

Self-referential theories

Samuel A. Alexander*

Department of Mathematics, the Ohio State University

February 9, 2016

Abstract

We study the structure of families of theories in the language of arithmetic extended to allow these families to refer to one another and to themselves. If a theory contains schemata expressing its own truth and expressing a specific Turing index for itself, and contains some other mild axioms, then that theory is untrue. We exhibit some families of true self-referential theories that barely avoid this forbidden pattern.

1 Introduction

This is a paper about families of r.e. theories, each capable of referring to itself and the others. Many of this paper's results first appeared in the author's dissertation [1]. There, they were stated in terms of families of interacting mechanical knowing agents. Here, we will speak instead of families of self-referential r.e. theories. We hope this will more directly expose the underlying mathematics.

In epistemology, it is well-known that a (suitably idealized) truthful knowing machine capable of arithmetic, logic, and self-reflection, cannot know its own truth and its own code. This is due, in various guises, to authors such as Lucas [7], Benacerraf [3], Reinhardt [10], Penrose [8], and Putnam [9]. In terms of self-referential theories, a true theory satisfying certain assumptions cannot contain schemata stating its own truth and its own Gödel number (if such a theory did exist, we could program a machine knower that knows precisely its consequences). Reinhardt conjectured, and Carlson proved [5], a truthful machine knower can know (in a local sense, i.e., expressed by infinite schemata rather than a single axiom) that it is truthful and has some code, without knowing which. A true self-referential theory can (in a local sense) state its own truth and recursive enumerability. We showed [2] that, alternatively, a truthful machine can (in a local sense) exactly know its own code, if not required to know its own truth. A true theory can state (in a local sense) its own Gödel number.

Our goal is to generalize the above consistency results to multiple theories. The paper contains four main findings. In the following list of promises, except where otherwise stated, \prec is an r.e. well-founded partial-order on ω , and *expresses* is meant in the local (infinite schema) sense.

1. There are true theories $(T_i)_{i \in \omega}$ such that T_i expresses a Gödel number of T_j (all i, j) and T_i expresses the truth of T_j (all $j \prec i$).
2. There are true theories $(T_i)_{i \in \omega}$ such that T_i expresses a Gödel number of T_j ($j \prec i$), the truth of T_j ($j \preceq i$), and the fact that T_j has some Gödel number (all i, j).
3. If \prec is ill-founded, and if we extend the base language to include a predicate for computable ordinals and require the theories to include rudimentary facts about them, then 1 and 2 fail.
4. Finally, if we do not extend the base language as in 3, then there do exist ill-founded r.e. partial orders \prec such that 1 and 2 hold.

Our proofs of 1 and 2 are constructive, but the proof of 4 is nonconstructive. In short, if 4 were false, either of 1 or 2 could be used to define the set WF of r.e. well-founded partial orders of ω using nothing but arithmetic and a truth predicate Tr for arithmetic. This is impossible since WF is Π_1^1 -complete and Tr is Δ_1^1 .

*Email: alexander@math.ohio-state.edu

2 Preliminaries

To us, *theory* and *schema* mean *set of sentences* (a *sentence* is a formula with no free variables).

Definition 1. (Standard Definitions)

1. When a first-order structure is clear from context, an *assignment* is a function s mapping first-order variables into the universe of that structure. If x is a variable and u is an element of the universe, $s(x|u)$ is the assignment that agrees with s except that it maps x to u .
2. We write $\mathcal{M} \models \phi[s]$ to indicate that the first-order structure \mathcal{M} satisfies the formula ϕ relative to the assignment s . We write $\mathcal{M} \models \phi$ just in case $\mathcal{M} \models \phi[s]$ for every assignment s . If T is a theory, $\mathcal{M} \models T$ means that $\mathcal{M} \models \phi$ for every $\phi \in T$.
3. We write $\text{FV}(\phi)$ for the set of free variables of ϕ .
4. We write $\phi(x|t)$ for the result of substituting term t for variable x in ϕ .
5. \mathcal{L}_{PA} is the language of Peano arithmetic, with constant symbol 0 and function symbols $S, +, \cdot$ with the usual arities. If \mathcal{L} extends \mathcal{L}_{PA} , an \mathcal{L} -structure *has standard first-order part* if it has universe \mathbb{N} and interprets $0, S, +$ and \cdot as intended.
6. We define \mathcal{L}_{PA} -terms \bar{n} ($n \in \mathbb{N}$), called *numerals*, so that $\bar{0} \equiv 0$ and $\overline{n+1} \equiv S(\bar{n})$.
7. We fix a computable bijection $\langle \bullet, \bullet, \bullet \rangle : \mathbb{N}^3 \rightarrow \mathbb{N}$. Being computable, this is \mathcal{L}_{PA} -definable, so we may freely act as if \mathcal{L}_{PA} contained a function symbol for this bijection. Similarly we may act as if \mathcal{L}_{PA} contained a binary predicate symbol $\bullet \in W_\bullet$ for membership in the n th r.e. set W_n .
8. Whenever a computable language is clear from context, $\phi \mapsto \ulcorner \phi \urcorner$ denotes Gödel numbering.
9. A *valid* formula is one that is true in every structure.
10. A *universal closure* of ϕ is a sentence $\forall x_1 \cdots \forall x_n \phi$ where $\text{FV}(\phi) \subseteq \{x_1, \dots, x_n\}$. We write $\text{ucl}(\phi)$ to denote a generic universal closure of ϕ .

Note that if \mathcal{M} is a structure and ψ is a universal closure of ϕ , in order to prove $\mathcal{M} \models \psi$ it suffices to let s be an arbitrary assignment and show $\mathcal{M} \models \phi[s]$.

To formalize self-referential theories, we employ an extension of first-order logic where languages may contain new unary connective symbols. This logic is borrowed from [5].

Definition 2. (The Base Logic) A language \mathcal{L} of the *base logic* is a first-order language \mathcal{L}_0 together with a class of symbols called *operators*. Formulas of \mathcal{L} are defined as usual, with the clause that $\mathbf{T}_i \models \phi$ is a formula whenever ϕ is a formula and $\mathbf{T}_i \models$ is an operator. Syntactic parts of Definition 1 extend to the base logic in obvious ways (we define $\text{FV}(\mathbf{T}_i \models \phi) = \text{FV}(\phi)$). An \mathcal{L} -structure \mathcal{M} is a first-order \mathcal{L}_0 -structure \mathcal{M}_0 together with a function that takes one operator $\mathbf{T}_i \models$, one \mathcal{L} -formula ϕ , and one assignment s , and outputs either True or False—in which case we write $\mathcal{M} \models \mathbf{T}_i \models \phi[s]$ or $\mathcal{M} \not\models \mathbf{T}_i \models \phi[s]$, respectively—satisfying the following three requirements.

1. Whether or not $\mathcal{M} \models \mathbf{T}_i \models \phi[s]$ does not depend on $s(x)$ if $x \notin \text{FV}(\phi)$.
2. If ϕ and ψ are *alphabetic variants* (meaning that one is obtained from the other by renaming bound variables so as to respect the binding of the quantifiers), then $\mathcal{M} \models \mathbf{T}_i \models \phi[s]$ if and only if $\mathcal{M} \models \mathbf{T}_i \models \psi[s]$.
3. For variables x and y such that y is substitutable for x in $\mathbf{T}_i \models \phi$, $\mathcal{M} \models \mathbf{T}_i \models \phi(x|y)[s]$ if and only if $\mathcal{M} \models \mathbf{T}_i \models \phi[s(x|s(y))]$.

The definition of $\mathcal{M} \models \phi[s]$ for arbitrary \mathcal{L} -formulas is obtained from this by induction. Semantic parts of Definition 1 extend to the base logic in obvious ways.

Traditionally the operator $\mathbf{T}_i \models$ would be written K_i , and the formula $K_i \phi$ would be read like “agent i knows ϕ ”. For the present paper, the added intuition would not be worth the philosophical distraction.

Theorem 3. (Completeness and compactness) Suppose \mathcal{L} is an r.e. language in the base logic.

1. The set of valid \mathcal{L} -formulas is r.e.
2. For any r.e. \mathcal{L} -theory Σ , $\{\phi : \Sigma \models \phi\}$ is r.e.
3. There is an effective procedure, given (a Gödel number of) an r.e. \mathcal{L} -theory Σ , to find (a Gödel number of) $\{\phi : \Sigma \models \phi\}$.
4. If Σ is an \mathcal{L} -theory and $\Sigma \models \phi$, there are $\sigma_1, \dots, \sigma_n \in \Sigma$ such that¹ $\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow \phi$ is valid.

Proof. By interpreting the base logic within first-order logic (for details see [1]). \square

Definition 4. If \mathcal{L} is a first-order language and I is an index set, let $\mathcal{L}(I)$ be the language (in the base logic) consisting of \mathcal{L} along with operators $\mathbf{T}_i \models$ for all $i \in I$.

In case I is a singleton, $\mathcal{L}_{\text{PA}}(I)$ is a form of Shapiro's [11] language of Epistemic Arithmetic.

Definition 5. For any $\mathcal{L}_{\text{PA}}(I)$ -formula ϕ with $\text{FV}(\phi) = \{x_1, \dots, x_n\}$, and for assignment s (into \mathbb{N}), let ϕ^s be the sentence

$$\phi^s \equiv \phi(x_1 | \overline{s(x_1)}) \cdots (x_n | \overline{s(x_n)})$$

obtained by replacing all free variables in ϕ by numerals for their s -values.

For example, if $s(x) = 0$ and $s(y) = 2$ then $(\forall z(x = y + z))^s \equiv \forall z(0 = S(S(0)) + z)$.

Definition 6. If $\mathbf{T} = (T_i)_{i \in I}$ is an I -indexed family of $\mathcal{L}_{\text{PA}}(I)$ -theories and \mathcal{N} is an $\mathcal{L}_{\text{PA}}(I)$ -structure, we say $\mathcal{N} \models \mathbf{T}$ if $\mathcal{N} \models T_i$ for all $i \in I$.

Definition 7. Suppose $\mathbf{T} = (T_i)_{i \in I}$ is an I -indexed family of $\mathcal{L}_{\text{PA}}(I)$ -theories. The *intended structure* for \mathbf{T} is the $\mathcal{L}_{\text{PA}}(I)$ -structure $\mathcal{M}_{\mathbf{T}}$ with standard first-order part, interpreting the operators $\mathbf{T}_i \models$ ($i \in I$) as follows:

$$\mathcal{M}_{\mathbf{T}} \models \mathbf{T}_i \models \phi[s] \text{ if and only if } T_i \models \phi^s.$$

If $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}$, we say \mathbf{T} is *true*.

Lemma 8. Let $\mathbf{T} = (T_i)_{i \in I}$ be a family of $\mathcal{L}_{\text{PA}}(I)$ -theories. For any $\mathcal{L}_{\text{PA}}(I)$ -formula ϕ and assignment s , $\mathcal{M}_{\mathbf{T}} \models \phi[s]$ if and only if $\mathcal{M}_{\mathbf{T}} \models \phi^s$.

Proof. By induction. \square

Definition 9. By the *axioms of Peano arithmetic* for $\mathcal{L}_{\text{PA}}(I)$ we mean the axioms of Peano arithmetic, with induction extended to $\mathcal{L}_{\text{PA}}(I)$.

Lemma 10. Any $\mathcal{L}_{\text{PA}}(I)$ -structure with standard first-order part and satisfying the conclusion of Lemma 8 satisfies the axioms of Peano arithmetic for $\mathcal{L}_{\text{PA}}(I)$.

