

WHEN BI-INTERPRETABILITY IMPLIES SYNONYMY

HARVEY M. FRIEDMAN^{†,‡} AND ALBERT VISSER[‡]

ABSTRACT. Two salient notions of sameness of theories are *synonymy*, aka *definitional equivalence*, and *bi-interpretability*. Of these two *definitional equivalence* is the strictest notion. In which cases can we infer synonymy from bi-interpretability? We study this question for the case of sequential theories. Our result is as follows. Suppose that two sequential theories are bi-interpretable and that the interpretations involved in the bi-interpretation are one-dimensional and identity preserving. Then, the theories are synonymous. We provide an example to show that this result is optimal. There are two finitely axiomatized sequential theories that are bi-interpretable but not synonymous, where precisely one of the interpretations involved in the bi-interpretation is not identity preserving.

The crucial ingredient of our proof is a version of the Schröder-Bernstein theorem under very weak conditions. We think this result has some independent interest.

1. INTRODUCTION

When are two theories the same? Are there reasonable ways of abstracting away from the precise choice of the signature? The notions of *synonymy* (or: *definitional equivalence*) and *bi-interpretability* provide two good answers to these questions.

The notion of synonymy was introduced by de Bouvere (see [dB65a] and [dB65b]) in 1965. It appears to be the strictest notion of sameness of theories except strict identity of signature and set of theorems. Two theories U and V are synonymous iff there is a theory W that is both a definitional extension of U and of V . Equivalently, U and V are synonymous iff there are interpretations $K : U \rightarrow V$ and $M : V \rightarrow U$, such that V proves that the composition $K \circ M$ is the identity interpretation on V and such that U proves that the composition $M \circ K$ is the identity interpretation on U . (Thus, synonymy is isomorphism in an appropriate category INT_0 of theories and interpretations.) For example, Peano Arithmetic, PA, is synonymous with an appropriate theory of strings.

The notion of *bi-interpretability* was introduced by Alhbrandt and Ziegler in 1986 (see [AZ86], see also [Hod93]). Two theories U and V are bi-interpretable iff there

2000 *Mathematics Subject Classification*. 03A05, 03F25, 03F35 .

Key words and phrases. Interpretations, Interpretability.

[†]Harvey M. Friedman, Distinguished University Professor of Mathematics, Philosophy, and Computer Science, Emeritus, Ohio State University, Columbus, Ohio 43235. This research was partially supported by the John Templeton Foundation grant ID #36297. The opinions expressed here are those of the author and do not necessarily reflect the views of the John Templeton Foundation.

[‡]We thank Tonny Hurkens who simplified our proof of the Schröder-Bernstein Theorem. We thank Allan van Hulst who verified our proof of the Schröder-Bernstein Theorem in Mizar. We are grateful to Peter Aczel, Ali Enayat and Wilfrid Hodges for stimulating conversations or e-mail correspondence.

are interpretations $K : U \rightarrow V$ and $M : V \rightarrow U$, such that there is a V -definable function F , such that V proves that F is an isomorphism between $K \circ M$ and the identity interpretation on V and such that there is a U -definable function G , such that U proves that G is an isomorphism between $M \circ K$ and the identity interpretation on U . (Thus, bi-interpretability is isomorphism in an appropriate category INT_1 of theories and interpretations.)

In terms of models the notion of bi-interpretability takes the following form. We note that an interpretation $K : U \rightarrow V$ gives us a construction of an internal model $\widetilde{K}(\mathcal{M})$ of U from a model \mathcal{M} of V . We find that U and V are bi-interpretable iff, there are interpretations $K : U \rightarrow V$ and $M : V \rightarrow U$ and formulas F and G , such that, for all models \mathcal{M} of V , the formula F defines an isomorphism between \mathcal{M} and $\widetilde{M}\widetilde{K}(\mathcal{M})$, and, for all models \mathcal{N} of U , the formula G defines an isomorphism between \mathcal{N} and $\widetilde{K}\widetilde{M}(\mathcal{N})$.

Bi-interpretability has a lot of good properties. E.g., it preserves automorphism groups, κ -categoricity, finite axiomatizability, etc. Still the stricter notion synonymy preserves more. For example, synonymy preserves the action of the automorphism group on the domain of the model. Bi-interpretability (without parameters) does preserve the automorphism group modulo isomorphism but does not necessarily preserve the action on the domain. (See Section 7 for an example illustrating this difference.)

