x_1, \dots, x_m \rangle \in R \equiv \langle a_1, \dots, a_{n-m}, x_1, \dots, x_m \rangle \in I_\phi.$$

In effect, a relation R is inductive just in case it is a ‘planar section’ of a fixed point or, again, if it is ‘parametrically definable’ in terms of a fixed point.

Definition. A relation whose complement is inductive is called *co-inductive*; a relation that is both inductive and co-inductive is said to be *hyper-elementary*. If S is actually defined by a formula $\phi(x_1, \dots, x_n)$ of $\mathcal{L}_{\mathcal{M}}$ itself, we say that it is *elementary*.

Here are a few basic facts about inductive relations, which we record without proof.

Proposition 10.

- (1) *Every relation elementary in $\mathcal{L}_{\mathcal{M}}$ is inductive over $\mathcal{L}_{\mathcal{M}}$.*
- (2) *Relations inductive over $\mathcal{L}_{\mathcal{M}}$ are closed under union and intersection.*
- (3) *Relations hyper-elementary over $\mathcal{L}_{\mathcal{M}}$ are closed under complementation.*
- (4) *The union and intersection of two inductive relations are inductive.*

Lemma 11 (Simultaneous Induction Lemma). *Let $\psi(y_1, \dots, y_m, \sigma^m, \tau^n)$ and $\phi(x_1, \dots, x_n, \sigma^m, \tau^n)$ be formulae of $\mathcal{L}_{\mathcal{M}}^*$ in which both σ^m and τ^n occur positively. Let $\Sigma_{\mathcal{M}, \psi}(S, T)$ and $\Sigma_{\mathcal{M}, \phi}(S, T)$ be the operators naturally associated with these formulae, as above. Then their minimal simultaneous fixed points $I_{\mathcal{M}, \psi}$ and $I_{\mathcal{M}, \phi}$ are both inductive over $\mathcal{L}_{\mathcal{M}}$.*

Theorem 12 (Transitivity Theorem). *Let $Q \subseteq D^n$ be a relation that is hyper-elementary over an interpreted language $\mathcal{L}_{\mathcal{M}}$ with domain \mathcal{D} . Let $\mathcal{L}_{\mathcal{M}}^+$ be the result of adding a new relation symbol χ^n to $\mathcal{L}_{\mathcal{M}}$, interpreted as having the extension Q , and let $R \subseteq D^m$ be a relation that is inductive over $\mathcal{L}_{\mathcal{M}}^+$. Then R is inductive over $\mathcal{L}_{\mathcal{M}}$.*

In short: Expanding an interpreted language by adding a symbol for a hyper-elementary relation does not allow one to define any new inductive relations.

4. TRUTH AND INDUCTIVE DEFINITIONS: TARSKI

In this section, we shall show how the consistency of Tarski’s theory of truth for a given language $\mathcal{L}_{\mathcal{M}}$ can be proven using the theory of inductive definitions.

As above, let $\mathcal{L}_{\mathcal{M}}$ be an interpreted language consisting of a language \mathcal{L} and an interpretation \mathcal{M} of \mathcal{L} with domain \mathcal{D} . Let $\mathcal{T}_{\mathcal{L}_{\mathcal{M}}}$ be its ‘diagram’, that is, the set of all true sentences of $\mathcal{L}_{\mathcal{M}}$. We assume that $\mathcal{T}_{\mathcal{L}_{\mathcal{M}}}$ interprets basic syntax: In particular, we can code finite sequences of elements of \mathcal{D} by means of elements of \mathcal{D} ; more precisely, we are assuming that notions such as atomic formula of \mathcal{L} , formula of \mathcal{L} , and the like, are $\mathcal{L}_{\mathcal{M}}$ -elementary, and that $\mathcal{T}_{\mathcal{L}_{\mathcal{M}}}$ proves their basic properties. We use ‘Seq(a, σ)’ to mean: σ is a sequence that assigns values to all variables free in the formula (with Gödel number) a ; ‘ \mathcal{L} -AtForm(a)’, that a is an atomic formula of \mathcal{L} ; ‘ \mathcal{L} -Form(a)’, that a is a formula of \mathcal{L} .

We will also assume, mostly for convenience, that satisfaction for atomic formulae, $\text{SatAt}_{\mathcal{M}}(a, \sigma)$, is $\mathcal{L}_{\mathcal{M}}$ -elementary, that is, that satisfaction for atomic formulae is definable by means of a formula of $\mathcal{L}_{\mathcal{M}}$.⁹

⁹In light of the transitivity theorem, it would actually be sufficient for our purposes to assume that $\text{SatAt}_{\mathcal{M}}(a, \sigma)$ is hyper-elementary.

Let me insert a note on the meanings of expressions such as ‘ $a = \neg b$ ’, used below. Of course, our formal language does not actually contain such expressions. By assumption, though, it does contain a formula $\text{neg}(a, b)$ meaning: a is (the code of a formula that is) a negation of (the formula with code) b . That is also what ‘ $a = \neg b$ ’ means, though in using this notation we are assuming that only formulae have negations; that every formula has a negation; and that no formula has more than one negation. We are thus assuming that $\mathcal{T}_{\mathcal{L}_{\mathcal{M}}}$ proves such basic syntactic facts.

Theorem 13 (Tarski). *Satisfaction for $\mathcal{L}_{\mathcal{M}}$ is not elementary, but it is hyper-elementary.*

Note what this means: That satisfaction is not elementary means that it is not defined by any formula of $\mathcal{L}_{\mathcal{M}}$. This follows from Tarski’s theorem, which will not be proven here.¹⁰ That it is hyper-elementary means that it is both inductive and co-inductive, which of course implies that both the set of truths and the set of falsehoods are inductive.

Note: *We are assuming we understand what satisfaction is.* Take it to be defined as Tarski defined it. Our goal is to prove facts *about* satisfaction, in particular, to prove that it can be defined by means of an inductive definition of the sort discussed in the previous section.

Proof of Theorem 13. Select arbitrary terms—which we shall abbreviate ‘ \perp ’ and ‘ \top ’—of $\mathcal{L}_{\mathcal{M}}$ such that $\top \neq \perp$ is true. Define $\text{Val}_{\mathcal{M}}(a, \sigma, t)$ as:

$$\mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{[t = \top \wedge \text{Sat}_{\mathcal{M}}(a, \sigma)] \vee [t = \perp \wedge \text{Sat}_{\mathcal{M}}(\neg a, \sigma)]\}.$$

(Thus, \top represents true and \perp represents false.) We shall prove that $\{ \langle a, \sigma, t \rangle : \text{Val}_{\mathcal{M}}(a, \sigma, t) \}$ is a fixed-point over $\mathcal{L}_{\mathcal{M}}$. Since $\text{Sat}_{\mathcal{M}}(a, \sigma) \equiv \text{Val}_{\mathcal{M}}(a, \sigma, \top)$, it follows that $\{ \langle a, \sigma \rangle : \text{SatAt}_{\mathcal{M}}(a, \sigma) \}$ is inductive. Since

$$\begin{aligned} \text{Sat}_{\mathcal{M}}(a, \sigma) &\equiv \mathcal{L}\text{-Form}(a) \wedge \neg \text{Val}_{\mathcal{M}}(a, \sigma, \perp) \\ &\equiv \neg[\neg \mathcal{L}\text{-Form}(a) \vee \text{Val}_{\mathcal{M}}(a, \sigma, \perp)] \end{aligned}$$

it will be co-inductive and so hyper-elementary.

Define $\phi(a, \sigma, t, \chi^3)$ as follows:

$$\begin{aligned} \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ & \\ & [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \top] \vee \\ & [\mathcal{L}\text{-AtForm}(a) \wedge \neg \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \perp] \vee \\ & \exists b[a = \neg b \wedge \chi(b, \sigma, \top) \wedge t = \perp] \vee \\ & \exists b[a = \neg b \wedge \chi(b, \sigma, \perp) \wedge t = \top] \vee \\ & \exists b \exists c[a = b \vee c \wedge (\chi(b, \sigma, \top) \vee \chi(c, \sigma, \top)) \wedge t = \top] \vee \\ & \exists b \exists c[a = b \vee c \wedge (\chi(b, \sigma, \perp) \wedge \chi(c, \sigma, \perp)) \wedge t = \perp] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge \chi(b, \tau, \top)) \wedge t = \top] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \forall \tau (\forall j \neq i (\sigma_j = \tau_j) \rightarrow \chi(b, \tau, \perp)) \wedge t = \perp] \} \end{aligned}$$

By inspection, ϕ is χ -positive. So it has a fixed point I_{χ} .