Proof. Let \mathcal{M} be any $\mathcal{L}_{\text{PA}}(I)$ -structure with standard first-order part and satisfying the conclusion of Lemma 8. Let σ be an axiom of Peano arithmetic for $\mathcal{L}_{\text{PA}}(I)$. If σ is not an instance of induction, then $\mathcal{M} \models \sigma$ since \mathcal{M} has standard first-order part. But suppose σ is $\text{ucl}(\phi(x|0) \rightarrow \forall x(\phi \rightarrow \phi(x|S(x))) \rightarrow \forall x\phi)$. To see $\mathcal{M} \models \sigma$, let s be an arbitrary assignment and assume $\mathcal{M} \models \phi(x|0)[s]$ and $\mathcal{M} \models \forall x(\phi \rightarrow \phi(x|S(x)))[s]$. By Lemma 8, $\mathcal{M} \models \phi^{s(x|0)}$ and $\forall m \in \mathbb{N}$, if $\mathcal{M} \models \phi^{s(x|m)}$ then $\mathcal{M} \models \phi(x|S(x))^{s(x|m)}$. Evidently $\phi(x|S(x))^{s(x|m)} \equiv \phi^{s(x|m+1)}$. By mathematical induction, $\forall m \in \mathbb{N}$, $\mathcal{M} \models \phi^{s(x|m)}$. By Lemma 8, $\mathcal{M} \models \forall x\phi[s]$. \square

Definition 11. Suppose $\mathbf{T} = (T_i)_{i \in I}$ is a family $\mathcal{L}_{\text{PA}}(I)$ -theories. If $\mathbf{T}^+ = (T_i^+)_{i \in I}$ is another such family, we say $\mathbf{T} \subseteq \mathbf{T}^+$ if $T_i \subseteq T_i^+$ for every $i \in I$. If T is a single $\mathcal{L}_{\text{PA}}(I)$ -theory, we say $T \subseteq \mathbf{T}$ if $T \subseteq T_i$ for all $i \in I$. If $\mathbf{T}^1 = (T_i^1)_{i \in I}$ and $\mathbf{T}^2 = (T_i^2)_{i \in I}$ are families of $\mathcal{L}_{\text{PA}}(I)$ -theories, $\mathbf{T}^1 \cup \mathbf{T}^2$ is the family $\mathbf{T}' = (T_i')_{i \in I}$ where each $T_i' = T_i^1 \cup T_i^2$. Arbitrary unions $\bigcup_{n \in \mathbb{N}} \mathbf{T}^n$ are defined similarly.

Definition 12. Suppose $\mathbf{T} = (T_i)_{i \in I}$ is a family of $\mathcal{L}_{\text{PA}}(I)$ -theories. For each $i \in I$, we say T_i is $\mathbf{T}_i \models$ -closed if $\mathbf{T}_i \models \phi \in T_i$ whenever $\phi \in T_i$. We say \mathbf{T} is *closed* if each T_i is $\mathbf{T}_i \models$ -closed.

Definition 13. If I is an r.e. index set, a family $\mathbf{T} = (T_i)_{i \in I}$ is *r.e. just in case* $\{(\phi, i) : \phi \in T_i\}$ is r.e.

¹We write $A \rightarrow B \rightarrow C$ for $A \rightarrow (B \rightarrow C)$, and likewise for longer chains.

3 Generic Axioms

If \mathbf{T} is a family of theories whose truth was in doubt, and if we state a theorem removing that doubt, we often state more: that $\mathbf{T} \cup \mathbf{S}$ is true, where \mathbf{S} is some background theory of provability, including non-controversial things like Peano arithmetic or the schema $\text{ucl}(\mathbf{T}_i \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_i \models \phi \rightarrow \mathbf{T}_i \models \psi)$. The choice of \mathbf{S} is somewhat arbitrary, or at best based on tradition. We will avoid this arbitrary choice by stating results in the form: “ \mathbf{T} is true together with any background theory of provability such that...”

Definition 14. A family \mathbf{T} of $\mathcal{L}_{\text{PA}}(\omega)$ -theories is *closed-r.e.-generic* if \mathbf{T} is r.e. and $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}$ for every closed r.e. family $\mathbf{T}' \supseteq \mathbf{T}$ of $\mathcal{L}_{\text{PA}}(\omega)$ -theories.

Lemma 15. If \mathbf{T} is a union of closed-r.e.-generic families and \mathbf{T} is r.e., then \mathbf{T} is closed-r.e.-generic.

Proof. Straightforward. □

Definition 16. For $j \in I$ and for T an $\mathcal{L}_{\text{PA}}(I)$ -theory, we write $[T]_j$ for the family $\mathbf{T} = (T_i)_{i \in I}$ where $T_j = T$ and $T_i = \emptyset$ for all $i \neq j$.

The following lemma provides building blocks that can be combined in diverse ways, via Lemma 15, to form background theories of provability. In the following lemma, part 4 is the whole reason for the “closed” in “closed-r.e.-generic”, and part 7 is part of the reason for the “r.e.”.

Lemma 17. For any $i, j \in \omega$, each of the following families is closed-r.e.-generic.

1. $[S]_i$ where S is: (j -Deduction) the schema $\text{ucl}(\mathbf{T}_j \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models \psi)$.
2. $[S]_i$ where S is: (Assigned Validity) the schema ϕ^s (ϕ valid, s an assignment).
3. $[\text{Assigned Validity}]_j \cup [S]_i$ where S is: (j -Validity) $\text{ucl}(\mathbf{T}_j \models \phi)$ for ϕ valid.
4. $[\text{Assigned Validity}]_j \cup [j\text{-Validity}]_j \cup [j\text{-Deduction}]_j \cup [S]_i$ where S is:
 $(j\text{-Introspection})$ the schema $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models \mathbf{T}_j \models \phi)$.
5. $[S]_i$ where S is the set of axioms of Peano arithmetic for $\mathcal{L}_{\text{PA}}(\omega)$.
6. $[S]_i$ where S is any r.e. set of true arithmetic sentences.
7. $[S]_i$ where S is: (j -SMT) (See [5] and [10]) $\text{ucl}(\exists e \forall x (\mathbf{T}_j \models \phi \leftrightarrow x \in W_e)), e \notin \text{FV}(\phi)$.
8. $\mathbf{T} \cup [S]_i$ where $\mathbf{T} = (T_k)_{k \in \omega}$ is closed-r.e.-generic and S is the schema $\mathbf{T}_j \models \phi$ ($\phi \in T_j$).

Proof.

(1) Let $\mathbf{T}' = (T'_k)_{k \in \omega}$ be any closed r.e. family of $\mathcal{L}_{\text{PA}}(\omega)$ -theories such that $\mathbf{T}' \supseteq [S]_i$ where S is j -Deduction. We must show $\mathcal{M}_{\mathbf{T}'} \models [S]_i$. In other words we must show $\mathcal{M}_{\mathbf{T}'} \models \text{ucl}(\mathbf{T}_j \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models \psi)$ for any ϕ, ψ . Let s be an assignment and assume $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models (\phi \rightarrow \psi)[s]$ and $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi[s]$, we must show $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \psi[s]$. By Definition of $\mathcal{M}_{\mathbf{T}'}$, $T'_j \models (\phi \rightarrow \psi)^s$ and $T'_j \models \phi^s$. Clearly $(\phi \rightarrow \psi)^s \equiv \phi^s \rightarrow \psi^s$ so by modus ponens $T'_j \models \psi^s$, that is, $\mathcal{M}_{\mathbf{T}'} \models \psi[s]$.

(2) Let $\mathbf{T}' = (T'_k)_{k \in \omega}$ be a closed r.e. superset of $[S]_i$ where S is Assigned Validity. We must show $\mathcal{M}_{\mathbf{T}'} \models [S]_i$. If $\phi \in [S]_i$ then ϕ is ϕ_0^s for some valid ϕ_0 and some assignment s . Since ϕ_0 is valid, $\mathcal{M}_{\mathbf{T}'} \models \phi_0[s]$. By Lemma 8, $\mathcal{M}_{\mathbf{T}'} \models \phi_0^s$.

(3) By Theorem 3, $[\text{Assigned Validity}]_j \cup [j\text{-Validity}]_i$ is r.e. Let $\mathbf{T}' = (T'_k)_{k \in \omega}$ be any closed r.e. family of $\mathcal{L}_{\text{PA}}(\omega)$ -theories such that T'_j contains Assigned Validity and T'_i contains j -Validity. We must show $\mathcal{M}_{\mathbf{T}'}$ satisfies Assigned Validity and j -Validity. For Assigned Validity, let ϕ be valid and s an assignment. Since ϕ is valid, $\mathcal{M}_{\mathbf{T}'} \models \phi[s]$, so by Lemma 8, $\mathcal{M}_{\mathbf{T}'} \models \phi^s$ as desired. For j -Validity, let ϕ be valid and s an assignment. Since T'_j contains Assigned Validity, $T'_j \models \phi^s$, so by definition of $\mathcal{M}_{\mathbf{T}'}$, $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi[s]$.

(4) Recursive enumerability is by Theorem 3. Let $\mathbf{T}' = (T'_k)_{k \in \omega}$ be any closed r.e. family of $\mathcal{L}_{\text{PA}}(\omega)$ -theories such that T'_j contains Assigned Validity, j -Validity and j -Deduction, and T'_i contains j -Introspection. That $\mathcal{M}_{\mathbf{T}'}$ satisfies Assigned Validity and j -Validity is as in (3). That $\mathcal{M}_{\mathbf{T}'}$ satisfies j -Deduction is straightforward. For j -Introspection, let s be an assignment and assume $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi[s]$, we will show $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \mathbf{T}_j \models \phi[s]$. Since $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi[s]$, $T'_j \models \phi^s$. By Theorem 3, there are $\sigma_1, \dots, \sigma_n \in T'_j$ such that $\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow \phi^s$ is valid. Since T'_j contains j -Validity, $T'_j \models \mathbf{T}_j \models (\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow \phi^s)$. By repeated applications of j -Deduction contained in T'_j , $T'_j \models \mathbf{T}_j \models \sigma_1 \rightarrow \dots \rightarrow \mathbf{T}_j \models \sigma_n \rightarrow \mathbf{T}_j \models (\phi^s)$. Since \mathbf{T}' is closed, T'_j is $\mathbf{T}_j \models$ -closed and so contains $\mathbf{T}_j \models \sigma_1, \dots, \mathbf{T}_j \models \sigma_n$. So $T'_j \models (\mathbf{T}_j \models \phi)^s$ and $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \mathbf{T}_j \models \phi[s]$.

(5) Follows from Lemma 10.

(6) Trivial.

(7) Suppose $\mathbf{T}' = (T'_i)_{i \in \omega}$ is a closed r.e. family of $\mathcal{L}_{\text{PA}}(\omega)$ -theories and $\mathbf{T}' \supseteq [S]_i$ where S is j -SMT. We must show $\mathcal{M}_{\mathbf{T}'} \models [S]_i$. That is, given ϕ with $e \notin \text{FV}(\phi)$, we must show $\mathcal{M}_{\mathbf{T}'} \models \text{ucl}(\exists e \forall x (\mathbf{T}_j \models \phi \leftrightarrow x \in W_e))$. Let s be an assignment and let $x_1, \dots, x_k = \overline{\text{FV}(\phi) \setminus \{x\}}$. Since T_j is r.e., by the S - m - n theorem there is some n such that $W_n = \{m : T_j \models \phi(x|\overline{m})(x_1|\overline{s(x_1)}) \dots (x_k|\overline{s(x_k)})\}$. Since $e \notin \text{FV}(\phi)$, and $\mathcal{M}_{\mathbf{T}'}$ has standard first-order part, it follows that $\mathcal{M}_{\mathbf{T}'} \models \forall x (\mathbf{T}_j \models \phi \leftrightarrow x \in W_e)[s(e|n)]$.

(8) Suppose $\mathbf{T}' = (T'_i)_{i \in \omega} \supseteq \mathbf{T} \cup [S]_i$ where $\mathbf{T} = (T_i)_{i \in \omega}$ is closed-r.e.-generic and S is the schema $\mathbf{T}_j \models \phi$ ($\phi \in T_j$). Right away $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}$ because \mathbf{T} is closed-r.e.-generic. It remains to show that $\mathcal{M}_{\mathbf{T}'} \models [S]_i$, i.e., that $\mathcal{M}_{\mathbf{T}'} \models S$. Fix $\phi \in T_j$ and let s be any assignment. Since ϕ is a sentence, $\phi \equiv \phi^s$ and thus $T_j \models \phi^s$. Since $T'_j \supseteq T_j$, $T'_j \models \phi^s$. By definition of $\mathcal{M}_{\mathbf{T}'}$, $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi[s]$. By arbitrariness of s , $\mathcal{M}_{\mathbf{T}'} \models \mathbf{T}_j \models \phi$. \square

4 First Consistency Result: Prioritizing Exact Codes

The following theorem fulfils the first promise from the introduction.

Theorem 18. Suppose \prec is an r.e. well-founded partial order on ω and $\mathbf{T}^0 = (T_i^0)_{i \in \omega}$ is closed-r.e.-generic. For each $n \in \mathbb{N}$, let $\mathbf{T}(n) = (T_i(n))_{i \in \omega}$ where each $T_i(n)$ is the smallest $\mathbf{T}_i \models$ -closed theory containing the following:

1. The axioms in T_i^0 .
2. $\forall x (\mathbf{T}_j \models \phi \leftrightarrow \langle \overline{\phi}, \bar{j}, x \rangle \in W_{\bar{n}})$ whenever $j \in \omega$, $\text{FV}(\phi) \subseteq \{x\}$.
3. $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$ whenever $j \prec i$.

There is some $n \in \mathbb{N}$ such that $\mathbf{T}(n)$ is true.