Surprisingly it is not easy to provide natural examples of pairs of theories that are bi-interpretable but not synonymous. For example PA is prima facie bi-interpretable with an appropriate theory of the hereditarily finite sets. However, on closer inspection, these theories are also synonymous. See Subsection 6.2. In Section 7, we give a verified example of two finitely axiomatized sequential theories that are bi-interpretable but not synonymous.

Our interest in this paper is in the relationship between synonymy and bi-interpretability for a special class of theories, *the sequential theories*. These are theories that have coding of sequences. Examples of sequential theories are Buss's theory S_2^1 , Elementary Arithmetic EA or EFA, $I\Sigma_1$, ZF, ZFC. We explain in detail what sequential theories are in Section 3.

We will show that, for *identity-preserving* interpretations between sequential theories, synonymy and bi-interpretability coincide (Section 5). In fact the proof works for a somewhat wider class: *the conceptual theories*.

A central ingredient of our proof is the Schröder-Bernstein Theorem that turns out to hold under surprisingly weak conditions. We give the proof of the Schröder-Bernstein Theorem under these weak conditions in Section 4.

2. BASIC NOTIONS

In this section, we formulate the basic notions employed in the paper. We keep the definitions here at an informal level. More detailed definitions are given in Appendix A.

2.1. Theories. The primary focus in this paper is on one-sorted theories of first order predicate logic of relational signature. 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. Many-sorted theories appear

as an auxiliary in the study of one-sorted theories. We will only consider theories with a finite number of sorts.

The results of the paper that are stated for one-sorted theories can be lifted to the many-sorted case in a fairly obvious way. We choose to restrict ourselves to the one-sorted case to keep the presentation reasonably light.

In this paper, no restriction is needed on the complexity of the set of axioms of a theory or on the size of the signature.

2.2. Interpretations. We describe the notion of an m -dimensional interpretation for a one-sorted language. An interpretation $K : U \rightarrow V$ is given by the theories U and V and a translation τ from the language of U to the language of V . The translation is given by a domain formula $\delta(\vec{x})$, where \vec{x} is a sequence of m variables, and a mapping from the predicates of U to formulas of V , where an n -ary predicate P is mapped to a formula $A(\vec{x}_0, \dots, \vec{x}_{n-1})$, where the \vec{x}_j are appropriately chosen pairwise disjoint sequences of m variables. We lift the translation to the full language in the obvious way making it commute with the propositional connectives and quantifiers, where we relativize the translated quantifiers to the domain δ . We demand that V proves all the translations of sentences of U .

We can compose interpretations in the obvious way. Note that the composition of an n -dimensional interpretation with an m -dimensional interpretation is $m \times n$ -dimensional.

An 1-dimensional interpretation is *identity preserving* if translates identity to identity. A 1-dimensional interpretation is *unrelativized* if its domain consists of all the objects of the interpreting theory. A 1-dimensional interpretation is *direct* if it is unrelativized and preserves identity. Note that all these properties are preserved by composition.

Each interpretation $K : U \rightarrow V$ gives us an inner model construction that builds a model $\tilde{K}(\mathcal{M})$ of U out of a model \mathcal{M} of V . Note that $\tilde{K}(\cdot)$ behaves contravariantly.

If we want to use interpretations to analyze sameness of theories, we will need, as we will see, to be able to say when two interpretations are ‘equal’. Strict identity of interpretations is too fine grained. It is too much dependent of arbitrary choices like which bound variables to use. We specify a first notion of equality between interpretations: two interpretations are *equal* when the *target theory thinks they are*. Modulo this identification, the operations identity and composition give rise to a category INT_0 , where the theories are objects and the interpretations arrows.¹

Let MOD be the category with as objects classes of models and as morphisms all functions between these classes. We define $\text{Mod}(U)$ as the class of all models of U . Suppose $K : U \rightarrow V$. Then, $\text{Mod}(K)$ is the function from $\text{Mod}(V)$ to $\text{Mod}(U)$ given by: $\mathcal{M} \mapsto \tilde{K}(\mathcal{M})$. It is clear that Mod is a *contravariant functor* from INT_0 to MOD .