Note. Think of the three-place predicate $\chi(a, \sigma, t)$ as meaning: t is the truth-value of a under the assignment σ of values to variables. One might have thought we should simply use a two-place formula $v(a, \sigma)$, meaning, intuitively: a is true

¹⁰See my “Formal Background for Theories of Truth” for a simple exposition.

under the assignment σ of values to variables (or just: σ satisfies a). We would then instead define $\phi(a, \sigma, v^2)$ as:

$$\begin{aligned} \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ \\ & [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\ & [\mathcal{L}\text{-AtForm}(a) \wedge \neg \text{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\ & \exists b[a = \perp b \wedge \neg v(b, \sigma)] \vee \\ & \exists b \exists c[a = b \vee c \wedge (v(b, \sigma) \vee v(c, \sigma))] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge v(b, \tau))] \} \end{aligned}$$

But then ϕ would not be v -positive—look at the clause for negation—and the theory of inductive definitions would not apply to it. In that respect, it is essential here to make use of truth-*values* and not just of a truth-*predicate*.

End of note.

We intend to show that $\text{Val}_{\mathcal{M}}(a, \sigma, t) \equiv \langle a, \sigma, t \rangle \in I_{\chi}$. To do this, we show by induction on α that

$$(i) \quad \langle a, \sigma, t \rangle \in I_{\chi}^{\alpha} \rightarrow \text{Val}_{\mathcal{M}}(a, \sigma, t)$$

Since $I_{\chi} = \bigcup_{\alpha} I_{\chi}^{\alpha}$, then if $\langle a, \sigma, t \rangle \in I_{\chi}$, also $\langle a, \sigma, t \rangle \in I_{\chi}^{\alpha}$, for some α . So we may conclude that $\langle a, \sigma, t \rangle \in I_{\chi} \rightarrow \text{Val}_{\mathcal{M}}(a, \sigma, t)$. We will then show that

$$(ii) \quad \text{Val}_{\mathcal{M}}(a, \sigma, t) \rightarrow \exists \alpha (\langle a, \sigma, t \rangle \in I_{\chi}^{\alpha})$$

and conclude that $\text{Val}_{\mathcal{M}}(a, \sigma, t) \rightarrow \langle a, \sigma, t \rangle \in I_{\chi}$.

We first prove (i). Obviously, this holds vacuously for $\alpha = 0$, since $I_{\chi}^0 = \emptyset$. So suppose it holds for $\xi < \alpha$; we want to show it also holds for α . Now, $I_{\chi}^{\alpha} = \Sigma_{\chi}^{\mathcal{M}}(\bigcup_{\xi < \alpha} I_{\chi}^{\xi})$, which—abusing notation slightly—is the set of $\langle a, \sigma, t \rangle$ such that $\phi(a, \sigma, t, \chi^3)$ when χ^3 has $\bigcup_{\xi < \alpha} I_{\chi}^{\xi}$ for its extension. So, if $\langle a, \sigma, t \rangle \in I_{\chi}^{\alpha}$, then $\mathcal{L}\text{-Form}(a)$ and $\text{Seq}(a, \sigma)$ and one of the disjuncts of ϕ must hold. Say it's the fourth, for illustration. Then $a = \perp b$ and $\chi(b, \sigma, \top)$ and $t = \perp$, and we must show that $\text{Val}_{\mathcal{M}}(a, \sigma, \perp)$. Since $\chi(b, \sigma, \top)$, we have $\langle b, \sigma, \top \rangle \in \bigcup_{\xi < \alpha} I_{\chi}^{\xi}$ and so, for some $\xi < \alpha$, $\langle b, \sigma, \top \rangle \in I_{\chi}^{\xi}$. So, by the induction hypothesis, $\text{Val}_{\mathcal{M}}(b, \sigma, \top)$, that is, $\text{Sat}_{\mathcal{M}}(b, \sigma)$. So, by the definition of satisfaction, $\neg \text{Sat}_{\mathcal{M}}(a, \sigma)$, whence $\text{Val}_{\mathcal{M}}(a, \sigma, \perp)$. The other cases are similar.

To prove (ii), we argue by induction on the complexity of the formula a that:

$$(*) \quad \forall \sigma \forall t [\text{Val}_{\mathcal{M}}(a, \sigma, t) \rightarrow \exists \alpha (\langle a, \sigma, t \rangle \in I_{\chi}^{\alpha})]$$

Certainly this holds for atomic formulae. If $t = \top$, then $\text{Sat}_{\mathcal{M}}(a, \sigma)$ because $\text{SatAt}_{\mathcal{M}}(a, \sigma)$, in which case $\langle a, \sigma, \top \rangle \in I_{\chi}^1$; if $t = \perp$, then $\neg \text{Sat}_{\mathcal{M}}(a, \sigma)$ because $\neg \text{SatAt}_{\mathcal{M}}(a, \sigma)$, in which case $\langle a, \sigma, \perp \rangle \in I_{\chi}^1$.

So suppose (*) holds for formulae of lesser complexity than a . By the induction hypothesis, for each such formula b and for each t and τ , if $\text{Val}_{\mathcal{M}}(b, \tau, t)$, then there is a ξ such that $\langle b, \tau, t \rangle \in I_{\chi}^{\xi}$. Let α be the supremum of all such ξ .¹¹ Then since X is monotone, $\langle b, \tau, t \rangle \in I_{\chi}^{\alpha}$. That is, for all b of complexity less than that of a and for all t and τ , if $\text{Val}_{\mathcal{M}}(b, \tau, t)$, then $\langle b, \tau, t \rangle \in I_{\chi}^{\alpha}$.

So we now need to check the various possibilities for the kind of formula a might be and for the value t might take. We'll do two cases. First, suppose $a = \perp b$ and $\text{Val}_{\mathcal{M}}(a, \sigma, t)$. If $t = \top$, then $\text{Sat}_{\mathcal{M}}(a, \sigma)$; so $\neg \text{Sat}_{\mathcal{M}}(b, \sigma)$, and so $\text{Val}_{\mathcal{M}}(b, \sigma, \perp)$.

¹¹The supremum—that is, the least α greater than or equal to all such ξ —exists because the $\langle b, \tau, t \rangle$ form a set.

By the results of the previous paragraph, $\langle b, \sigma, \perp \rangle \in I_\chi^\alpha$. But then $\phi(a, \sigma, \top, \chi^3)$ will hold when χ^3 takes I_χ^α for its extension, since the third disjunct in the definition of ϕ will hold. So $\langle a, \sigma, \top \rangle \in I_\chi^{\alpha+1}$. And if $t = \perp$, just swap \perp and \top in the foregoing.

Suppose that $a = \exists v_1 b$ and that $\text{Val}_\mathcal{M}(a, \sigma, t)$. Let $t = \top$. So $\text{Sat}_\mathcal{M}(a, \sigma)$ and hence, for some τ such that $\forall j \neq i(\sigma_j = \tau_j)$, $\text{Sat}_\mathcal{M}(b, \sigma)$ and so $\text{Val}_\mathcal{M}(b, \tau, \top)$. So $\langle b, \tau, \top \rangle \in I_\chi^\alpha$. But then we have $\phi(a, \sigma, \top, \chi^3)$ when χ^3 takes I_χ^α for its extension, since the seventh of the disjuncts above will hold. It follows that $\langle a, \sigma, \top \rangle \in I_\chi^{\alpha+1}$. If $t = \perp$, then $\neg \text{Sat}_\mathcal{M}(a, \sigma)$, so for all τ such that $\forall j \neq i(\sigma_j = \tau_j)$, $\neg \text{Sat}_\mathcal{M}(b, \tau)$ and so $\text{Val}_\mathcal{M}(b, \tau, \perp)$. But then $\langle b, \tau, \perp \rangle \in I_\chi^\alpha$ and so $\phi(a, \sigma, \perp, \chi^3)$ when χ^3 takes I_χ^α for its extension, since the eighth disjunct holds, so $\langle a, \sigma, \perp \rangle \in I_\chi^{\alpha+1}$. \square