Proof. By the S - m - n Theorem, there is a total computable $f : \mathbb{N} \rightarrow \mathbb{N}$ such that $\forall n \in \mathbb{N}$,

$$W_{f(n)} = \{\langle \overline{\phi}, \bar{i}, m \rangle : \text{FV}(\phi) \subseteq \{x\} \text{ and } T_i(n) \models \phi(x|\overline{m})\}.$$

Using the Recursion Theorem, fix $n \in \mathbb{N}$ such that $W_{f(n)} = W_n$. For brevity write \mathbf{T} for $\mathbf{T}(n)$ and T_i for $T_i(n)$. We will show $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}$. This is a self-referential statement: to show T_i is true includes showing $\mathcal{M}_{\mathbf{T}} \models \text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$, which is essentially the statement that T_j is true. Hence the restriction $j \prec i$, which allows induction since \prec is well founded. We will show, by \prec -induction on i , that $\mathcal{M}_{\mathbf{T}} \models T_i$ for every $i \in \omega$. Fix $i \in \omega$ and assume $\mathcal{M}_{\mathbf{T}} \models T_j$ for all $j \prec i$. Suppose $\sigma \in T_i$, we will show $\mathcal{M}_{\mathbf{T}} \models \sigma$.

Case 1: $\sigma \in T_i^0$. Then $\mathcal{M}_{\mathbf{T}} \models \sigma$ because \mathbf{T}^0 is closed-r.e.-generic and $\mathbf{T} \supseteq \mathbf{T}^0$ is closed r.e.

Case 2: σ is $\forall x (\mathbf{T}_j \models \phi \leftrightarrow \langle \overline{\phi}, \bar{j}, x \rangle \in W_{\bar{n}})$ for some $j \in \omega$, $\text{FV}(\phi) \subseteq \{x\}$. Let s be an assignment, $m \in \mathbb{N}$. The following are equivalent.

$$\begin{aligned} \mathcal{M}_{\mathbf{T}} \models \mathbf{T}_j \models \phi[s(x|m)] & & & \\ T_j \models \phi^{s(x|m)} & & \text{(Definition of } \mathcal{M}_{\mathbf{T}}) & \\ T_j \models \phi(x|\overline{m}) & & \text{(Since } \text{FV}(\phi) \subseteq \{x\}) & \\ \langle \overline{\phi}, \bar{j}, m \rangle \in W_n & & \text{(By definition of } n) & \\ \mathcal{M}_{\mathbf{T}} \models \langle \overline{\phi}, \bar{j}, \overline{m} \rangle \in W_{\bar{n}} & & (\mathcal{M}_{\mathbf{T}} \text{ has standard first-order part)} & \\ \mathcal{M}_{\mathbf{T}} \models \langle \overline{\phi}, \bar{j}, x \rangle \in W_{\bar{n}}[s(x|m)]. & & \text{(Lemma 8)} & \end{aligned}$$

Case 3: σ is $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$ for some $j \prec i$. Let s be an assignment and assume $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}_j \models \phi[s]$. This means $T_j \models \phi^s$. By our \prec -induction hypothesis, $\mathcal{M}_{\mathbf{T}} \models T_j$, so $\mathcal{M}_{\mathbf{T}} \models \phi^s$. By Lemma 8, $\mathcal{M}_{\mathbf{T}} \models \phi[s]$.

Case 4: σ is only present in T_i because of the clause that T_i is $\mathbf{T}_i \models$ -closed. Then σ is $\mathbf{T}_i \models \sigma_0$ for some $\sigma_0 \in T_i$. Being in T_i , σ_0 is a sentence, so for any assignment s , $\sigma_0 \equiv \sigma_0^s$, $T_i \models \sigma_0^s$, and finally $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}_i \models \sigma_0[s]$.

By \prec -induction, $\mathcal{M}_{\mathbf{T}} \models T_i$ for all $i \in \omega$. This shows $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}$, that is, \mathbf{T} is true. \square

The first promise from the introduction is met: for any r.e. well-founded partial order \prec on ω , there are theories $(T_n)_{n \in \omega}$ such that $\forall i, j, k \in \omega$ with $j \prec i$, T_i expresses the truth of T_j , and T_i expresses a Gödel number of T_k . In order to fulfil the second promise we will extend Carlson's notion of *stratification* to the case of multiple operators, and introduce *stratifiers*, a tool used to deal with subtleties that arise when multiple self-referential theories refer to one another.

In [2] the technique behind Theorem 18 was used to exhibit a machine that knows its own code.

5 Stratification

For the second promise from the introduction, we need to prove a result like Theorem 18 where T_i includes $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$ for all $j \preceq i$, not just $j \prec i$. This rules out the direct induction used above. Instead we will induct on formula complexity in a specific way. Naive formula complexity will not work: we would need to show all of T_i consistent just to show $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}_i \models (1 = 0) \rightarrow (1 = 0)$. So we will “stratify” theories T_i into increasing families T_i^α ($\alpha \in \text{Ord}$), and show these stratified theories are true. At large enough ordinals, the theories will repeat themselves (or *collapse*) in a technical sense, allowing recovery of T_i .

Definition 19. We define a binary relation \leq_1 on Ord by transfinite recursion so that for all $\alpha, \beta \in \text{Ord}$, $\alpha \leq_1 \beta$ if and only if $\alpha \leq \beta$ and (α, \leq, \leq_1) is a Σ_1 -elementary substructure of (β, \leq, \leq_1) .

The following theorem is based on calculations from [4]. It was used by Carlson to prove Reinhardt's conjecture [5]. We state it here without proof.

Theorem 20.

1. The binary relation \leq_1 is a recursive partial ordering on $\epsilon_0 \cdot \omega$.
2. For all positive integers $m \leq n$, $\epsilon_0 \cdot m \leq_1 \epsilon_0 \cdot n$.
3. (See Figure 1) For any $\alpha \leq \beta \in \text{Ord}$, $\alpha \leq_1 \beta$ if and only if the following statement is true. For every finite set $X \subseteq \alpha$ and every finite set $Y \subseteq [\alpha, \beta)$, there is a set $X < \tilde{Y} < \alpha$ such that $X \cup \tilde{Y} \cong_{(\leq, \leq_1)} X \cup Y$.

(Figure 1 here)

The usefulness of Theorem 20 will first appear in Theorem 31, but first we need some machinery.

Definition 21. Let $\mathcal{I} = ((\epsilon_0 \cdot \omega) \times \omega) \sqcup \omega$. Thus $\mathcal{L}_{\text{PA}}(\mathcal{I})$ contains operators $\mathbf{T}_{(\alpha, i)} \models$ for all $\alpha \in \epsilon_0 \cdot \omega$, $i \in \omega$, along with operators $\mathbf{T}_i \models$ for all $i \in \omega$. As abbreviation, we write $\mathbf{T}_i^\alpha \models$ for $\mathbf{T}_{(\alpha, i)} \models$, and refer to α as its *exponent*.

Definition 22. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , $\text{On}(\phi) \subseteq \epsilon_0 \cdot \omega$ denotes the set of exponents appearing in ϕ .

Definition 23. Suppose $i \in \omega$. The *i-stratified* formulas of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ are defined as follows (where ϕ ranges over $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formulas).

1. If ϕ is $\mathbf{T}_j \models \phi_0$ for some $j \neq i$, then ϕ is *i-stratified* if and only if ϕ is an $\mathcal{L}_{\text{PA}}(\omega)$ -formula.
2. If ϕ is $\mathbf{T}_j^\alpha \models \phi_0$ for some $j \neq i$, then ϕ is not *i-stratified*.
3. If ϕ is $\mathbf{T}_i \models \phi_0$, then ϕ is not *i-stratified*.
4. If ϕ is $\mathbf{T}_i^\alpha \models \phi_0$, then ϕ is *i-stratified* if and only if ϕ_0 is *i-stratified* and $\alpha > \text{On}(\phi_0)$.
5. If ϕ is $\neg \phi_0$, $\phi_1 \rightarrow \phi_2$, or $\forall x \phi_0$, then ϕ is *i-stratified* if and only if its immediate subformula(s) are.

6. If ϕ is atomic, then ϕ is i -stratified.

An $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory T is i -stratified if ϕ is i -stratified whenever $\phi \in T$. An $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ is *very i -stratified* if ϕ is i -stratified and $\text{On}(\phi) \subseteq \{\epsilon_0 \cdot 1, \epsilon_0 \cdot 2, \dots\}$.

For example:

- $\mathbf{T}_7^\omega \models \mathbf{T}_7^5 \models (1 = 0) \rightarrow \mathbf{T}_8 \models (1 = 0)$ is 7-stratified but not 6- or 8-stratified.
- $\mathbf{T}_7^5 \models \mathbf{T}_7^\omega \models (1 = 0)$ is not 7-stratified, nor is $\mathbf{T}_7^5 \models \mathbf{T}_7 \models (1 = 0)$.
- $\mathbf{T}_7^5 \models \mathbf{T}_8 \models \mathbf{T}_7 \models (1 = 0)$ is 7-stratified but $\mathbf{T}_7^5 \models \mathbf{T}_8 \models \mathbf{T}_7^4 \models (1 = 0)$ is not.

Definition 24. Suppose $X \subseteq \epsilon_0 \cdot \omega$ and $h : X \rightarrow \epsilon_0 \cdot \omega$ is order preserving. For each $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , let $h(\phi)$ be the $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula obtained by applying h to every superscript in ϕ that is in X .

For example if $X = \{1, \omega\}$, $h(1) = 0$, and $h(\omega) = \omega \cdot 2 + 1$, then

$$h(\mathbf{T}_i^0 \models (1 = 0) \rightarrow \mathbf{T}_i^1 \models (1 = 0) \rightarrow \mathbf{T}_i^\omega \models (1 = 0)) \equiv \mathbf{T}_i^0 \models (1 = 0) \rightarrow \mathbf{T}_i^0 \models (1 = 0) \rightarrow \mathbf{T}_i^{\omega \cdot 2 + 1} \models (1 = 0).$$

Definition 25. Suppose $X \subseteq \epsilon_0 \cdot \omega$ and $h : X \rightarrow \epsilon_0 \cdot \omega$ is order preserving. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{N} , we define an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure $h(\mathcal{N})$ that has the same universe as \mathcal{N} , agrees with \mathcal{N} on $\mathcal{L}_{\text{PA}}(\omega)$, and interprets $\mathcal{L}_{\text{PA}}(\mathcal{I}) \setminus \mathcal{L}_{\text{PA}}(\omega)$ so that

$$h(\mathcal{N}) \models \mathbf{T}_i^\alpha \models \phi[s] \text{ if and only if } \mathcal{N} \models h(\mathbf{T}_i^\alpha \models \phi)[s].$$

Lemma 26. Suppose $X \subseteq \epsilon_0 \cdot \omega$, $h : X \rightarrow \epsilon_0 \cdot \omega$ is order preserving, and \mathcal{N} is an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , $h(\mathcal{N}) \models \phi[s]$ if and only if $\mathcal{N} \models h(\phi)[s]$.

Proof. By induction. □

Corollary 27. Suppose $X \subseteq \epsilon_0 \cdot \omega$ and $h : X \rightarrow \epsilon_0 \cdot \omega$ is order preserving. For any valid $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , $h(\phi)$ is valid.

Proof. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{N} and assignment s , $h(\mathcal{N}) \models \phi[s]$ by validity, so $\mathcal{N} \models h(\phi)[s]$ by Lemma 26. □

Definition 28. If $X \subseteq \text{Ord}$ and $h : X \rightarrow \text{Ord}$, we call h a *covering* if h is order preserving and whenever $x, y \in X$ and $x \leq_1 y$, $h(x) \leq_1 h(y)$.

Definition 29. Suppose $i \in \omega$. A set T of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentences is i -uniform if the following conditions hold:

1. T is i -stratified.
2. Whenever $\phi \in T$, $X \subseteq \epsilon_0 \cdot \omega$, $\text{On}(\phi) \subseteq X$, and $h : X \rightarrow \epsilon_0 \cdot \omega$ is a covering, then $h(\phi) \in T$.

Definition 30. If T is an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory and $\alpha \in \epsilon_0 \cdot \omega$, let $T \cap \alpha$ be the set $\{\phi \in T : \text{On}(\phi) \subseteq \alpha\}$ of sentences in T that do not contain any superscripts $\geq \alpha$.

Theorem 31. (The Collapse Theorem) Suppose T is an i -uniform $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory.

1. If n is a positive integer and $\text{On}(\phi) \subseteq \epsilon_0 \cdot n$, then $T \models \phi$ if and only if $T \cap (\epsilon_0 \cdot n) \models \phi$.
2. If $\alpha \leq_1 \beta$ and $\text{On}(\phi) \subseteq \alpha$, then $T \cap \alpha \models \phi$ if and only if $T \cap \beta \models \phi$.

Proof. Note that since T is i -uniform, in particular T is i -stratified. We will prove (1), the proof of (2) is similar.

