

THE ARITHMETICS OF A THEORY

ALBERT VISSER

ABSTRACT. In this paper we study the interpretations of a weak arithmetic, like Buss' theory S_2^1 , in a given theory U . We call these interpretations *the arithmetics of U* . We develop the basics of the structure of the arithmetics of U . We study the provability logic(s) of U from the standpoint of the framework of the arithmetics of U . Finally, we provide a deeper study of the arithmetics of a finitely axiomatized sequential theory.

1. Introduction	2
About this Paper	2
Prerequisites	3
2. Basics	3
2.1. Theories	3
2.2. Translations and Interpretations	4
2.3. i-morphisms	5
2.4. Sequential Theories	6
2.5. Complexity and Satisfaction	7
3. The Arithmetics of a Theory	7
4. Introduction to Provability Logic	12
5. Solovay's Theorem	15
5.1. The Guaspari-Solovay System R^-	15
5.2. Arithmetical Interpretations of R^-	16
5.3. The Basic Proof	17
5.4. Application to Arithmetic	20
5.5. All Arithmetics of a Theory	20
5.6. Theories with a Σ_1 -sound Arithmetic	21
6. Deep Arithmetics	22
7. An Example	28
References	31
Appendix A. More Details on the Basics	33
A.1. Translations and Interpretations	33
A.2. i-morphisms	35
A.3. Piecewise Interpretations	36
A.4. Parameters	37
A.5. Complexity Measures	38
Appendix B. On Faithful Interpretability	39

Contents

Date: March 8, 2012.

I am grateful to Lev Beklemishev, Emil Jeřábek and Vincent van Oostrom for helpful discussions.

1. INTRODUCTION

In this paper, we propose and expose a particular way of viewing theories. We look at theories as a class of interpretations of a given weak arithmetical theory. Consider a theory U . We view the interpretations of the given weak arithmetical theory in U as ‘occurrences’ of that given theory in U . Thus, U appears as a class of copies of the given weak theory. If we consider a model \mathcal{M} of U , the versions of the weak theory sitting inside U take the form of the set of internal models of the weak theory in \mathcal{M} .

We will call an interpretation N of the given weak theory in U *an arithmetic of U* . The arithmetics of U have a natural ordering, the (definable) initial embedding ordering \preceq . We study basic facts concerning the arithmetics of U and the ordering \preceq in Section 3.

From the perspective of theories as containers of (possibly) lots of arithmetics, we study the provability logics of theories. We fully characterize the propositional modal principles for provability that hold in all arithmetics in any theory U . The only assumption being a constraint on the complexity of the set of axioms of U . The comparatively easy success of this characterization contrasts with the remaining great open questions of provability logic concerning the provability logics of theories like S_2^1 or $I\Delta_0 + \Omega_1$.

Section 4 briefly reviews some basic ideas concerning provability logic.

In Section 5, we study Solovay’s Theorem in various settings. In Subsections 5.1, 5.2, 5.3, and 5.4, we present a proof of Solovay’s Completeness Theorem for Löb’s Logic via a wonderful version of the proof given by Dick de Jongh, Marc Jumelet and Franco Montagna. The main part of the proof is itself formulated in a richer modal logic which was formulated and studied by David Guaspari and Robert Solovay. The advantage of the de Jongh, Jumelet, Montagna proof is that it allows us to see clearly what arithmetical principles are involved in Solovay’s proof. In Subsection 5.5 we prove our characterization of the provability logic of all arithmetics of a given theory. In Subsection 5.6, we give a sufficient condition for when the provability logic of all theories is assumed at a single arithmetic N in U .

In Section 7, we provide an example of a theory U where the provability logic of U is not assumed at any arithmetic N in U .

In Section 6 we study the wondrous world of the arithmetics of a finitely axiomatized sequential theory U . In the sequential case we have many extra properties of our structure of arithmetics to work with. In this section we strengthen certain results due to Harvey Friedman and, indepently, to Jan Krajíček. We use the methods of Section 6 to construct the example of Section 7.

The reader interested in Provability Logic could very well choose to read Sections 2,3,4,5. The reader who is interested in the fine structure of the arithmetics of a theory could study Sections 2,3,6,(7). More details on the basics are provided in Appendix A.

About this Paper. The present paper is, in a sense, a remake of my paper [Vis91a]. It is the result of my reflection on what the earlier paper is saying.

We strengthen the results of that paper by presenting them in a better framework and we add new results relevant to the framework.

Prerequisites. We will presuppose some knowledge of weak arithmetics. See e.g. [Bus86] or [HP91, chapter V]. Some basic knowledge of provability logic will help to understand the paper. At present there are many expositions: [Smo85b], [BS91], [Boo93], [Lin96], [JdJ98], [Šve00], [AB04]. The most comprehensive source concerning the provability logic of weak theories is [Ver93].

2. BASICS

In this section we sketch the framework in which our discussion will take place. One problem of sections with basic notions and facts is that they are so long and so boring that the reader gets stuck in them and never arrives at the real stuff. So what I did is to make the present section rather sketchy. At the end of the paper, in Appendix A more details are provided. Regrettably, even without the details this section is rather long, so the reader is advised to go over it lightly and come back to it or Appendix A when needed.

2.1. Theories. Theories are, in this paper, one-sorted theories of first-order predicate logic, that have a finite signature and that are axiomatized by an axiom set that is represented by a Δ_1^b -formula.¹

Remark 2.1. The demand for Δ_1^b -axiomatization seems to be rather restrictive. However, it seems to me, that every real-life theory is given by an axiomatization that is p-time decidable.

Because in S_2^1 , we have the Σ_1^b -replacement axiom, we can relax our demand to the consideration of theories which are Σ_1^b -axiomatized. In this case, the witnesses of $\Box_U A$ would not be really a code of a proof but a somewhat modified object.

Note that, by a version of Craig's trick, every RE theory in extension can be given a Δ_1^b -axiomatization. Of course, a weak theory will not be able to see that both axiomatizations prove the same theorems, so for the eyes of the weak theory the *craigified* theory will be a different theory. We need a theory like EA, aka $I\Delta_0 + \exp$, plus Σ_1 -replacement to make this construction work in a verifiable way. \square

The formula specifying the axiom set is part of the data for the theory. Thus, we treat theories *intensionally* and not as mere sets of theorems. We will explain why this is important for our purposes in Section 4.

We say that a theory is *finitely axiomatized* if its axiomatization has the form $\bigvee_{i < n} x = \ulcorner A_i \urcorner$.² Note that S_2^1 may prove that a theory has an axiom-set of, say, less than two axioms, without being able to prove the equivalence of the formula defining the axiom set with any formula of the prescribed form.

¹See [Bus86] or [HP91] for an explanation of the relevant formula classes.

²The function $\ulcorner \cdot \urcorner$ sends a syntactical object to its Gödel number. The function (\cdot) sends a number to its numeral. We will employ efficient numerals that reflect binary notation.

We take identity to be a logical constant. Our official signatures are relational, however, via the term-unwinding algorithm, we can also accommodate signatures with functions.

2.2. Translations and Interpretations. The notion of interpretation that we will employ in this paper will be *m-dimensional interpretation without parameters*. There are two extensions of this notion: we can consider piecewise interpretations and we can add parameters. We refrain from considering piecewise interpretations. We explain why in Subsection A.3 of Appendix A. We sketch a few basic ingredients of adding parameters in Subsection A.4 of Appendix A. We explain why, in the sequential case, addition of parameters makes no difference for the provability logic of all arithmetics of a given theory in Remark 3.9.

Consider two signatures Σ and Θ . An m -dimensional translation $\tau : \Sigma \rightarrow \Theta$ is a quadruple $\langle \Sigma, \delta, \mathcal{F}, \Theta \rangle$, where $\delta(v_0, \dots, v_{m-1})$ is a Θ -formula and where for any n -ary predicate P of Σ , $\mathcal{F}(P)$ is a formula $A(\vec{v}_0, \dots, \vec{v}_{n-1})$ in the language of signature Θ , where $\vec{v}_i = v_{i0}, \dots, v_{i(m-1)}$. Both in the case of δ and A all free variables are among the variables shown. Moreover, if $i \neq j$ or $k \neq \ell$, then v_{ik} is syntactically different from $v_{j\ell}$.

We demand that we have $\vdash \mathcal{F}(P)(\vec{v}_0, \dots, \vec{v}_{n-1}) \rightarrow \bigwedge_{i < n} \delta(\vec{v}_i)$. Here \vdash is provability in predicate logic. This demand is inessential, but it is convenient to have.

We define B^τ as follows:

- $(P(x_0, \dots, x_{n-1}))^\tau := \mathcal{F}(P)(\vec{x}_0, \dots, \vec{x}_{n-1})$.
- $(\cdot)^\tau$ commutes with the propositional connectives.
- $(\forall x A)^\tau := \forall \vec{x} (\delta(\vec{x}) \rightarrow A^\tau)$.
- $(\exists x A)^\tau := \exists \vec{x} (\delta(\vec{x}) \wedge A^\tau)$.

We allow identity to be translated to a formula that is not identity. We can define the identity translation id_Σ on Σ , the composition $\rho \circ \tau$ of translations τ and ρ , and the disjunctive translation $\tau \langle A \rangle \rho$, that is τ if A and ρ if $\neg A$. We refer the reader to Appendix A for details.

A translation relates signatures; an interpretation relates theories. An interpretation $K : U \rightarrow V$ is a triple $\langle U, \tau, V \rangle$, where U and V are theories and $\tau : \Sigma_U \rightarrow \Sigma_V$. We demand: for all axioms A of U , we have $V \vdash A^\tau$.

In the context of the formalization of interpretability, we have to distinguish between *axioms-interpretability*, which is the notion we just introduced and *theorems-interpretability*, where we demand that: for all theorems A of U , we have $V \vdash A^\tau$. In the real world these notions are equivalent, but we need a principle like Σ_1 -collection to prove that, so, for example Buss' theory S_2^1 does not 'know' this equivalence. See [Vis91b] for more information about this matter.

Here are some further definitions and conventions.

- Suppose $K : U \rightarrow V$. We often write A^K for A^{τ_K} , in the context of a theory W that extends V .

- We write \overline{U} for the set of theorems of U . Suppose $K : U \rightarrow V$. We write $\overline{K} := \{A \mid V \vdash A^K\}$. We note that $\overline{U} \subseteq \overline{K}$. If $\overline{K} = \overline{U}$, we will say that K is *faithful*.
- $\text{ID}_U : U \rightarrow U$ is the interpretation $\langle U, \text{id}_{\Sigma_U}, U \rangle$.
- Suppose $\overline{U} \subseteq \overline{V}$. Then, $\mathcal{E}_{UV} : U \rightarrow V$ is $\langle U, \text{id}_{\Sigma_U}, V \rangle$.
- Suppose $K : U \rightarrow V$ and $M : V \rightarrow W$. Then, $KM := M \circ K : U \rightarrow W$ is $\langle U, \tau_M \circ \tau_K, W \rangle$.
- Suppose $K : W \rightarrow U$ and $U \subseteq V$. We write $K \uparrow V$ for $\mathcal{E}_{UV} \circ K$.
- Suppose $M : V \rightarrow Z$ and $U \subseteq V$. We write $U \downarrow M$ for $M \circ \mathcal{E}_{UV}$.
- Suppose $K : U \rightarrow (V + A)$ and $M : U \rightarrow (V + \neg A)$. Then $K\langle A \rangle M : U \rightarrow V$ is the interpretation $\langle U, \tau_K\langle A \rangle\tau_M, V \rangle$. In an appropriate category $K\langle A \rangle M$ is a special case of a product.

The notation $K : U \rightarrow V$ is inspired by the idea of interpretations as arrows in a category. There is also an intuition of interpretability as a generalization of provability. The traditional notations and notions associated to this intuition are:

- $K : U \triangleleft V$ stands for $K : U \rightarrow V$.
- $K : V \triangleright U$ stands for $K : U \rightarrow V$.
- $U \triangleleft V$ stands for $\exists K K : U \triangleleft V$. We say: U is *interpretable* in V .
- $V \triangleright U$ stands for $\exists K K : V \triangleright U$. We say: V *interprets* U .
- $U \triangleleft_{\text{loc}} V$ means: all finitely axiomatized subtheories U_0 of U are interpretable in V . We say that U is *locally interpretable* in V .
- $U \triangleleft_{\text{mod}} V$ means that, for every $\mathcal{M} \models V$, there is a translation τ such that $\tau(\mathcal{M}) \models U$. We say that U is *model-interpretable* in V .

2.3. i-morphisms. Consider an interpretation $K : U \rightarrow V$. We can view this interpretation as a uniform way of constructing internal models $\tau_K(\mathcal{M})$ of U from models \mathcal{M} of V . This construction gives us the contravariant model functor as soon as we have defined an appropriate category of interpretations.

Now consider two interpretations $K, M : U \rightarrow V$. Between the inner models $\tau_K(\mathcal{M})$ and $\tau_M(\mathcal{M})$ we have the usual structural morphisms of models. We are interested in the case where these morphisms are V -definable and uniform over models. This idea leads to the notion of i-morphism. An i-morphism $F : K \rightarrow M$ is a triple $\langle K, F(\vec{u}, \vec{v}), M \rangle$, where $F(\vec{u}, \vec{v})$ is a V -formula and where in all models of V , F represents a morphism of models from $\tau_K(\mathcal{M})$ to $\tau_M(\mathcal{M})$.

Two i-morphisms $F, G : K \rightarrow M$ are *i-equal*, when:

$$V \vdash \forall \vec{u}, \vec{v} (F(\vec{u}, \vec{v}) \leftrightarrow G(\vec{u}, \vec{v})).$$

We will think about i-morphisms modulo i-equality without dividing this equivalence relation out.

In the obvious way, we can define the identity i-morphism $\text{id}_K : K \rightarrow K$, composition of i-morphisms, i-isomorphisms, etc. All these operations preserve i-equality.

We easily see that i-isomorphisms really are isomorphisms in the category given by these operations.

We will say that two interpretations K, M are *i-equivalent* when there is an i-isomorphism between them, i.e. they are i-isomorphic. The notion of i-equivalence is our intended notion of sameness of interpretations. We will, however, *not* divide out i-equivalence. This enables us to use the notation τ_M meaningfully, to speak about the dimension of an interpretation, etc. Of course, we demand that operations on interpretations preserve i-equivalence. One may show that operations like identity, composition, $(\cdot)(\cdot)(\cdot)$ do indeed preserve i-equivalence. Moreover, if K is i-equivalent to M , then $\overline{K} = \overline{M}$.

The category INT_1 is the category of theories (as objects) and interpretations modulo i-equivalence (as arrows). One may show that we have indeed defined a category. Two theories U and V are *bi-interpretable* if they are isomorphic in INT_1 . Wilfrid Hodges calls this notion: *homotopy*. See [Hod93], p222.

Thus, U and V are bi-interpretable if there are interpretations $K : U \rightarrow V$ and $M : V \rightarrow U$, so that $M \circ K$ is i-isomorphic to ID_U and $K \circ M$ is i-isomorphic to ID_V . We call the pair K, M a *bi-interpretation* between U and V . One can show that the components of a bi-interpretation are faithful interpretations. Many good properties of theories like finite axiomatizability, decidability, κ -categoricity are preserved by bi-interpretations.

2.4. Sequential Theories. The sequential theories form an important class of theories in this paper. A sequential theory provides an interpretation N of a weak number theory, say S_2^1 , and sequences of all objects of the domain of the theories with projections in N . We can use these sequences to develop partial satisfaction predicates. Using these we can prove restricted consistency statements of U in U . See Subsection 2.5 for more about this.

The notion of sequential theory has a very simple definition discovered by Pavel Pudlák. We first need the definition of a very weak set theory. The theory Adjunctive Set Theory or AS is a one-sorted theory with a binary relation \in .

$$\text{AS1} \vdash \exists x \forall y (y \notin x),$$

$$\text{AS2} \vdash \forall x, y \exists z \forall u (u \in z \leftrightarrow (u \in x \vee u = y)).$$

We note that we do not demand extensionality. E.g., in AS we could have lots of ‘empty sets’.

An interpretation is *direct* iff it is one-dimensional, unrelativised (i.e. it has the trivial domain) and identity preserving (i.e. it translates identity to identity).

A theory U is sequential iff it directly interprets AS. By a substantial bootstrap, we can define, in a sequential theory U , an interpretation N of a weak number theory, sequences of all objects, etc.

For details see, e.g., [Pud83], [Pud85], [MPS90], [HP91], [Vis09] and [Vis10].

We can generalize the notion of sequentiality a bit to *poly-sequentiality* by replacing *direct interpretation* in the definition by its obvious generalization to the

m -dimensional case. All results in this paper that we prove for sequential theories also hold for poly-sequential theories.

2.5. Complexity and Satisfaction. *Restricted provability* plays an important role in this paper. An n -proof is a proof from axioms with Gödel number smaller or equal than n only involving formulas of complexity smaller or equal than n . To work conveniently with this notion, a good complexity measure ρ is needed. This should satisfy three conditions. (i) Eliminating terms in favour of a relational formulation should raise the complexity only by a fixed standard number. (ii) Translation of a formula via the translation corresponding to an interpretation K should raise the complexity of the formula by a fixed standard number depending only on K . (iii) The tower of exponents involved in cut-elimination should be of height linear in the complexity of the formulas involved in the proof.

Such a good measure of complexity together with a verification of desideratum (iii)—a form of nesting degree of quantifier alternations—is supplied in the work of Philipp Gerhardy. See [Ger03] and [Ger05]. It is also provided by Samuel Buss in his preliminary draft [Bus11]. Buss also proves that (iii) is fulfilled. We give some details about these measures in Appendix A.

We will use $\text{proof}_{U,n}$ for the proof predicate where only U -axioms with Gödel numbers $\leq n$ are allowed and where the formulas occurring in the proof are in the complexity class Γ_n of all formulas of complexity $\leq n$. Similarly we use $U \vdash_n A$, $\text{con}_n(U)$, $\Box_{U,m}A$, etc.