Corollary 14. *For every formula a and sequence σ , either $\langle a, \sigma, \top \rangle \in I_\chi^\mathcal{M}$ or $\langle a, \sigma, \perp \rangle \in I_\chi^\mathcal{M}$, but not both.*

Proof. This follows immediately from the preceding, since either $\text{Sat}_\mathcal{M}(a, \sigma)$ or $\neg \text{Sat}_\mathcal{M}(a, \sigma)$. But it can also be proven by induction on the complexity of a . \square

Theorem 15 (Tarski). *Let \mathcal{T} be any theory in \mathcal{L} meeting these two assumptions:*

- (1) \mathcal{T} interprets basic syntax;
- (2) *There is a formula $\text{SatAt}_\mathcal{M}(a, \sigma)$ of \mathcal{L} such that, for every atomic formula ϕ , \mathcal{T} proves $\text{SatAt}_\mathcal{M}(\ulcorner \phi(x_1, \dots, x_n \urcorner, \sigma) \equiv \phi(\sigma_1, \dots, \sigma_n)$.*

Let \mathcal{T}_{sat} be the theory in $\mathcal{L} + \text{'sat'}$ that is the result of adding to \mathcal{T} the single axiom:

$$\begin{aligned} \text{sat}(a, \sigma) \equiv & \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ \\ & [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_\mathcal{M}(a, \sigma)] \vee \\ & \exists b[a = \neg b \wedge \neg \text{sat}(b, \sigma)] \vee \\ & \exists b \exists c[a = b \vee c \wedge (\text{sat}(b, \sigma) \vee \text{sat}(c, \sigma))] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \exists t \tau (\forall j \neq i(\sigma_j = \tau_j) \wedge \text{sat}(b, \tau))] \}, \end{aligned}$$

Then \mathcal{T}_{sat} is a conservative extension of \mathcal{T} and so is consistent if \mathcal{T} is.

Proof. Let \mathcal{M} be any model of \mathcal{T} .

We show first that condition (2) implies that $\text{SatAt}_\mathcal{M}(a, \sigma)$ defines satisfaction in \mathcal{M} for atomic formulae of \mathcal{L} . Let $\phi(x_1, \dots, x_n)$ be an atomic formula of \mathcal{L} . Since \mathcal{M} is a model of \mathcal{T} , $\mathcal{M} \models \forall \sigma [\text{SatAt}_\mathcal{M}(\ulcorner \phi(x_1, \dots, x_n \urcorner, \sigma) \equiv \phi(\sigma_1, \dots, \sigma_n)]$. Now let s be a finite sequence of elements of the domain of \mathcal{M} and let σ be a closed term coding s . Either $\mathcal{M}, s \models \phi(x_1, \dots, x_n)$ or not. If so, then $\mathcal{M} \models \text{SatAt}_\mathcal{M}(\ulcorner \phi(x_1, \dots, x_n \urcorner, \sigma)$; if not, then $\mathcal{M} \models \neg \text{SatAt}_\mathcal{M}(\ulcorner \phi(x_1, \dots, x_n \urcorner, \sigma)$. Hence, $\mathcal{M}, s \models \phi(x_1, \dots, x_n)$ iff $\mathcal{M} \models \text{SatAt}_\mathcal{M}(\ulcorner \phi(x_1, \dots, x_n \urcorner, \sigma)$, and $\text{SatAt}_\mathcal{M}(a, \sigma)$ defines \mathcal{M} -satisfaction for atomic formulae of \mathcal{L} .

Let $I_{\mathcal{M}, \phi}$ be the fixed point over \mathcal{M} of the formula $\phi(a, \sigma, t, \chi^3)$ defined in the proof of Theorem 13. Expand \mathcal{M} to a model \mathcal{M}_χ of $\mathcal{L} + \chi$ by assigning $I_{\mathcal{M}, \phi}$ as the extension of ' χ '; rewrite ' χ ' as ' Val ' for clarity.

Now consider the result of adding to \mathcal{T} the following three axioms:

$$\begin{aligned} & \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \rightarrow \text{Val}_\mathcal{M}(a, \sigma, \top) \vee \text{Val}_\mathcal{M}(a, \sigma, 0) \\ & \neg[\text{Val}_\mathcal{M}(a, \sigma, \top) \wedge \text{Val}_\mathcal{M}(a, \sigma, 0)] \end{aligned}$$

$$\begin{aligned}
\text{Val}_{\mathcal{M}}(a, \sigma, t) \equiv & \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ \\
& [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \top] \vee \\
& [\mathcal{L}\text{-AtForm}(a) \wedge \neg \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \perp] \vee \\
& \exists b[a = \neg b \wedge \text{Val}_{\mathcal{M}}(b, \sigma, 1) \wedge t = \perp] \vee \\
& \exists b[a = \neg b \wedge \text{Val}_{\mathcal{M}}(b, \sigma, 0) \wedge t = \top] \vee \\
& \exists b \exists c[a = b \vee c \wedge (\text{Val}_{\mathcal{M}}(b, \sigma, \top) \vee \text{Val}_{\mathcal{M}}(c, \sigma, \top)) \wedge t = \top] \vee \\
& \exists b \exists c[a = b \vee c \wedge (\text{Val}_{\mathcal{M}}(b, \sigma, \perp) \wedge \text{Val}_{\mathcal{M}}(c, \sigma, \perp)) \wedge t = \perp] \vee \\
& \exists b \exists i[a = \exists v_i b \wedge \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge \text{Val}_{\mathcal{M}}(b, \tau, \top)) \wedge t = \top] \vee \\
& \exists b \exists i[a = \exists v_i b \wedge \forall \tau (\forall j \neq i (\sigma_j = \tau_j) \rightarrow \text{Val}_{\mathcal{M}}(b, \tau, \perp)) \wedge t = \perp] \}
\end{aligned}$$

Call this theory \mathcal{T}_{Val} . It is a conservative extension of \mathcal{T} , by theorem 5 and corollary 14.

Simple logical manipulations applied to the third axiom give us:

$$\begin{aligned}
\text{Val}_{\mathcal{M}}(a, \sigma, 1) \equiv & \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ \\
& [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \top] \vee \\
& \exists b[a = \neg b \wedge \text{Val}_{\mathcal{M}}(b, \sigma, \perp) \wedge t = \top] \vee \\
& \exists b \exists c[a = b \vee c \wedge (\text{Val}_{\mathcal{M}}(b, \sigma, \top) \vee \text{Val}_{\mathcal{M}}(c, \sigma, \top)) \wedge t = \top] \vee \\
& \exists b \exists i[a = \exists v_i b \wedge \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge \text{Val}_{\mathcal{M}}(b, \tau, \top)) \wedge t = \top] \}
\end{aligned}$$

Now define $\text{sat}(a, \sigma)$ as: $\text{Val}_{\mathcal{M}}(a, \sigma, 1)$. The first two axioms then yield:

$$\mathcal{L}\text{-Form}(x) \wedge \text{Seq}(a, \sigma) \rightarrow [\neg \text{sat}(x, \sigma) \equiv \neg \text{Val}_{\mathcal{M}}(x, \sigma, 0)]$$

Standard logical manipulations then show that the axiom mentioned in the theorem follows logically from the preceding and the definition of ‘ sat ’. So \mathcal{T}_{sat} is definitionally equivalent to \mathcal{T}_{Val} , which is itself a conservative extension of \mathcal{T} . So \mathcal{T}_{sat} is a conservative extension of \mathcal{T} . \square

Corollary 16. *Let \mathcal{T} be any theory meeting the assumptions stated in Theorem 15. Then the theory \mathcal{T}_{Tarski} in $\mathcal{L} + \text{‘sat’}$, the result of adding to \mathcal{T} the axioms:*

- (1) $\text{sat}(a, \sigma) \rightarrow \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma)$
- (2) $\mathcal{L}\text{-AtForm}(a) \rightarrow [\text{sat}(a, \sigma) \equiv \text{SatAt}_{\mathcal{M}}(a, \sigma)]$
- (3) $\text{sat}(\ulcorner \neg A \urcorner, \sigma) \equiv \neg \text{sat}(A, \sigma)$
- (4) $\text{sat}(\ulcorner A \vee B \urcorner, \sigma) \equiv \text{sat}(A, \sigma) \vee \text{sat}(B, \sigma)$
- (5) $\text{sat}(\ulcorner \exists v_i A \urcorner, \sigma) \equiv \exists \tau [\forall j \neq i (\sigma_j = \tau_j) \wedge \text{sat}(A, \tau)]$

is a conservative extension of \mathcal{T} and so is consistent if \mathcal{T} is.