(\Leftarrow) Immediate since $T \cap (\epsilon_0 \cdot n) \subseteq T$.

(\Rightarrow) Assume $T \models \phi$. By Theorem 3 there are $\sigma_1, \dots, \sigma_k \in T$ such that

$$\Phi \equiv \sigma_1 \rightarrow \dots \rightarrow \sigma_k \rightarrow \phi$$

is valid. Let $X = \text{On}(\Phi) \cap (\epsilon_0 \cdot n)$, $Y = \text{On}(\Phi) \cap [\epsilon_0 \cdot n, \infty)$, note $|X|, |Y| < \infty$.

Since Y is finite, there is some integer $n' \succ n$ such that $Y \subseteq \epsilon_0 \cdot n'$. By Theorem 20 part 2, $\epsilon_0 \cdot n \leq_1 \epsilon_0 \cdot n'$. By Theorem 20 part 3, there is some $X < \tilde{Y} < \epsilon_0 \cdot n$ such that $X \cup \tilde{Y} \cong_{(\leq, \leq_1)} X \cup Y$.

Let $h : X \cup Y \rightarrow X \cup \tilde{Y}$ be a (\leq, \leq_1) -isomorphism. Since $\text{On}(\phi) \subseteq \epsilon_0 \cdot n$, $h(\phi) = \phi$. By Corollary 27,

$$h(\Phi) \equiv h(\sigma_1) \rightarrow \dots \rightarrow h(\sigma_k) \rightarrow \phi$$

is valid. Since T is i -uniform, $h(\sigma_1), \dots, h(\sigma_k) \in T$. Finally since $\text{range}(h) < \epsilon_0 \cdot n$, $h(\sigma_1), \dots, h(\sigma_k) \in T \cap (\epsilon_0 \cdot n)$, showing $T \cap (\epsilon_0 \cdot n) \models \phi$. \square

Loosely speaking, what we have done in Theorem 31 is we have taken a proof of ϕ and we have *collapsed* the proof, shrinking its ordinals by using Theorem 20 part 3.

Definition 32. For every $i \in \omega$ we define the following $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -schema:

- (i -Collapse) $\text{ucl}(\mathbf{T}_i^\alpha \models \phi \leftrightarrow \mathbf{T}_i^\beta \models \phi)$ whenever $\mathbf{T}_i^\alpha \models \phi$ is i -stratified and $\alpha \leq_1 \beta$.

Definition 33. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , ϕ^- is the result of erasing all superscripts from ϕ . If T is an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory, $T^- = \{\sigma^- : \sigma \in T\}$.

For example, if ϕ is $\mathbf{T}_5^\omega \models (1 = 0) \rightarrow \mathbf{T}_5^{\omega+1} \models \mathbf{T}_5^\omega \models (1 = 0)$, then ϕ^- is $\mathbf{T}_5 \models (1 = 0) \rightarrow \mathbf{T}_5 \models \mathbf{T}_5 \models (1 = 0)$.

Lemma 34. If T is i -uniform then for every $\phi \in T$ there is some $\psi \in T$ such that ψ is very i -stratified and $\psi^- \equiv \phi^-$.

Proof. Let $X = \text{On}(\phi) = \{\alpha_1 < \dots < \alpha_n\}$, $Y = \{\epsilon_0 \cdot 1, \dots, \epsilon_0 \cdot n\}$, and define $h : X \rightarrow Y$ by $h(\alpha_j) = \epsilon_0 \cdot j$. Clearly h is injective and order preserving; by Theorem 20 part 2, h is a covering. Since T is i -uniform, T contains $\psi \equiv h(\phi)$. Clearly ψ is very i -stratified and $\psi^- \equiv \phi^-$. \square

Definition 35. For any $\mathcal{L}_{\text{PA}}(\omega)$ -structure \mathcal{N} , we define an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{N}^- that has the same universe as \mathcal{N} , agrees with \mathcal{N} on $\mathcal{L}_{\text{PA}}(\omega)$, and interprets $\mathcal{L}_{\text{PA}}(\mathcal{I}) \setminus \mathcal{L}_{\text{PA}}(\omega)$ as follows. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , $\alpha \in \epsilon_0 \cdot \omega$, $i \in \mathbb{N}$, and assignment s ,

$$\mathcal{N}^- \models \mathbf{T}_i^\alpha \models \phi[s] \text{ if and only if } \mathcal{N} \models (\mathbf{T}_i^\alpha \models \phi)^-[s].$$

Lemma 36. Suppose \mathcal{N} is an $\mathcal{L}_{\text{PA}}(\omega)$ -structure. For every $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , $\mathcal{N}^- \models \phi[s]$ if and only if $\mathcal{N} \models \phi^-[s]$.

Proof. By induction. \square

Corollary 37. If ϕ is a valid $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula, then ϕ^- is a valid $\mathcal{L}_{\text{PA}}(\omega)$ -formula.

Proof. Similar to the proof of Corollary 27. \square

A converse-like statement holds for Corollary 37 as well.

Lemma 38. For any valid $\mathcal{L}_{\text{PA}}(\omega)$ -sentence ϕ and $i \in \omega$, there is a valid very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentence ψ such that $\psi^- \equiv \phi$.

We postpone the proof of Lemma 38 until Section 6.

Definition 39. Let $i \in \omega$. We define the following $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -schemas.

- (i -Stratificality) $\text{ucl}(\mathbf{T}_i^\alpha \models \phi)$ whenever ϕ is a valid $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula and $\mathbf{T}_i^\alpha \models \phi$ is i -stratified.
- (i -Stratideduction) $\text{ucl}(\mathbf{T}_i^\alpha \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_i^\alpha \models \phi \rightarrow \mathbf{T}_i^\alpha \models \psi)$ whenever this formula is i -stratified.

An $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory T is *i -stratifiable* if the following conditions hold:

1. T is i -uniform.
2. T includes i -Strativalidity, i -Stratideduction and i -Collapse.
3. For every $\phi \in T$, if $\mathbf{T}_i^\alpha \models \phi$ is i -stratified then $\mathbf{T}_i^\alpha \models \phi \in T$.

A family $\mathbf{T} = (T_i)_{i \in \omega}$ is *stratifiable* if each T_i is i -stratifiable.

The reason for the word *stratifiable* in Definition 39 will become clearer later on when we arrive at Definition 41.

The following theorem serves as an omnibus of results from Section 5 of [5].

Theorem 40. (Proof Stratification) Suppose T is an i -stratifiable $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory. Then:

1. Whenever $T \cap \alpha \models \phi$, $\mathbf{T}_i^\alpha \models \phi$ is an i -stratified sentence, and $\beta > \alpha$, then $T \cap \beta \models \mathbf{T}_i^\alpha \models \phi$.
2. For any very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentences ρ and σ , if $\rho^- \equiv \sigma^-$ then $T \models \rho \leftrightarrow \sigma$.
3. For any very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentence ϕ , $T \models \phi$ if and only if $T^- \models \phi^-$.

Proof. Note that since T is i -stratifiable, in particular T is i -uniform and hence, i -stratified.

Claim 0: Any time $T \models \mathbf{T}_i^\alpha \models (\rho \leftrightarrow \sigma)$ and this is i -stratified, $T \models \mathbf{T}_i^\alpha \models \rho \leftrightarrow \mathbf{T}_i^\alpha \models \sigma$.

Assume the hypotheses. By i -Strativalidity, $T \models \mathbf{T}_i^\alpha \models ((\rho \leftrightarrow \sigma) \rightarrow (\rho \rightarrow \sigma))$. By i -Stratideduction,

$$\begin{aligned} T &\models \mathbf{T}_i^\alpha \models ((\rho \leftrightarrow \sigma) \rightarrow (\rho \rightarrow \sigma)) \rightarrow \mathbf{T}_i^\alpha \models (\rho \leftrightarrow \sigma) \rightarrow \mathbf{T}_i^\alpha \models (\rho \rightarrow \sigma) \\ \text{and } T &\models \mathbf{T}_i^\alpha \models (\rho \rightarrow \sigma) \rightarrow \mathbf{T}_i^\alpha \models \rho \rightarrow \mathbf{T}_i^\alpha \models \sigma. \end{aligned}$$

It follows that $T \models \mathbf{T}_i^\alpha \models \rho \rightarrow \mathbf{T}_i^\alpha \models \sigma$. The reverse implication is similar.

Claim 1: If $T \cap \alpha \models \phi$, $\mathbf{T}_i^\alpha \models \phi$ is an i -stratified sentence, and $\beta > \alpha$, then $T \cap \beta \models \mathbf{T}_i^\alpha \models \phi$.

Given $T \cap \alpha \models \phi$, there are $\sigma_1, \dots, \sigma_n \in T \cap \alpha$ such that $\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow \phi$ is valid. By instances of i -Strativalidity and i -Stratideduction contained in $T \cap \beta$, and by the last hypothesis of the theorem, $T \cap \beta \models \mathbf{T}_i^\alpha \models \phi$.

Claim 2: If ρ and σ are very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentences and $\rho^- \equiv \sigma^-$, then $T \models \rho \leftrightarrow \sigma$.

By induction on ρ . Note that ρ is not of the form $\mathbf{T}_j^\alpha \models \rho_0$ (with $j \neq i$), as that is not i -stratified. If ρ is $\mathbf{T}_j^\alpha \models \rho_0$ then $\rho \equiv \rho^- \equiv \sigma^- \equiv \sigma$ and the claim is immediate.

The only nontrivial remaining case is when ρ is $\mathbf{T}_i^\alpha \models \rho_0$. Since ρ is very i -stratified, this implies $\alpha = \epsilon_0 \cdot n$ (some positive integer n) and ρ_0 is very i -stratified. Since $\sigma^- \equiv \rho^-$ and σ is very stratified, this implies $\sigma \equiv \mathbf{T}_i^{\epsilon_0 \cdot m} \models \sigma_0$ for some positive integer m and very i -stratified σ_0 with $\sigma_0^- \equiv \rho_0^-$. Assume $m \leq n$, the other case is similar.

By induction, $T \models \rho_0 \leftrightarrow \sigma_0$. By compactness, there is a natural $\ell \geq n$ such that $T \cap (\epsilon_0 \cdot \ell) \models \rho_0 \leftrightarrow \sigma_0$. By Claim 1, $T \models \mathbf{T}_i^{\epsilon_0 \cdot \ell} \models (\rho_0 \leftrightarrow \sigma_0)$; Claim 0 then gives $T \models \mathbf{T}_i^{\epsilon_0 \cdot \ell} \models \rho_0 \leftrightarrow \mathbf{T}_i^{\epsilon_0 \cdot \ell} \models \sigma_0$. The claim now follows since T contains i -Collapse and $\epsilon_0 \cdot m \leq \epsilon_0 \cdot n \leq \epsilon_0 \cdot \ell$ (Theorem 20 part 2).

Claim 3: If ϕ is a very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentence and $T \models \phi$, then $T^- \models \phi^-$.

By compactness, find $\sigma_1, \dots, \sigma_n \in T$ such that $\sigma_1 \rightarrow \dots \rightarrow \sigma_n \rightarrow \phi$ is valid. By Corollary 37, so is $\sigma_1^- \rightarrow \dots \rightarrow \sigma_n^- \rightarrow \phi^-$, witnessing $T^- \models \phi^-$.

Claim 4: If ϕ is a very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentence and $T^- \models \phi^-$, then $T \models \phi$.

By compactness, there is a valid sentence

$$\Phi \equiv \sigma_1^- \rightarrow \dots \rightarrow \sigma_n^- \rightarrow \phi^-$$

where each $\sigma_j \in T$. By Lemma 38, there is a valid very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sentence Ψ such that $\Psi^- \equiv \Phi$. And because $\Psi^- \equiv \Phi$, this implies

$$\Psi \equiv \sigma_1^* \rightarrow \dots \rightarrow \sigma_n^* \rightarrow \phi^*$$

where each $(\sigma_j^*)^- \equiv \sigma_j^-$, $(\phi^*)^- \equiv \phi^-$, and $\sigma_1^*, \dots, \sigma_n^*, \phi^*$ are very i -stratified.

By Lemma 34, there are very i -stratified $\sigma_1^{**}, \dots, \sigma_n^{**} \in T$ with each $(\sigma_j^{**})^- \equiv \sigma_j^- \equiv (\sigma_j^*)^-$. By Claim 2, $T \models \phi^* \leftrightarrow \phi$, and for $j = 1, \dots, n$, $T \models \sigma_j^{**} \leftrightarrow \sigma_j^*$. Thus

$$T \models (\sigma_1^{**} \rightarrow \dots \rightarrow \sigma_n^{**} \rightarrow \phi) \leftrightarrow \Psi,$$

and since Ψ is valid and the $\sigma_j^{**} \in T$, this shows $T \models \phi$. \square

Definition 41. If $\mathbf{T} = (T_i)_{i \in \omega}$ is a stratifiable family of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theories, its *stratification*, written $\text{Str}(\mathbf{T})$, is the family $\text{Str}(\mathbf{T}) = (S_i)_{i \in \mathcal{I}}$, where for every $i \in \omega$, $S_i = T_i^-$ and $\forall \alpha \in \epsilon_0 \cdot \omega$, $S_{(\alpha, i)} = T_i \cap \alpha$.