In sequential theories we can define partial satisfaction predicates for formulas with complexity below n , for any n . The presence of these predicates has as a consequence that for any sequential theory U and for any n , we can find an interpretation N of a weak arithmetic like Buss' S_2^1 in U such that $U \vdash \text{con}_n^N(U)$. See e.g. [Vis93] for more details.

3. THE ARITHMETICS OF A THEORY

There are many heuristic ways to look at interpretations. For example, an interpretation is a uniform internal model construction. In the case of definitional extensions an interpretation is an enrichment of the interpreting or target theory. In this paper, we opt for a third heuristic: we view an interpretation as *the interpreted theory in the context of the interpreting theory*.

We will say that an interpretation $N : S_2^1 \rightarrow U$ is an *arithmetic in U* or an *arithmetic of U* . The theory S_2^1 is Buss' theory of p-time computability. See [Bus86]. We stipulate we work with a version of S_2^1 that is formulated in the language of arithmetic with (the relational versions of) 0 , S , $+$ and \times .

Remark 3.1. In our definition of arithmetic we are both rewarded and punished for having a strict typing regime on interpretations. The reward is that the target theory or interpreting theory or context is part of the data for an arithmetic. So we can speak about an arithmetic N without mentioning the context. The punishment is that, e.g., an interpretation $K : EA \rightarrow U$ is *not* an arithmetic. The associated arithmetic is $S_2^1 \downarrow K$. See also Remark A.1. \square

There are several reasons for choosing S_2^1 . First it is a sequential theory. Secondly, the usual metamathematics leading up to Gödel's Incompleteness Theorems can be formalized in S_2^1 without the use of any extraneous tricks. Moreover it is a reasonably weak choice among theories in which this can be done. Thirdly, S_2^1 is finitely axiomatizable. On the other hand the results of the present paper are rather robust w.r.t. to different choices of the basic arithmetic. E.g., T_2^1 or $I\Delta_0 + \Omega_1$ would have worked as well. (But often some extra care is needed for $I\Delta_0 + \Omega_1$, since it is not known whether this theory is finitely axiomatizable.)

The main structure between arithmetics that we will consider is the initial embedding ordering \preceq . Consider two arithmetics N and N' in U . An initial embedding $F : N \rightarrow N'$ is an i-morphism satisfying the following additional property:

$$\bullet \ U \vdash (F(\vec{x}_0, \vec{y}_0) \wedge \vec{y}_1 <_{N'} \vec{y}_0) \rightarrow \exists \vec{x}_1 <_N \vec{x}_0 \ F(\vec{x}_1, \vec{y}_1).$$

We write $N \preceq N'$ for: there is an initial embedding $\langle N, F, N' \rangle$ of N in N' . We note that \preceq is preserved by i-equivalence. I.o.w., i-equivalence is a congruence relation for the arithmetics of U with \preceq . So, i-equivalence is a subrelation of the induced equivalence relation of \preceq .

We call N a *cut* of N' iff $\text{emb} : N \preceq N'$, where emb is the identical embedding.

The most salient fact about \preceq is upwards preservation of Σ_1 -sentences and downward preservation of Π_1 -sentences. We formulate this as a theorem.

Theorem 3.2. *Suppose N and N' are arithmetics in U and $N \preceq N'$. Let S be any Σ_1 -statement and let P be any Π_1 statement. We have: $U \vdash S^N \rightarrow S^{N'}$ and $U \vdash P^{N'} \rightarrow P^N$.*

We leave the trivial proof to the reader. Arithmetics commute in all the right ways with bi-interpretations, as is shown in the next theorem.

Theorem 3.3. *Suppose $K : U \rightarrow V$ and $M : V \rightarrow U$ are a bi-interpretation between U and V . Then the mapping $\Phi : N \mapsto NK$ is a bijection between the arithmetics of U and the arithmetics of V modulo i-equivalence, that is an isomorphism w.r.t. \preceq . Moreover, Φ is an isomorphism w.r.t. \preceq and $\overline{N} = \overline{\Phi(N)}$.*

Proof. Let $\Psi : N' \mapsto N'M$ be a mapping between the arithmetics of V and the arithmetics of U . It is easy to see that Ψ is the inverse of Φ , modulo i-equivalence. Clearly Φ and Ψ preserve \preceq , so it easily follows that Φ is an isomorphism w.r.t. \preceq .

Since the interpretations of a bi-interpretation are faithful, we find $\overline{N} = \overline{\Phi(N)}$. \square

Arithmetics of a sequential theory can always be assumed to be one-dimensional, as is formulated in the following theorem.

Theorem 3.4. *Suppose U is sequential and N is an arithmetic in U . Then there is a 1-dimensional arithmetic M in U that is i-equivalent to N .³*

³Similarly, if U is polysequential via an m -dimensional interpretation of AS, any arithmetic N in U is i-equivalent to an m -dimensional arithmetic N .

We leave the trivial proof to the reader.

A fundamental fact about the arithmetics of a sequential theory follows by Pavel Pudlák's adaptation of Dedekind's proof of his Categoricity Theorem for second order arithmetic.

Theorem 3.5 (Púdlak). *Consider a sequential theory U . Let N_0 and N_1 be arithmetics in U . Then, there is an arithmetic M in U such that $M \preceq N_0$ and $M \preceq N_1$.*

For a proof see e.g. [Pud85].

In a sequential theory we have a convenient reflection principle. We write Γ_n for the class of all formulas of complexity $\leq n$.

Theorem 3.6. *Consider any sequential theory U and let N be an arithmetic in U . For any n , we can find an arithmetic $M \preceq N$ such that, verifiably in S_2^1 , we have, for all sentences A in Γ_n , that $U \vdash \Box_{U,n}^M A \rightarrow A$.*

Sketch of the proof. The idea of the proof is that, in U , we can define a satisfaction predicate for Γ_n , using the N -numbers, and prove Γ_n -reflection by replacing induction over proof-length by the use of a definable cut M of N . For details see the proof of Fact 2.4.5(ii) of [Vis93]. In [Vis93] only verifiability of this fact in EA was claimed. However, the big disjunctions and conjunctions of exponential size used there are not needed, since for each proof p we only need the truth of the axioms occurring in p . So the disjunctions we really need are polynomial in p . \square

As far as we can ascertain, this theorem was known (or versions of it were known), at an early stage, to, independently, Pavel Pudlák, Robert Solovay, Alex Wilkie & Jeff Paris and Harvey Friedman. The paper [Pud85] contains a version.

The previous theorem shows that, for any n , we can 'improve' a given N to obtain n -reflection. In contrast, if U is finitely axiomatized, for any N , we can find an n such that, for any $m \geq n$, we have anti- m -reflection, i.o.w., a version of Löb's theorem for N and m . We first need Lemma 4.1. of [Vis93]. For the convenience of the reader we reproduce it here.

Lemma 3.1. *The following fact is verifiable in S_1^2 . Suppose A is any finitely axiomatized theory, $\rho(A) \leq m$, $\rho(B) \leq m$ and $A \vdash B$. Then, $S_2^1 \vdash \Box_{A,m} B$.*

We note that without the verifiability clause we could conclude $A \vdash_m B$ from $A \vdash B$. Since this step uses superexponentiation, it is not available in the context of S_2^1 .

Proof. We can prove the lemma in two ways.

The first uses the insight of [Pud86, Lemma 2.2] that, in S_2^1 , we have, for all x and y , that $S_2^1 \vdash \text{itexp}(x, |y|)$ exists. Here we define:

- $\text{itexp}(x, 0) := x$,
- $\text{itexp}(x, z + 1) := 2^{\text{itexp}(x, z)}$.

Suppose $A \vdash B$. By $\exists\Sigma_1^b$ -completeness, we find, for some p , $S_2^1 \vdash \text{proof}_A(p, B)$. Since we have, $S_2^1 \vdash \text{itexp}(kp, k'|p|)$ exists, for standard k, k' , we can apply cut-elimination to p inside S_2^1 .

The second way is to note that, in $S_2^1 + \text{con}_m(A + \neg B)$, we can build a Henkin interpretation H of $A + \neg B$. It follows that $\text{con}(S_2^1 + \text{con}_m(A + \neg B))$ implies $\text{con}(A + \neg B)$. We find the desired result by contraposition. \square

Theorem 3.7. *Suppose A is a finitely axiomatized sequential theory and that N is an arithmetic in A . We can find an n , such that $m \geq n$, we have, for all $B \in \Gamma_m$, if $A \vdash \Box_{A,m}^N B \rightarrow B$, then $A \vdash B$. This fact is verifiable in S_2^1 .*

This theorem is a weaker version of Theorem 4.1 of [Vis93]. We sketch the proof since it is easier to read without the ballast of the stronger version of [Vis93].

Proof. The proof is just the usual proof of Löb's Theorem with some checks that all the complexities are correct and one step involving Lemma 3.1 added. We choose:

$$n := \max(\rho(\text{prov}_{A,y}^N(z)) + 1, \rho(\text{sub}^N(x, y)) + 1, \rho(A)).$$

Here sub is the formula defining the Gödel substitution function. It follows that $n \geq \Box_{A,m}^N B$, for any B and m , since both B and n appear as numerals and, thus, only add a non-alternating block of quantifiers.

Let C be Gödel fixed point with $A \vdash C \leftrightarrow (\Box_{A,m}^N C \rightarrow B)$. The complexity of C is again m as can be seen by inspecting the construction.

Note that e.g. $\rho(\text{prov}_{A,y}^N(z))$ is polynomial in the data for N .

Suppose $A \vdash \Box_{A,m}^N B \rightarrow B$. We have:

$$\begin{aligned} A \vdash \Box_{A,m}^N C &\rightarrow (\Box_{A,m}^N \Box_{A,m}^N C \wedge \Box_{A,m}^N (\Box_{A,m}^N C \rightarrow B)) \\ &\rightarrow \Box_{A,m} B \\ &\rightarrow B \end{aligned}$$

So, (a) $A \vdash \Box_{A,m}^N C \rightarrow B$, and, hence, $A \vdash C$. By Lemma 3.1, we find that (b) $A \vdash \Box_{A,m}^N C$. Combining (a) and (b), we may conclude that $A \vdash B$. \square

Theorem 3.8. *Consider a theory U and an arithmetic N in U . Then, there is an arithmetic $N' \preceq N$ and a U -formula TRUE such that, for Σ_1 -sentences S , we have $U \vdash \text{TRUE}(S) \leftrightarrow S^{N'}$. (Here S is coded in N' .)*

Proof. We first work S_2^1 . Let $\text{sat}(v)$ be a Δ_0 -satisfaction predicate for formulas with just one designated variable v free. The main ingredients for the construction of such a predicate can be found in [HP91, Chapter V(5)].⁴ We will use the following two properties of sat : for $D(v)$ in Δ_0 ,

S1. $S_2^1 \vdash \forall x (\text{sat}(x, D(v)) \rightarrow D(x))$,

S2. $S_2^1 \vdash \forall x ((2^{2^x} \text{ exists} \wedge D(x)) \rightarrow \text{sat}(x, D(v)))$.

⁴The two classical works on this subject are [Les78] and [PD82].

Let J be an S_2^1 -cut such that $S_2^1 \vdash x \in J \rightarrow 2^{2^x}$ exists. We suppose Σ_1 -sentences are written in the form $\exists x S_0(x)$ where S_0 is Δ_0 . (If not, we add an algorithm that rewrites the Σ_1 -sentence to this normal form.) We define:

- $\text{true}(\exists x S_0(x)) := \exists x \in J \text{ sat}(x, S_0(v))$,
- $N' := N \circ J$,
- $\text{TRUE}(x) := \text{true}^N(x)$.

We easily verify that TRUE has the desired property. \square

Inspection of the proof of Theorem 3.8 shows that we can obtain reasonable commutation properties for TRUE in addition to mere Tarskian disquotation.

Remark 3.9. Suppose U is sequential. Let N be an arithmetic with parameters in U . In a model \mathcal{M} of U we can view N as a definable class of internal models parametrized by models of U . Theorem 3.5 tells us how to construct an parameter-free arithmetic below two parameter-free arithmetics. With some care we can generalize the construction to produce one parameter-free arithmetic below N viewed as a class of internal models. For details on such a construction, see [Vis10], the second proof of Theorem 5.2. As a result of this observation, the provability logic of all parameter-free arithmetics of U is the same as the provability logic of all arithmetics of U with parameters. \square

We end this section with a tentative discussion of what it means that \preceq has a minimal element.

Theorem 3.10. *Consider any theory U . Suppose N is a \preceq -minimal arithmetic in U , i.e., for any arithmetic M in U with $M \preceq N$, we have $N \preceq M$. Then, we have:*

- i. *For any Σ_1 -sentence S , and any $M \preceq N$, that $U \vdash S^N \rightarrow S^M$.*
- ii. *U proves parameter-free Π_1 -induction for the N -numbers. In other words, we have $U \vdash (\text{III}_1^-)^N$. As a consequence, we have sentential Σ_1 -completeness in N .*
- iii. *We have a Σ_1 -truth predicate TRUE satisfying Tarskian disquotation for Σ_1 -sentences on N .*

Proof. We have (i) simply because if $M \preceq N$, then $M \preceq N$, and Σ_1 -sentences are upwards preserved.

Ad (ii): As is easily seen parameter-free Π_1 -induction is equivalent to the parameter-free Σ_1 -minimum principle (over PA^-).⁵ We prove the parameter-free Σ_1 -minimum principle for N . We reason in U . Suppose $\exists x \in N S(x)$, where S is Σ_1 . Consider the virtual class $X := \{x \in N \mid \forall y < x \neg S(y)\}$. Clearly $0 \in X$. If X is not closed under successor, there is a $z \in N$ such that $z \in X$ and $Sz \notin X$. By elementary reasoning we find that z is the minimal N -number such that Sz . If X is not closed under successor, we can shorten X to a cut J that satisfies S_2^1 . Thus J is an arithmetic below N . It follows that on J we have both $\neg \exists x Sx$ and $\exists x Sx$. A contradiction.

⁵See [KPD88] and [CFL11] for details on III_1^- .

The theory III_1^- proves sentential Σ_1 -completeness since EA is conservative over III_1^- w.r.t. Σ_2 -sentences as was proved in [KPD88].

Ad (iii): The existence of the desired truth-predicate is immediate from Theorem 3.8. \square

Theorem 3.11. *Consider any sequential theory U . Suppose N is a \preceq -minimal arithmetic in U . It follows that:*

- i. N is a \preceq -minimum in U , i.o.w., for all arithmetics M in U , $N \preceq M$.*
- ii. U is parameter-free essentially reflexive for N , i.e., for any n , and any sentence $B \in \Gamma_n$, we have $U \vdash \Box_{U,n}^N B \rightarrow B$.*
- iii. U is not finitely axiomatizable.*

Proof. (i) is immediate by Theorem 3.5. (ii) follows from Theorem 3.6 in combination with Theorem 3.10(i). Claim (iii) is immediate from (ii) in combination with Theorem 3.7. \square

Remark 3.12. Consider a sequential theory U and suppose that N is \preceq -minimal in U . It follows that the interpretability logic of U (for sentential substitutions), w.r.t. arithmetization in N , is ILM. See [BV05] for most ingredients of the proof. \square

Open Question 3.13. Suppose U is sequential and has a \preceq -minimal arithmetic N . Can we get a precise estimate what this implies? E.g., can one show that we do not get full induction for N ? Is any $M \preceq N$ i-equivalent to N ? Such questions are both interesting in general and in the sequential case. \square

4. INTRODUCTION TO PROVABILITY LOGIC

We start with the basics concerning Löb's Logic GL. We define the language \mathcal{L}_{mod} of propositional modal logic by:

- $\alpha ::= p_0 \mid p_1 \mid \dots$
- $\phi ::= \alpha \mid \perp \mid \top \mid \neg \phi \mid \Box \phi \mid (\phi \wedge \phi) \mid (\phi \vee \phi) \mid (\phi \rightarrow \phi)$

The logic GL is axiomatized by the following axioms and rules.

GL1. We have all substitution instances of propositional tautologies.

GL2. $\vdash \Box(\phi \rightarrow \psi) \rightarrow (\Box \phi \rightarrow \Box \psi)$.

GL3. $\vdash \Box \phi \rightarrow \Box \Box \phi$.

GL4. $\vdash \Box(\Box \phi \rightarrow \phi) \rightarrow \Box \phi$.

GL5. If $\vdash \phi \rightarrow \psi$ and $\vdash \phi$, then $\vdash \psi$.

GL6. If $\vdash \phi$, then $\vdash \Box \phi$.

We have a completeness theorem for GL in finite, transitive, irreflexive Kripke models.

We define arithmetical interpretations of the modal language as follows. Let U be a theory and let N be an arithmetic in U . We define an N -translation σ as a mapping of the formulas of \mathcal{L}_{mod} to the sentences of the language of U , where σ commutes with the propositional connectives and where:

$$\sigma(\Box\phi) := \Box_U(\sigma(\phi))^N := \text{prov}_U(\ulcorner \sigma(\phi) \urcorner^N).$$

The variable ‘ a ’ will range over $0, 1, \dots, \infty$. We define:

- $\Box^0\phi := \phi$, $\Box^{n+1}\phi := \Box\Box^n\phi$, $\Box^\infty\phi := \top$.
- $\phi \in \text{prl}(N)$ iff, for all N -translations σ , $U \vdash \sigma(\phi)$.
- $\phi \in \text{prl}_{\text{all}}(U)$ iff, for all arithmetics N in U , $\phi \in \text{prl}(N)$.
- $\deg(N) := \min(\{a \mid \Box^a \perp \in \text{prl}(N)\})$.
- $\deg_{\text{all}}(U) := \min(\{a \mid \Box^a \perp \in \text{prl}_{\text{all}}(U)\})$.
- In case U is an extension of S_2^1 in the language of arithmetic, we write $\text{prl}(U)$ for $\text{prl}(\mathcal{E}_{S_2^1} U)$ and $\deg(U)$ for $\deg(\mathcal{E}_{S_2^1} U)$.

We note that $\deg_{\text{all}}(U) := \sup(\{\deg(N) \mid N \text{ is an arithmetic in } U\})$.

We have the following two major insights. Let exp be the axiom stating that exponentiation is total.