Proof. Each of (1)–(5) is provable in \mathcal{T}_{sat} , so \mathcal{T}_{Tarski} is a sub-theory of \mathcal{T}_{sat} . \square

Note that the fact that \mathcal{T}_{sat} and \mathcal{T}_{Tarski} are conservative extensions of \mathcal{T} depends essentially upon the fact that ‘ sat ’ is not in the language of \mathcal{T} and so does not occur in any axiom of the theory \mathcal{T} . In particular, it may not appear in formulae instantiating \mathcal{T} ’s axiom schemes, should it have any.¹² If we do allow ‘ sat ’ to be used in that way, then \mathcal{T}_{sat} and \mathcal{T}_{Tarski} are not conservative extensions of \mathcal{T} but will often prove the consistency of \mathcal{T} .

¹²That is, we are not thinking of \mathcal{T} as a *schematic theory* in Feferman’s sense.

5. TRUTH AND INDUCTIVE DEFINITIONS: KRIPKE

In this section, we shall show how the consistency of Kripke's theory of truth, for an interpreted language $\mathcal{L}_{\mathcal{M}}$, can be proven using the theory of inductive definitions. In fact, the changes needed to get Kripke's theory rather than Tarski's are surprisingly minimal.

Let $\mathcal{L}_{\mathcal{M}}$ be an interpreted language sufficient for syntax; let $\mathcal{L}_{\mathcal{M}}^T$ be the result of adding a single one-place predicate T intended to mean: is true, to $\mathcal{L}_{\mathcal{M}}$. We want to produce a truth-definition for $\mathcal{L}_{\mathcal{M}}^T$. We know, basically, what to do with the vocabulary of $\mathcal{L}_{\mathcal{M}}$. But what should the clause for the predicate T be like? Well, we want ' $T(\mathbf{n})$ ' to be true just in case \mathbf{n} is the Gödel number of a true sentence. We can think of this as part of an inductive specification of the extension of the predicate T : Once a sentence \mathbf{n} has been deemed true, at the next stage of the induction, the sentence ' $T(\mathbf{n})$ ' will be deemed true. On the other hand, we do *not* want ' $T(\mathbf{n})$ ' to be deemed *false* simply because the sentence \mathbf{n} has not yet been deemed true. We want ' $T(\mathbf{n})$ ' to be deemed false only when the sentence \mathbf{n} has itself been deemed false.

We can incorporate this idea into the inductive definition given earlier, though we shall do so for the general case of satisfaction, not just for truth. As before, let $\text{SatAt}_{\mathcal{M}}(a, \sigma)$ be a formula of \mathcal{L} that defines satisfaction in \mathcal{M} for atomic formulae of \mathcal{L} ; $\mathcal{L}\text{-AtForm}(a)$, the notion of an atomic formula of \mathcal{L} ; and $\mathcal{L}\text{-Form}(a)$, that of a formula. We assume further that there are formulae $\text{Term}(t)$ and $\text{den}_{\mathcal{M}}(t, \sigma, n)$ in \mathcal{L} that define: a is a formula of \mathcal{L} ; t is a term of \mathcal{L} , and: t denotes n under σ , respectively.¹³ Further still, we assume that the theories we shall be considering below prove the basic facts about all of these notions. In particular, we assume that in all such theories we can prove:

$$\text{SatAt}_{\mathcal{M}}(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) \equiv \phi(\sigma_1, \dots, \sigma_n)$$

for atomic formula ϕ and:

$$\text{den}_{\mathcal{M}}(\ulcorner t(x_1, \dots, x_n) \urcorner, \sigma, n) \equiv t(\sigma_1, \dots, \sigma_n) = n$$

for each term t . Such a theory might be called 'adequate'.

For convenience, we shall let $\text{den}_{\mathcal{M}}(t, \sigma)$ abbreviate: $\ulcorner \text{den}_{\mathcal{M}}(t, \sigma, n) \urcorner$, the description being proper since every term provably has a unique denotation (in adequate theories).¹⁴

As before, we now expand the language $\mathcal{L}_{\mathcal{M}}$ by adding a three-place predicate χ . Now let $\phi(a, \sigma, t, \chi^3)$ be the formula in Figure 5.1 on page 14. Note that this is the same definition of ϕ as used above, except for the last two disjuncts, which govern sentences of the form: $\chi(a, \sigma, t)$. Here, $a = \ulcorner \chi(t, u, 1) \urcorner$ means: a is (the code of) the formula: $\chi(\mathbf{t}, \mathbf{u}, 1)$, where \mathbf{t} and \mathbf{u} are the numerals denoting the numbers t and u , respectively.

What does this definition have to do with Kripke's? Kripke takes the truth-predicate to have both an *extension* and an *anti-extension*. As in the above discussion of Tarski, we are using a three-place predicate $\chi(a, \sigma, t)$, which we can think of as meaning: t is the truth-value of a under the assignment σ of values to variables. We may thus regard $\chi(a, \sigma, \top)$ as meaning: $\langle a, \sigma \rangle$ is in the extension of χ ; and $\chi(a, \sigma, \perp)$ as meaning: $\langle a, \sigma \rangle$ is in the anti-extension of χ . The fact that the logic is three-valued is reflected in the fact that not every formula will

¹³One could take $\text{den}_{\mathcal{M}}(t, \sigma, n)$ to abbreviate: $\text{SatAt}_{\mathcal{M}}(\ulcorner t = n \urcorner, \sigma)$.

¹⁴That is, each formula $\dots \text{den}_{\mathcal{M}}(t, \sigma) \dots$ abbreviates the corresponding formula: $\exists n[\text{den}_{\mathcal{M}}(t, \sigma, n) \wedge \dots n \dots]$.

$$\begin{aligned}
& \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \{ \\
& \quad [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \top] \vee \\
& \quad [\mathcal{L}\text{-AtForm}(a) \wedge \neg \text{SatAt}_{\mathcal{M}}(a, \sigma) \wedge t = \perp] \vee \\
& \quad \exists b[a = \neg b \wedge \chi(b, \sigma, \top) \wedge t = \perp] \vee \\
& \quad \exists b[a = \neg b \wedge \chi(b, \sigma, 0) \wedge t = \top] \vee \\
& \quad \exists b \exists c[a = b \vee c \wedge (\chi(b, \sigma, \top) \vee \chi(c, \sigma, \top)) \wedge t = \top] \vee \\
& \quad \exists b \exists c[a = b \vee c \wedge (\chi(b, \sigma, 0) \wedge \chi(c, \sigma, \perp)) \wedge t = \perp] \vee \\
& \quad \exists b \exists i[a = \exists v_i b \wedge \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge \chi(b, \tau, \top)) \wedge t = \top] \vee \\
& \quad \exists b \exists i[a = \exists v_i b \wedge \forall \tau (\forall j \neq i (\sigma_j = \tau_j) \rightarrow \chi(b, \tau, \perp)) \wedge t = \perp] \vee \\
& \quad \exists t \exists u [\text{Term}(t) \wedge \text{Term}(u) \wedge a = \ulcorner \chi(t, u, \top) \urcorner \wedge \\
& \quad \quad \chi(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma), 1) \wedge t = \top] \vee \\
& \quad \exists t \exists u [\text{Term}(t) \wedge \text{Term}(u) \wedge a = \ulcorner \chi(t, u, \perp) \urcorner \wedge \\
& \quad \quad \chi(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma), \perp) \wedge t = \perp] \}
\end{aligned}$$

FIGURE 5.1. Kripke's Definition of Satisfaction

turn up in the 'extension' or 'anti-extension' in the minimal fixed point: We shall not, that is to say, be able to prove an analogue of Corollary 14. The fact that we are working with the *Strong* Kleene scheme is reflected in the character of the clauses for disjunction. To get the Weak Kleene scheme, define:

$$TV(x, \sigma) \equiv \exists y(\chi(x, \sigma, y)),$$

meaning: x has a truth-value under σ , and replace the first clause for disjunction with this one:

$$\exists b \exists c[a = b \vee c \wedge TV(b, \sigma) \wedge TV(c, \sigma) \wedge (\chi(b, \sigma, 1) \vee \chi(c, \sigma, 1)) \wedge t = 1]$$

A similar change needs to be made to the clause for the existential quantifier.