Theorem 42. (The Stratification Theorem) Suppose $\mathbf{T} = (T_i)_{i \in \omega}$ is a stratifiable family of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theories. For any $i \in \omega$, any very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , and any assignment s , $\mathcal{M}_{\text{Str}(\mathbf{T})} \models \phi[s]$ if and only if $\mathcal{M}_{\text{Str}(\mathbf{T})} \models \phi^-[s]$.

Proof. By induction on ϕ . The only nontrivial case is when ϕ is $\mathbf{T}_i^\alpha \models \psi$. Since ϕ is very i -stratified, ψ is very i -stratified and we may write $\alpha = \epsilon_0 \cdot n$ for some positive integer n , $\text{On}(\psi) \subseteq \epsilon_0 \cdot n$. The following are equivalent.

$$\begin{aligned} \mathcal{M}_{\text{Str}(\mathbf{T})} \models \mathbf{T}_i^{\epsilon_0 \cdot n} \models \psi[s] & & (\text{Definition of } \mathcal{M}_{\text{Str}(\mathbf{T})}) \\ T_i \cap (\epsilon_0 \cdot n) \models \psi^s & & (\text{Theorem 31}) \\ T_i \models \psi^s & & (\text{Theorem 40}) \\ T_i^- \models (\psi^s)^- & & (\text{Theorem 40}) \\ T_i^- \models (\psi^-)^s & & (\text{Clearly } (\psi^s)^- \equiv (\psi^-)^s) \\ \mathcal{M}_{\text{Str}(\mathbf{T})} \models \mathbf{T}_i \models \psi^-[s] & & (\text{Definition of } \mathcal{M}_{\text{Str}(\mathbf{T})}) \end{aligned}$$

\square

6 Stratifiers

In order to apply theorems from the previous section, it is necessary to work with families $\mathbf{T} = (T_i)_{i \in \omega}$ where each T_i is i -stratified. If we want T_i^- to (locally) express the truthfulness of T_j^- , we cannot simply add a schema like $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$ to T_i , because this is not necessarily i -stratified: for example, the particular instance $\mathbf{T}_j \models \mathbf{T}_i \models (1 = 0) \rightarrow \mathbf{T}_i \models (1 = 0)$ is not i -stratified. But neither is, say, $\mathbf{T}_j \models \mathbf{T}_i^\alpha \models (1 = 0) \rightarrow \mathbf{T}_i^\alpha \models (1 = 0)$, where $\mathbf{T}_i^\alpha \models$ occurs within the scope of $\mathbf{T}_j \models$. We will use a schema $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi^+)$, where \bullet^+ varies over what we call i -stratifiers.

Definition 43. Suppose $X \subseteq \epsilon_0 \cdot \omega$, $|X| = \infty$, and $i \in \omega$. The i -stratifier given by X is the function $\phi \mapsto \phi^+$ taking $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formulas to $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formulas as follows.

1. If ϕ is atomic or of the form $\mathbf{T}_j \models \phi_0$ with $j \neq i$, then $\phi^+ \equiv \phi$.
2. If ϕ is $\mathbf{T}_i \models \phi_0$ ($\phi_0 \in \mathcal{L}_{\text{PA}}(\omega)$) then $\phi^+ \equiv \mathbf{T}_i^\alpha \models \phi_0^+$ where $\alpha = \min\{x \in X : x > \text{On}(\phi_0^+)\}$.
3. If ϕ is $\mathbf{T}_i \models \phi_0$ ($\phi_0 \notin \mathcal{L}_{\text{PA}}(\omega)$) then $\phi^+ \equiv \mathbf{T}_i \models \phi_0$.
4. If ϕ is $\neg\psi$, $\psi \rightarrow \rho$, or $\forall x\psi$, then ϕ^+ is $\neg\psi^+$, $\psi^+ \rightarrow \rho^+$ or $\forall x\psi^+$, respectively.

By an i -stratifier we mean an i -stratifier given by some X . By the i -veristratifier we mean the i -stratifier given by $X = \{\epsilon_0 \cdot 1, \epsilon_0 \cdot 2, \dots\}$.

For example, if \bullet^+ is the i -veristratifier and $j \neq i$ then

$$(\mathbf{T}_j \models \mathbf{T}_i \models (1 = 0) \rightarrow \mathbf{T}_i \models \mathbf{T}_i \models (1 = 0))^+ \equiv \mathbf{T}_j \models \mathbf{T}_i \models (1 = 0) \rightarrow \mathbf{T}_i^{\epsilon_0 \cdot 2} \models \mathbf{T}_i^{\epsilon_0} \models (1 = 0).$$

On the other hand, if \bullet^+ is the i -veristratifier and $\alpha \in \epsilon_0 \cdot \omega$, then

$$(\mathbf{T}_i \models \mathbf{T}_i^\alpha \models (1 = 0))^+ \equiv \mathbf{T}_i \models \mathbf{T}_i^\alpha \models (1 = 0)$$

because the formula $\mathbf{T}_i^\alpha \models (1 = 0)$ is $\notin \mathcal{L}_{\text{PA}}(\omega)$. Where Definition 43 is concerned we are really interested in formulas where $\mathbf{T}_i \models$ occurs but not $\mathbf{T}_i^\alpha \models$ for any α ; or where $\mathbf{T}_i^\alpha \models$ occurs for various α but $\mathbf{T}_i \models$ does not occur. It is technically simpler to use the single language $\mathcal{L}_{\text{PA}}(\mathcal{I})$ for both, while bearing in mind that it also contains formulas such as $\mathbf{T}_i \models \mathbf{T}_i^\alpha \models (1 = 0)$ that we don't really care about.

Lemma 44. Suppose $Z \subseteq \epsilon_0 \cdot \omega$, $h : Z \rightarrow \epsilon_0 \cdot \omega$ is order preserving, $i \in \omega$, and \bullet^+ is an i -stratifier. For any $\mathcal{L}_{\text{PA}}(\omega)$ -formula θ with $\text{On}(\theta^+) \subseteq Z$, there is a computable i -stratifier \bullet^* with $\theta^* \equiv h(\theta^+)$.

Proof. Let $X_0 = \{h(\alpha) : \alpha \in \text{On}(\theta^+)\}$, let $X = X_0 \cup \{\alpha \in \epsilon_0 \cdot \omega : \alpha > X_0\}$, and let \bullet^* be the i -stratifier given by X . By induction, for every subformula θ_0 of θ , $\theta_0^* \equiv h(\theta_0^+)$. \square

Definition 45. Suppose \mathcal{N} is an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure and \bullet^+ is an i -stratifier. We define an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{N}^+ as follows. The universe and interpretation of arithmetic of \mathcal{N}^+ agree with those of \mathcal{N} , as do the interpretations of $\mathbf{T}_j \models$ ($j \neq i$) and $\mathbf{T}_j^{\alpha} \models$ (any α, j). As for $\mathbf{T}_i \models$, for any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s ,

$$\mathcal{N}^+ \models \mathbf{T}_i \models \phi[s] \text{ if and only if } \mathcal{N} \models (\mathbf{T}_i \models \phi)^+[s].$$

Lemma 46. (Compare Lemma 36) Suppose \mathcal{N} is an $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure, $i \in \omega$, and \bullet^+ is an i -stratifier. For every $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , $\mathcal{N}^+ \models \phi[s]$ if and only if $\mathcal{N} \models \phi^+[s]$.

Proof. By induction. Note that since $\phi \in \mathcal{L}_{\text{PA}}(\omega)$, clause 3 of Definition 43 plays no part in the definition of $\mathcal{N} \models \phi^+[s]$. \square

Lemma 47. For any $\mathcal{L}_{\text{PA}}(\omega)$ -formula ϕ , any $i \in \omega$, and any i -stratifier \bullet^+ , ϕ is valid if and only if ϕ^+ is valid.

Proof.

(\Rightarrow) Assume ϕ is valid. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{N} and assignment s , $\mathcal{N}^+ \models \phi[s]$ by validity, so $\mathcal{N} \models \phi^+[s]$ by Lemma 46.

(\Leftarrow) By Corollary 37. \square

Given Lemma 47, Lemma 38 (which we promised to prove) is trivial.

Proof of Lemma 38. By Lemma 47 with \bullet^+ taken to be the i -veristratifier. \square

For the remainder of the section, fix a strict r.e. well-founded partial-order \prec on ω .

The next definition needs motivation. We seek a true stratified family $\mathbf{T} = (T_i)_{i \in \omega}$ where T_i declares the truth of T_j ($j \preceq i$). By collapsing \mathbf{T} we hope to fulfil the second promise from the introduction. For $j \prec i$, T_i cannot contain $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$: it is not i -stratified (e.g. if ϕ is $\mathbf{T}_i \models (1 = 0)$). An alternative is $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi^+)$ where \bullet^+ ranges over i -stratifiers. But imagine proving $\mathcal{M}_{\mathbf{T}} \models \text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi^+)$. We would assume $T_j \models \phi^s$ and try to show $\mathcal{M}_{\mathbf{T}} \models \phi^+[s]$. Since $j \prec i$ we could assume $\mathcal{M}_{\mathbf{T}} \models T_j$, so $\mathcal{M}_{\mathbf{T}} \models \phi^s$, so $\mathcal{M}_{\mathbf{T}} \models \phi[s]$ —not what we want. If only $\mathcal{M}_{\mathbf{T}}^+ \models \phi[s]$, Lemma 46 would give $\mathcal{M}_{\mathbf{T}} \models \phi^+[s]$. In a sense, \bullet^+ is invisible to T_j (T_j can't see T_i because $j \prec i$). So our proof that $\mathcal{M}_{\mathbf{T}} \models T_j$ could've gone through for $\mathcal{M}_{\mathbf{T}}^+$. The following definition will allow the proof to go through for $\mathcal{M}_{\mathbf{T}}^+$ via strong induction hypothesis.

Definition 48. Suppose $\mathcal{M}, \mathcal{M}'$ are $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structures, $i \in \omega$. We say $\mathcal{M} \Rightarrow_i \mathcal{M}'$ if there is a sequence

$$i_1, \bullet^1, i_2, \bullet^2, \dots, i_n, \bullet^n,$$

each $i_k \succ i$, each \bullet^k a computable i_k -stratifier, such that $\mathcal{M}' = (\dots(\mathcal{M}^1)^2)\dots)^n$.

Lemma 49. Suppose $\mathcal{M}, \mathcal{M}', \mathcal{M}''$ are $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structures and $i, j, k \in \omega$.

1. $\mathcal{M} \Rightarrow_i \mathcal{M}$.
2. If $\mathcal{M} \Rightarrow_i \mathcal{M}'$ then \mathcal{M} and \mathcal{M}' have the same universe and agree on all symbols of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ except possibly for some symbols $\mathbf{T}_j \models$ where $j \succ i$.
3. If $\mathcal{M} \Rightarrow_i \mathcal{M}'$ then for any $j \in \omega$, either \mathcal{M} and \mathcal{M}' agree on $\mathbf{T}_j \models$, or there is some j -stratifier \bullet^+ such that \mathcal{M}^+ and \mathcal{M}' agree on $\mathbf{T}_j \models$.
4. If $\mathcal{M} \Rightarrow_i \mathcal{M}'$, $\mathcal{M}' \Rightarrow_j \mathcal{M}''$, and $k \preceq i, j$, then $\mathcal{M} \Rightarrow_k \mathcal{M}''$.
5. If $\mathcal{M} \Rightarrow_i \mathcal{M}'$ and $j \prec i$ then $\mathcal{M} \Rightarrow_j \mathcal{M}'$.

6. If $j < i$ and \bullet^+ is a computable i -stratifier, then $\mathcal{M}^+ \Rightarrow_j \mathcal{M}$.

Proof. Straightforward. \square

Lemma 50. Suppose $\mathcal{M}, \mathcal{M}'$ are $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structures and $\mathcal{M} \Rightarrow_i \mathcal{M}'$ for some $i \in \omega$. Further suppose \mathcal{M} has the property that for every $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , $\mathcal{M} \models \phi[s]$ if and only if $\mathcal{M} \models \phi^s$. Then for all such ϕ and s , $\mathcal{M}' \models \phi[s]$ if and only if $\mathcal{M}' \models \phi^s$.