Theorem 4.1. *Consider any theory U . Let N be an arithmetic in U . We have:*

- I. $\text{prl}(N)$ extends GL (and is closed under the rules of GL).
- II. If $U \vdash \text{exp}^N$, then $\text{prl}(N) = \text{GL} + \Box^{\deg(N)} \perp$.

In essence the proof of (I) is given in [Bus86]. Most of the ideas are also in [WP87]. Robert Solovay proved (II) for theories like PA which are reasonably strong and Σ_1 -sound. The extension to the case of Σ_1 -unsound theories extending PA was proved in [Vis81]. The fact that EA was needed on the designated interpretation of arithmetic slowly emerged. See [Ver93]. In Section 5 we give a sharper formulation of Theorem 4.1(II).

The gap between (I) and (II) provides the great open problem of provability logic. What happens in the gap? For an extensive discussion of this problem, see [BV06].

Is the provability logic of an arithmetic a good property of arithmetics? It should at least be preserved under our chosen notion of sameness of arithmetics. We note that, if N is i -equivalent to N' , then S_2^1 verifies this i -equivalence. It follows that:

Theorem 4.2. *Consider any theory U and suppose that N and N' are arithmetics in U and that N is i -equivalent to N' . Then, for any N -translation σ and N' -translation σ' , we have: if $\sigma(p) = \sigma'(p)$ for all atoms p , then, for all ϕ , we have $U \vdash \sigma(\phi) \leftrightarrow \sigma'(\phi)$. It follows that $\text{prl}(N) = \text{prl}(N')$ and $\deg(N) = \deg(N')$. We have the same result on the weaker assumption that $N \preceq N'$ and $N' \preceq N$.*

So, in this sense, the provability logic of an arithmetic is a good property. On the other hand the provability logic of a theory-in-extension is dependent on the specification of the axiom set. The provability logic of a theory is *intensional*.⁶

Example 4.3. Consider the theory $U := \text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp$ with the obvious axiomatization. Clearly $U \vdash \Box_U\Box_U\perp$. On the other hand, suppose $U \vdash \Box_U\perp$. Then,

$$\text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp \vdash \Box_{\text{PA}+\Box_{\text{PA}}\Box_{\text{PA}}\perp}\perp.$$

So,

$$\text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp \vdash \Box_{\text{PA}}\neg\Box_{\text{PA}}\Box_{\text{PA}}\perp.$$

And, hence,

$$\text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp \vdash \Box_{\text{PA}}\neg\Box_{\text{PA}}\perp.$$

By applying the formalized Second Incompleteness Theorem to the conclusion, we get: $\text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp \vdash \Box_{\text{PA}}\perp$. By Löb's Theorem, we obtain $\text{PA} \vdash \Box_{\text{PA}}\perp$. Quod non. So $U \not\vdash \Box_U\perp$.

Now we modify the formula defining U , thus obtaining the theory \tilde{U} , by taking something to be an axiom if it is an axiom of U or it is of the form $\underline{p} \neq \underline{p}$, where p is a PA -proof of $\Box_{\text{PA}}\perp$. Clearly, U and \tilde{U} are extensionally equal.

Clearly $\text{PA} + \Box_{\text{PA}}\Box_{\text{PA}}\perp \vdash \Box_{\tilde{U}}\perp$, and, hence, $\tilde{U} \vdash \Box_{\tilde{U}}\perp$.

We conclude that $\deg_{\text{all}}(U) = \deg(U) = 2$ and $\deg_{\text{all}}(\tilde{U}) = \deg(\tilde{U}) = 1$.

We note that the arithmetics in our example are Σ_1 -unsound. It is unknown whether we can find two extensionally equal theories V and V' and two arithmetics $N := \langle S_2^1, \tau, V \rangle$ and $N' := \langle S_2^1, \tau, V' \rangle$ such that N and, *a fortiori*, N' are Σ_1 -sound that give rise to different provability logics. In case $V \vdash \text{exp}^N$, where exp is the axiom stating that exponentiation is total, we will see that $\text{prl}_{\text{all}}(U) = \text{prl}(N) = \text{GL} = \text{prl}(N') = \text{prl}_{\text{all}}(V')$. So any counterexamples for $\text{prl}(N)$ and $\text{prl}(N')$ should fail to prove the totality of exponentiation for N and any counterexamples for $\text{prl}_{\text{all}}(V)$ and $\text{prl}_{\text{all}}(V')$ should not contain any Σ_1 -sound arithmetic M that proves the principle exp^M . \square

Because of intensionality, provability logics and degrees need not to be preserved by bi-interpretation. To get the appropriate notion of sameness that preserves provability logics and degrees we consider bi-interpretability for S_2^1 -verifiable interpretations.

- $K : U \triangleright V$ is S_2^1 -verifiably an interpretation, if $S_2^1 \vdash \forall A (\Box_V A \rightarrow \Box_U A^K)$.

We have chosen the formalized version of the *theorems* formulation of interpretability. This is convenient but not really necessary. As Emil Jeřábek pointed out to me Buss' witnessing theorem implies that S_2^1 -verifiable axioms-interpretability implies S_2^1 -verifiable theorems-interpretability.

Theorem 4.4. *Suppose that $K : U \rightarrow V$ and $M : V \rightarrow U$ form a bi-interpretation. Suppose further that both K and M are S_2^1 -verifiably interpretations. Let N be an arithmetic in U . Then, $K \circ N$ is an arithmetic in V . We have: $\text{prl}(K \circ N) = \text{prl}(N)$. It follows that $\text{prl}(U) = \text{prl}(V)$. Similarly for the \deg .*

⁶This fact is folklore. I learned it first from Sergei Artemov around 1984.

Proof. We first note that (\dagger) U proves, for some F that “ F is an isomorphism between ID_U and $M \circ K$ ”. It follows that S_2^1 verifies the formalization of (\dagger) . Similarly, for the isomorphism between ID_V and $K \circ M$. Thus, we may conclude that (\dagger) $S_2^1 \vdash \Box_U A \leftrightarrow \Box_V A^K$, etcetera.

Suppose σ is an N -translation. We prove by induction on ϕ , that, for al ϕ , we have $V \vdash (\sigma(\phi))^K \leftrightarrow (K \circ \sigma)(\phi)$. The only interesting case is when $\phi = \Box\psi$. We have:

$$\begin{aligned}
 (1) \quad & V \vdash (\sigma(\Box\psi))^K \leftrightarrow (\Box_U^N \sigma(\psi))^K \\
 (2) \quad & \leftrightarrow \Box_U^{K \circ N} \sigma(\psi) \\
 (3) \quad & \leftrightarrow \Box_V^{K \circ N} (\sigma(\psi))^K \\
 (4) \quad & \leftrightarrow \Box_V^{K \circ N} (K \circ \sigma)(\psi)
 \end{aligned}$$

We note that step (3) uses (\dagger) and that step (4) uses the induction hypothesis.

Suppose that ϕ is in $\text{prl}(K \circ N)$, then, for any N -translation σ , we have $V \vdash (K \circ \sigma)(\phi)$. Ergo, $V \vdash \sigma(\phi)^K$. Hence, $U \vdash \sigma(\phi)$. It follows that $\phi \in \text{prl}(N)$.

Conversely, suppose ϕ is in $\text{prl}(N)$. Then, by Theorem 4.2, it follows that ϕ is in $\text{prl}(M \circ K \circ N)$. By the above argument applied with V and U and K and M interchanged and with $K \circ N$ in the role of N , we find: If ϕ is in $\text{prl}(K \circ N)$, then ϕ is in $\text{prl}(M \circ K \circ N)$. Hence, if ϕ is in $\text{prl}(K \circ N)$, then ϕ is in $\text{prl}(N)$. \square

Thus, if U and V are bi-interpretable via S_2^1 -verifiable interpretations, then the interpretations provide an isomorphism between their arithmetics that preserves \preceq and deg and prl .

5. SOLOVAY'S THEOREM

In this section, we study the forms that Solovay's Theorem takes in various settings.

5.1. The Guaspari-Solovay System R^- . In this subsection we give a careful analysis of the proof of Solovay's Theorem. We follow the modal presentation of the proof due Dick de Jongh, Marc Jumelet and Franco Montagna in their paper [dJJM91].

We introduce the logic R^- of Guaspari and Solovay and some subsystems of this logic. See [GS79]. The language of R^- is given by:

- $\alpha ::= p_0 \mid p_1 \mid \dots$
- $\phi ::= \alpha \mid \perp \mid \top \mid \neg\phi \mid \Box\phi \mid (\phi \wedge \phi) \mid (\phi \vee \phi) \mid (\phi \rightarrow \phi) \mid \Box\phi < \Box\phi \mid \Box\phi \leq \Box\phi$

The logic R^- is axiomatized by the axioms and rules of **GL** (for the extended language) plus the following axioms.

- $R^-1. \vdash \Box\phi \leq \Box\psi \rightarrow \Box\phi$
- $R^-2. \vdash (\Box\phi \leq \Box\psi \wedge \Box\psi \leq \Box\chi) \rightarrow \Box\phi \leq \Box\chi$
- $R^-3. \vdash \Box\phi < \Box\psi \leftrightarrow (\Box\phi \leq \Box\psi \wedge \neg\Box\psi \leq \Box\phi)$
- $R^-4. \vdash \Box\phi \rightarrow (\Box\phi \leq \Box\psi \vee \Box\psi \leq \Box\phi)$
- $R^-5. \vdash \Box\phi \leq \Box\psi \rightarrow \Box(\Box\phi \leq \Box\psi)$

R⁻6. $\vdash \Box\phi < \Box\psi \rightarrow \Box(\Box\phi < \Box\psi)$

We can split Axiom R⁻4 into two parts that are jointly equivalent to R⁻4:

R⁻4a. $\vdash \Box\phi \leq \Box\psi \rightarrow (\Box\phi \leq \Box\psi \vee \Box\psi \leq \Box\phi)$

R⁻4b. $\vdash \Box\phi \rightarrow \Box\phi \leq \Box\phi$

We will consider two subsystems of R⁻ to wit R₀⁻ and R₁⁻. R₀⁻ is given by the axioms and rules of GL (for the extended language) plus R⁻1, R⁻2, R⁻3, R⁻4a and R₁⁻ is R₀⁻ plus R⁻4b, i.o.w., R₁⁻ is given by R⁻1, R⁻2, R⁻3, R⁻4.

5.2. Arithmetical Interpretations of R⁻. Consider any arithmetical theory U and any arithmetic N in U . We specify what it is to be an interpretation of the language of R⁻ for U, N .

We remind the reader of the witness comparison notation. We define, for any $C = \exists x C_0(x)$ and $D = \exists y D_0(y)$:

- $C \leq D := \exists x (C_0(x) \wedge \forall y < x \neg D_0(y))$,
- $C < D := \exists x (C_0(x) \wedge \forall y \leq x \neg D_0(y))$,
- $(C \leq D)^\perp := (D < C)$.

We interpret the language of R⁻ as follows. An N -translation σ sends the propositional variables to U -sentences, commutes with the propositional connectives and satisfies:

- $\sigma(\Box\phi) = \text{prov}_U^N(\ulcorner \sigma(\phi) \urcorner)$,
- $\sigma(\Box\phi \leq \Box\psi) = \text{prov}_U^N(\ulcorner \sigma(\phi) \urcorner) \leq \text{prov}_U^N(\ulcorner \sigma(\psi) \urcorner)$,
- $\sigma(\Box\phi < \Box\psi) = \text{prov}_U^N(\ulcorner \sigma(\phi) \urcorner) < \text{prov}_U^N(\ulcorner \sigma(\psi) \urcorner)$.

We assume that we are employing an ordinary single conclusion proof predicate. A modal formula ϕ is N -valid if for all N -translations σ , we have $U \vdash \sigma(\phi)$.

It is easily seen that the theory R₀⁻ is N -valid, for any N . The principles R⁻4b, on the one hand, and R⁻5 and R⁻6, on the other, are not known to be N -valid. e.g. in case $N = \text{ID}_{S_2^1}$.

The principle R⁻4b is a modal articulation of a special case of the minimum-principle. It tells us that if a certain sentence has a proof than it has a minimal proof. Since the proof-predicate is $\Delta_1^b(S_2^1)$, a reasonable principle to ask for is the Σ_1^b -minimum principle, i.e., Buss' minimization axiom Σ_1^b -MIN. By the results of [Bus86, Section 2.9], this principle is equivalent over S_2^1 with Σ_1^b -induction, i.e., Buss' principle Σ_1^b -IND. This means that a salient not-too-strong theory in which we have R⁻4b is the theory T_2^1 . Thus, if $U \vdash (T_2^1)^N$, then R₁⁻ is N -valid.

The principles R⁻5 and R⁻6 are instances of sentential $\exists\Pi_1^b$ -completeness. Is there a natural theory extending T_2^1 that is as weak as possible that delivers this principle? Of course, $B_0 := T_2^1 + \{S \rightarrow \Box_{S_2^1} S \mid S \in \exists\Pi_1^b\text{-sent}\}$ does the trick but this principle involves coding. Let $\mathcal{J} := \{x \mid 2^x \downarrow\}$. We define:

$$B := T_2^1 + \{S \in \exists\Pi_1^b\text{-sent} \mid S \rightarrow S^{\mathcal{J}}\}.$$

Since the proof of completeness for $\exists\Pi_1^b$ -sentences only needs one exponent, we find $B \vdash S \rightarrow \Box_{S_2^1} S$, for S a $\exists\Pi_1^b$ -sentence. So B extends B_0 .⁷ Thus, if $U \vdash B^N$, then R^- is N -valid.

The two main theorems concerning provability logic of this paper, to wit Theorems 5.5 and 5.9, will employ T_2^1 to ensure the principle R^-4b . In contrast, we will not use B to obtain R^-5 & R^-6 . This last theory is, in a sense, still too strong. The theory T_2^1 is interpretable in S_2^1 on a cut, but B is locally but not globally interpretable in S_2^1 . See Remark 5.1.

Remark 5.1. The two minimal *salient* theories in the literature in which we have B_0 are EA and Π_1^{-1} . Since Π_1^{-1} does not fit our framework, we will consider $\Pi_1^{-1} + \Omega_1$ instead. The theory B is both a subtheory of EA and of $\Pi_1^{-1} + \Omega_1$. The theories EA and Π_1^{-1} have the same $B\Sigma_1$ -consequences. See [KPD88] or [CFL11].

In [CFL11, Theorem 1.3(2)] it is shown that the $B\Sigma_1$ -consequences of Π_0^{-1} and, hence, of EA can be axiomatized by the theory:

$$CFL := I\Delta_0 + \{S \in \Sigma_1\text{-sent} \mid S \rightarrow S^{\mathcal{J}}\},^8$$

Clearly, $CFL + \Omega_1$ extends B . The theories $\Pi_1^{-1} + \Omega_1$ and, *a fortiori*, $CFL + \Omega_1$ and B are locally interpretable in S_2^1 . The proof that $\Pi_1^{-1} + \Omega_1$ is locally interpretable in S_2^1 can be found in [Vis12]. Thus, they are *locally weak*. One can show that S_2^1 does not interpret B , so B is not a weak theory and, *a fortiori*, neither are $CFL + \Omega_1$ and $\Pi_1^{-1} + \Omega_1$. To prove this one shows that $B \vdash \text{con}_n(S_2^1)$, for every n . By the results of [Pud85], S_2^1 cannot interpret $S_2^1 + \{\text{con}_n(S_2^1) \mid n \in \omega\}$. See also e.g. [Vis11]. \square

5.3. The Basic Proof. In this subsection we present the version of Solovay's proof that is due to de Jongh, Jumelet & Montagna.

Our first aim is to embed a finite Kripke frame for ordinary modal logic in the logic R^- extended with a finite set of constants and a finite set of axioms concerning these constants. Via the arithmetical validity of our modal theory this embedding subsequently induces an embedding in an arithmetic.

Let \mathcal{F} be a finite, irreflexive, transitive Kripke frame on worlds $\{0, \dots, n-1\}$. Our frame need not be rooted.

We write $i \parallel j$ for: i and j are *incomparable*, i.e., $i \not\leq j$ and $j \not\leq i$.

For $i = 0, \dots, n-1$, we add constants ℓ_i to the language of R^- . Consider the following axioms.

$$\mathcal{F}1. \vdash \ell_i \leftrightarrow (\Box \neg \ell_i \wedge \bigwedge_{j \succ i} \Diamond \ell_j \wedge \bigwedge_{j \parallel i} \bigvee_{k \leq i, k \parallel j} \Box \neg \ell_k < \Box \neg \ell_j).$$

$\mathcal{F}2.$ For $i \neq j$, we have:

$$\vdash \Box \neg \ell_i \leq \Box \neg \ell_j \rightarrow \Box \neg \ell_i < \Box \neg \ell_j.$$

⁷The theory B as defined here seems to suffice. However, I am not sure that a definition using $\exists^*\Pi_1^b$ -sentences is not more natural.

⁸It takes a little argument to prove the equivalence of our formulation of CFL and the formulation in [CFL11]. For the definition it does not matter, modulo provable equivalence over T_2^1 , whether the Σ_1 sentences are written in the form $\exists\Delta_0$ or $\exists^*\Delta_0$. We may also consider sentences in $\exists^*\Delta_0$ and just relativize only the first of the block of existential quantifiers to \mathcal{J} .

We add these axioms to R_i^- and to R^- , thus obtaining $R_{i,\mathcal{F}}^-$ and $R_{\mathcal{F}}^-$. (Here we let the axioms and rules of $R_{(i)}^-$ apply to the extended language with the new axioms.)

We adhere to the usual convention that the empty conjunction is \top and the empty disjunction is \perp .

The N -interpretation of these principles is given as follows. By the Multiple Fixed Point Lemma we find sentences L_i such that:

$$S_2^1 \vdash L_i \leftrightarrow (\Box_U \neg L_i^N \wedge \bigwedge_{j \succ i} \Diamond L_j^N \wedge \bigwedge_{j \parallel i} \bigvee_{k \preceq i, k \parallel j} \Box_U \neg L_k^N < \Box_U \neg L_j^N).$$

We will assume that, for $i \neq j$, we have $L_i \neq L_j$.⁹

We demand that $\sigma(\ell_i) := L_i^N$. Thus we treat the ℓ_i as constants.¹⁰ For an arbitrary arithmetic N this gives us the validity of $R_{0,\mathcal{F}}^-$.