One can make the definition look more like Kripke's by adding *two* two-place predicates τ and φ , intended to express truth and falsity or, more generally, satisfaction and non-satisfaction. Consider, then, the two formulae in Figure 5.2 on page 15, $TSat(a, \sigma, \tau^2, \varphi^2)$ and $FSat(a, \sigma, \tau^2, \varphi^2)$. Here τ represents satisfaction; φ , non-satisfaction. It follows from Lemma 11 that there are simultaneous fixed points I_{TSat} and I_{FSat} of $TSat$ and $FSat$: These will be the extension and anti-extension of the truth-predicate in Kripke's construction.

The above definition of the formula $\phi(a, \sigma, t, \chi^3)$ is essentially what one gets by applying the procedure described in the proof of Theorem 9 to $TSat(a, \sigma, \tau^2, \varphi^2)$ and $FSat(a, \sigma, \tau^2, \varphi^2)$. One could debate whether the pair of formulae or the single formula yields a more natural theory. Formally, however, it is a bit easier to work with the pair of definitions, so we shall use them. But it is important to appreciate the mathematical relationship between Kripke's treatment and Tarski's: Though their definitions of truth *look* very different, they are, both mathematically and conceptually, very similar.

There is another difference between this definition of satisfaction and Kripke's. In Kripke's construction, one begins by putting all true sentences of the base language $\mathcal{L}_{\mathcal{M}}$ into the extension. To get a definition more like Kripke's, one can replace the clauses of the $TSat$ and $FSat$ concerning atomic formulae of \mathcal{L} with

$$\begin{aligned}
& \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \\
& \quad \{ [\mathcal{L}\text{-AtForm}(a) \wedge \text{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\
& \quad \exists b [a = \neg b \wedge \varphi(b, \sigma)] \vee \\
& \quad \exists b \exists c [a = b \vee c \wedge (\tau(b, \sigma) \vee \tau(c, \sigma))] \vee \\
& \quad \exists b \exists i [a = \exists v_i b \wedge \exists v (\forall j \neq i (\sigma_j = v_j) \wedge \tau(b, v))] \vee \\
& \quad \exists t \exists u [\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner \tau(t, u) \urcorner \wedge \\
& \quad \quad \tau(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))] \vee \\
& \quad \exists t \exists u [\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner \varphi(t, u) \urcorner \wedge \\
& \quad \quad \varphi(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))] \}
\end{aligned}$$

$$\begin{aligned}
& \mathcal{L}\text{-Form}(a) \wedge \text{Seq}(a, \sigma) \wedge \\
& \quad \{ [\mathcal{L}\text{-AtForm}(a) \wedge \neg \text{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\
& \quad \exists b [a = \neg b \wedge \tau(b, \sigma)] \vee \\
& \quad \exists b \exists c [a = b \vee c \wedge \varphi(b, \sigma) \wedge \varphi(c, \sigma)] \vee \\
& \quad \exists b \exists i [a = \exists v_i b \wedge \forall v (\forall j \neq i (\sigma_j = v_j) \rightarrow \varphi(b, v))] \vee \\
& \quad \exists t \exists u [\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner \varphi(t, u) \urcorner \wedge \\
& \quad \quad \tau(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))] \vee \\
& \quad \exists t \exists u [\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner \tau(t, u) \urcorner \wedge \\
& \quad \quad \varphi(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))] \}
\end{aligned}$$

FIGURE 5.2. Kripke's Definition, Again

these clauses:

$$\begin{aligned}
& \mathcal{L}\text{-Form}(a) \wedge \text{sat}_{\mathcal{L}_{\mathcal{M}}}(a, \sigma) \\
& \mathcal{L}\text{-Form}(a) \wedge \neg \text{sat}_{\mathcal{L}_{\mathcal{M}}}(a, \sigma)
\end{aligned}$$

where $\text{sat}_{\mathcal{L}_{\mathcal{M}}}$ is a satisfaction predicate for $\mathcal{L}_{\mathcal{M}}$, defined as above. Since satisfaction for $\mathcal{L}_{\mathcal{M}}$ is hyper elementary over $\mathcal{L}_{\mathcal{M}}$, there is a sense in which this difference makes no difference: By Theorem 12, any relation inductive over $\mathcal{L}_{\mathcal{M}} + \text{sat}_{\mathcal{M}}$ is also inductive over $\mathcal{L}_{\mathcal{M}}$.

On the other hand, this difference does affect the 'levels' at which sentences get truth-values. On Kripke's treatment, the level of every (code of a) sentence of the base language is 1. This will not be so if we use the inductive definition given above: Only *atomic* sentences of $\mathcal{L}_{\mathcal{M}}$ will get truth-values at stage 1. Perhaps Kripke's way of doing things is preferable if one is interested only in the semantics of sentences containing the truth-predicate. But the definition above is more natural if we are not. This definition does not treat 'true' in any way peculiarly; 'true' is treated pretty much the same way any other atomic predicate is treated (a fact that is partly hidden by our use of ' $\text{SatAt}_{\mathcal{M}}(a, \sigma)$ ' rather than individual clauses for the various atomic expressions of $\mathcal{L}_{\mathcal{M}}$). That is to say: We are not thinking of *adding* a truth-predicate to a language for which truth has already been defined; we are defining satisfaction for a language that already contains a truth-predicate. We shall therefore stick with our present definition.

At the moment, then, what we have is an inductive definition of a pair of notions, satisfaction and non-satisfaction. We have thus now reached a point in our exposition comparable the one we had reached in our discussion of Tarski after the proof of Theorem 13. We now want to consider what *formal theories* of truth we can prove to be consistent on the basis of these definitions. We begin by simply converting the formulae $TSat$ and $FSat$ into a pair of axioms.

Theorem 17. *Let \mathcal{T} be a theory in a language \mathcal{L} . Let $\mathcal{T}_{TSat,FSat}$ be a theory in $\mathcal{L} + TSat + FSat$ containing the axioms of \mathcal{T} plus the three additional axioms:*

$$\neg \exists x [TSat(x, \sigma) \wedge FSat(x, \sigma)]$$

$$\begin{aligned} TSat(a, \sigma) \equiv & \mathcal{L}\text{-Form}(a) \wedge \mathbf{Seq}(a, \sigma) \wedge \\ & \{[\mathcal{L}\text{-AtForm}(a) \wedge \mathbf{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\ & \exists b[a = \neg b \wedge FSat(b, \sigma)] \vee \\ & \exists b \exists c[a = b \vee c \wedge (TSat(b, \sigma) \vee TSat(c, \sigma))] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \exists v(\forall j \neq i(\sigma_j = v_j) \wedge TSat(b, v))] \vee \\ & \exists t \exists u[\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner TSat(t, u) \urcorner \wedge \\ & \quad TSat(\mathbf{den}_{\mathcal{M}}(t, \sigma), \mathbf{den}_{\mathcal{M}}(u, \sigma))] \vee \\ & \exists t \exists u[\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a = \ulcorner FSat(t, u) \urcorner \wedge \\ & \quad FSat(\mathbf{den}_{\mathcal{M}}(t, \sigma), \mathbf{den}_{\mathcal{M}}(u, \sigma))]\} \end{aligned}$$

$$\begin{aligned} FSat(a, \sigma) \equiv & Fmla(a) \wedge \mathbf{Seq}(a, \sigma) \wedge \\ & \{[\mathcal{L}\text{-AtForm}(a) \wedge \neg \mathbf{SatAt}_{\mathcal{M}}(a, \sigma)] \vee \\ & \exists b[a = \neg b \wedge TSat(b, \sigma)] \vee \\ & \exists b \exists c[a = b \vee c \wedge FSat(b, \sigma) \wedge FSat(c, \sigma)] \vee \\ & \exists b \exists i[a = \exists v_i b \wedge \forall v(\forall j \neq i(\sigma_j = v_j) \rightarrow FSat(b, v))] \vee \\ & \exists t \exists u[\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a \wedge \ulcorner FSat(t, u) \urcorner \wedge \\ & \quad TSat(\mathbf{den}_{\mathcal{M}}(t, \sigma), \mathbf{den}_{\mathcal{M}}(u, \sigma))] \vee \\ & \exists t \exists u[\mathbf{Term}(t) \wedge \mathbf{Term}(u) \wedge a \wedge \ulcorner TSat(t, u) \urcorner \wedge \\ & \quad FSat(\mathbf{den}_{\mathcal{M}}(t, \sigma), \mathbf{den}_{\mathcal{M}}(u, \sigma))]\} \end{aligned}$$

Then $\mathcal{T}_{TSat,FSat}$ is a conservative extension of \mathcal{T} and so is consistent if \mathcal{T} is. Moreover, the formulae:

$$\begin{aligned} TSat(\neg t, \sigma) & \equiv FSat(t, \sigma) \\ FSat(\neg t, \sigma) & \equiv TSat(t, \sigma) \end{aligned}$$

are theorems of $\mathcal{T}_{TSat,FSat}$.