2.3. Sameness of Interpretations. For each sufficiently good notion of sameness of interpretations there is an associated category of theories and interpretations: the category of interpretations modulo that notion of sameness. Any such a category gives us a notion of isomorphism of theories which can function as a notion of sameness.

¹For many reasons, the choice for the reverse direction of the arrows would be more natural. However, our present choice coheres with the extensive tradition in degrees of interpretability. So, we opted to adhere to the usual choice.

We present a *basic list* of salient notions of sameness. For all items in the list it is easily seen that sameness is preserved by composition.

2.3.1. *Equality.* The interpretations $K, K' : U \rightarrow V$ are equal when V ‘thinks’ K and K' are identical. By the Completeness Theorem, this is equivalent to saying that, for all V -models \mathcal{M} , $\widetilde{K}(\mathcal{M}) = \widetilde{K}'(\mathcal{M})$. This notion gives rise to the category INT_0 . Isomorphism in INT_0 is *synonymy* or *definitional equivalence*.

2.3.2. *i-Isomorphism.* An i-isomorphism between interpretations $K, M : U \rightarrow V$ is given by a V -formula F . We demand that V verifies that “ F is an isomorphism between K and M ”, or, equivalently, that, for each model \mathcal{M} of V , the function $F^{\mathcal{M}}$ is an isomorphism between $\widetilde{K}(\mathcal{M})$ and $\widetilde{M}(\mathcal{M})$.

Two interpretations $K, K' : U \rightarrow V$, are *i-isomorphic* iff there is an i-isomorphism between K and K' . Wilfrid Hodges calls this notion: *homotopy*. See [Hod93], p222.

We can also define the notion of being i-isomorphic semantically. The interpretations $K, K' : U \rightarrow V$, are *i-isomorphic* iff there is V -formula F such that for all V -models \mathcal{M} , the relation $F^{\mathcal{M}}$ is an isomorphism between $\widetilde{K}(\mathcal{M})$ and $\widetilde{K}'(\mathcal{M})$.

In case the signature of U is finite, being i-isomorphic has a third characterization. The interpretations $K, K' : U \rightarrow V$, are *i-isomorphic* iff, for every V -model \mathcal{M} , there is an \mathcal{M} -definable isomorphism between $\widetilde{K}(\mathcal{M})$ and $\widetilde{K}'(\mathcal{M})$. (See Theorem A.1.)

Clearly, if K, K' are equal in the sense of the previous subsection, they will be i-isomorphic. The notion of i-isomorphism give rise to a category of interpretations modulo i-isomorphism. We call this category INT_1 . Isomorphism in INT_1 is *bi-interpretability*.

2.3.3. *Isomorphism.* Our third notion of sameness of the basic list is that K and K' are the same if, for all models \mathcal{M} of V , the internal models $\widetilde{K}(\mathcal{M})$ and $\widetilde{K}'(\mathcal{M})$ are isomorphic. We will simply say that K and K' are isomorphic. Clearly, i-isomorphism implies isomorphism. We call the associated category INT_2 . Isomorphism in INT_2 is *iso-congruence*.

2.3.4. *Elementary Equivalence.* The fourth notion is to say that two interpretations are the same if, for each \mathcal{M} , the internal models $\widetilde{K}(\mathcal{M})$ and $\widetilde{K}'(\mathcal{M})$ are elementary equivalent. We will say that K and K' are elementary equivalent. By the Completeness Theorem, we easily see that this notion can be alternatively defined by saying that K is the same as K' iff, for all U -sentences A , we have $V \vdash A^K \leftrightarrow A^{K'}$. It is easy to see that isomorphism implies elementary equivalence. We call the associated category INT_3 . Isomorphism in INT_3 is *elementary congruence* or *sentential congruence*.

2.3.5. *Identity of Source and Target.* Finally, we have the option of abstracting away from the specific identity of interpretations completely, simply counting any two interpretations $K, K' : U \rightarrow V$ the same. The associated category is INT_4 . This is simply the structure of degrees of one-dimensional interpretability. Isomorphism in INT_4 is *mutual interpretability*.