Below we want to reason in an informal way in the theory $R_{(1),\mathcal{F}}^-$. We want to reason as if we have predicate logic available, so that we can talk about statements like $i \preceq j$ and so that we can quantify over the nodes of F . These problems can be solved as follows. A statement like $i \preceq j$ in the context of $R_{(1),\mathcal{F}}^-$ stands for \top when it is true and for \perp when it is false. Quantification over our finite domain is handled by translating it to iterated conjunctions and disjunctions.

We define:

$$\bullet h_i := (\Box \neg \ell_i \wedge \bigwedge_{j \parallel i} \bigvee_{k \preceq i, k \parallel j} \Box \neg \ell_k < \Box \neg \ell_j).$$

Lemma 5.1 ($R_{1,\mathcal{F}}^-$). *Suppose $i \parallel j$, then $\neg(h_i \wedge h_j)$. I.o.w., h_i and h_j , implies $i \preceq j$ or $j \preceq i$.*

Proof. Reason in $R_{1,\mathcal{F}}^-$. Suppose $i \parallel j$ and h_i and h_j . Consider the i' , such that $i' \preceq i$, $i' \parallel j$, and $\Box \neg \ell_{i'}$. We note that there is such an i' , to wit i , because $i \preceq i$, $i \parallel j$ and $\Box \neg \ell_i$. The $\Box \neg \ell_{i'}$ are linearly ordered in the witness comparison ordering $<$. Suppose $\Box \neg \ell_{i^*}$ is the $<$ -smallest such element. Consider the j' , such that $j' \preceq j$, $j' \parallel i$, and $\Box \neg \ell_{j'}$. The node j is an example of such a j' . The $\Box \neg \ell_{j'}$ are linearly ordered in the witness comparison ordering $<$. Suppose $\Box \neg \ell_{j^*}$ is the \leq -smallest such element.

By the second conjunct of h_i applied to j^* , we find $\Box \neg \ell_{i^*} < \Box \neg \ell_{j^*}$. By the second conjunct of h_j applied to i^* , we find $\Box \neg \ell_{j^*} < \Box \neg \ell_{i^*}$. A contradiction. \square

Lemma 5.2 ($R_{1,\mathcal{F}}^-$). *Suppose $i \neq j$, then $\neg(\ell_i \wedge \ell_j)$.*

⁹I was aware of two essentially different constructions for the Multiple Fixed Point Lemma. Vincent van Oostrom, after I asked him, provided a third construction. All three constructions automatically guarantee the desired property even in the presence of non-trivial automorphisms of the frame. One reason that this happens is that the choice of substitution-variables is explicitly arithmetically coded in the construction.

¹⁰Note that the L_i are not necessarily uniquely determined by the fixed point equations. Thus, we are looking at some choice of the L_i .

Proof. Reason in $R_{1,\mathcal{F}}^-$. In case i and j are incomparable, this is immediate by the previous lemma. Suppose, e.g., $i \prec j$. Suppose ℓ_i and ℓ_j . From ℓ_i , we have $\Diamond \ell_j$, and, from ℓ_j , we have $\Box \neg \ell_j$. A contradiction. \square

Lemma 5.3 ($R_{1,\mathcal{F}}^-$). *Suppose h_i and $\neg \ell_i$. Then, for some $j \succ i$, we have h_j .*

Proof. Reason in $R_{1,\mathcal{F}}^-$. Suppose h_i and $\neg \ell_i$. Then for some $j' \succ i$, we have $\Box \neg \ell_{j'}$. The $\Box \neg \ell_{j'}$ with $j' \succ i$ can be linearly ordered by the witness comparison ordering $<$. Let $\Box \neg \ell_{j^*}$ be the $<$ -minimal element among these j' .

Consider any $m \parallel j^*$. If $m \parallel i$, by h_i , we can find such that $k \preceq i \prec j^*$, and $k \parallel m$ and $\Box \neg \ell_k < \Box \neg \ell_m$. If not $m \parallel i$, we must have $i \prec m$. In case $\Box \neg \ell_m$, by the choice of j^* , we find $\Box \neg \ell_{j^*} < \Box \neg \ell_m$. In case $\neg \Box \neg \ell_m$, the axioms of R^- imply $\Box \neg \ell_{j^*} < \Box \neg \ell_m$. So in all cases we can find a k , such that $k \preceq j^*$ and $k \parallel m$ and $\Box \neg \ell_k < \Box \neg \ell_m$. We may conclude h_{j^*} . \square

Lemma 5.4 ($R_{1,\mathcal{F}}^-$). *Suppose h_i , then, for some $j \succeq i$, we have ℓ_j .*

Proof. Reason in $R_{1,\mathcal{F}}^-$. Suppose h_i . If ℓ_i , we are done. If not, by Lemma 5.3, there is a $i' \succ i$ such that $h_{i'}$. If $\ell_{i'}$, we are done. By repeating this procedure, we eventually find a $j \succeq i$, such that ℓ_j . \square

Lemma 5.5 ($R_{1,\mathcal{F}}^-$). *Suppose $\Box \neg \ell_i$. Then, for some j , we have ℓ_j .*

Proof. Reason in $R_{1,\mathcal{F}}^-$. Suppose $\Box \neg \ell_i$. Consider all j' such that $\Box \neg \ell_{j'}$. There is one such j' , to wit i . The $\Box \neg \ell_{j'}$ are linearly ordered by the witness comparison ordering $<$. Let j^* be the minimal such. It is easy to see that h_{j^*} . By Lemma 5.4, we find a $j \succeq j^*$ such that ℓ_j . \square

We define the theory $R_{2,\mathcal{F}}^-$ as $R_{1,\mathcal{F}}^-$ plus the following axioms:

$$\mathcal{F}3. \vdash \bigwedge_{i < n} (h_i \rightarrow \Box h_i)$$

Let \mathcal{K} be any Kripke model on the frame \mathcal{F} . We define an interpretation σ^* from \mathcal{L}_{mod} to the closed formulas of the language of $R_{\mathcal{F}}^-$, by setting $\sigma^*(p) := \bigvee_{j \Vdash p} \ell_j$, where σ^* commutes with the propositional connectives and \Box . We have:

Theorem 5.2. *We have, for every formula ϕ of the modal language:*

- if $i \Vdash \phi$, then $R_{2,\mathcal{F}}^- \vdash \ell_i \rightarrow \sigma^*(\phi)$;
- if $i \not\Vdash \phi$, then $R_{2,\mathcal{F}}^- \vdash \ell_i \rightarrow \neg \sigma^*(\phi)$;

Proof. The proof is by induction on ϕ . The cases of the atoms and of the propositional connectives are trivial using Lemma 5.2

Suppose $\phi = \Box \psi$.

Suppose $i \Vdash \Box \psi$. We reason in $R_{2,\mathcal{F}}^-$. Suppose ℓ_i . Then, h_i and, hence $\Box h_i$. By Lemma 5.4, in combination with $\Box \neg \ell_i$, we find $\Box \bigvee_{j \succ i} \ell_j$. So, by the induction hypothesis, we have $\Box \sigma^*(\psi)$.

Suppose $i \not\models \Box\psi$. Then, for some $j \succ i$, $j \not\models \psi$. We reason in $R_{2,\mathcal{F}}^-$. Suppose ℓ_i . It follows that $\Diamond\ell_j$. Ergo, by the induction hypothesis, $\Diamond\neg\sigma^*(\psi)$. Hence, $\neg\Box\sigma^*(\psi)$. \square

Remark 5.3. The most naive attempt to avoid the use of $\vdash h_i \rightarrow \Box h_i$, is to replace $\Box\neg\ell_i$ by $\Box\bigvee_{j \succ i} \ell_j$ (where we read \Diamond as $\neg\Box\neg$) everywhere in the above definitions and arguments. This certainly will give us the “ $i \Vdash \Box\psi$ ”-part in the proof of Theorem 5.2 for free. It may amuse the reader to try this and to see where and why precisely it goes wrong. \square

5.4. Application to Arithmetic. In this subsection, we articulate what Theorem 5.2 tells us about a theory U with arithmetic N .

Consider a theory U and an arithmetic N in U . Suppose that $\deg(N) = \alpha$ and $U \vdash (\mathsf{T}_2^1)^N$. Suppose $\mathsf{GL}_\alpha \not\models \phi$. Let \mathcal{K} be a finite counter Kripke model with frame \mathcal{F} . We choose \mathcal{K} in such a way that the set of worlds is $\{0, \dots, n-1\}$, that the root is 0 and $0 \not\models \phi$. Note that the depth of the root must be $k \leq \alpha$.

Let τ be the N -interpretation of the language of $R_{\mathcal{F}}^-$ that is generated by $\ell_i \mapsto L_i^N$. (We are only interested in τ on the closed fragment of $R_{\mathcal{F}}^-$.) Clearly, $\tau(h_i)$ is of the form H_i^N , where H_i is S_2^1 -provably equivalent to an $\exists\Pi_1^b$ -sentence. Let σ^* be the interpretation of the language of provability logic in the closed fragment of $R_{2,\mathcal{F}}^-$ generated by $p \mapsto \bigvee_{i \Vdash p} \ell_i$. We take the interpretation ν of the language of provability logic into the language of U to be $\tau \circ \sigma^*$. Thus, ν is the N -interpretation generated by $p \mapsto \bigvee_{i \Vdash p} L_i^N$. We assume that $U \vdash \bigwedge (H_i \rightarrow \Box_U H_i)$, in other words, that $R_{2,\mathcal{F}}^-$ is N -valid. We show that $U \not\models \nu(\phi)$.

Suppose $U \vdash \nu(\phi)$. Since $U \vdash L_0^N \rightarrow \neg\nu(\phi)$, we find that $U \vdash \neg L_0^N$. It follows that $U \vdash \Box_U^N \neg L_0^N$, and, hence $U \vdash \bigvee_{j < n} L_j^N$. Since $U \vdash \neg L_0^N$, we find $U \vdash \bigvee_{0 < j < n} L_j^N$. Thus, since each $j > 0$ satisfies $\Box^{k-1} \perp$, we find $U \vdash \Box_U^{N,k-1} \perp$, quod non. We may conclude that $U \not\models \nu(\phi)$.

The following theorem is an immediate consequence of these considerations.

Theorem 5.4. *Consider a theory U and an arithmetic N in U . We suppose that $U \vdash (\mathsf{T}_2^1)^N$ and $U \vdash S^N \rightarrow \Box_U^N S^N$, for all sentences S in $\exists\Pi_1^b$. Then $\mathsf{prl}(N) = \mathsf{GL} + \Box^{\deg(N)} \perp$.*

5.5. All Arithmetics of a Theory. In this subsection, we apply the Solovay argument to all arithmetics of a theory U . We first show how we can improve our local arithmetic.

Consider any set of Σ_1 -sentences \mathcal{S} with n elements. Let $C_{\mathcal{S}} := \bigwedge_{S \in \mathcal{S}} (S \rightarrow \Box_{\mathsf{S}_2^1} S)$. Let J be an S_2^1 -cut such that $\mathsf{S}_2^1 \vdash \forall x \in J \exists y \ 2^{2^x} = y$. We define: $J_0 := \text{ID}$, $J_{k+1} := \text{ID}\langle C_{\mathcal{S}} \rangle(J \circ J_k)$, $J_{\mathcal{S}} := J_n$.

The following argument is taken from [Vis91a].

Lemma 5.6. *We have: $\mathsf{S}_2^1 \vdash C_{\mathcal{S}}^{J_{\mathcal{S}}}$.*

Proof. Reason in S_2^1 . If we have $C_{\mathcal{S}}$, clearly $J_n = \text{ID}$ and we are done. Otherwise, for some S in \mathcal{S} , we have S and $\neg\Box_{\mathsf{S}_2^1} S$. So, inside J , the sentence S will be false.

It follows that, inside J , the number of true S 's from \mathcal{S} is at least one less than inside ID . Now the game repeats itself inside J for J_{n-1} . Each time we have $\neg C_{\mathcal{S}}$, we move inside J and loose at least one true S . If at some point, we have $C_{\mathcal{S}}$, we are done. Otherwise we end up with zero true S 's and we have $C_{\mathcal{S}}$ in J_n . (Since n is standard the whole argument can be spelled out with big disjunctions, etc.) \square

We have the following theorem.

Theorem 5.5. *Consider any theory U . We have $\text{prov}_{\text{all}}(U) = \text{GL} + \Box^{\text{deg}(U)} \perp$.*

We note that the result also is valid for the case that $\text{deg}_{\text{all}}(U) = 0$, i.e. when either U is inconsistent or \mathbf{S}_2^1 is not interpretable in U .

Proof. Consider any theory U and suppose $\text{deg}_{\text{all}}(U) = \alpha$. It is easily seen that $\text{GL} + \Box^\alpha \perp \subseteq \text{prl}(U)$.

Suppose $\text{GL}_\alpha \not\vdash \phi$. Then, there is a finite Kripke model \mathcal{K} with nodes $\{0, \dots, n-1\}$ and with root 0, such that $0 \not\models \phi$ and $d(0) \leq \alpha$.

Since $\text{deg}_{\text{all}}(U) := \sup(\{\text{deg}(M) \mid M \text{ is an arithmetic in } U\})$ and $d(0)$ is finite, we can find an arithmetic N_0 with $d(0) \leq \text{deg}(N_0) \leq \alpha$. We can shorten N_0 to an arithmetic $N_1 \preceq N_0$ in which we have \mathbf{T}_2^1 . (See e.g. [HP91]. In fact, we can shorten N_0 to a cut on which we have $\text{ID}_0 + \Omega_1$.) We note that $\text{deg}(N_1) \geq \text{deg}(N_0) \geq d(0)$.

We simultaneously construct a cut N in N_0 and the L_i using the Gödel Fixed point Lemma. We find L_i such that:

$$\mathbf{S}_2^1 \vdash L_i \leftrightarrow (\Box_U \neg L_i^{N_1 \circ J_{\mathcal{H}}} \wedge \bigwedge_{j \succ i} \Diamond L_j^{N_1 \circ J_{\mathcal{H}}} \wedge \bigwedge_{j \parallel i} \bigvee_{k \preceq i, k \parallel j} \Box_U \neg L_k^{N_1 \circ J_{\mathcal{H}}} < \Box_U \neg L_j^{N_1 \circ J_{\mathcal{H}}}).$$

Here:

- $H_i := (\Box_U \neg L_i^{N_1 \circ J_{\mathcal{H}}} \wedge \bigwedge_{j \parallel i} \bigvee_{k \preceq i, k \parallel j} \Box_U \neg L_k^{N_1 \circ J_{\mathcal{H}}} < \Box_U \neg L_j^{N_1 \circ J_{\mathcal{H}}})$.
- $\mathcal{H} := \{H_0, \dots, H_{n-1}\}$.

We note that we have indeed a valid application of the Fixed Point Lemma since \mathcal{H} occurs ‘inside the box’.

We take $N := N_1 \circ J_{\mathcal{H}}$. We note that $\mathbf{S}_2^1 \vdash \Box_{\mathbf{S}_2^1} H_i \rightarrow \Box_U H_i^N$. Hence, we find that $U \vdash H_i^N \rightarrow \Box_U^N H_i^N$. Moreover, $U \vdash (\mathbf{T}_2^1)^N$, since \mathbf{T}_2^1 is downwards preserved over \preceq . Finally, $\text{deg}(N) \geq \text{deg}(N_1) \geq d(0)$.

We now employ the interpretation ν of Subsection 5.4 using the L_i constructed above. We find: $U \not\models \nu(\phi)$. \square

5.6. Theories with a Σ_1 -sound Arithmetic. In this subsection we provide a sufficient condition for a theory to contain an arithmetic N with $\text{prl}(N) = \text{GL}$.

We will use the following facts.

Fact 5.6. *Suppose N is an arithmetic in U . Then $U \triangleright (U + \Box_U^N \perp)$.*

This insight was first due Solomon Feferman in his classical paper [Fef60]. The simple proof below was discovered independently by Per Lindström (see [Fef97]) and the author (see [Vis90]).

Proof. Suppose N is an arithmetic in U . We have $U + \neg\Box_U^N \perp \vdash \Diamond_U^N(\Box_U^N \perp)$, by Löb's Theorem. Hence, it follows that $(U + \neg\Box_U^N \perp) \triangleright (U + \Diamond_U^N \Box_U^N \perp)$. So, using a Henkin interpretation, we may conclude that $(U + \neg\Box_U^N \perp) \triangleright (U + \Box_U^N \perp)$. On the other hand, we trivially have $(U + \Box_U^N \perp) \triangleright (U + \Box_U^N \perp)$. Thus, using a disjunctive interpretation, we find that $U \triangleright (U + \Box_U^N \perp)$. \square

Fact 5.7. *Suppose $U \triangleright V$. Let N be an arithmetic in V . Then $U \triangleright (V + \Box_U^N \perp)$.*

Proof. Suppose $M : U \triangleright V$. We apply Fact 5.6 to the arithmetic $M \circ N$ to find the desired result. \square

We note that Facts 5.6 and 5.7 can be considered as nice and general formulations of the Second Incompleteness Theorem. Suppose that, for some arithmetic N in U , we would have $U \vdash \text{con}^N(U)$. Since, by Fact 5.6, we have $U \triangleright (U + \Box_U^N \perp)$, it follows that U is inconsistent.

Fact 5.8. *Suppose $U \triangleright V$ and suppose that U contains a Σ_1 -sound arithmetic N , i.e., for all Σ_1 -sentences S , if $U \vdash S^N$, then S is true. Then $U \triangleright_{\text{faith}} V$.*

This fact is a direct consequence of Theorem B.4 of Appendix B. It was first proved in [Vis05]. The basic idea of the proof is due to Per Lindström.

We prove the following theorem.