Proof. Let \mathcal{M} be a model of \mathcal{T} . As mentioned, Lemma 11 guarantees that there are simultaneous fixed points $I_{TSat}^{\mathcal{M}}$ and $I_{FSat}^{\mathcal{M}}$ of the formulae $TSat$ and $FSat$ in Figure 5.2 on page 15. Expand \mathcal{M} to an interpretation \mathcal{M}^* of $\mathcal{L} + TSat + FSat$ by taking the extension of $TSat$ to be $I_{TSat}^{\mathcal{M}}$ and that of $FSat$ to be $I_{FSat}^{\mathcal{M}}$. Then \mathcal{M}^* is a model of $\mathcal{T}_{TSat,FSat}$.

That $\neg \exists x [TSat(x, \sigma) \wedge FSat(x, \sigma)]$ holds in \mathcal{M}^* is proven by induction on the stages of the inductive definition, that is, on α in I_{TSat}^{α} and I_{FSat}^{α} . The details are left as an exercise.

The final claim follows from the clauses for negation. □

Remark. We do not have $\neg T\text{sat}(x, \sigma) \rightarrow F\text{sat}(x, \sigma)$ or, equivalently, $T\text{sat}(x, \sigma) \vee F\text{sat}(x, \sigma)$.

We can simplify $\mathcal{T}_{T\text{sat}, F\text{sat}}$ substantially. We begin by stating an obvious corollary of the preceding theorem.

Corollary 18 (Kripke). *Let \mathcal{T} be a theory in a language \mathcal{L} . Let \mathcal{T}_{TF} be the theory in $\mathcal{L}+T\text{sat}+F\text{sat}$ containing the axioms of \mathcal{T} plus the universal closures of:*

- (1) $T\text{sat}(x, \sigma) \rightarrow \mathcal{L}\text{-Form}(x) \wedge \text{Seq}(a, \sigma)$
 $F\text{sat}(x, \sigma) \rightarrow \mathcal{L}\text{-Form}(x) \wedge \text{Seq}(a, \sigma)$
- (2) $T\text{sat}(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) \equiv \phi(\sigma_1, \dots, \sigma_n)$ and
 $F\text{sat}(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) \equiv \neg \phi(\sigma_1, \dots, \sigma_n)$, for ϕ an atomic formula of \mathcal{L}
- (3) $T\text{sat}(t \vee u, \sigma) \equiv T\text{sat}(t, \sigma) \vee T\text{sat}(u, \sigma)$
 $F\text{sat}(t \vee u, \sigma) \equiv F\text{sat}(t, \sigma) \wedge F\text{sat}(u, \sigma)$
- (4) $T\text{sat}(\exists v_i t, \sigma) \equiv \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge T\text{sat}(t, \tau))$
 $F\text{sat}(\exists v_i t, \sigma) \equiv \forall \tau (\forall j \neq i (\sigma_j = \tau_j) \rightarrow F\text{sat}(t, \tau))$
- (5) $T\text{sat}(\ulcorner T\text{sat}(t, u) \urcorner, \sigma) \equiv T\text{sat}(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$
 $T\text{sat}(\ulcorner F\text{sat}(t, u) \urcorner, \sigma) \equiv F\text{sat}(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$
 $F\text{sat}(\ulcorner T\text{sat}(t, u) \urcorner, \sigma) \equiv F\text{sat}(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$
 $F\text{sat}(\ulcorner F\text{sat}(t, u) \urcorner, \sigma) \equiv T\text{sat}(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$
- (6) $\neg \exists x [T\text{sat}(x, \sigma) \wedge F\text{sat}(x, \sigma)]$
- (7) $T\text{sat}(\neg t, \sigma) \equiv F\text{sat}(t, \sigma)$
 $F\text{sat}(\neg t, \sigma) \equiv T\text{sat}(t, \sigma)$

\mathcal{T}_{TF} has the same theorems as $\mathcal{T}_{T\text{sat}, F\text{sat}}$.

Note what the various clauses in (5) mean: The first means that a formula of the form $T\text{sat}(t, u)$ is satisfied by a sequence σ if, and only if, the formula whose Gödel number is the denotation of t with respect to σ is satisfied by the sequence whose Gödel number is the denotation of u with respect to σ . Thus, $T\text{sat}(t, u)$ ‘says’ that the formula (with Gödel number) t is ‘truly’ satisfied by the sequence (with Gödel number) u .

Proof. Axioms (1)–(5) of \mathcal{T}_{TF} are easily provable in $\mathcal{T}_{T\text{sat}, F\text{sat}}$: They all follow immediately from appropriate clauses of its second and third axioms. Axioms (6) and (7) are already axioms of $\mathcal{T}_{T\text{sat}, F\text{sat}}$. So \mathcal{T}_{TF} is a sub-theory of $\mathcal{T}_{T\text{sat}, F\text{sat}}$.

Conversely, the three axioms of $\mathcal{T}_{T\text{sat}, F\text{sat}}$ are provable in \mathcal{T}_{TF} . We prove just the second.

If we suppose that $T\text{sat}(a, \sigma)$, then by axiom (1) of \mathcal{T}_{TF} , a is a formula and σ is a sequence. So a is either an atomic formula of \mathcal{L} , a negation, a disjunction, an existential quantification, or is of the form $T\text{sat}(b, \tau)$ or $F\text{sat}(b, \tau)$. The parts of axioms (2), (3), (4), (5), and (7) concerned with $T\text{sat}$ then entail the corresponding disjuncts on the right-hand side of axiom (1). Conversely, if we suppose the right-hand side of the second axiom of $\mathcal{T}_{T\text{sat}, F\text{sat}}$, then, again, a is a formula and σ is a sequence; a must be of one of the specified forms; and the appropriate disjunct together with the corresponding axiom of $\mathcal{T}_{T\text{sat}, F\text{sat}}$ will entail $T\text{sat}(a, \sigma)$. \square

Even a casual examination of the axioms of \mathcal{T}_{TF} reveals a massive and dissatisfying repetitiveness that is due to the fact that we are using distinct satisfaction and non-satisfaction axioms. But we can easily eliminate this redundancy: In virtue of (7), we can replace $F\text{sat}(t, \sigma)$ throughout by $T\text{sat}(\neg t, \sigma)$. This gives us the following.