2.4. The Many-sorted Case. Interpretability can be extended to interpretability between many-sorted theories. However to do that properly, we would need to develop the notion of piece-wise interpretation. Since this notion is not needed in the present paper, we just describe interpretations of many-sorted theories in one-sorted theories. These are precisely what one would expect: the interpretation K does not specify just one domain, but, for each sort \mathfrak{a} , a domain $\delta_{\mathfrak{a}}$. We allow a different dimension for each sort. The translation of a quantifier $\forall x^{\mathfrak{a}}$ is $\forall \vec{x} (\delta_{\mathfrak{a}}(\vec{x}) \rightarrow \dots)$. We translate a predicate P of type $\mathfrak{a}_0, \dots, \mathfrak{a}_{n-1}$ to a formula $A(\vec{x}_0, \dots, \vec{x}_{n-1})$, where the target theory verifies, for $i < n$, the formula $A(\vec{x}_0, \dots, \vec{x}_{n-1}) \rightarrow \delta_{\mathfrak{a}_i}(\vec{x}_i)$.²

We will consider theories with a designated sort \mathfrak{o} of objects. An interpretation of such a theory into a one-sorted theory is \mathfrak{o} -direct iff it is one-dimensional for sort \mathfrak{o} , and has $\delta_{\mathfrak{o}}(x) := (x = x)$ and translates identity on \mathfrak{o} to identity simpliciter. In other words, the interpretation is direct when we restrict our attention to the single sort \mathfrak{o} .

2.5. Parameters. We can extend our notion of interpretation to *interpretation with parameters* as follows. Say our interpretation is $K : U \rightarrow V$. In the target theory, we have a parameter domain $\alpha(\vec{z})$, which is V -provably non-empty. The definition of interpretation remains the same but for the fact that the parameters \vec{z} . Our condition for K to be an interpretation becomes:

$$U \vdash A \Rightarrow V \vdash \forall \vec{z} (\alpha(\vec{z}) \rightarrow A^{K, \vec{z}}).$$

We note that an interpretation $K : U \rightarrow V$ with parameters provides a parametrized set of inner models of U inside a model of V .

3. SEQUENTIALITY AND CONCEPTUALITY

We are interested in theories *with coding*. There are several ‘degrees’ of coding, like pairing, sequences, etcetera. We want a notion that allows us to build arbitrary sequences of all objects of our domain. The relevant notion is *sequentiality*. We also define a wider notion *conceptuality*. This last notion is proof-generated: it gives us the most natural class of theories for which our proof works. All sequential theories are conceptual, but not vice versa.

We have a simple and elegant definition of sequentiality. A theory U is *sequential* iff it directly interprets *adjunctive set theory* AS. Here AS is the following theory in the language with only one binary relation symbol.

AS1. $\vdash \exists x \forall y y \notin x$,

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

So the basic idea is that we can define a predicate \in^* in U such that \in^* satisfies a very weak set-theory involving *all* the objects of U . Given this weak set theory, we can develop a theory of sequences for all the objects in U , which again gives us partial truth-predicates, etc. In short, the notion of sequentiality explicates the idea of a *theory with coding*.

Remark 3.1. To develop the notion of sequentiality in a proper way for many-sorted theories we would need the idea of a *piecewise* interpretation. We do not develop the idea of piecewise interpretation here. Fortunately one can forget the framework and give the definition in a theory-free way. It looks like this. Let U

²Note that the sequence \vec{x}_i has as length the dimension associated to the sort \mathfrak{a}_i .

be a theory with sorts \mathcal{S} . The theory U is sequential when we can define, for each $\mathbf{a}, \mathbf{b} \in \mathcal{S}$, a binary predicate $\in^{\mathbf{a}\mathbf{b}}$ of type $\mathbf{a}\mathbf{b}$ such that:

- a. $U \vdash \bigvee_{\mathbf{a} \in \mathcal{S}} \exists x^{\mathbf{a}} \bigwedge_{\mathbf{b} \in \mathcal{S}} \forall y^{\mathbf{b}} y^{\mathbf{b}} \notin^{\mathbf{b}\mathbf{a}} x^{\mathbf{a}}$,
- b. $U \vdash \bigwedge_{\mathbf{a}, \mathbf{b} \in \mathcal{S}} \forall x^{\mathbf{a}}, y^{\mathbf{b}} \bigvee_{\mathbf{c} \in \mathcal{S}} \exists z^{\mathbf{c}} \bigwedge_{\mathbf{d} \in \mathcal{S}} \forall u^{\mathbf{d}} (u^{\mathbf{d}} \in^{\mathbf{d}\mathbf{c}} z^{\mathbf{c}} \leftrightarrow (u^{\mathbf{d}} \in^{\mathbf{d}\mathbf{a}} x^{\mathbf{a}} \vee u^{\mathbf{d}} \in^{\mathbf{d}\mathbf{b}} y^{\mathbf{b}}))$.