Theorem 5.9. *Suppose U contains a Σ_1 -sound arithmetic N_0 . Then there is an N in U such that $\text{prl}(N) = \text{GL}$.*

Proof. Suppose that U contains a Σ_1 -sound arithmetic N_0 . We can find an interpretation of T_2^1 by shortening N_0 . By Fact 5.7, we find $U \triangleright (\mathsf{T}_2^1 + \Box_U \perp)$. Let W be T_2^1 plus sentential Σ_1 -completeness for U . Since $\mathsf{T}_2^1 + \Box_U \perp$ extends W we find: $U \triangleright W$. By Fact 5.8, we can find an K such that $K : U \triangleright_{\text{faith}} W$. Finally, let $N := \mathsf{S}_2^1 \downarrow K$. Since W is a true theory, N is a Σ_1^0 -sound arithmetic in U . Hence $\text{deg}(N) = \infty$. By Theorem 5.4 we find that $\text{prl}(N) = \text{GL}$. \square

6. DEEP ARITHMETICS

In this section we study the fine structure of the arithmetics of a finitely axiomatized sequential theory. Finitely axiomatized sequential theories have many surprising properties. The present section builds on and extends a line earlier work, to wit: [Smo85a], [Pud85], [Kra87], [Vis93] and [Vis05].

We have the following definition. Suppose A is a finitely axiomatized sequential theory. (We confuse these theories with their single axiom.) Let N be an arithmetic in A .

- N is Σ_1 -veracious in A iff

$$S_2^1 \vdash \forall S \in \Sigma_1\text{-sent} (\Box_A S^N \rightarrow \Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S)).$$

Thus, we see that Σ_1 -veracity is the S_2^1 -verifiable Σ_1 -conservativity of N over $\mathcal{E}_{S_2^1(S_2^1 + \text{con}_{\rho(A)}(A))}$.

- N is *strong* in A iff $A \vdash \text{con}_{\rho(A)}^N(A)$.
- N is *deep* in A iff N is both Σ_1 -veracious and strong in A .

Σ_1 -veracity is connected to Σ_1 -soundness: this is elucidated by the following theorem.

Theorem 6.1. *Suppose that A is a finitely axiomatized sequential theory and N is Σ_1 -veracious in A . Then,*

$$\text{I}\Delta_0 + \text{supexp} + \text{con}(A) \vdash \forall S \in \Sigma_1\text{-sent} (\Box_A S^N \rightarrow \text{true}(S)).$$

Here true is a Σ_1 truth predicate.

Proof. By [WP87], the theory EA , aka $\text{I}\Delta_0 + \text{exp}$, proves uniform Π_2 -reflection for cutfree-provability in S_2^1 . Hence, $\text{I}\Delta_0 + \text{supexp}$ proves uniform Π_2 -reflection for ordinary provability in S_2^1 . Our theorem is immediate from this. \square

In the definition of Σ_1 -veracious theory we may replace $\rho(A)$ by any $m \geq \rho(A)$, in the light of the following theorem.

Theorem 6.2. *Let A be a finitely axiomatized sequential theory. Suppose that $m \geq \rho(A)$. We have:*

$$S_2^1 \vdash \forall S \in \Sigma_1\text{-sent} (\Box_{S_2^1}(\text{con}_m(A) \rightarrow S) \leftrightarrow \Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S)).$$

Proof. The right-to-left direction of our theorem is trivial.

To go from m -provability to $\rho(A)$ -provability we need to eliminate standard (proof-theoretical) cuts. So we only need a multi-exponential transformation.¹¹ Thus, there is an S_2^1 -cut J , such that $S_2^1 \vdash \text{con}_{\rho(A)}(A) \rightarrow \text{con}_m^J(A)$.

Reason in S_2^1 . Consider any Σ_1 -sentence S . Suppose that $\Box_{S_2^1}(\text{con}_m(A) \rightarrow S)$. So, $\Box_{S_2^1}(\text{con}_m(A) \rightarrow S)^J$, and hence $\Box_{S_2^1}(\text{con}_m^J(A) \rightarrow S)$. Thus, we may conclude $\Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S)$. \square

If an arithmetic is deep, we can strengthen the implication in Σ_1 -veracity to a bi-implication.

Theorem 6.3. *Let A be a finitely axiomatized sequential theory. Suppose N is a deep arithmetic in A . We have:*

$$S_2^1 \vdash \forall S \in \Sigma_1\text{-sent} (\Box_A S^N \leftrightarrow \Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S)).$$

We leave the simple proof to the reader.

¹¹We use the work of [Ger03], [Ger05] and [Bus11] here. See Subsection A.5 of Appendix A.

Theorem 6.4. *Let A be a finitely axiomatized sequential theory. Each of the following classes is, S_2^1 -verifiably downwards closed under \preceq : the Σ_1 -veracious theories, the strong theories and the deep theories.*

We leave the simple proof to the reader.

Iterated inconsistencies take a simple form for Σ_1 -veracious arithmetics, as will be proved in the next theorem.

Theorem 6.5. *Suppose A is sequential and N is a Σ_1 -veracious arithmetic in A . We have:*

$$(\dagger_n) \quad S_2^1 \vdash \Box_A \Box_A^{N,n} \perp \leftrightarrow \Box_{S_2^1}^n \Box_A \perp.$$

Proof. We prove then left-to-right direction by induction on n . The case $n = 0$ is trivial.

Suppose we have (\dagger_n) . We prove (\dagger_{n+1}) . We have in S_2^1 ,

$$\begin{aligned} (5) \quad & \Box_A \Box_A^{N,n+1} \perp \rightarrow \Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow \Box_A \Box_A^{N,n} \perp) \\ (6) \quad & \rightarrow \Box_{S_2^1}(\Box_{A,\rho(A)} \perp \vee \Box_A \Box_A^{N,n} \perp) \\ (7) \quad & \rightarrow \Box_{S_2^1} \Box_A \Box_A^{N,n} \perp \\ (8) \quad & \rightarrow \Box_{S_2^1} \Box_{S_2^1}^n \Box_A \perp \\ (9) \quad & \rightarrow \Box_{S_2^1}^{n+1} \Box_A \perp \end{aligned}$$

Here step (8) uses the induction hypothesis.

The right-to-left direction is proved by a trivial induction. \square

Corollary 6.6. *Suppose A is a finitely axiomatized sequential theory and N is a Σ_1 -veracious arithmetic in A . We have:*

- i. $A \vdash \Box_A^{N,n+1} \perp \leftrightarrow (\Box_{S_2^1}^n \Box_A \perp)^N$.
- ii. $\text{I}\Delta_0 + \text{supexp} \vdash \Box_A \Box_A^{N,n} \perp \leftrightarrow \Box_A \perp$.

In the next theorem, we establish the existence of lots of deep arithmetics in a finitely axiomatized sequential theory. The proof of the theorem employs a form of the Friedman-Goldfarb-Harrington fixed point. See [Vis05] for a discussion of this fixed point.

Theorem 6.7. *For every finitely axiomatized sequential theory A , and, for every arithmetic N_0 in A , there is a deep arithmetic N in A with $N \preceq N_0$. This theorem is verifiable in S_2^1 .*

Proof. Let A be a finitely axiomatized sequential theory and let N_0 be an arithmetic in A .

Let true be a Σ_1 -truth predicate. For the construction of such a truth predicate, see [HP91, Chapter V(5)]. The two classical works on this subject are [Les78] and [PD82]. We will use the following two properties of true : for S in Σ_1 ,

$$\text{T1. } S_2^1 \vdash \text{true}(S) \rightarrow S,$$

T2. Suppose J is an S_2^1 -cut such that $S_2^1 \vdash x \in J \rightarrow 2^{2^x}$ exists. Then, we have $S_2^1 \vdash S^J \rightarrow \text{true}(S)$.

We remind the reader of the witness comparison ordering. We define, for any $C = \exists x C_0(x)$ and $D = \exists y D_0(y)$:

- $C \leq D := \exists x (C_0(x) \wedge \forall y < x \neg D_0(y))$,
- $C < D := \exists x (C_0(x) \wedge \forall y \leq x \neg D_0(y))$,
- $(C \leq D)^\perp := (D < C)$, and $(C < D)^\perp := (D \leq C)$.

By the Gödel Fixed Point Lemma, we find R such that, for a suitably chosen m :

$$S_2^1 \vdash R \leftrightarrow \text{true}(S) \leq \square_{A,m} R^{N_0}.$$

We note that the complexity $\rho(R)$ of R is *not* dependent on S and m , since numerals do not change the complexity of a formula even if the numeral is given a relational representation. Moreover, for any B and K , $\rho(B^K)$ is linear in $\rho(B)$. Hence, we may choose m so large that $\max(\rho(A), \rho(R^{N_0})) \leq m$.

We choose N_1 to be an initial segment of N_0 such that:

- U1. $A \vdash \square_{A,m}^{N_1} B \rightarrow B$, for any B with $\rho(B) \leq m$.
- U2. $A \vdash (\forall S \in \Sigma_1\text{-sent} (\text{true}(S) \rightarrow \text{true}(S) \leq \text{true}(S)))^{N_1}$,
i.o.w., A proves that, in N_1 , if $\text{true}(S)$ is witnessed, then $\text{true}(S)$ has a minimal witness.

We can always find such an N_1 since (i) we have a truth predicate for formulas of complexity $\leq m$ and since (ii) we can interpret $\text{ID}_0 + \Omega_1$ in S_2^1 .

Let J be an S_2^1 -cut such that $S_2^1 \vdash x \in J \rightarrow 2^{2^x}$ exists. We take $N := N_1 \circ J$.

We note that N_1 is strong, and that, hence, N is strong. We show that N is Σ_1 -veracious.

We reason, for the rest of the proof, in S_2^1 . Consider any Σ_1 -sentence S . Suppose $\square_A S^N$. It follows, by (T2) that $\square_A (\text{true}(S))^{N_1}$. Ergo, by (U2), $\square_A (R \vee R^\perp)^{N_1}$. Thus, $\square_A (R \vee \square_{U,m} R^{N_0})^{N_1}$. Hence, by (U1), $\square_A (R^{N_1} \vee R^{N_0})$. Since N_1 is a cut of N_0 , we get $\square_A R^{N_0}$.

By Lemma 3.1, we may conclude $(\dagger) \square_{S_2^1} \square_{A,m} R^{N_0}$.

We can find an S_2^1 -cut J^* , on which we have T_2^1 , so that if something is A -provable with a proof in J^* , then there is a minimal proof. We can arrange that J^* is so small that $S_2^1 \vdash x \in J^* \rightarrow 2^{2^x}$ exists. We find from (\dagger) : $\square_{S_2^1} (\square_{A,m} R^{N_0})^{J^*}$. Hence, $\square_{S_2^1} (R \vee R^\perp)^{J^*}$. We may conclude:

$$\square_{S_2^1} (\text{true}(S) \vee (R^\perp \wedge \square_{A,m} R^{N_0}))^{J^*}.$$

Since we have Σ_1 -completeness in the presence of double exponentiation, it follows that:

$$\square_{S_2^1} (\text{true}(S) \vee \square_{A,m} (R^\perp \wedge R)^{N_0}).$$

Hence, $\square_{S_2^1} (S \vee \square_{A,m} \perp)$, or, in a different formulation: $\square_{S_2^1} (\text{con}_m(A) \rightarrow S)$.

By Theorem 6.2, we may conclude that $\square_{S_2^1} (\text{con}_{\rho(A)}(A) \rightarrow S)$.

We have proved that, for every finitely axiomatized sequential theory A , and, for every arithmetic N_0 in A , there is a deep arithmetic N in A with $N \preceq N_0$. To see that this argument is verifiable in S_2^1 , we have to see that the construction of N from N_0 is feasible. We note that m in our argument remains standard even if S is non-standard. As a consequence e.g., $\delta_N = \phi(\delta_{N_0}, Z_{N_0}, S_{N_0}, \dots)$, where ϕ is a fixed standard context. Thus N will be p-time in N_0 . \square

Discussion 6.8. Clearly, the Second Incompleteness Theorem implies that adding $\text{con}(U)$ to a consistent theory U that contains an arithmetic, gives us a stronger theory, a theory that is, so to say, *one gödel* stronger. However, it is clear that we have to ask: to what arithmetic in U are we adding the consistency statement?

Consider GB and let neumann be the interpretation of S_2^1 in the finite von Neumann ordinals. Clearly,

$$\text{PA} \not\vdash \text{con}(\text{GB}) \rightarrow \text{con}(\text{GB} + \text{con}^{\text{neumann}}(\text{GB})).$$

In fact, by the Second Incompleteness Theorem, GB cannot prove this statement w.r.t. the neumann -interpretation. However, for a Σ_1 -veracious arithmetic N in GB, we have:

$$\text{I}\Delta_0 + \text{supexp} \vdash \text{con}(\text{GB}) \rightarrow \text{con}(\text{GB} + \text{con}^N(\text{GB})).$$

Thus, in which theories the relative consistency of a theory plus its consistency statement w.r.t. that theory can be verified is dependent on the chosen arithmetic. Adding $\text{con}^N(\text{GB})$ adds less strength to GB than adding $\text{con}^{\text{neumann}}(\text{GB})$ does.

So *the gödel* is not such a good unit when we define it as *how much stronger a theory becomes when we add its consistency statement*. My proposal would be to take as the theory that is *one gödel* stronger: $S_2^1 + \text{con}(U)$. Note that the strength of $S_2^1 + \text{con}(U)$ still depends on the chosen axiomatization of U .

In the case of a finitely axiomatized sequential theory A and a Σ_1 -veracious arithmetic N in A , we have:

$$S_2^1 \vdash \text{con}(A + \text{con}^N(A)) \leftrightarrow \text{con}(S_2^1 + \text{con}(A)).$$

So, by the measure of S_2^1 -verifiable relative consistency, adding the consistency statement for a Σ_1 -veracious arithmetic in A , is adding one gödel. Note that there are no arithmetics in, say, PA with the same property. \square

The next theorem shows that under a verifiability condition, Theorem 6.7 can be strengthened to theories that are mutually interpretable with a finitely axiomatized sequential theory.

Theorem 6.9. *Suppose A is a consistent finitely axiomatized sequential theory and U is any theory. Suppose $K : A \triangleright U$ and $M : U \triangleright A$. Then, there is an arithmetic N in U , such that:*

$$S_2^1 \vdash K : A \triangleright_{\text{thm}} U \rightarrow \forall S \in \Sigma_1\text{-sent} (\Box_U S^N \rightarrow \Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S)).$$

Here $\triangleright_{\text{thm}}$ stands for theorems-interpretability, where we demand that the interpreting theory proves the translations of the theorems of the interpreted theory. In the context of arithmetics without Σ_1 -collection this notion is not provably equivalent to the usual notion of axioms-interpretability.

Proof. Suppose $K : A \triangleright U$ and $M : U \triangleright A$. We find $\Box_{S_2^1}(M : U \triangleright A)$. Consider any arithmetic N_0 in A . We note that $N_1 := N_0MK$ is also an arithmetic in A .¹² Let N_2 be an arithmetic in A such that $N_2 \preceq N_0$ and $N_2 \preceq N_1$. We may assume that in N_2 we have $I\Delta_0 + \Omega_1$, or a sufficiently large, finitely axiomatized part of $I\Delta_0 + \Omega_1$.

Let k be the complexity of $(\text{true}(x))^{N_0}$, where true is the Σ_1 -truth predicate. By Theorem 3.7, we can choose m so large that, S_2^1 -verifiably,

$$(10) \quad \forall B \in \Gamma_k (\Box_A(\Box_{A,m}^{N_2} B \rightarrow B) \rightarrow \Box_{A,m} B).$$

Let $N_3 \leq N_0$ be an arithmetic in A such that, verifiably in S_2^1 ,

$$(11) \quad \forall B \in \Gamma_m \Box_A(\Box_{A,m}^{N_3} B \rightarrow B).$$

Let N_4 be a cut of N_3 such that

$$(12) \quad \Box_A \forall x \in N_4 \exists y \in N_3 2^{2^x} = y.$$

Finally, we take $N := N_4M$. So N is an arithmetic in U .

By the Gödel Fixed Point Lemma, we find R such that:

$$(13) \quad \Box_{S_2^1}(R \leftrightarrow \text{true}(S) \leq \Box_{A,m} R^{N_0}).$$

We reason in S_2^1 .

We have, for all Σ_1 -sentences S ,

$$(14) \quad \Box_A (S^{N_4} \rightarrow (\text{true}(S))^{N_3})$$

$$(15) \quad \rightarrow (R \vee R^\perp)^{N_3}$$

$$(16) \quad \rightarrow (R^{N_0} \vee \Box_{A,m}^{N_3} R^{N_0})$$

$$(17) \quad \rightarrow R^{N_0}$$

Suppose $K : A \triangleright_{\text{thm}} U$. We also have: $M : U \triangleright A$. Consider any S and suppose $\Box_U S^N$. It follows that $\Box_A S^{N_4MK}$. By the previous result, we may conclude $\Box_A R^{N_0MK}$, i.e., $\Box_A R^{N_1}$. From this, we have:

$$(18) \quad \Box_A (\Box_{A,m}^{N_2} R^{N_0} \rightarrow R^{N_1} \wedge \Box_{A,m}^{N_2} R^{N_0})$$

$$(19) \quad \rightarrow R^{N_2}$$

$$(20) \quad \rightarrow R^{N_0}$$

So, by Equation 10, we have: $\Box_A R^{N_0}$. It follows that $\Box_{S_2^1} \Box_{A,m} R^{N_0}$. We now may repeat the reasoning of the proof of Theorem 6.7. So we get

$$\Box_{S_2^1}(\text{con}_{\rho(A)}(A) \rightarrow S).$$

And we are done. □

¹²We remind the reader that N_0MK stands for $K \circ M \circ N_0$.

7. AN EXAMPLE

In this section, we provide an example of a sequential theory W , such that the degrees of the arithmetics in W are finite and cofinal in ω . So, for every n there is an arithmetic in W with degree $k \geq n$, but there is no arithmetic in W with degree ∞ .

We start with a consistent finitely axiomatized sequential theory A . Pick any arithmetic N in A . Let Σ be the signature of A and let Θ be the signature of arithmetic.