Proof. By induction on sequence length, we may assume $\mathcal{M}' = \mathcal{M}^+$ for some computable j -stratifier \bullet^+ , $j > i$. For any $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , the following are equivalent.

$$\begin{aligned} \mathcal{M}^+ &\models \phi[s] \\ \mathcal{M} &\models \phi^+[s] && \text{(Lemma 46)} \\ \mathcal{M} &\models (\phi^+)^s && \text{(Hypothesis)} \\ \mathcal{M} &\models (\phi^s)^+ && \text{(Clearly } (\phi^+)^s \equiv (\phi^s)^+ \text{)} \\ \mathcal{M}^+ &\models \phi^s. && \text{(Lemma 46)} \end{aligned}$$

\square

Lemma 51. Suppose the $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -structure \mathcal{M} is an instance of Definition 7. For any $i \in \omega$, $\mathcal{M}' \Rightarrow_i \mathcal{M}$, $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ and assignment s , $\mathcal{M}' \models \phi[s]$ if and only if $\mathcal{M}' \models \phi^s$.

Proof. By Lemmas 8 and 50. \square

7 Second Consistency Result: Prioritizing Self-Truth

In this section we fulfil the second promise from the introduction. Throughout, $<$ is an r.e. well-founded partial-order of ω .

Definition 52. (Compare Definition 14) Suppose $\mathbf{T} = (T_i)_{i \in \omega}$ is an r.e. family of i -uniform $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theories. We say \mathbf{T} is *stratified-r.e.-generic* if for every stratifiable r.e. family $\mathbf{U} \supseteq \mathbf{T}$, every $i \in \omega$, and every $\mathcal{M} \Rightarrow_i \mathcal{M}_{\text{Str}(\mathbf{U})}$, $\mathcal{M} \models T_i$.

Lemma 53. If the family $\mathbf{T} = (T_i)_{i \in \omega}$ of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -sets is r.e. and is a union of stratified-r.e.-generic families, then \mathbf{T} is stratified-r.e.-generic.

Proof. Straightforward. \square

Lemma 54. (Compare Lemma 17) For any $i, j \in \mathbb{N}$, each of the following families is stratified-r.e.-generic.

1. $[i\text{-Stratideduction}]_i$.
2. $[j\text{-Deduction}]_i$ (if $j \not\prec i$).
3. $[S]_i$ (if $j \neq i$) where S is: (Weak j -Deduction) $\text{ucl}(\mathbf{T}_j \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models (\psi \wedge (\phi \vee \neg \phi)))$.
4. $[S]_i$ where S is: (i -Assigned Strativalidity) the schema ϕ^s (ϕ valid and i -stratified, s an assignment).
5. $[i\text{-Assigned Strativalidity}]_i \cup [i\text{-Strativalidity}]_i$.
6. $[i\text{-Assigned Strativalidity}]_i \cup [i\text{-Validity}]_j$ (if $j \neq i$).
7. $[i\text{-Assigned Strativalidity}]_i \cup [i\text{-Strativalidity}]_i \cup [i\text{-Stratideduction}]_i \cup [i\text{-Introspection}]_j$ ($j \neq i$).
8. $[i\text{-Assigned Strativalidity}]_i \cup [i\text{-Strativalidity}]_i \cup [i\text{-Stratideduction}]_i \cup [S]_i$ where S is:

$$(i\text{-Stratrospection}) \text{ ucl}(\mathbf{T}_i^\alpha \models \phi \rightarrow \mathbf{T}_i^\beta \models \mathbf{T}_i^\alpha \models \phi) \text{ whenever this is } i\text{-stratified.}$$

9. $[S]_i$ where S is the set of those axioms of Peano arithmetic for $\mathcal{L}_{\text{PA}}(\mathcal{I})$ that are i -stratified.
10. $[S]_i$ where S is any r.e. set of true arithmetic sentences.
11. $[j\text{-SMT}]_i$ ($j \neq i$).
12. $[S]_i$, where S is: (i -Strati-SMT) $\text{ucl}(\exists e \forall x (\mathbf{T}_i^\alpha \models \phi \leftrightarrow x \in W_e))$ when this is i -stratified, $e \notin \text{FV}(\phi)$.
13. $\mathbf{T} \cup [S]_i$ where $\mathbf{T} = (T_k)_{k \in \omega}$ is stratified-r.e.-generic and S is the schema $\mathbf{T}_i^\alpha \models \phi$ ($\phi \in T_i$ such that this is i -stratified).

Proof. For uniformity of 4–8, use Corollary 27. Uniformity of the other families is clear. Recursive enumerability follows from the fact that \prec is r.e. In each case below, let $\mathbf{U} = (U_k)_{k \in \omega}$ be a stratifiable r.e. family extending the family in question. For brevity let $\hat{\mathbf{U}} = \text{Str}(\mathbf{U})$.

(1) Similar to part 1 of Lemma 17.

(2) Let $\mathcal{M} \Rightarrow_i \mathcal{M}_{\hat{\mathbf{U}}}$, we must show $\mathcal{M} \models \text{ucl}(\mathbf{T}_j \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models \psi)$. Let s be an assignment and assume $\mathcal{M} \models \mathbf{T}_j \models (\phi \rightarrow \psi)[s]$ and $\mathcal{M} \models \mathbf{T}_j \models \phi[s]$. Since $j \neq i$, Lemma 49 says \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}$ agree on $\mathbf{T}_j \models$. By definition of $\mathcal{M}_{\hat{\mathbf{U}}}$, $U_j^- \models \phi^s \rightarrow \psi^s$ and $U_j^- \models \phi^s$, thus $U_j^- \models \psi^s$, so $\mathcal{M}_{\hat{\mathbf{U}}} \models \mathbf{T}_j \models \psi[s]$ and so does \mathcal{M} .

(3) Let $\mathcal{M} \Rightarrow_i \mathcal{M}_{\hat{\mathbf{U}}}$, we must show $\mathcal{M} \models \text{ucl}(\mathbf{T}_j \models (\phi \rightarrow \psi) \rightarrow \mathbf{T}_j \models \phi \rightarrow \mathbf{T}_j \models (\psi \wedge (\phi \vee \neg \phi)))$. Let s be an assignment and assume $\mathcal{M} \models \mathbf{T}_j \models (\phi \rightarrow \psi)[s]$ and $\mathcal{M} \models \mathbf{T}_j \models \phi[s]$. If \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}$ agree on $\mathbf{T}_j \models$, reason as in (2) above. If not, Lemma 49 says there is some j -stratifier \bullet^+ such that \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}^+$ agree on $\mathbf{T}_j \models$. By definition of $\mathcal{M}_{\hat{\mathbf{U}}}^+$, $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models (\mathbf{T}_j \models (\phi \rightarrow \psi))^+[s]$ and $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models (\mathbf{T}_j \models \phi)^+[s]$. Let $\alpha, \beta \in \epsilon_0 \cdot \omega$ be such that $(\mathbf{T}_j \models (\phi \rightarrow \psi))^+ \equiv \mathbf{T}_j^\alpha \models (\phi^+ \rightarrow \psi^+)$ and $(\mathbf{T}_j \models \phi)^+ \equiv \mathbf{T}_j^\beta \models \phi^+$. Then $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models \mathbf{T}_j^\alpha \models (\phi^+ \rightarrow \psi^+)[s]$ and $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models \mathbf{T}_j^\beta \models \phi^+[s]$. This means $U_j \cap \alpha \models (\phi^+ \rightarrow \psi^+)^s$ and $U_j \cap \beta \models (\phi^+)^s$. Since ϕ is a subformula of $\phi \rightarrow \psi$, it follows $\beta \leq \alpha$, thus $U_j \cap \alpha \models (\psi^+)^s$, and by tautology, $U_j \cap \alpha \models (\psi^+ \wedge (\phi^+ \vee \neg \phi^+))^s$. So $\mathcal{M}_{\hat{\mathbf{U}}} \models \mathbf{T}_j^\alpha \models (\psi^+ \wedge (\phi^+ \vee \neg \phi^+))[s]$. By Definition 43,

$$\mathbf{T}_j^\alpha \models (\psi^+ \wedge (\phi^+ \vee \neg \phi^+)) \equiv (\mathbf{T}_j \models (\psi \wedge (\phi \vee \neg \phi)))^+$$

(this is the reason for the $\phi \vee \neg \phi$ clause) and finally $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models \mathbf{T}_j \models (\psi \wedge (\phi \vee \neg \phi))[s]$.

(4) Similar to part 2 of Lemma 17.

(5) Similar to part 3 of Lemma 17.

(6) By (4), $\widetilde{\mathcal{M}} \models i$ -Assigned Strativalidity whenever $\widetilde{\mathcal{M}} \Rightarrow_i \mathcal{M}_{\hat{\mathbf{U}}}$. Let $\mathcal{M} \Rightarrow_j \mathcal{M}_{\hat{\mathbf{U}}}$, we will show $\mathcal{M} \models i$ -Validity. Let ϕ be a valid $\mathcal{L}_{\text{PA}}(\omega)$ -formula, s an assignment.

Case 1: \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}$ agree on $\mathbf{T}_i \models$. Let \bullet^+ be an i -stratifier. Since ϕ is valid, so is ϕ^+ (by Lemma 47), so $(\phi^+)^s \in U_i$ (since $[i$ -Assigned Strativalidity] $_i$ is part of line 6). Clearly $((\phi^+)^s)^- \equiv \phi^s$, so $\phi^s \in U_i^-$, thus $\mathcal{M}_{\hat{\mathbf{U}}} \models \mathbf{T}_i \models \phi[s]$, and so does \mathcal{M} .

Case 2: \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}$ disagree on $\mathbf{T}_i \models$. By Lemma 49 there is an i -stratifier \bullet^+ such that \mathcal{M} and $\mathcal{M}_{\hat{\mathbf{U}}}^+$ agree on $\mathbf{T}_i \models$. Let $\alpha \in \epsilon_0 \cdot \omega$ be such that $(\mathbf{T}_i \models \phi)^+ \equiv \mathbf{T}_i^\alpha \models \phi^+$. As in Case 1, $(\phi^+)^s \in U_i$. In fact by choice of α , $(\phi^+)^s \in U_i \cap \alpha$, so $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models \mathbf{T}_i^\alpha \models \phi^+[s]$, that is, $\mathcal{M}_{\hat{\mathbf{U}}} \models (\mathbf{T}_i \models \phi)^+[s]$. By Lemma 46, $\mathcal{M}_{\hat{\mathbf{U}}}^+ \models \mathbf{T}_i \models \phi[s]$.

(7–8) Similar to part 4 of Lemma 17. For 8, use part 3 of Definition 39.

(9) Similar to Lemma 10.

(10) Trivial.

(11–12) Similar to part 7 of Lemma 17. Use the fact that the stratifiers in Definition 48 are computable.

(13) Similar to part 8 of Lemma 17; use part 3 of Definition 39. \square

Definition 55. If $\mathbf{T}^0 = (T_i^0)_{i \in \omega}$ where each T_i^0 is an $\mathcal{L}_{\text{PA}}(\omega)$ -theory, we say \mathbf{T}^0 is *stratifiable-r.e.-generic* if there is some stratified-r.e.-generic family $\mathbf{T} = (T_i)_{i \in \omega}$ of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theories such that each $T_i^- = T_i^0$.

Theorem 56. Let $\mathbf{T}^0 = (T_i^0)_{i \in \omega}$ be any stratifiable-r.e.-generic family of $\mathcal{L}_{\text{PA}}(\omega)$ -theories. For every $i \in \omega$ and $n \in \mathbb{N}$, let $T_i(n)$ be the smallest $\mathbf{T}_i \models$ -closed $\mathcal{L}_{\text{PA}}(\omega)$ -theory containing the following axioms.

1. The axioms contained in T_i^0 .
2. Assigned Validity, i -Validity and i -Deduction.
3. $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi)$ whenever $j \preceq i$.
4. $\forall x(\mathbf{T}_j \models \phi \leftrightarrow \langle \overline{\Gamma \phi}, \bar{j}, x \rangle \in W_{\bar{n}})$ whenever $j \prec i$, $\text{FV}(\phi) \subseteq \{x\}$.

Let each $\mathbf{T}(n) = (T_i(n))_{i \in \omega}$. There is some $n \in \mathbb{N}$ such that $\mathbf{T}(n)$ is true.

Proof. By the S - m - n Theorem, there is a total computable $f : \mathbb{N} \rightarrow \mathbb{N}$ such that $\forall n \in \mathbb{N}$,

$$W_{f(n)} = \{ \langle \overline{\Gamma \phi}, j, m \rangle \in \mathbb{N} : \phi \text{ is an } \mathcal{L}_{\text{PA}}(\omega)\text{-formula, } \text{FV}(\phi) \subseteq \{x\}, \text{ and } T_j(n) \models \phi(x|\bar{m}) \}.$$

By the Recursion Theorem, there is an $n \in \mathbb{N}$ such that $W_n = W_{f(n)}$. Fix this n for the rest of the proof and write \mathbf{T} for $\mathbf{T}(n)$, T_i for $T_i(n)$.