Here ‘ $\in^{\mathbf{d}\mathbf{a}}$ ’ is not really in the language if $\mathbf{d} \neq \mathbf{a}$. In this case we read $u^{\mathbf{d}} \in^{\mathbf{d}\mathbf{a}} y^{\mathbf{b}}$ simply as \perp .

It’s a nice exercise to show that e.g. ACA_0 and GB are sequential. \square

Closely related to AS is *adjunctive class theory* ac . We define this theory as follows. The theory ac is two-sorted with sorts \mathbf{o} (of objects) and \mathbf{c} (of classes). We have identity for every sort and one relation symbol \in between objects and classes, i.e. of type \mathbf{oc} . We let x, y, \dots range over objects and X, Y, \dots range over classes. We have the following axioms

- $\text{ac1. } \vdash \exists X \forall x x \notin X$,
- $\text{ac2. } \vdash \forall Y, y \exists X \forall x (x \in X \leftrightarrow (x \in Y \vee x = y))$,
- $\text{ac3. } \vdash X = Y \leftrightarrow \forall z (z \in X \leftrightarrow z \in Y)$.

Note that extensionality is cheap since we could treat identity on classes as *defined* by the relation of extensional sameness. The theory ac is much weaker than AS , since it admits finite models. The following theorem is easy to see.

Theorem 3.2. *A theory U is sequential iff there is an \mathbf{o} -direct interpretation of ac in U that is one-dimensional in the interpretation of classes.*

A theory U is *conceptual* iff there is an \mathbf{o} -direct interpretation of ac in U . We note that there are conceptual theories that are not sequential. For example, sequential theories always have infinite domain, but there are conceptual theories with finite models. Note also that AS is sequential, but ac is *not* conceptual (not even in the appropriate many-sorted formulation).

For more information on sequentiality and conceptuality, see Appendix B.

4. THE SCHRÖDER-BERNSTEIN THEOREM

We start with a brief story of the genesis of the theorem. The first step was taken by Harvey Friedman who saw that sequential theories should satisfy a version of the Schröder-Bernstein Theorem.³ Albert Visser subsequently wrote down a proof, discovering that one needs even less than sequentiality: the thing to use is adjunctive class theory. Allan van Hulst verified Visser’s version of the proof in Mizar as part of his master’s project under Freek Wiedijk in 2009. After hearing a presentation by Allan van Hulst, Tonny Hurkens found a simplification of the proof. Hurkens proof is shorter and conceptually simpler. In our presentation here we include Hurkens’ simplification. We thank Tonny for his gracious permission to do so.

We work in the theory SB which is ac extended with two unary predicates on objects: \mathbf{A} and \mathbf{B} and four binary predicates on objects: $\mathbf{E}_\mathbf{A}$, $\mathbf{E}_\mathbf{B}$, \mathbf{F} , \mathbf{G} , plus axioms expressing that $\mathbf{E}_\mathbf{A}$ is an equivalence relation on \mathbf{A} , $\mathbf{E}_\mathbf{B}$ is an equivalence relation on \mathbf{B} , \mathbf{F} is an injection from $\mathbf{A}/\mathbf{E}_\mathbf{A}$ to $\mathbf{B}/\mathbf{E}_\mathbf{B}$, \mathbf{G} is an injection from $\mathbf{B}/\mathbf{E}_\mathbf{B}$ to $\mathbf{A}/\mathbf{E}_\mathbf{A}$. We construct a formula \mathbf{H} that SB -provably defines a bijection between $\mathbf{A}/\mathbf{E}_\mathbf{A}$ and $\mathbf{B}/\mathbf{E}_\mathbf{B}$.