Let $\tau : \Theta \rightarrow \Sigma$ be a translation. We define $\tilde{\tau} := \langle S_2^1, \tau \langle (S_2^1)^\tau \rangle_{\tau_N}, A \rangle$. Here we assume that the axioms of identity are an explicit part of the axiomatization of S_2^1 . It is easily seen that $\tilde{\tau}$ is an arithmetic in A . We assign to any translation $\tau : \Theta \rightarrow \Sigma$ a Gödel number $\text{gn}(\tau)$. We define:

$$W := A + \{(\Box_{S_2^1}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}} \mid \tau : \Theta \rightarrow \Sigma\}.^{13}$$

We note that there is a p-time algorithm to decide whether a sentence is of the form $(\Box_{S_2^1}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}$. So, W is Δ_1^b -axiomatized.¹⁴

Consider any arithmetic K in W . We have $W \vdash (\Box_{S_2^1}^{\text{gn}(\tau_K)} \Box_A \perp)^{\tilde{\tau}_K}$. Clearly, in W , the interpretations K and $\tilde{\tau}_K \upharpoonright W$ coincide, hence, by an easy induction, $W \vdash \Box_W^{K, \text{gn}(\tau_K)+1} \perp$. So, each arithmetic K in W has a finite degree.

We show that, for any n , the theory W contains an arithmetic with degree $\geq n$. Consider any number n .

Let N^* be a strong arithmetic in A that has an initial embedding in all arithmetics $\tilde{\tau}$ in A with $\text{gn}(\tau) \leq n$. Let $N^\circ \preceq (N^* \upharpoonright (A + \Box_A^{N^*} \perp))$ be a deep arithmetic in $A + \Box_A^{N^*} \perp$. Let $N_\circ := \tilde{\tau}_{N^\circ}$. We note that N_\circ is an arithmetic in A .

We want to show that $W \not\vdash \Box_W^{N_\circ \upharpoonright W, n} \perp$. This will be a direct consequence of the following claim.

Claim: We have, $S_2^1 \vdash \Box_W \Box_W^{N_\circ \upharpoonright W, n} \perp \rightarrow \Box_{S_2^1}^n \Box_A \perp$.

We first show how our desired result follows from the claim. Suppose $W \vdash \Box_W^{N_\circ \upharpoonright W, n} \perp$. Since, S_2^1 is a true theory, the claim gives us, by applying reflection a number of times, $\Box_A \perp$. Quod non. Note that this argument can be formalized in $\text{ID}_0 + \text{supexp}$.

Proof of the claim: By our conventions, we may write $\Box_W^{N_\circ \upharpoonright W, n} \perp$ as $\Box_W^{N_\circ, n} \perp$. We will apply this convention to increase readability. We prove by induction that, for each $j \leq n$,

$$(\$_j) \quad S_2^1 \vdash \Box_W \Box_W^{N_\circ, j} \perp \rightarrow \Box_{S_2^1}^j \Box_A \perp.$$

¹³In stead of $(\Box_{S_2^1}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}$, we could also have used $\Box_A^{\tilde{\tau}, \text{gn}(\tau)+1} \perp$, but it seems to me that the argument is a bit shorter under the present choice.

¹⁴Alternatively, we could have constructed W using a version of Craig's trick, taking as axioms $(\underline{p} = \underline{p}) \wedge (\Box_{S_2^1}^{\text{gn}(N)} \Box_A \perp)^N$, where p is an A -proof that N is an arithmetic.

For the case $j = 0$, we have to prove: $S_2^1 \vdash \Box_W \perp \rightarrow \Box_A \perp$. We reason in S_2^1 . Suppose $\Box_W \perp$. Consider any W -proof p of \perp . If p only employs the axiom A , we are done. Suppose p employs at least one axiom of the form $(\Box_{S_2^1}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}$. Let X be the set of such axioms employed in p . By our assumption X is not empty.

We construct an arithmetic M in A , such that $M \preceq \tilde{\tau}$, for all $\tau \in X$. Suppose $\tau_0 \in X$. For each τ we construct an initial embedding F_τ of an initial segment J_τ of $\tilde{\tau}_0$ in $\tilde{\tau}$. This construction is uniform in τ and $|F_\tau|$ is linear in $|\tau|$. We take M to be the intersection of the J_τ .

The definition of M involves e.g. a conjunction $\delta_M(x) : \leftrightarrow \bigwedge_{\tau \in X} J_\tau(x)$. Why is this conjunction not too big? Under reasonable assumptions, we have:

$$\text{gn}(\tau) \leq \text{gn}((\Box_{S_2^1}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}) \leq p.$$

Moreover, $|J_\tau|$ is linear in $|\text{gn}(\tau)|$ and, hence, linear in $|p|$. Moreover the size of X is $< |p|$. Hence $|M|$ is bounded by $L(|p|) \times |p|$, where L is a standard linear term. So M is below $a \cdot (\omega_1(p))^b$, for some standard a and b .

Clearly, each axiom in X is implied by $\Box_A^M \perp$. So, we have $\Box_{A+\Box_A^M} \perp$, and, hence $\Box_A \perp$, by the Second Incompleteness Theorem.

An alternative way, to prove this, is to adapt the proof of the second incompleteness theorem as follows. We still reason inside S_2^1 . By the Gödel Fixed Point Lemma, we construct G such that:

$$\Box_{S_1^1}(G \leftrightarrow \neg \Box_A \bigvee_{\tau \in X} G^{\tilde{\tau}}).$$

Note that the big disjunction exists inside S_2^1 , since the set X is derived from p . We have:

$$(21) \quad \Box_{S_2^1}(\neg G \rightarrow \Box_A \bigvee_{\tau \in X} G^{\tilde{\tau}})$$

$$(22) \quad \rightarrow (\Box_A \bigvee_{\tau \in X} G^{\tilde{\tau}} \wedge \Box_A \bigwedge_{\nu \in X} \Box_A^{\tilde{\nu}} \bigvee_{\tau \in X} G^{\tilde{\tau}})$$

$$(23) \quad \rightarrow (\Box_A \bigvee_{\tau \in X} G^{\tilde{\tau}} \wedge \Box_A \bigwedge_{\nu \in X} \neg G^{\tilde{\nu}})$$

$$(24) \quad \rightarrow \Box_A \perp$$

Step (22), uses the fact that we have $\exists \Sigma_1^b$ -completeness for every $\tilde{\tau}$.

From our assumption on X , it clearly follows that $\Box_{A+\bigwedge_{\tau \in X} \Box_A^{\tilde{\tau}} \perp}$. Hence, we find $\Box_A \bigvee_{\tau \in X} \text{con}^{\tilde{\tau}}(A)$, and, so, by (24), $\Box_A \bigvee_{\tau \in X} G^{\tilde{\tau}}$. We may conclude $\Box_A \perp$.

The nice feature of this second argument is that it does not use sequentiality.

We stop reasoning in S_2^1 .

We now prove $(\$_{j+1})$, for $j+1 \leq n$, where we use the induction hypothesis $(\$_j)$.

We reason again in S_2^1 .

Suppose $\Box_W \Box_W^{N_0 \cdot j+1} \perp$. Our induction hypothesis, $\$ _j$, gives us: $\Box_W (\Box_{S_2^1}^j \Box_A \perp)^{N_0}$.

Let p be a proof witnessing $\Box_W(\Box_{S_2}^j \Box_A \perp)^{N^\circ}$. Let X be the set of all τ such that $(\Box_{S_2}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}$ occurs as an axiom in p . We clearly have:

$$(25) \quad \text{proof}_{A+\{(\Box_{S_2}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}} \mid \tau \in X\}} (p, (\Box_{S_2}^j \Box_A \perp)^{N^\circ}).$$

Let X_0 be the set of elements τ of X with $\text{gn}(\tau) \leq j$, and let X_1 be the set of elements of X with $\text{gn}(\tau) > j$. Since $j+1 \leq n$, we find that $N^* \preceq \tilde{\tau}$, for any τ with $\text{gn}(\tau) \leq j$. It follows that, inside \Box_{S_2} , $\Box_A^{N^*} \perp$ implies $(\Box_{S_2}^{\text{gn}(\tau)} \Box_A \perp)^{\tilde{\tau}}$, for each τ in X_0 .

Reasoning as in the case $j = 0$, we can find an arithmetic M^* in A , such that $M^* \preceq \tilde{\tau}$, for all $\tau \in X_1$. Moreover, we can choose M^* in such a way that it is deep.

We have to move in a careful way, at this point, to compensate for our lack of induction. Clearly, we can find a standard cut J , such that:

$$(26) \quad \Box_{S_2} \forall z \in J (\Box_{S_2}^j \Box_A \perp \rightarrow \Box_{S_2}^{j+z} \Box_A \perp).$$

It follows that for p-time computable f with standard code:

$$(27) \quad \forall \tau \in X_1 \exists q < f(\text{gn}(\tau)) \text{proof}_{S_2^1}(q, (\Box_{S_2}^j \Box_A \perp \rightarrow \Box_{S_2}^{\text{gn}(\tau)} \Box_A \perp)).$$

We note that, we can find a p-time computable g with standard code such that:

$$(28) \quad \forall \tau \in X_1 \exists r < g(\text{gn}(\tau)) \text{proof}_A(r, (\Box_{S_2}^{M^*, \text{gn}(\tau)} \Box_A \perp \rightarrow \Box_{S_2}^{\tilde{\tau}, \text{gn}(\tau)} \Box_A \perp)).$$

Using Equations (27) and (28), we can transform the proof p of Equation (25), to a proof s witnessing the following provability:

$$(29) \quad \Box_{A+\Box_A^{N^*} \perp + (\Box_{S_2}^{j+1} \Box_A \perp)^{M^*}} (\Box_{S_2}^j \Box_A \perp)^{N^\circ}.$$

Hence, using $\boxplus B$ for $B \wedge \Box B$:

$$(30) \quad \boxplus_{S_2^1} (\Box_A (\Box_{S_2}^{j+1} \Box_A \perp)^{M^*} \rightarrow \Box_{A+\Box_A^{N^*} \perp} (\Box_{S_2}^j \Box_A \perp)^{N^\circ}).$$

Since, $N^\circ \uparrow (A + \Box_A^{N^*} \perp) = N^\circ$, we have:

$$(31) \quad \boxplus_{S_2^1} (\Box_A (\Box_{S_2}^{j+1} \Box_A \perp)^{M^*} \rightarrow \Box_{A+\Box_A^{N^*} \perp} (\Box_{S_2}^j \Box_A \perp)^{N^\circ}).$$

Since M^* is deep, we find:

$$(32) \quad \boxplus_{S_2^1} (\Box_{S_2^1} (\text{con}_{\rho(A)}(A) \rightarrow \Box_{S_2}^{j+1} \Box_A \perp) \rightarrow \Box_{A+\Box_A^{N^*} \perp} (\Box_{S_2}^j \Box_A \perp)^{N^\circ}).$$

Hence,

$$(33) \quad \boxplus_{S_2^1} (\Box_{S_2^1}^{j+2} \Box_A \perp \rightarrow \Box_{A+\Box_A^{N^*} \perp} (\Box_{S_2}^j \Box_A \perp)^{N^\circ}).$$

Since N° is deep, we get:

$$(34) \quad \boxplus_{S_2^1} (\Box_{S_2^1}^{j+2} \Box_A \perp \rightarrow \Box_{S_2^1} (\text{con}_{\rho(A+\Box_A^{N^*} \perp)}(A + \Box_A^{N^*} \perp) \rightarrow \Box_{S_2}^j \Box_A \perp)).$$

It follows that:

$$(35) \quad \boxplus_{S_2^1} (\Box_{S_2^1}^{j+2} \Box_A \perp \rightarrow \Box_{S_2^1} (\Box_A \neg \Box_A^{N^*} \perp \vee \Box_{S_2}^j \Box_A \perp)).$$

By the Second Incompleteness Theorem, we have:

$$(36) \quad \boxplus_{S_2^1} (\Box_{S_2^1}^{j+2} \Box_A \perp \rightarrow \Box_{S_2^1} (\Box_A \perp \vee \Box_{S_2}^j \Box_A \perp)).$$

Ergo:

$$(37) \quad \boxplus_{S_2^1}(\boxdot_{S_2^1}^{j+2}\Box_A\perp \rightarrow \boxdot_{S_2^1}^{j+1}\Box_A\perp).$$

So, by Löb's Theorem,

$$(38) \quad (\boxdot_{S_2^1}^{j+2}\Box_A\perp \rightarrow \boxdot_{S_2^1}^{j+1}\Box_A\perp) \wedge \boxdot_{S_2^1}^{j+2}\Box_A\perp$$

We may conclude:

$$(39) \quad \boxdot_{S_2^1}^{j+1}\Box_A\perp$$

This is what we wanted to prove. \square

We end this section by providing a bit of information on the relationship between A and W . By Theorem 6.9, it follows that W is not interpretable in A . On the other hand, it turns out that W is model-interpretable in A . Consider any model \mathcal{M} of A . In case for all arithmetics N in \mathcal{M} we have $\Box_A^N\perp$, we find that $\mathcal{M} \models W$. So we can take the identity translation to provide an inner model of M . Suppose, for some N , we have $\mathcal{M} \models \text{con}^N(A)$. We note that the proof of $(\$_0)$ works for an arbitrary arithmetic. So we have $\mathcal{M} \models \text{con}^N(W)$. We use the Henkin interpretation to provide an inner model of W .

REFERENCES

- [AB04] S.N. Artemov and L.D. Beklemishev. Provability logic. In D. Gabbay and F. Guenther, editors, *Handbook of Philosophical Logic, 2nd ed.*, volume 13, pages 229–403. Springer, Dordrecht, 2004.
- [Boo93] G. Boolos. *The logic of provability*. Cambridge University Press, Cambridge, 1993.
- [BS91] G. Boolos and G. Sambin. Provability: the emergence of a mathematical modality. *Studia Logica*, 50:1–23, 1991.
- [Bus86] S.R. Buss. *Bounded Arithmetic*. Bibliopolis, Napoli, 1986.
- [Bus11] S.R. Buss. Cut elimination *in situ*. <http://math.ucsd.edu/~sbuss/>, 2011.
- [BV05] L.D. Beklemishev and A. Visser. On the limit existence principles in elementary arithmetic and Σ_n^0 -consequences of theories. *Annals of Pure and Applied Logic*, 136:56–74, 2005.
- [BV06] L.D. Beklemishev and A. Visser. Problems in the Logic of Provability. In Dov M. Gabbay, Sergei S. Concharov, and Michael Zakharyashev, editors, *Mathematical Problems from Applied Logic I, Logics for the XXIst Century*, volume 4 of *International Mathematical Series*, pages 77–136. Springer, New York., 2006.
- [CFL11] A. Cordon-Franco, A. Fernández-Margarit, and F.F. Lara-Martín. A note on parameter free Π_1 -induction and restricted exponentiation. *Mathematical Logic Quarterly*, 57(5):444–455, 2011.
- [dJMM91] D.H.J. de Jongh, M. Jumelet, and F. Montagna. On the proof of Solovay's theorem. *Studia Logica*, 50:51–70, 1991.
- [Fef60] S. Feferman. Arithmetization of metamathematics in a general setting. *Fundamenta Mathematicae*, 49:35–92, 1960.
- [Fef97] S. Feferman. My route to arithmetization. *Theoria*, 63:168–181, 1997.
- [Ger03] Philipp Gerhardy. Refined Complexity Analysis of Cut Elimination. In Matthias Baaz and Johann Makovsky, editors, *Proceedings of the 17th International Workshop CSL 2003*, volume 2803 of *LNCS*, pages 212–225. Springer-Verlag, Berlin, 2003.
- [Ger05] Philipp Gerhardy. The Role of Quantifier Alternations in Cut Elimination. *Notre Dame Journal of Formal Logic*, 46, no. 2:165–171, 2005.
- [GS79] D. Guaspari and R.M. Solovay. Rosser sentences. *Annals of Mathematical Logic*, 16:81–99, 1979.
- [Hod93] W. Hodges. *Model theory*. Encyclopedia of Mathematics and its Applications, vol. 42. Cambridge University Press, Cambridge, 1993.