Theorem 19 (Kripke, Feferman). *Let \mathcal{T} be a theory in a language \mathcal{L} . Let \mathcal{T}_{KF} be the theory in $\mathcal{L}+Sat$ containing the axioms of \mathcal{T} plus the universal closures of:*

- (1) $Sat(x, \sigma) \rightarrow \mathcal{L}\text{-Form}(x) \wedge \text{Seq}(a, \sigma)$
- (2) $Sat(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) \equiv \phi(\sigma_1, \dots, \sigma_n)$, for ϕ an atomic or negated atomic formula of \mathcal{L}
- (3) a. $Sat(t \ulcorner u \urcorner, \sigma) \equiv Sat(t, \sigma) \vee Sat(u, \sigma)$
b. $Sat(\ulcorner t \ulcorner u \urcorner \urcorner, \sigma) \equiv Sat(\ulcorner t \urcorner, \sigma) \wedge Sat(\ulcorner u \urcorner, \sigma)$
- (4) a. $Sat(\ulcorner \exists v_i t \urcorner, \sigma) \equiv \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge Sat(t, \tau))$
b. $Sat(\ulcorner \neg \exists v_i t \urcorner, \sigma) \equiv \forall \tau (\forall j \neq i (\sigma_j = \tau_j) \rightarrow Sat(\ulcorner t \urcorner, \tau))$
- (5) a. $Sat(\ulcorner Sat(t, u) \urcorner, \sigma) \equiv Sat(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$
b. $Sat(\ulcorner \neg Sat(t, u) \urcorner, \sigma) \equiv Sat(\ulcorner \neg \text{den}_{\mathcal{M}}(t, \sigma) \urcorner, \text{den}_{\mathcal{M}}(u, \sigma))$
- (6) $\neg [Sat(x, \sigma) \wedge Sat(\ulcorner \neg x \urcorner, \sigma)]$
- (7) $Sat(\ulcorner \neg \ulcorner t \urcorner \urcorner, \sigma) \equiv Sat(t, \sigma)$

\mathcal{T}_{KF} is a definitional extension of \mathcal{T}_{TF} . It is therefore a conservative extension of \mathcal{T} and so is consistent if \mathcal{T} is.

The axioms just mentioned are called the *Kripke-Feferman satisfaction axioms*.

Proof. Define $Sat(a, \sigma)$ as: $Tsat(a, \sigma)$. Axioms (1)–(7) of \mathcal{T}_{KF} then follow immediately from the corresponding axioms of \mathcal{T}_{TF} . \square

The result just proven has a converse.

Theorem 20. \mathcal{T}_{TF} is a definitional extension of \mathcal{T}_{KF} .

Proof. We need only define $Tsat(x, y)$ as: $Sat(x, y)$, and $Fsat(x, y)$ as: $Sat(\ulcorner \neg x \urcorner, y)$. \square

The theory \mathcal{T}_{KF} has a very nice property: It delivers half of Tarski's T-scheme.

Theorem 21. *For every formula ϕ of \mathcal{L} , \mathcal{T}_{KF} proves:*

$$Sat(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) \rightarrow \phi(\sigma_1, \dots, \sigma_n)$$

Proof. The proof is, of course, by induction on the complexity of formulas. In fact, the induction proceeds by showing that, for every formula ϕ , \mathcal{T}_{KF} proves both of the followign:

$$\begin{aligned} Sat(\ulcorner \phi(x_1, \dots, x_n) \urcorner, \sigma) &\rightarrow \phi(\sigma_1, \dots, \sigma_n) \\ Sat(\ulcorner \neg \phi(x_1, \dots, x_n) \urcorner, \sigma) &\rightarrow \neg \phi(\sigma_1, \dots, \sigma_n) \end{aligned}$$

We'll see the reason we need so to proceed shortly.

If ϕ is atomic, it is either an atomic formula of \mathcal{L} or is of the form $Sat(t, u)$. If the former, the result follows from axiom (2) of \mathcal{T}_{KF} . So suppose ϕ is of the form $Sat(t, u)$. Then axiom (5a) gives us:

$$Sat(\ulcorner Sat(t, u) \urcorner, \sigma) \equiv Sat(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma))$$

But we are assuming that \mathcal{T} proves: $\text{den}_{\mathcal{M}}(t, \sigma) = t(\sigma_1, \dots, \sigma_n)$. So

$$Sat(\text{den}_{\mathcal{M}}(t, \sigma), \text{den}_{\mathcal{M}}(u, \sigma)) \equiv Sat(t(\sigma_1, \dots, \sigma_n), u(\sigma_1, \dots, \sigma_n))$$

and similarly for the case of $\neg Sat(t, u)$.

If ϕ is a disjunction $\psi \vee \chi$, then we have:

$$\begin{aligned} & Sat(\ulcorner \psi(x_1, \dots, x_n) \vee \chi(x_1, \dots, x_n) \urcorner, \sigma) \rightarrow \\ & Sat(\ulcorner \psi(x_1, \dots, x_n) \urcorner, \sigma) \vee Sat(\ulcorner \chi(x_1, \dots, x_n) \urcorner, \sigma) \rightarrow \\ & \psi(\sigma_1, \dots, \sigma_n) \vee \chi(\sigma_1, \dots, \sigma_n) \end{aligned}$$

The first inference is by axiom (3a); the second, by the induction hypothesis. The case of negated disjunctions is similar.

So suppose ϕ is $\exists v_i \psi$; for convenience we assume $i = 1$. Then we have:

$$\begin{aligned} & Sat(\ulcorner \exists v_1 \phi(v_1, \dots, v_n) \urcorner, \sigma) \rightarrow \\ & \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge Sat(\ulcorner \phi(v_1, \dots, v_n) \urcorner, \tau)) \rightarrow \\ & \exists \tau (\forall j \neq i (\sigma_j = \tau_j) \wedge \phi(\tau_1, \dots, \tau_n)) \rightarrow \\ & \exists \tau \phi(\tau_1, \sigma_2, \dots, \sigma_n) \rightarrow \\ & \exists v_1 \phi(v_1, \sigma_2, \dots, \sigma_n) \rightarrow \end{aligned}$$

The first inference is by (4a); the second is by the induction hypothesis; the last two are simple logical manipulations. And again, the case of negated existentials is similar. \square

Corollary 22. In \mathcal{T}_{KF} , define a truth-predicate thus:

$$Tr(t) \equiv Sent(t) \wedge Sat(t, \langle \rangle)$$

where $\langle \rangle$ is the empty sequence. Then for every sentence ϕ , \mathcal{T}_{KF} proves: $Tr(\ulcorner \phi \urcorner) \rightarrow \phi$. Moreover, \mathcal{T}_{KF} proves: $\neg \exists t [Tr(t) \wedge Tr(\ulcorner \neg t \urcorner)]$.

Proof. The first claim is immediate from the preceding corollary. The second follows from axiom (6) of \mathcal{T}_{KF} . \square

Proposition 23. \mathcal{T}_{KF} is inconsistent with the *Tr*-introduction rule: $A \vdash Tr(\ulcorner A \urcorner)$.

This rule is similar to the rule of necessitation in modal logic. What it says, in effect, is that all *theorems* of \mathcal{T}_{KF} are true. Indeed, as we shall see, the only instance of this rule we need is one for a theorem of pure logic, not of \mathcal{T}_{KF} .

Proof. Let λ be a liar sentence. So $\mathcal{T}_{KF} \vdash \lambda \equiv \neg Tr(\ulcorner \lambda \urcorner)$. Then \mathcal{T}_{KF} proves:

- | | | |
|------|---|---------------------------|
| (1) | $Tr(\ulcorner \lambda \urcorner) \rightarrow \lambda$ | by Corollary 22 |
| (2) | $\lambda \rightarrow \neg Tr(\ulcorner \lambda \urcorner)$ | since λ is a liar |
| (3) | $\neg Tr(\ulcorner \lambda \urcorner)$ | from (1) and (2) |
| (4) | $Tr(\ulcorner \neg \lambda \urcorner) \rightarrow \neg \lambda$ | by Corollary 22 |
| (5) | $\neg \lambda \rightarrow Tr(\ulcorner \lambda \urcorner)$ | since λ is a liar |
| (6) | $Tr(\ulcorner \neg \lambda \urcorner) \rightarrow Tr(\ulcorner \neg \lambda \urcorner) \wedge Tr(\ulcorner \lambda \urcorner)$ | from (4) and (5) |
| (7) | $\neg [Tr(\ulcorner \neg \lambda \urcorner) \wedge Tr(\ulcorner \lambda \urcorner)]$ | by Corollary 22 |
| (8) | $\neg Tr(\ulcorner \neg \lambda \urcorner)$ | from (6) and (7) |
| (9) | $Tr(\ulcorner \lambda \vee \neg \lambda \urcorner) \rightarrow Tr(\ulcorner \lambda \urcorner) \vee Tr(\ulcorner \neg \lambda \urcorner)$ | by axiom (5) of KF |
| (10) | $\neg Tr(\ulcorner \lambda \vee \neg \lambda \urcorner)$ | from (3), (8), and (9) |

But \mathcal{T}_{KF} proves: $\lambda \vee \neg \lambda$, trivially, as a truth of logic. So *Tr*-introduction would yield $Tr(\ulcorner \lambda \vee \neg \lambda \urcorner)$, and \mathcal{T}_{KF} would be inconsistent. \square

Corollary 24. Let (E, A) be a fixed point of the Kripke construction over an interpreted language $\mathcal{L}_{\mathcal{M}}$ (i.e., what we get from the proof of Theorem 17). Then all axioms of \mathcal{T}_{KF} are true when ‘*Sat*’ is interpreted by E . Conversely, if all axioms of \mathcal{T}_{KF} are true when ‘*Sat*’ is interpreted by E , and A is defined as: $\langle t, \sigma \rangle \in A$ iff $\langle \ulcorner t \urcorner, \sigma \rangle \in E$, then (E, A) is a fixed point of the Kripke construction over $\mathcal{L}_{\mathcal{M}}$.