Since \mathbf{T}^0 is stratifiable-r.e.-generic, there is a stratified-r.e.-generic family $\mathbf{V} = (V_i)_{i \in \omega}$ of $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theories such that each $V_i^- = T_i^0$. For every $i \in \mathbb{N}$, let U_i be the smallest i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -theory such that the following hold.

1. U_i contains V_i .
2. U_i contains i -Assigned Strativalidity, i -Strativalidity, i -Stratideduction and i -Collapse.
3. U_i contains $\text{ucl}(\mathbf{T}_i^{\alpha} \models \phi \rightarrow \phi)$ whenever $\mathbf{T}_i^{\alpha} \models \phi$ is i -stratified.
4. U_i contains $\text{ucl}(\mathbf{T}_j \models \phi \rightarrow \phi^+)$ for every $\mathcal{L}_{\text{PA}}(\omega)$ -formula ϕ , $j \prec i$, and i -stratifier \bullet^+ .
5. U_i contains $\forall x(\mathbf{T}_j \models \phi \leftrightarrow \langle \overline{\Gamma \phi}, \bar{j}, x \rangle \in W_{\bar{n}})$ whenever $j \prec i$, $\text{FV}(\phi) \subseteq \{x\}$ and ϕ is an $\mathcal{L}_{\text{PA}}(\omega)$ -formula.
6. Whenever $\phi \in U_i$ and $\mathbf{T}_i^{\alpha} \models \phi$ is i -stratified, $\mathbf{T}_i^{\alpha} \models \phi \in U_i$.

Let $\mathbf{U} = (U_i)_{i \in \omega}$. Observe that \mathbf{U} is stratifiable r.e. (to see \mathbf{U} is uniform, use Lemma 44, to see \mathbf{U} is r.e., use Theorem 20 part 1); $\mathbf{U} \supseteq \mathbf{V}$; and for each $i \in \omega$, $U_i^- = T_i$.

Let $\mathbf{S} = (S_i)_{i \in \mathcal{I}} = \text{Str}(\mathbf{U})$. By definition this means that for all $i \in \omega$ and $\alpha \in \epsilon_0 \cdot \omega$,

$$S_i = U_i^- = T_i \text{ and } S_{(\alpha, i)} = U_i \cap \alpha.$$

In order to show $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}$ and thus prove the theorem, we will show $\mathcal{M}_{\mathbf{S}} \models \mathbf{S}$, which is more than sufficient, since $\mathcal{M}_{\mathbf{S}}$ and $\mathcal{M}_{\mathbf{T}}$ agree on $\mathcal{L}_{\text{PA}}(\omega)$. But for sake of a stronger induction hypothesis, we will prove more. We will prove that for every $i \in \omega$, every $j \preceq i$, every $\mathcal{M} \Rightarrow_i \mathcal{M}_{\mathbf{S}}$, and every $\alpha \in \epsilon_0 \cdot \omega$, $\mathcal{M} \models S_j \cup S_{(\alpha, j)}$.

Fix $i \in \omega$. Since \prec is well-founded, we may assume the following:

(*) For every $k \preceq j \prec i$, every $\mathcal{M} \Rightarrow_j \mathcal{M}_{\mathbf{S}}$, and every $\alpha \in \epsilon_0 \cdot \omega$, $\mathcal{M} \models S_k \cup S_{(\alpha, k)}$.

Fix $\mathcal{M} \Rightarrow_i \mathcal{M}_{\mathbf{S}}$. For all $j \prec i$, Lemma 49 says $\mathcal{M} \Rightarrow_j \mathcal{M}_{\mathbf{S}}$ and therefore by (*) we already have $\mathcal{M} \models S_j \cup S_{(\alpha, j)}$. It remains to show $\forall \alpha \in \epsilon_0 \cdot \omega$, $\mathcal{M} \models S_i \cup S_{(\alpha, i)}$.

Claim 1: $\forall \alpha \in \epsilon_0 \cdot \omega$, $\mathcal{M} \models S_{(\alpha, i)}$.

By induction on α . Let $\sigma \in S_{(\alpha, i)}$. This means $\sigma \in U_i \cap \alpha$.

Case 1: $\sigma \in V_i$. Then $\mathcal{M} \models \sigma$ because \mathbf{V} is stratified-r.e.-generic and $\mathbf{U} \supseteq \mathbf{V}$ is stratifiable r.e.

Case 2: σ is an instance of i -Assigned Strativalidity, i -Strativalidity, or i -Stratideduction. Then $\mathcal{M} \models \sigma$ by Lemma 54.

Case 3: σ is $\text{ucl}(\mathbf{T}_i^{\alpha_0} \models \phi \rightarrow \phi)$ for some i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ such that $\mathbf{T}_i^{\alpha_0} \models \phi$ is i -stratified. Since $\sigma \in U_i \cap \alpha$, this forces $\alpha_0 < \alpha$. Let s be an assignment and assume $\mathcal{M} \models \mathbf{T}_i^{\alpha_0} \models \phi[s]$, then:

$$\begin{aligned} \mathcal{M} &\models \mathbf{T}_i^{\alpha_0} \models \phi[s] && \text{(Assumption)} \\ \mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_i^{\alpha_0} \models \phi[s] && (\mathcal{M} \text{ and } \mathcal{M}_{\mathbf{S}} \text{ agree on } \mathbf{T}_i^{\alpha_0} \models \text{ by Lemma 49}) \\ S_{(\alpha_0, i)} &\models \phi^s && \text{(Definition of } \mathcal{M}_{\mathbf{S}}) \\ \mathcal{M} &\models \phi^s && \text{(By induction, } \mathcal{M} \models S_{(\alpha_0, i)}) \\ \mathcal{M} &\models \phi[s]. && \text{(By Lemma 51)} \end{aligned}$$

Case 4: σ is $\text{ucl}(\mathbf{T}_j \models \phi \leftrightarrow \phi^+)$ for some $\mathcal{L}_{\text{PA}}(\omega)$ -formula ϕ , $j \prec i$, and i -stratifier \bullet^+ . By Lemma 44 we may assume \bullet^+ is computable. Let s be an assignment and assume $\mathcal{M} \models \mathbf{T}_j \models \phi[s]$, then:

$$\begin{aligned}
\mathcal{M} &\models \mathbf{T}_j \models \phi[s] && \text{(Assumption)} \\
\mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_j \models \phi[s] && \text{(Since } j \prec i, \mathcal{M} \text{ and } \mathcal{M}_{\mathbf{S}} \text{ agree on } \mathbf{T}_j \models \text{ by Lemma 49)} \\
T_j &\models \phi^s && \text{(Definition of } \mathcal{M}_{\mathbf{S}}) \\
\mathcal{M}^+ &\models \phi^s && \text{(By Lemma 49, } \mathcal{M}^+ \Rightarrow_j \mathcal{M}_{\mathbf{S}}, \text{ so by } (*), \mathcal{M}^+ \models T_j) \\
\mathcal{M} &\models (\phi^s)^+ && \text{(Lemma 46)} \\
\mathcal{M} &\models (\phi^+)^s && \text{(Clearly } (\phi^s)^+ \equiv (\phi^+)^s) \\
\mathcal{M} &\models \phi^+[s]. && \text{(Lemma 51)}
\end{aligned}$$

Case 5: σ is $\forall x(\mathbf{T}_j \models \phi \leftrightarrow \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, \bar{j}, x \rangle \in W_{\bar{n}})$ for some $\mathcal{L}_{\text{PA}}(\omega)$ -formula ϕ with $\text{FV}(\phi) \subseteq \{x\}$ and $j \prec i$. Let s be any assignment, say $s(x) = m$. The following biconditionals are equivalent:

$$\begin{aligned}
\mathcal{M} &\models \mathbf{T}_j \models \phi \leftrightarrow \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, \bar{j}, x \rangle \in W_{\bar{n}}[s] \\
\mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_j \models \phi \leftrightarrow \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, \bar{j}, x \rangle \in W_{\bar{n}}[s] && (\mathcal{M} \text{ and } \mathcal{M}_{\mathbf{S}} \text{ agree on the symbols in question)} \\
\mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_j \models \phi[s] \text{ iff } \mathcal{M}_{\mathbf{S}} \models \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, \bar{j}, \bar{m} \rangle \in W_{\bar{n}} && \text{(Lemma 51)} \\
\mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_j \models \phi[s] \text{ iff } \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, j, m \rangle \in W_n && (\mathcal{M}_{\mathbf{S}} \text{ has standard first-order part)} \\
T_j &\models \phi^s \text{ iff } \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, j, m \rangle \in W_n && \text{(Definition of } \mathcal{M}_{\mathbf{S}}) \\
T_j &\models \phi(x|\bar{m}) \text{ iff } \langle \bar{\Gamma}\phi^{\bar{\Gamma}}, j, m \rangle \in W_n. && \text{(Since } \text{FV}(\phi) \subseteq \{x\})
\end{aligned}$$

The latter is true by definition of n .

Case 6: σ is an instance $\mathbf{T}_i^{\beta} \models \phi \leftrightarrow \mathbf{T}_i^{\gamma} \models \phi$ of i -Collapse (so $\beta \leq_1 \gamma$ and $\mathbf{T}_i^{\beta} \models \phi \leftrightarrow \mathbf{T}_i^{\gamma} \models \phi$ is i -stratified). Let s be an assignment, since \mathcal{M} and $\mathcal{M}_{\mathbf{S}}$ agree on $\mathbf{T}_i^{\beta} \models$ and $\mathbf{T}_i^{\gamma} \models$, we need only show $\mathcal{M}_{\mathbf{S}} \models \mathbf{T}_i^{\beta} \models \phi \leftrightarrow \mathbf{T}_i^{\gamma} \models \phi[s]$. In other words we must show $U_i \cap \beta \models \phi^s$ if and only if $U_i \cap \gamma \models \phi^s$. This is by Theorem 31.

Case 7: σ is $\mathbf{T}_i^{\alpha_0} \models \phi$ for some $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ such that $\mathbf{T}_i^{\alpha_0} \models \phi$ is i -stratified and $\phi \in U_i$. Since $\mathbf{T}_i^{\alpha_0} \models \phi$ is i -stratified, $\text{On}(\phi) \subseteq \alpha_0$, so $\phi \in U_i \cap \alpha_0$. Thus $\mathcal{M}_{\mathbf{S}} \models \mathbf{T}_i^{\alpha_0} \models \phi$, so $\mathcal{M} \models \mathbf{T}_i^{\alpha_0} \models \phi$ since \mathcal{M} and $\mathcal{M}_{\mathbf{S}}$ agree on $\mathbf{T}_i^{\alpha_0} \models$.

Cases 1–7 establish $\mathcal{M} \models S_{(\alpha, i)}$. By arbitrariness of α , Claim 1 is proved.

Claim 2: For any assignment s and any very i -stratified $\mathcal{L}_{\text{PA}}(\mathcal{I})$ -formula ϕ , $\mathcal{M} \models \phi[s]$ if and only if $\mathcal{M} \models \phi^-[s]$.

By induction on ϕ . The only interesting cases are the following.

Case 1: ϕ is $\mathbf{T}_j \models \psi$ for some j . Then $\phi^- \equiv \phi$ and the claim is trivial.

Case 2: ϕ is $\mathbf{T}_j^{\alpha} \models \psi$ for some $j \neq i$. Impossible, this is not i -stratified.

Case 3: ϕ is $\mathbf{T}_i^{\alpha} \models \psi$. The following are equivalent:

$$\begin{aligned}
\mathcal{M} &\models (\mathbf{T}_i^{\alpha} \models \psi)^-[s] \\
\mathcal{M}_{\mathbf{S}} &\models (\mathbf{T}_i^{\alpha} \models \psi)^-[s] && (\mathcal{M} \text{ and } \mathcal{M}_{\mathbf{S}} \text{ agree on } \mathbf{T}_i \models) \\
\mathcal{M}_{\mathbf{S}} &\models \mathbf{T}_i^{\alpha} \models \psi[s] && \text{(Theorem 42)} \\
\mathcal{M} &\models \mathbf{T}_i^{\alpha} \models \psi[s]. && (\mathcal{M} \text{ and } \mathcal{M}_{\mathbf{S}} \text{ agree on } \mathbf{T}_i^{\alpha} \models)
\end{aligned}$$

Claim 3: $\mathcal{M} \models S_i$.