- [HP91] P. Hájek and P. Pudlák. *Metamathematics of First-Order Arithmetic*. Perspectives in Mathematical Logic. Springer, Berlin, 1991.
- [JdJ98] G. Japaridze and D. de Jongh. The logic of provability. In S. Buss, editor, *Handbook of proof theory*, pages 475–546. North-Holland Publishing Co., Amsterdam, 1998.
- [KPD88] R. Kaye, J. Paris, and C. Dimitracopoulos. On parameter-free induction schemas. *Journal of Symbolic Logic*, 53(4):1082–1097, 1988.
- [Kra87] J. Krajčček. A note on proofs of falsehood. *Archiv für Mathematische Logik und Grundlagenforschung*, 26:169–176, 1987.
- [Les78] H. Lesan. *Models of arithmetic*. Dissertation. University of Manchester, Manchester, 1978.
- [Lin96] P. Lindström. Provability logic – a short introduction. *Theoria*, 62(1-2):19–61, 1996.
- [MPS90] J. Mycielski, P. Pudlák, and A.S. Stern. *A lattice of chapters of mathematics (interpretations between theorems)*, volume 426 of *Memoirs of the American Mathematical Society*. AMS, Providence, Rhode Island, 1990.
- [PD82] J. B. Paris and C. Dimitracopoulos. Truth definitions and Δ_0 formulae. In *Logic and algorithmic*, Monographie de L’Enseignement Mathématique 30, pages 317–329, Geneve, 1982.
- [Pud83] P. Pudlák. Some prime elements in the lattice of interpretability types. *Transactions of the American Mathematical Society*, 280:255–275, 1983.
- [Pud85] P. Pudlák. Cuts, consistency statements and interpretations. *The Journal of Symbolic Logic*, 50:423–441, 1985.
- [Pud86] P. Pudlák. On the length of proofs of finitistic consistency statements in finitistic theories. In J.B. Paris, A.J. Wilkie, and G.M. Wilmers, editors, *Logic Colloquium ’84*, pages 165–196. North-Holland, 1986.
- [Smo85a] C. Smoryński. Nonstandard models and related developments. In L.A. Harrington, M.D. Morley, A. Scedrov, and S.G. Simpson, editors, *Harvey Friedman’s Research on the Foundations of Mathematics*, pages 179–229. North Holland, Amsterdam, 1985.
- [Smo85b] C. Smoryński. *Self-Reference and Modal Logic*. Universitext. Springer, New York, 1985.
- [Šve00] V. Švejdar. On provability logic. *Nordic Journal of Philosophical Logic*, 4(2):95–116, 2000.
- [Ver93] L.C. Verbrugge. *Efficient metamathematics*. ILLC-dissertation series 1993-3, Amsterdam, 1993.
- [Vis81] A. Visser. *Aspects of diagonalization and provability*. Ph.D. Thesis, Department of Philosophy, Utrecht University, 1981.
- [Vis90] A. Visser. Interpretability logic. In P.P. Petkov, editor, *Mathematical logic, Proceedings of the Heyting 1988 summer school in Varna, Bulgaria*, pages 175–209. Plenum Press, Boston, 1990.
- [Vis91a] A. Visser. The Σ_1^0 -conservativity of Σ_1^0 -completeness. *Notre Dame Journal of Formal Logic*, 32:554–561, 1991.
- [Vis91b] A. Visser. The formalization of interpretability. *Studia Logica*, 51:81–105, 1991.
- [Vis92] A. Visser. An inside view of EXP. *The Journal of Symbolic Logic*, 57:131–165, 1992.
- [Vis93] A. Visser. The unprovability of small inconsistency. *Archive for Mathematical Logic*, 32:275–298, 1993.
- [Vis05] A. Visser. Faith & Falsity: a study of faithful interpretations and false Σ_1^0 -sentences. *Annals of Pure and Applied Logic*, 131:103–131, 2005.
- [Vis09] A. Visser. Cardinal arithmetic in the style of baron von Münchhausen. *Review of Symbolic Logic*, 2(3):570–589, 2009. doi: 10.1017/S1755020309090261.
- [Vis10] A. Visser. What is the right notion of sequentiality? Logic Group Preprint Series 288, Department of Philosophy, Utrecht University, Heidelberglaan 8, 3584 CS Utrecht, (<http://www.phil.uu.nl/preprints/lgps/>), 2010.
- [Vis11] A. Visser. Can we make the Second Incompleteness Theorem coordinate free. *Journal of Logic and Computation*, 21(4):543–560, 2011. First published online August 12, 2009, doi: 10.1093/logcom/exp048.
- [Vis12] A. Visser. III_1^- is locally interpretable in PA^- . Unpublished note, 2012.
- [WP87] A. Wilkie and J.B. Paris. On the scheme of induction for bounded arithmetic formulas. *Annals of Pure and Applied Logic*, 35:261–302, 1987.

APPENDIX A. MORE DETAILS ON THE BASICS

In this appendix we explain some basic notions in somewhat more detail.

A.1. Translations and Interpretations. The notion of interpretation that we will employ in this paper will be *m-dimensional interpretation without parameters*. There are two extensions of this notion: we can consider piecewise interpretations and we can add parameters. We refrain from considering piecewise interpretations. We explain why in Subsection A.3. We sketch a few basic ingredients of adding parameters in Subsection A.4. We explain why, in the sequential case, addition of parameters makes no difference for the provability logic of all arithmetics of a given theory in Remark 3.9.

Consider two signatures Σ and Θ . An *m*-dimensional translation $\tau : \Sigma \rightarrow \Theta$ is a quadruple $\langle \Sigma, \delta, \mathcal{F}, \Theta \rangle$, where $\delta(v_0, \dots, v_{m-1})$ is a Θ -formula and where for any *n*-ary predicate *P* of Σ , $\mathcal{F}(P)$ is a formula $A(\vec{v}_0, \dots, \vec{v}_{n-1})$ in the language of signature Θ , where $\vec{v}_i = v_{i_0}, \dots, v_{i_{(m-1)}}$. Both in the case of δ and *A* all free variables are among the variables shown. Moreover, if $i \neq j$ and $k \neq \ell$, then v_{ik} is syntactically different from $v_{j\ell}$.

We demand that we have $\vdash \mathcal{F}(P)(\vec{v}_0, \dots, \vec{v}_{n-1}) \rightarrow \bigwedge_{i < n} \delta(\vec{v}_i)$. Here \vdash is provability in predicate logic. This demand is inessential, but it is convenient to have.

We define B^τ as follows:

- $(P(x_0, \dots, x_{n-1}))^\tau := \mathcal{F}(P)(\vec{x}_0, \dots, \vec{x}_{n-1})$.
- $(\cdot)^\tau$ commutes with the propositional connectives.¹⁵
- $(\forall x A)^\tau := \forall \vec{x} (\delta(\vec{x}) \rightarrow A^\tau)$.
- $(\exists x A)^\tau := \exists \vec{x} (\delta(\vec{x}) \wedge A^\tau)$.

There are two worries about this definition. First, what variables \vec{x}_i on the side of the translation A^τ correspond with x_i in the original formula *A*? The second worry is that substitution of variables in δ and $\mathcal{F}(P)$ may cause variable clashes. These worries are never important in practice: we choose ‘suitable’ sequences \vec{x} to correspond to variables *x*, and we avoid clashes by α -conversions. However, if we want to give precise definitions of translations and e.g. of composition of translations these problems come into play. We will address these problems elsewhere.

We allow identity to be translated to a formula that is not identity. There is some tension between this choice and the treatment of identity as a logical constant. The reason is that the notion of logical constant can do several kinds of work. It may be *obligatory in the language* and it may be *preserved under translation*. For identity we only ask that it is obligatory.

¹⁵If we have a complex formula *A*, the translation A^τ could be satisfied in a model even if the sequences of values of the variables corresponding to the free variables in *A* are not in the domain of the translation in that model. One alternative option for the definition is to add a conjunction that stipulates that these sequences are in the domain. Thus, we would always have $\vdash A^\tau \rightarrow \delta_K(\vec{x})$, where \vec{x} is a sequence corresponding to a free variable in *A*. We will refrain from doing this. The cost is that, e.g., the definition of composition of translations becomes more complicated.

There are several important operations on translations.

- id_Σ is the identity translation. We take $\delta_{\text{id}_\Sigma}(v) := v = v$ and $\mathcal{F}(P) := P(\vec{v})$.
- We can compose translations. Suppose $\tau : \Sigma \rightarrow \Theta$ and $\nu : \Theta \rightarrow \Lambda$. Then $\nu \circ \tau$ or $\tau\nu$ is a translation from Σ to Λ . We define:

$$\begin{aligned} - \delta_{\tau\nu}(\vec{v}_0, \dots, \vec{v}_{m_\tau-1}) &:= \bigwedge_{i < m_\tau} \delta_\nu(\vec{v}_i) \wedge (\delta_\tau(v_0, \dots, v_{m_\tau-1}))^\nu. \\ - P_{\tau\nu}(\vec{v}_{0,0}, \dots, \vec{v}_{0,m_\tau-1}, \dots, \vec{v}_{n-1,0}, \dots, \vec{v}_{n-1,m_\tau-1}) &:= \\ &\bigwedge_{i < n, j < m_\tau} \delta_\nu(\vec{v}_{i,j}) \wedge (P(v_0, \dots, v_{n-1})^\tau)^\nu. \end{aligned}$$

- Let $\tau, \nu : \Sigma \rightarrow \Theta$ and let A be a sentence of signature Θ . We define the disjunctive translation $\sigma := \tau\langle A \rangle\nu : \Sigma \rightarrow \Theta$ as follows. We take $m_\sigma := \max(m_\tau, m_\nu)$. We write $\vec{v} \upharpoonright n$, for the restriction of \vec{v} to the first n variables, where $n \leq \text{length}(\vec{v})$.

$$\begin{aligned} - \delta_\sigma(\vec{v}) &:= (A \wedge \delta_\tau(\vec{v} \upharpoonright m_\tau)) \vee (\neg A \wedge \delta_\nu(\vec{v} \upharpoonright m_\nu)). \\ - P_\sigma(\vec{v}_0, \dots, \vec{v}_{n-1}) &:= (A \wedge P_\tau(\vec{v}_0 \upharpoonright m_\tau, \dots, \vec{v}_{n-1} \upharpoonright m_\tau)) \vee \\ &\quad (\neg A \wedge P_\nu(\vec{v}_0 \upharpoonright m_\nu, \dots, \vec{v}_{n-1} \upharpoonright m_\nu)) \end{aligned}$$

Note that in the definition of $\tau\langle A \rangle\nu$ we used a padding mechanism. In case e.g. $m_\tau < m_\nu$, the variables $v_{m_\tau}, \dots, v_{m_\nu-1}$ are used ‘vacuously’ when we have A . If we had piecewise interpretations, where domains are built up from pieces with possibly different dimensions, we could avoid padding by building the domain of disjoint pieces with different dimensions.

A translation relates signatures; an interpretation relates theories. An interpretation $K : U \rightarrow V$ is a triple $\langle U, \tau, V \rangle$, where U and V are theories and $\tau : \Sigma_U \rightarrow \Sigma_V$. We demand: for all axioms A of U , we have $V \vdash A^\tau$.

In the context of the formalization of interpretability, we have to distinguish between *axioms-interpretability*, which is the notion we just introduced and *theorems-interpretability*, where we demand that: for all theorems A of U , we have $V \vdash A^\tau$. In the real world these notions are equivalent, but we need a principle like Σ_1 -collection to prove that, so, for example Buss’ theory S_2^1 does not ‘know’ this equivalence. See [Vis91b] for more information about this matter.

Remark A.1. The design choice to make interpretations a transition between theories has many advantages. It allows us to build various categories of theories and interpretations; it allows us to have a decent model-functor on categories of theories and interpretations; in various arguments, it reminds us *where we are*, etc. However, in some cases, the typing regime is somewhat stifling. E.g., if you have an interpretation $K : U \rightarrow V$ and an extension W of V , then it would seem that K is also an interpretation of U in W . The typing regime forces us to say that it is a lifting $K \upharpoonright W : U \rightarrow W$, that is the interpretation based on τ_K , etc. In this paper we will remain faithful to the typing regime, but we will alleviate it a bit by the convention below. \square

- Suppose $K : U \rightarrow V$. We often write A^K for A^{τ_K} , in the context of a theory W that extends V .

Here are some further definitions.

- We write \overline{U} for the set of theorems of U . Suppose $K : U \rightarrow V$. We write $\overline{K} := \{A \mid V \vdash A^K\}$. We note that $\overline{U} \subseteq \overline{K}$. If $\overline{K} = \overline{U}$, we will say that K is *faithful*.
- $\text{ID}_U : U \rightarrow U$ is the interpretation $\langle U, \text{id}_{\Sigma_U}, U \rangle$.
- Suppose $\overline{U} \subseteq \overline{V}$. Then, $\mathcal{E}_{UV} : U \rightarrow V$ is $\langle U, \text{id}_{\Sigma_U}, V \rangle$.
- Suppose $K : U \rightarrow V$ and $M : V \rightarrow W$. Then, $KM := M \circ K : U \rightarrow W$ is $\langle U, \tau_M \circ \tau_K, W \rangle$.
- Suppose $K : W \rightarrow U$ and $U \subseteq V$. We write $K \uparrow V$ for $\mathcal{E}_{UV} \circ K$.
- Suppose $M : V \rightarrow Z$ and $U \subseteq V$. We write $U \downarrow M$ for $M \circ \mathcal{E}_{UV}$.
- Suppose $K : U \rightarrow (V + A)$ and $M : U \rightarrow (V + \neg A)$. Then $K \langle A \rangle M : U \rightarrow V$ is the interpretation $\langle U, \tau_K \langle A \rangle \tau_M, V \rangle$. In an appropriate category $K \langle A \rangle M$ is a special case of a product.

The notation $K : U \rightarrow V$ is inspired by the idea of interpretations as arrows in a category. There is also an intuition of interpretability as a generalization of provability. The traditional notations and notions associated to this intuition are:

- $K : U \triangleleft V$ stands for $K : U \rightarrow V$.
- $K : V \triangleright U$ stands for $K : U \rightarrow V$.
- $U \triangleleft V$ stands for $\exists K K : U \triangleleft V$. We say: U is *interpretable* in V .
- $V \triangleright U$ stands for $\exists K K : V \triangleright U$. We say: V *interprets* U .
- $U \triangleleft_{\text{loc}} V$ means: all finitely axiomatized subtheories U_0 of U are interpretable in V . We say that U is *locally interpretable* in V .
- $U \triangleleft_{\text{mod}} V$ means that, for every $\mathcal{M} \models V$, there is a translation τ such that $\tau(\mathcal{M}) \models U$. We say that U is *model-interpretable* in V .

A.2. i-morphisms. Consider an interpretation $K : U \rightarrow V$. We can view this interpretation as a uniform way of constructing internal models $\tau_K(\mathcal{M})$ of U from models \mathcal{M} of V . This construction gives us the contravariant model functor as soon as we have defined an appropriate category of interpretations.

Now consider two interpretations $K, M : U \rightarrow V$. Between the inner models $\tau_K(\mathcal{M})$ and $\tau_M(\mathcal{M})$ we have the usual structural morphisms of models. We are interested in the case where these morphisms are V -definable and uniform over models. This idea leads to the following definition. An i-morphism $M : K \rightarrow M$ is a triple $\langle K, F(\vec{u}, \vec{v}), M \rangle$, where $F(\vec{u}, \vec{v})$ is a V -formula and where \vec{u} has length m_K and \vec{v} has length m_M . We demand:

- $V \vdash F(\vec{u}, \vec{v}) \rightarrow (\delta_K(\vec{u}) \wedge \delta_M(\vec{v}))$,
- $V \vdash \delta_K(\vec{u}) \rightarrow \exists \vec{v} (\delta_M(\vec{v}) \wedge F(\vec{u}, \vec{v}))$,
- $V \vdash (\vec{u}_0 =_K \vec{u}_1 \wedge F(\vec{u}_0, \vec{v}_0) \wedge F(\vec{u}_1, \vec{v}_1)) \rightarrow \vec{v}_0 =_M \vec{v}_1$,
- $V \vdash (\vec{u}_0 =_K \vec{u}_1 \wedge \vec{v}_0 =_M \vec{v}_1 \wedge F(\vec{u}_0, \vec{v}_0)) \rightarrow F(\vec{u}_1, \vec{v}_1)$,
- $V \vdash (P_K(\vec{u}_0, \dots, \vec{u}_{n-1}) \wedge \bigwedge_{i < n} F(\vec{u}_i, \vec{v}_i)) \rightarrow P_M(\vec{v}_0, \dots, \vec{v}_{n-1})$.

Clearly, $F : K \rightarrow M$ is an i-morphism iff, for all models \mathcal{M} of V , $F^{\mathcal{M}}$ represents a morphism of models from $\tau_K(\mathcal{M})$ to $\tau_M(\mathcal{M})$.

Two i-morphisms $F, G : K \rightarrow M$ are *i-equal*, when $V \vdash \forall \vec{u}, \vec{v} (F(\vec{u}, \vec{v}) \leftrightarrow G(\vec{u}, \vec{v}))$.

In the obvious way, we can define the identity i-morphism $\text{id}_K : K \rightarrow K$, composition of i-morphisms, i-isomorphisms, etc. One can show that these operations preserve i-equality. Moreover, i-isomorphisms really are isomorphisms in the categories given by these operations.

We will say that two interpretations K, M are *i-equivalent* when there is an i-isomorphism between them, i.e. they are i-isomorphic.

We will *not* divide out i-equivalence of interpretations. This enables us to use the notation τ_M meaningfully, to speak about the dimension of an interpretation, etc. However, we demand that operations on interpretations preserve i-equivalence. It is easy to see that e.g. the operation $K, M \mapsto K\langle A \rangle M$ preserves i-equivalence. Moreover, if K and M are i-equivalent, then $\bar{K} = \bar{M}$.

One can show, by a simple compactness argument, that K and M are i-isomorphic iff, for every $\mathcal{M} \models V$, there is an F such that $F^{\mathcal{M}}$ represents an isomorphism between $\tau_K(\mathcal{M})$ and $\tau_M(\mathcal{M})$.

The category INT_1 is the category of theories (as objects) and interpretations modulo i-equivalence (as arrows). One may show that we have indeed defined a category. The relation of i-equivalence is preserved by composition, etcetera. Two theories U and V are *bi-interpretable* if they are isomorphic in INT_1 . Wilfrid Hodges calls this notion: *homotopy*. See [Hod93], p222.

Thus, U and V are bi-interpretable if there are interpretations $K : U \rightarrow V$ and $M : V \rightarrow U$, so that $M \circ K$ is i-isomorphic to id_U and $K \circ M$ is i-isomorphic to id_V . We call the pair K, M a *bi-interpretation* between U and V . One can show that the components of a bi-interpretation are faithful interpretations. Many good properties of theories like finite axiomatizability, decidability, κ -categoricity are preserved by bi-interpretations.

A.3. Piecewise Interpretations. There is a notion of *piecewise interpretability* where we allow the domain of the interpretation to be built up from finitely many pieces with possibly different dimensions. An example of this is the construction where we add points at infinity to a points-only version of plane geometry. We could have a piece with the original points and strict identity and a piece with pairs of distinct points with the following equivalence relation (x, y) is equivalent to (u, v) is there is no point w that is both collinear with x and y and with u and v . Of course, we can replace this interpretation by a one-piece interpretation that is isomorphic (in an appropriate sense) to it in various obvious ways. E.g. a pair (x, y) could represent a point at infinity if $x \neq y$ and the point x if $x = y$.

One can show that, if we have $V \vdash \exists x, y \ x \neq y$, then any piecewise interpretation is isomorphic (in an appropriate sense) to an interpretation without pieces (but in general with higher dimension). It follows that a theory piecewise interprets a weak arithmetic if and only if it interprets this arithmetic non-piecewise via an interpretation that is in a relevant sense *i-equivalent* to the original one.