Proof. Left-to-right: Theorem 17 shows that, if (E, A) is a fixed point of the Kripke construction, then all axioms of $\mathcal{T}_{Tsat, Fsat}$ are true when ‘ $Tsat$ ’ is interpreted by E and ‘ $Fsat$ ’ is interpreted by A . But by Corollary 18 and Theorem 19, all axioms of \mathcal{T}_{KF} are provable in $\mathcal{T}_{Tsat, Fsat}$ if ‘ $Tsat$ ’ is re-written as ‘ Sat ’, so all axioms of \mathcal{T}_{KF} must be true when ‘ Sat ’ is interpreted by E .

Right-to-left: Suppose all axioms of \mathcal{T}_{KF} are true when ‘ Sat ’ is interpreted by E and define A as stated. By Theorem 20, if we re-write ‘ Sat ’ as ‘ $Tsat$ ’ and treat ‘ $Fsat(a, \sigma)$ ’ as an abbreviation of: $Tsat(\neg a, \sigma)$, then all axioms of \mathcal{T}_{TF} are theorems of \mathcal{T}_{KF} . So all axioms of $\mathcal{T}_{Tsat, Fsat}$ must be true when ‘ $Tsat$ ’ is interpreted by E and ‘ $Fsat$ ’ is treated as an abbreviation. \square

6. SATISFACTION-PREDICATES AND THREE-VALUED LOGIC

The theories whose consistency we just proved were all *classical* theories: That is, the underlying logic of each of these theories was classical, two-valued logic. The theory thus does not yet look much like Kripke’s theory of truth, which employs a three-valued logic, though there is an obvious sort of correspondence between the satisfaction and non-satisfaction predicates that appear above and the extension and anti-extension of the truth-predicate in Kripke’s treatment, a correspondence we just exploited. We now show how to use it to recover Kripke’s original theory.

Given an interpreted language $\mathcal{L}_{\mathcal{M}}$, define an interpreted language $\mathcal{L}_{\mathcal{M}}^K$ as follows. The atomic expressions of $\mathcal{L}_{\mathcal{M}}^K$ are to be those of $\mathcal{L}_{\mathcal{M}}$ plus a predicate $TSat_K(a, \sigma)$ intended to express satisfaction for $\mathcal{L}_{\mathcal{M}}^K$. Let the extension and anti-extension of $TSat_K(x, y)$ be the fixed points $I_{TSat, \mathcal{M}}$ and $I_{FSat, \mathcal{M}}$ and extend this to an interpretation of the whole of $\mathcal{L}_{\mathcal{M}}^K$ by means of the Strong Kleene scheme. That is, assume we have a definition of satisfaction for atomic formulae of $\mathcal{L}_{\mathcal{M}}$ and simultaneously define satisfaction and anti-satisfaction as follows:

σ satisfies ϕ if, and only if, either:

- (1) ϕ is an atomic formula of $\mathcal{L}_{\mathcal{M}}$, and $\text{SatAt}_{\mathcal{M}}(\phi, \sigma)$;
- (2) ϕ is a conjunction $\psi \wedge \chi$ and σ satisfies both ψ and χ ;
- (3) ϕ is a negation $\neg\psi$ and σ anti-satisfies ψ ;
- (4) ϕ is an existential quantification $\exists v_i\psi$ and there is some sequence τ agreeing with σ everywhere except, possibly, at i such that σ satisfies ψ ;
- (5) ϕ is $TSat_K(t, u)$, where t is the Gödel number of some formula ψ and u is the code of a sequence τ and τ satisfies ψ .

σ anti-satisfies ϕ if, and only if, either:

- (1) ϕ is an atomic formula of $\mathcal{L}_{\mathcal{M}}$, and $\neg\text{SatAt}_{\mathcal{M}}(\phi, \sigma)$;
- (2) ϕ is a conjunction $\psi \wedge \chi$ and σ anti-satisfies at least one of ψ and χ ;
- (3) ϕ is a negation $\neg\psi$ and σ satisfies ψ ;
- (4) ϕ is an existential quantification $\exists v_i\psi$ and for every sequence τ agreeing with σ everywhere except, possibly, at i , τ anti-satisfies ψ ;
- (5) ϕ is $TSat_K(t, u)$, where t is the Gödel number of some formula ψ and u is the code of a sequence τ and τ anti-satisfies ψ .

We say that ϕ is a true sentence of $\mathcal{L}_{\mathcal{M}}^K$ if it is satisfied by every sequence; a false sentence, if it is anti-satisfied by every sequence.

Theorem 25 (Kripke). $TSat_K(a, \sigma)$ is a satisfaction predicate for $\mathcal{L}_{\mathcal{M}}^K$. That is: $TSat_K(a, \sigma)$ is a true sentence of $\mathcal{L}_{\mathcal{M}}^K$ just in case σ is (the code of) a sequence that satisfies the sentence (whose code is) a , and $TSat_K(a, \sigma)$ is a false sentence of $\mathcal{L}_{\mathcal{M}}^K$ just in case σ is (the code of) a sequence that anti-satisfies the sentence (whose code is) a .¹⁵

Proof. The crucial observation is simply that the definition of satisfaction and anti-satisfaction just given transcribes the clauses of the simultaneous inductive definition of $ISat_{\mathcal{M}}$ and $ISat_{\mathcal{M}}$. The only difference is that we have omitted any reference to a formula $FSat_K(x, y)$, which might have been taken to express anti-satisfaction. The theorem can thus be proven by a straightforward inductive argument. \square

To get an actual *formal theory* of truth for $\mathcal{L}_{\mathcal{M}}^K$, one must axiomatize the underlying three-valued logic of the theory. There are different ways to do this: See Kremer [3] and McGee [5] for a couple of different approaches. Once one has done that, however, the principles governing the truth- and satisfaction-predicates are easily stated. These will be the so-called T-rules:

$$\begin{aligned} A &\vdash T(\ulcorner A \urcorner) \\ \neg A &\vdash \neg T(\ulcorner A \urcorner) \\ T(\ulcorner A \urcorner) &\vdash A \\ \neg T(\ulcorner A \urcorner) &\vdash \neg A \end{aligned}$$

with corresponding rules for satisfaction.

7. EXERCISES

Exercise 1. Show that the converse of Lemma 8 also holds: If (σ) is true when σ^n is taken to have S as its extension, then S is a fixed point of $\Sigma_{\phi}^{\mathcal{M}}$.

Exercise 2. Complete the proof of Theorem 13 by doing the cases we omitted.

REFERENCES

- [1] Solomon Feferman. Toward useful type-free theories I. *Journal of Symbolic Logic*, 49:75–111, 1984.
- [2] Alexander Kechris and Yiannis Moschovakis. Recursion in higher types. In John Barwise, editor, *A Handbook of Mathematical Logic*, pages 681–739. North-Holland Publishing, New York, 1977.
- [3] Michael Kremer. Kripke and the logic of truth. *Journal of Philosophical Logic*, 17:225–78, 1988.
- [4] Saul Kripke. Outline of a theory of truth. *Journal of Philosophy*, 72:690–716, 1975.
- [5] Vann McGee. *Truth, Vagueness, and Paradox: An Essay on the Logic of Truth*. Hackett, 1990.
- [6] Yiannis Moschovakis. *Elementary Induction on Abstract Structures*. North Holland Publishing, Amsterdam, 1974.

¹⁵One can state this more generally, to allow for the possibility that variables occur in the ‘sentence’ $TSat(a, \sigma)$, but we’ll leave doing so as an exercise. (In the case of arithmetic, this adds nothing, since every element of the domain has a name.)