³In this paper we use the version of the Schröder-Bernstein Theorem that is formulated in terms of injections.

We will employ the usual notations like: \emptyset , $\{x_0, \dots, x_{n-1}\}$, \subseteq .

Our definition of what it means that F is a function includes: if $x \in A$, then $x \in E_A x' F y' \in E_B y$, then $x F y$. Similarly for G . We will treat A as a virtual class and write $x \in A$, etc. We use pairing of classes (X, Y) in a virtual way. E.g. $(X, Y) \subseteq (X', Y')$ is defined as a quaternary relation. Here are some definitions.

- $(X, Y) \subseteq (X', Y') : \leftrightarrow X \subseteq X' \wedge Y \subseteq Y'$,
- A pair of classes (X, Y) is *downwards closed* if,
 - (1) $(X, Y) \subseteq (A, B)$,
 - (2) if $u G v \in X$, then there is a u' with $Y \ni u' G v$,
 - (3) if $u F v \in Y$, then there is a u' with $X \ni u' F v$.
- We say that (X, Y) is an *x-switch* if (i) (X, Y) is downwards closed; (ii) x is a member of X ; (iii) each member of X is in the range of G .
- $x H y$ iff (there is no x -switch and $x F y$) or (there is an x -switch and $y G x$).

Lemma 4.1. H is a function from A/E_A to B/E_B .

Proof. We prove that H preserves equivalences. Suppose $x \in A$, $y \in B$ and $x H y$. Suppose there is an x -switch (X, Y) . It is easy to see that $(X \cup \{x'\}, Y)$ is an x' -switch. It is now immediate that $x' H y'$. Similarly, if we are given an x' -switch, we may conclude that there is an x -switch. The remaining case where there is neither an x -switch nor an x' -switch is again immediate.

Suppose $x H y$ and $x H y'$. If there is an x -switch, we have $y G x$ and $y' G x$. So we are done by the injectivity of G . If there is no x -switch, we have $x F y$ and $x F y'$. So we are done by the functionality of F .

We prove that H is total. If there is no x -switch, we are done. If there is an x -switch then x is in the range of G , and we are again done. \square

Lemma 4.2. H is injective from A/E_A to B/E_B .

Proof. Suppose $x H y$ and $x' H y$. If in both cases the same clauses in the definition of H are active, we are easily done.

Suppose there is no x -switch and there is an x' -switch, say (X', Y') . By the definition of H , we have $x F y G x'$. It follows that $x \in E_A x''$, for some x'' in X' . Hence, $(X' \cup \{x\}, Y')$ is an x -switch. A contradiction. \square

Lemma 4.3. H is surjective.

Proof. Consider any $y \in B$. First suppose y is not in the range of F . Let $y G x$. Then $(\{x\}, \{y\})$ is an x -switch, and we have $x H y$.

Next suppose $x F y$ and there is no x -switch. In this case $x H y$.

Suppose $x F y$ and there is an x -switch (X, Y) . Let $y G x_1$. Then $(X \cup \{x_1\}, Y \cup \{y\})$ is an x_1 -switch. Hence, $x_1 H y$.

Thus, in all cases y is in the image of H . \square

We have proved the following theorem.

Theorem 4.4. In SB we can construct a function H that is a bijection between A/E_A and B/E_B .

Example 4.5. Consider a model of SB . We write A for the interpretation of A , etc. We note that the function H constructed by the proof of the theorem depends on our choice of classes. Suppose, e.g., that our objects are the integers, A is the set of even integers, B is the set of odd integers, our equivalence relations are identity

on the given virtual class, F is the successor function domain-restricted to A , and G is the successor function domain-restricted to B . In the case that our classes are all possible classes of numbers, the pair of the class of all even numbers and all odd numbers is an a -switch, for each even a . So $H = G^{-1}$, i.e. the predecessor function domain-restricted to A . In the case that our classes are the finite classes, there is no a -switch for any even a . So, $H = F$. \square

For us the following obvious corollary is relevant.