A.4. Parameters. In general interpretations are allowed to have parameters. We will briefly sketch how to add parameters to our framework. We first define a translation with parameters. The parameters of the translation are given by a fixed sequence of variables \vec{w} that we keep apart from all other variables. A translation is defined as before, but for the fact that now the variables \vec{w} are allowed to occur in the domain and in the translations of the predicate symbols in addition to the variables that correspond to the argument places. Officially, we represent a translation $\tau_{\vec{w}}$ with parameters \vec{w} as a quintuple $\langle \Sigma, \delta, \vec{w}, F, \Theta \rangle$. The parameter sequence may be empty: in this case our interpretation is parameter-free.

An interpretation with parameters $K : U \rightarrow V$ is a quadruple $\langle U, \alpha, E, \tau_{\vec{w}}, V \rangle$, where $\tau_{\vec{w}} : \Sigma_U \rightarrow \Sigma_V$ is a translation and α is a V -formula containing at most \vec{w} free. The formula α represents the parameter domain. E.g., if we interpret the Hyperbolic Plain in the Euclidean Plain via the Poincaré interpretation, we need two distinct points to define a circular disk. These points are parameters of the construction, the parameter domain is $\alpha(w_0, w_1) = (w_0 \neq w_1)$. (For this specific example, we can also find a parameter-free interpretation.) The formula E represents an equivalence relation on the parameter domain. In practice this is always pointwise identity for parameter sequences, but for reasons of theory one must admit other equivalence relations too. We demand:

- $\vdash \delta_{\tau, \vec{w}}(\vec{v}) \rightarrow \alpha(\vec{w})$,
- $\vdash P_{\tau, \vec{w}}(\vec{v}_0, \dots, \vec{v}_{n-1}) \rightarrow \alpha(\vec{w})$.
- $V \vdash \exists \vec{w} \alpha(\vec{w})$;
- $V \vdash E(\vec{w}, \vec{z}) \rightarrow (\alpha(\vec{w}) \wedge \alpha(\vec{z}))$;
- V proves that E represents an equivalence relation on the sequences forming the parameter domain;
- $\vdash E(\vec{w}, \vec{z}) \rightarrow \forall \vec{x} (\delta_{\tau, \vec{w}}(\vec{x}) \leftrightarrow \delta_{\tau, \vec{z}}(\vec{x}))$;
- $\vdash E(\vec{w}, \vec{z}) \rightarrow \forall \vec{x}_0, \dots, \vec{x}_{n-1} (P_{\tau, \vec{w}}(\vec{x}_0, \dots, \vec{x}_{n-1}) \leftrightarrow P_{\tau, \vec{z}}(\vec{x}_0, \dots, \vec{x}_{n-1}))$;
- for all U -axioms A , $V \vdash \forall \vec{w} (\alpha(\vec{w}) \rightarrow A^{\tau, \vec{w}})$.

We can lift the various operations in the obvious way. Note that the parameter domain of $N := M \circ K$ and the corresponding equivalence relation should be:

- $\alpha_N(\vec{w}, \vec{u}_0, \dots, \vec{u}_{k-1}) := \alpha_M(\vec{w}) \wedge \bigwedge_{i < k} \delta_{\tau_M}(\vec{w}, \vec{u}_i) \wedge (\alpha_K(\vec{u}))^{\tau_M, \vec{w}}$.
- $E_N(\vec{w}, \vec{u}_0, \dots, \vec{u}_{k-1}, \vec{z}, \vec{v}_0, \dots, \vec{v}_{k-1}) := E_M(\vec{w}, \vec{z}) \wedge \bigwedge_{i < k} \delta_{\tau_M}(\vec{w}, \vec{u}_i) \wedge \bigwedge_{i < k} \delta_{\tau_M}(\vec{w}, \vec{v}_i) \wedge (E_K(\vec{u}, \vec{v}))^{\tau_M, \vec{w}}$.

Consider interpretations $K, M : U \rightarrow V$. An i-morphism $\phi : K \rightarrow M$ is a triple $\langle K, G, F, M \rangle$, where $G(\vec{u}, \vec{w})$ and $F(\vec{u}, \vec{w}, \vec{x}, \vec{y})$ are V -formulas.¹⁶ We write $F^{\vec{u}; \vec{w}}(\vec{x}, \vec{y})$ for F . We demand that:

- V proves that G is a surjective relation between α_K/E_K and α_M/E_M ;¹⁷

¹⁶In G and F we could allow extra parameters, \vec{z} , the *eigenparameters* of G and F . We will refrain from doing that here to unburden the presentation a bit.

¹⁷It seems a more logical choice to demand that G represents a function from α_K/E_K to α_M/E_M . There are also sound theoretical reasons for that choice. However, the definition of

- $V \vdash F^{\vec{u};\vec{w}}(\vec{x}, \vec{y}) \rightarrow G(\vec{u}, \vec{w})$;
- V proves that, if $G(\vec{u}, \vec{w})$, then $F^{\vec{u};\vec{w}}$ is a function from $\delta_K/=_K$ to $\delta_M/=_M$.
- V proves that if $E_K(\vec{u}_0, \vec{u}_1)$ and $E_M(\vec{w}_0, \vec{w}_1)$, then $F^{\vec{u}_0, \vec{w}_0}$ is the same function as $F^{\vec{u}_1, \vec{w}_1}$.

Finally, we say that two i-maps ϕ_0 and ϕ_1 are *i-equal* if V proves that G_{ϕ_0} and G_{ϕ_1} and F_{ϕ_0} and F_{ϕ_1} are the same.

The definitions of the identity i-morphism and of composition of i-morphisms are as is to be expected. We can compute what an i-isomorphism is: G is, V -verifiably, a bijection between α_K/E_K and α_M/E_M , and V proves that, if $G(\vec{u}, \vec{w})$, then $F^{\vec{u};\vec{w}}$ is a bijection between $\delta_K/=_K$ and $\delta_M/=_M$.

A.5. Complexity Measures. *Restricted provability* plays an important role in this paper. An n -proof is a proof from axioms with Gödel number smaller or equal than n only involving formulas of complexity smaller or equal than n . To work conveniently with this notion, a good complexity measure is needed. This should satisfy three conditions. (i) Eliminating terms in favour of a relational formulation should raise the complexity only by a fixed standard number. (ii) Translation of a formula via the translation corresponding to an interpretation K should raise the complexity of the formula by a fixed standard number depending only on K . (iii) The tower of exponents involved in cut-elimination should be of height linear in the complexity of the formulas involved in the proof.

Such a good measure of complexity together with a verification of desideratum (iii)—a form of nesting degree of quantifier alternations—is supplied in the work of Philipp Gerhardy. See [Ger03] and [Ger05]. It is also provided by Samuel Buss in his preliminary draft [Bus11]. Buss also proves that (iii) is fulfilled.

Gerhardy's measure corresponds to the following formula classes:

- AT is the class of atomic formulas.
- $\mathbf{N}_{-1}^* = \Sigma_{-1}^* = \Pi_{-1}^* := \emptyset$.
- $\mathbf{N}_n^* ::= \text{AT} \mid \neg \mathbf{N}_n^* \mid (\mathbf{N}_n^* \wedge \mathbf{N}_n^*) \mid (\mathbf{N}_n^* \vee \mathbf{N}_n^*) \mid (\mathbf{N}_n^* \rightarrow \mathbf{N}_n^*) \mid \forall \Pi_n^* \mid \exists \Sigma_n^*$.
- $\Sigma_n^* ::= \text{AT} \mid \neg \Pi_n^* \mid (\Sigma_{n-1}^* \wedge \Sigma_{n-1}^*) \mid (\Sigma_{n-1}^* \vee \Sigma_{n-1}^*) \mid (\Sigma_{n-1}^* \rightarrow \Sigma_{n-1}^*) \mid \forall \Pi_{n-1}^* \mid \exists \Sigma_n^*$.
- $\Pi_n^* ::= \text{AT} \mid \neg \Sigma_n^* \mid (\Pi_n^* \wedge \Pi_n^*) \mid (\Pi_n^* \vee \Pi_n^*) \mid (\Pi_n^* \rightarrow \Pi_n^*) \mid \forall \Pi_n^* \mid \exists \Sigma_{n-1}^*$.

We may define $\rho(A)$ as the minimal n such that A is in \mathbf{N}_n^* .¹⁸

Samuel Buss gives the following formula classes.

- $\Sigma_0^* = \Pi_0^* =$ the class of quantifier-free formulas.
- $\Sigma_n^* ::= \Sigma_{n-1}^* \mid \Pi_{n-1}^* \mid \neg \Pi_n^* \mid (\Sigma_n^* \wedge \Sigma_n^*) \mid (\Sigma_n^* \vee \Sigma_n^*) \mid (\Pi_n^* \rightarrow \Sigma_n^*) \mid \exists \Sigma_n^*$.
- $\Pi_n^* ::= \Sigma_{n-1}^* \mid \Pi_{n-1}^* \mid \neg \Sigma_n^* \mid (\Pi_n^* \wedge \Pi_n^*) \mid (\Pi_n^* \vee \Pi_n^*) \mid (\Sigma_n^* \rightarrow \Pi_n^*) \mid \forall \Pi_n^*$.

initial embedding that we need in Section 3 does not work under this second choice. So for the purposes of at least this paper we seem to need the definition given in the main text.

¹⁸Vincent van Oostrom gave a variant of this formulation of Gerhardy's measure in conversation.

We may define $\rho(A)$ as the smallest n such that A is in Σ_n^* . This is the same measure, as was employed in [Vis93]. For our purposes it does not matter whether we use Gerhardy's or Buss' definition.

APPENDIX B. ON FAITHFUL INTERPRETABILITY

We assume that the formalization of syntax is standard, so that the code of a subformula C of B is smaller than the code of B , etc. We also assume that the proof-predicate is standard, so that every proof p has a single conclusion C with $C < p$, etc.

Theorem B.1. *Consider a theory T and suppose that N is an arithmetic in T . Let Γ be any T -definable class of T -sentences for which T contains a definable truth predicate, say TRUE for sentences coded in N . We only need that TRUE satisfies Tarski's convention. We assume that the set of codes of elements of Γ has a fixed binumeration in T (which, par abus de langage, we call also Γ). So we assume: if $A \in \Gamma$, then $T \vdash A \in \Gamma$ and, if $A \notin \Gamma$, then $T \vdash A \notin \Gamma$. Then, there is a unary predicate $A(x)$, such that:*

- i. $T \vdash A(x) \rightarrow x \in N$.
- ii. $T \vdash (A(x) \wedge x =_N y) \rightarrow A(y)$.
- iii. $T \vdash (A(x) \wedge A(y)) \rightarrow x =_N y$.
- iv. For any n , $T + A(\underline{n})$ is Γ -conservative over T .
Here \underline{n} is the N -numeral of n .

Proof. We define, for $p \in N$ and B an N -code of a formula with at most one designated variable v_0 free:

$$\bullet \mathfrak{S}(p, x, B) :\leftrightarrow p, x \in N \wedge \exists C \in \Gamma (\text{proof}_T^N(p, B(x) \rightarrow C) \wedge \neg \text{TRUE}(C)).$$

Here $B(x)$ in the context of proof means the code of the result of substituting the numeral of x for v_0 in B . We find, using the Gödel Fixed Point Lemma, a formula A with the following property.

$$T \vdash A(x) \leftrightarrow \exists p (\mathfrak{S}(p, x, A) \wedge \forall q <_N p \forall y \in N \neg \mathfrak{S}(q, y, A)).$$

Clearly, we have (i) and (ii). We prove the uniqueness clause (iii).

Reason in T . Suppose that $x \neq_N y$ and $A(x)$ and $A(y)$. Let p be a witness for $A(x)$ and let q be a witness of $A(y)$. By our assumption about the proof predicate, it find that $p \neq_N q$, and, hence, $p <_N q$ or $q <_N p$. By the specification of A , this is impossible.

We move to the metatheory again. We prove (iv). We write $r : T \vdash E$ for: r is (a code of) a T -proof of E .

We assume, to get a contradiction, that, for some n , $A(n)$ is not Γ -conservative over T . Let p be the smallest proof such that, for some n and some $C \in \Gamma$, we have $p : T \vdash A(\underline{n}) \rightarrow C$ and $T \not\vdash C$. It follows that, for all $q < p$, and all m and all $D \in \Gamma$, if $q : T \vdash A(\underline{m}) \rightarrow D$, then $T \vdash D$. It follows that: $T \vdash \forall q <_N p \forall y \in N \neg \mathfrak{S}(q, y, A)$.

We also find: $T \vdash \neg C \rightarrow \exists(\underline{p}, \underline{n}, A)$. Ergo $T + \neg C \vdash A(\underline{n})$. In other words, we have both $T + A(\underline{n}) \vdash C$ and $T + \neg A(\underline{n}) \vdash C$. Ergo $T \vdash C$. A contradiction. \square

By the specification of the above formula A it functions as a closed partial N -numerical term in T . For this reason we will write $\tau \simeq_N x$ for $A(x)$.

Theorem B.2. *Let T be a theory and suppose that N is an arithmetic in T . Let Σ be a finite signature for predicate logic. We call predicate logic of signature Σ : FOL_Σ . Let $\alpha(x)$ be a formula in the language of T such that T proves that all elements of $\{x \mid \alpha(x)\}$ are N -codes of Σ -sentences. We write \Box_α for provability from the sentences coded by the elements of $\{x \mid \alpha(x)\}$. We write $\text{con}(\alpha)$ for $\neg \Box_\alpha \perp$.*

There is an interpretation $H : (T + \text{con}(\alpha)) \triangleright \text{FOL}_\Sigma$ such that, for any Σ -sentence A , we have $T + \text{con}(\alpha) + \Box_\alpha A \vdash A^H$. We say that H is a Henkin interpretation of α .

Proof. We can see this by inspection of the usual proof of the Interpretation Existence Lemma. The basic idea is that we formalize the Henkin construction, employing definable cuts whenever we would have used induction in PA. See e.g. [Vis91b] or [Vis92]. \square

We proceed with our upperbound result.

Theorem B.3. *Let T be any theory. Suppose $K : T \triangleright U$. Let A be any T -sentence and let N be an arithmetic in $T + A$. Then there is an interpretation $M : T \triangleright U$ such that, for any U -sentence B , $T \vdash B^M \Rightarrow T + A \vdash \Box_U^N B$.*

Proof. Consider $T + A$. We first show that we may assume without loss of generality that we have a Σ_1 -truth predicate for N .

By Theorem 3.8, we may shorten N to a $T + A$ -definable cut N' such that $T + A$ contains a truth predicate, say **TRUE**, for the Σ_1 -sentences of N' , i.e., for every S in Σ_1 , $U \vdash S^{N'} \leftrightarrow \text{TRUE}(S)$, where S inside the truth predicate is coded in N' .

Note that:

$$T + A \vdash \Box_U^{N'} B \Rightarrow T + A \vdash \Box_U^N B.$$

It follows that it is sufficient to prove our theorem for N' .

Thus, we may assume that T contains a truth predicate, say **TRUE**, for the Σ_1 -sentences of N .

Let τ be the partial closed term promised by Theorem B.1 for N and Σ_1 . We fix some standard enumeration C_x of the U -sentences in such a way that T verifies its elementary properties w.r.t. N . We specify M by cases. In case we have $\neg A$, we take M equal to K . Suppose we have A . We may now work in $T + A$. Let $U^* := U + \{C_x \mid \tau \simeq_N x\}$. Note that (i) U^* is not Δ_1^b -axiomatized, and that (ii) in talking about U^* we are really talking about the formula defining the axiom set and that (iii) the definition of U^* only makes sense in the presence of A . In case $\text{incon}^N(U^*)$,¹⁹ we take M again equal to K . If $\text{con}(U^*)$, we take M equal to the

¹⁹In writing $(\text{in})\text{con}^N(U^*)$, we intend no relativization of the formula defining the axiom set, only of the proofs.

Henkin-interpretation H of U^* . In other words, we take

$$M := H\langle A \wedge \text{con}^N(U^*) \rangle K.^{20}$$

Clearly, $M : T \triangleright U$. Suppose $T \vdash B^M$. Let $\neg B = C_n$. We have:

$$T + A + \tau \simeq_N \underline{n} \vdash "(U + \neg B) = U^*".$$

Here “=” stands for extensional identity. Hence,

$$T + A + (\tau = \underline{n}) + \text{con}^N(U + \neg B) \vdash \neg B^M.$$

Thus, $T + A \vdash (\tau \simeq_N \underline{n}) \rightarrow \Box_U^N B$. By the Σ_1 -conservativity of $\tau \simeq_N \underline{n}$, we find $T + A \vdash \Box_U^N B$. \square

From Theorem B.3 we can derive basic result about interpretability. We say that a theory U is *trustworthy* if, whenever $U \triangleright V$, then $U \triangleright_{\text{faith}} V$.

Theorem B.4. *The following are equivalent.*

- i. U is trustworthy.
- ii. U faithfully interprets predicate logic with a binary predicate R .
- iii. For some A , $U + A$ contains a Σ_1 -sound arithmetic N .

Proof. Trivially (ii) follows from (i).

Suppose (ii). Let B be the single axiom of adjunctive set theory AS. In AS we can provide a Σ_1 -sound interpretation M of S_2^1 . Suppose K is the promised faithful interpretations of predicate logic with a binary relation symbol R in U . Then, as is easily seen, $A := B^K$ and $N := K \circ M$ satisfy the desiderata of (iii).

Finally, we assume (iii). Suppose $K : U \triangleright V$. Let M be the interpretations of V provided by Theorem B.3, such that $M : U \triangleright V$ and $U \vdash B^M$ implies $U + A \vdash \Box_V^N B$. By the Σ_1 -soundness of N , we may conclude: that $U \vdash B^M$ implies $V \vdash B$, and we are done. \square

DEPARTMENT OF PHILOSOPHY, UTRECHT UNIVERSITY, JANSKERKHOF 13A, 3512BL UTRECHT, THE NETHERLANDS

E-mail address: `albert.visser@phil.uu.nl`

²⁰Strictly speaking we should not have K here but $K \upharpoonright (T + (\neg A \vee \text{incon}(U^*)))$.