For any $\sigma \in S_i$, there is some $\tau \in U_i$ such that $\tau^- \equiv \sigma$; since U_i is i -uniform, we may take τ to be very i -stratified (Lemma 34). By Claim 1, $\mathcal{M} \models U_i$, so $\mathcal{M} \models \tau$. By Claim 2, $\mathcal{M} \models \sigma$. \square

This satisfies the second promise from the introduction: given a well-founded r.e. partial order \prec on ω , we have exhibited true theories $(T_i)_{i \in \omega}$ such that T_i expresses a Gödel number of T_j ($j \prec i$) and the truth of T_j ($j \preceq i$). These theories can further be taken so that T_i expresses the fact that T_j has some Gödel number (all i, j), by Lemma 54 parts 11–12.

8 Well-Foundation and Ill-Foundation

The following is a variation on Kleene's \mathcal{O} .

Definition 57. Simultaneously define $\mathcal{O} \subseteq \mathbb{N}$ and $|\bullet| : \mathcal{O} \rightarrow \text{Ord}$ so that $\mathcal{O} \subseteq \mathbb{N}$ is the smallest set such that:

1. $0 \in \mathcal{O}$ (it represents the ordinal $|0| = 0$).
2. $\forall n \in \mathcal{O}, 2^n \in \mathcal{O}$ (it represents the ordinal $|2^n| = |n| + 1$).
3. If φ_e (the e th partial recursive function) is total and $\text{range}(\varphi_e) \subseteq \mathcal{O}$, then $3 \cdot 5^e \in \mathcal{O}$ (it represents the ordinal $|3 \cdot 5^e| = \sup\{|\varphi_e(0)|, |\varphi_e(1)|, \dots\}$).

To avoid technical complications, we have differed from the usual Kleene's \mathcal{O} in the following way: in the usual definition, in order for $3 \cdot 5^e$ to lie in \mathcal{O} , it is also required that $|\varphi_e(0)| < |\varphi_e(1)| < \dots$.

Definition 58. $\mathcal{L}_{\text{PA}}^{\mathcal{O}}$ is the language of Peano arithmetic extended by a unary predicate \mathcal{O} . The following notions are defined by analogy with Section 2:

1. For any assignment s and $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -formula ϕ with $\text{FV}(\phi) = \{x_1, \dots, x_n\}$, $\phi^s \equiv \phi(x_1 | \overline{s(x_1)}) \cdots (x_n | \overline{s(x_n)})$.
2. If $\mathbf{T} = (T_i)_{i \in I}$ is an I -indexed family of $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -theories, the *intended structure* for \mathbf{T} is the $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -structure $\mathcal{M}_{\mathbf{T}}$ with universe \mathbb{N} , interpreting symbols of PA as usual and interpreting \mathcal{O} as \mathcal{O} , and interpreting $\mathbf{T}_i \models (i \in I)$ as in Definition 7. For any $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -structure \mathcal{N} , we write $\mathcal{N} \models \mathbf{T}$ if $\forall i \in I, \mathcal{N} \models T_i$. We say \mathbf{T} is *true* if $\mathcal{M}_{\mathbf{T}} \models \mathbf{T}$.

Definition 59. If I is an index set and $\mathbf{T} = (T_i)_{i \in I}$ is a family of $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -theories, then for any $i \in I$ such that $\mathcal{M}_{\mathbf{T}} \models T_i$, we define the ordinal $\|\mathbf{T}_i\| = \sup\{|m| + 1 : T_i \models \mathcal{O}(\overline{m})\}$.

The above definition makes sense: since $\mathcal{M}_{\mathbf{T}} \models T_i$ and $\mathcal{O}^{\mathcal{M}_{\mathbf{T}}} = \mathcal{O}$, the supremands are defined.

Definition 60. The *basic axioms* of \mathcal{O} are the following $\mathcal{L}_{\text{PA}}^{\mathcal{O}}$ -axioms.

1. $\mathcal{O}(0)$.
2. $\mathcal{O}(\overline{n}) \rightarrow \mathcal{O}(\overline{2^n})$, for every $n \in \mathbb{N}$.
3. $\forall x(\varphi_{\overline{n}}(x) \downarrow \ \& \ \mathcal{O}(\varphi_{\overline{n}}(x))) \rightarrow \mathcal{O}(\overline{3 \cdot 5^n})$, for every $n \in \mathbb{N}$.

We have written the last two lines using infinite schemata to strengthen the following result.

Theorem 61. Let I be an index set, \prec a binary relation on I . Suppose $\mathbf{T} = (T_i)_{i \in I}$ is a family of $\mathcal{L}_{\text{PA}}^{\mathcal{O}}(I)$ -theories with the following properties:

1. $\forall i \in I, T_i$ contains the axioms of Peano arithmetic.
2. $\forall i \in I, T_i$ contains the basic axioms of \mathcal{O} .
3. $\forall i \in I, \forall j \prec i, \exists n \in \mathbb{N}$ such that $T_i \models \forall x(\mathbf{T}_j \models \mathcal{O}(x) \leftrightarrow x \in W_{\overline{n}})$.
4. $\forall i \in I, \forall j \prec i, T_i \models \forall x(\mathbf{T}_j \models \mathcal{O}(x) \rightarrow \mathcal{O}(x))$.

If $\mathcal{M}_{\mathbf{T}} \models T_i \cup T_j$ (in particular if \mathbf{T} is true) and $j \prec i$, then $\|\mathbf{T}_j\| < \|\mathbf{T}_i\|$.

Proof. Assume $\mathcal{M}_{\mathbf{T}} \models T_i \cup T_j$ and $j \prec i$. By hypothesis there is some $n \in \mathbb{N}$ such that $T_i \models \forall x(\mathbf{T}_j \models \mathcal{O}(x) \leftrightarrow x \in W_{\overline{n}})$ and $T_i \models \forall x(\mathbf{T}_j \models \mathcal{O}(x) \rightarrow \mathcal{O}(x))$. From these, $T_i \models \forall x(x \in W_{\overline{n}} \rightarrow \mathcal{O}(x))$.

Since $\mathcal{M}_{\mathbf{T}} \models T_i$, in particular $\mathcal{M}_{\mathbf{T}} \models \forall x(\mathbf{T}_j \models \mathcal{O}(x) \leftrightarrow x \in W_{\overline{n}})$. This means $W_n = \{m \in \mathbb{N} : T_j \models \mathcal{O}(\overline{m})\}$. Since T_j includes the axiom $\mathcal{O}(0)$, $W_n \neq \emptyset$.

Since $W_n \neq \emptyset$, by computability theory there is some $k \in \mathbb{N}$ such that

$$\text{PA} \models (\text{domain}(\varphi_{\overline{k}}) = \mathbb{N}) \wedge (\text{range}(\varphi_{\overline{k}}) = W_{\overline{n}}).$$

Since T_i includes PA, T_i also implies as much. Combined with $T_i \models \forall x(x \in W_{\bar{n}} \rightarrow \mathcal{O}(x))$, it follows that $T_i \models \forall x(\varphi_{\bar{k}}(x) \downarrow \ \& \ \mathcal{O}(\varphi_{\bar{k}}(x)))$. Since T_i contains the basic axiom $\forall x(\varphi_{\bar{k}}(x) \downarrow \ \& \ \mathcal{O}(\varphi_{\bar{k}}(x))) \rightarrow \mathcal{O}(\overline{3 \cdot 5^k})$, $T_i \models \mathcal{O}(\overline{3 \cdot 5^k})$.

To finish the proof, calculate

$$\begin{aligned}
\|T_j\| &= \sup\{|m| + 1 : T_j \models \mathcal{O}(\overline{m})\} \\
&= \sup\{|m| : T_j \models \mathcal{O}(\overline{m})\} && \text{(Since } T_j \text{ contains } \mathcal{O}(\overline{n}) \rightarrow \mathcal{O}(\overline{2^n}) \text{ for all } n \in \mathbb{N}) \\
&= \sup\{|m| : m \in W_n\} && \text{(Since } W_n = \{m \in \mathbb{N} : T_j \models \mathcal{O}(\overline{m})\}) \\
&= \sup\{|\varphi_k(0)|, |\varphi_k(1)|, \dots\} && \text{(By choice of } k) \\
&= |3 \cdot 5^k| && \text{(Definition 57)} \\
&< \sup\{|m| + 1 : T_i \models \mathcal{O}(\overline{m})\} && \text{(Since } T_i \models \mathcal{O}(\overline{3 \cdot 5^k})) \\
&= \|T_i\|.
\end{aligned}$$

□

Corollary 62. (Well-Foundedness of True Self-Referential Theories) Let I, \mathbf{T}, \prec be as in Theorem 61. If \mathbf{T} is true then \prec is well founded, by which we mean there is no infinite descending sequence $i_0 \succ i_1 \succ \dots$.

In particular Corollary 62 says that if I, \mathbf{T}, \prec are as in Theorem 61 and \mathbf{T} is true then \prec is strict: there is no i with $i \prec i$. This gives a new form (under the additional new assumption of containing/knowing basic rudiments of computable ordinals) of the Lucas–Penrose–Reinhardt argument that a truthful theory (or machine) cannot state (or know) its own truth and its own Gödel number.

We could remove Peano arithmetic from Theorem 61 if we further departed from Kleene and changed line 3 of Definition 57 to read:

3. If $W_e \subseteq \mathcal{O}$, then $3 \cdot 5^e \in \mathcal{O}$ (and $|3 \cdot 5^e| = \sup\{|n| : n \in W_e\}$, or $|3 \cdot 5^e| = 0$ if $W_e = \emptyset$)

(and altered Definition 60 accordingly). The previous paragraph would still stand, in fact giving a version of the Lucas–Penrose–Reinhardt argument in which the theory (machine) is not required to contain (know) arithmetic.

We close the paper by showing that Corollary 62 fails without \mathcal{O} . Let WF be the set of all r.e. well-founded partial orders on ω and let Tr be the set of all true \mathcal{L}_{PA} -sentences. It is well-known that WF is computability theoretically Π_1^1 -complete and Tr is Δ_1^1 , so WF cannot be defined in $\mathcal{L}_{\text{PA}} \cup \{\text{Tr}\}$.

Theorem 63. (Ill-Foundedness of True Self-Referential Theories)

1. For any closed-r.e.-generic $\mathbf{T}^0 = (T_i^0)_{i \in \omega}$, there is an r.e., ill-founded partial order \prec on ω and an $n \in \mathbb{N}$ such that $\mathbf{T}(n)$ is true, where $\mathbf{T}(n)$ is as in Theorem 18.
2. For any stratifiable-r.e.-generic $\mathbf{T}^0 = (T_i^0)_{i \in \omega}$, there is an r.e., ill-founded partial order \prec on ω and an $n \in \mathbb{N}$ such that $\mathbf{T}(n)$ is true, where $\mathbf{T}(n)$ is as in Theorem 56.

Proof. We prove (1), (2) is similar. Assume $\neg(1)$. If \prec is any r.e. partial order on ω , then, combining $\neg(1)$ with Theorem 18, \prec is well founded if and only if the conclusion of Theorem 18 holds for \prec . Thus it is possible to define WF in $\mathcal{L}_{\text{PA}} \cup \{\text{Tr}\}$. Absurd. □

References

- [1] Alexander, S. (2013). The Theory of Several Knowing Machines. Doctoral dissertation, the Ohio State University.
- [2] Alexander, S. (preprint). A machine that knows its own code. To appear in *Studia Logica*.
arXiv: <http://arxiv.org/abs/1305.6080>
- [3] Benacerraf, P. (1967). God, the Devil, and Gödel. *The Monist*, **51**, 9–32.

- [4] Carlson, T.J. (1999). Ordinal arithmetic and Σ_1 -elementarity. *Archive for Mathematical Logic*, **38**, 449–460.
- [5] Carlson, T.J. (2000). Knowledge, machines, and the consistency of Reinhardt’s strong mechanistic thesis. *Annals of Pure and Applied Logic*, **105**, 51–82.
- [6] Carlson, T.J. (2001). Elementary patterns of resemblance. *Annals of Pure and Applied Logic*, **108**, 19–77.
- [7] Lucas, J.R. (1961). Minds, machines, and Gödel. *Philosophy*, **36**, 112–127.
- [8] Penrose, R. (1989). *The Emperor’s New Mind: Concerning Computers, Minds, and the Laws of Physics*. Oxford University Press.
- [9] Putnam, H. (2006). After Gödel. *Logic Journal of the IGPL*, **14**, 745–754.
- [10] Reinhardt, W. (1985). Absolute versions of incompleteness theorems. *Nous*, **19**, 317–346.
- [11] Shapiro, S. (1985). Epistemic and Intuitionistic Arithmetic. In: S. Shapiro (ed.), *Intensional Mathematics* (North-Holland, Amsterdam), pp. 11–46.