Corollary 4.6. *Let T be a conceptual theory. Suppose we have formulas Ax , By , E_A , E_B , F , G , where T proves that E_A is an equivalence relation on $\{x \mid Ax\}$, that E_B is an equivalence relation on $\{y \mid By\}$, that F is an injection from $\{x \mid Ax\}/E_A$ to $\{y \mid By\}/E_B$, and that G is an injection from $\{y \mid By\}/E_B$ to $\{x \mid Ax\}/E_A$. Then we can find a formula H that T -provably defines a bijection between the virtual classes $\{x \mid Ax\}/E_A$ and $\{y \mid By\}/E_B$.*

Our corollary can be rephrased as follows. Suppose that T is conceptual and we have one-dimensional interpretations $K, M : \text{EQ} \rightarrow T$, where EQ is the theory of equality. Suppose further that $F : K \rightarrow M$ and $G : M \rightarrow K$ are injections. Then we can find a bijection $H : K \rightarrow M$, and, thus, K and M are i -isomorphic.

We note that if T is sequential, we can drop the demand that K and M are one-dimensional, since every interpretation in a sequential theory is i -isomorphic with a one-dimensional interpretation.

5. FROM BI-INTERPRETABILITY TO SYNONYMY

In this section we prove our main result. If two theories are bi-interpretable via identity-preserving interpretations, then they are synonymous.

Theorem 5.1. *Let U and V be any theories and suppose $K : U \rightarrow V$ and $M : V \rightarrow U$. Suppose that, for any model \mathcal{M} of V , we have $\widetilde{M}\widetilde{K}(\mathcal{M}) = \mathcal{M}$ (in other words, $K \circ M = \text{id}_V$ in INT_0). Suppose further that, for any model \mathcal{N} of U , the model $\widetilde{K}\widetilde{M}(\mathcal{N})$ is elementary equivalent to \mathcal{N} (in other words, $M \circ K = \text{id}_U$ in INT_3). Then U and V are synonymous.*

Here is a different formulation: if K, M witness that V is an INT_0 -retract of U and that U is an INT_3 -retract of V , then U and V are synonymous.

Proof. Consider any model \mathcal{N} of U . We have $\widetilde{K}\widetilde{M}\widetilde{K}\widetilde{M}(\mathcal{N}) = \widetilde{K}\widetilde{M}(\mathcal{N})$. Let $\mathcal{P} := \widetilde{K}\widetilde{M}(\mathcal{N})$. So, we have $\widetilde{K}\widetilde{M}(\mathcal{P}) = \mathcal{P}$. We note that the identity of $\widetilde{K}\widetilde{M}(\mathcal{P})$ and \mathcal{P} is witnessed by such statements as $\forall x \delta_{M \circ K}(x)$ and $\forall \vec{x} (P_{M \circ K} \vec{x} \leftrightarrow P \vec{x})$. Since \mathcal{P} is elementary equivalent to \mathcal{N} , we have $\widetilde{K}\widetilde{M}(\mathcal{N}) = \mathcal{N}$. So $M \circ K = \text{id}_U$ in INT_0 . \square

There is, of course, also a model-free proof of the result.

Theorem 5.2. *Suppose that K, M witness that V is an INT_1 -retract of U , and U is an INT_3 -retract of V . Suppose further that M is direct. Then U and V are synonymous.*

Proof. Since M is direct, it follows that, in V , we have $\delta_{K \circ M} = \delta_K$. We replace K by a definably isomorphic direct interpretation K' . Suppose F is the promised isomorphism between $K \circ M$ and id_V . We take, for P of arity n , in the signature of U :

$$\bullet P_{K'}(v_0, \dots, v_{n-1}) := \exists \vec{u}_0 \in \delta_K, \dots, \exists \vec{u}_{n-1} \in \delta_K \\ (\bigwedge_{i < n} \vec{u}_i F v_i \wedge P_K(\vec{u}_0, \dots, \vec{u}_{n-1})).$$

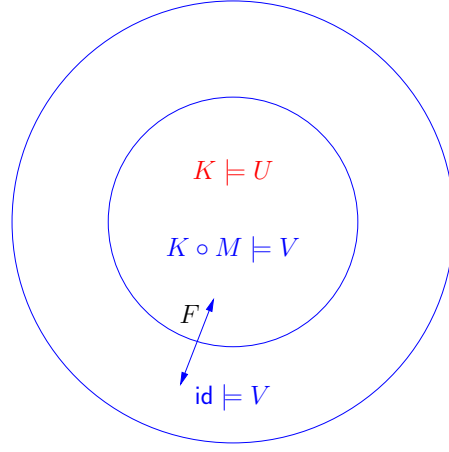


FIGURE 1. Illustration of the Proof of Theorem 5.2

Clearly, we have an isomorphism $F' : K \rightarrow K'$, based on the same underlying formula as F . Hence K', M witness that U is an INT_3 -retract of V .

We note that $K' \circ M$ is direct. Suppose R is an m -ary predicate of V . We have:

$$\begin{aligned}
V \vdash R_{K' \circ M}(x_0, \dots, x_{m-1}) &\leftrightarrow (R_M(x_0, \dots, x_{m-1}))^{K'} \\
&\leftrightarrow \exists \vec{y}_0 \in \delta_K, \dots, \exists \vec{y}_{m-1} \in \delta_K \\
&\quad \left(\bigwedge_{j < m} \vec{y}_j F x_j \wedge (R_M(\vec{y}_0, \dots, \vec{y}_{m-1}))^K \right) \\
&\leftrightarrow \exists \vec{y}_0 \in \delta_K, \dots, \exists \vec{y}_{m-1} \in \delta_K \\
&\quad \left(\bigwedge_{j < m} \vec{y}_j F x_j \wedge R_{K \circ M}(\vec{y}_0, \dots, \vec{y}_{m-1}) \right) \\
&\leftrightarrow R(x_0, \dots, x_{m-1})
\end{aligned}$$

Thus, we find: $K' \circ M = \text{id}_V$ in INT_0 , in other words, V is an INT_0 -retract of U . We apply Theorem 5.1 to K', M to obtain the desired result: U and V are synonymous. \square

We are mainly interested in the following corollary.

Corollary 5.3. *Suppose that $K : U \rightarrow V$ and $M : V \rightarrow U$ form a bi-interpretation and that M is direct. Then U and V are synonymous.*

We now prove our main theorem.

Theorem 5.4. *Suppose V is conceptual. Suppose that K, M witness that U is an INT_3 -retract of V and V is an INT_1 -retract of U . Suppose finally that K and M are both identity-preserving. Then, U and V are synonymous.*

Proof. We note that in V , δ_K is a (virtual) subclass of the full domain. Hence we have an definable injection from δ_K to the full domain.

Again $\delta_{K \circ M} = \delta_M^K \cap \delta_K$ is a (virtual) subclass of δ_K . Moreover, we have a definable bijection F between the full domain and $\delta_{K \circ M}$. Hence, we have a definable injection from the full domain into δ_K .

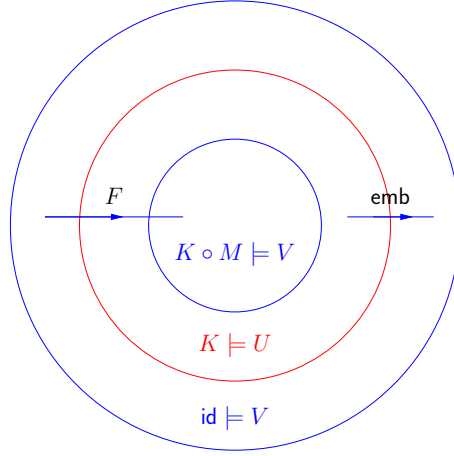


FIGURE 2. Illustration of the Proof of Theorem 5.4

We apply the Schröder-Bernstein Theorem to the full domain and δ_K providing us with a bijection G between the full domain and δ_K . We define a new interpretation $K' : U \rightarrow V$, by setting:

$$\bullet P_{K'}(v_0, \dots, v_{n-1}) := \leftrightarrow \exists w_0 \in \delta_K, \dots, w_{n-1} \in \delta_K \\ (\bigwedge_{i < n} v_i G w_i \wedge P_K(w_0, \dots, w_{n-1})).$$

Clearly K' is direct and isomorphic to K . We apply Theorem 5.2 and conclude that U and V are synonymous. \square

We note that in the circumstances of Theorem 5.4, it follows that U is also conceptual. Here is the salient corollary.

Corollary 5.5. *Suppose V is conceptual. Suppose $K : U \rightarrow V$ and $M : V \rightarrow U$ form a bi-interpretation and are both identity-preserving. Then, U and V are synonymous. A fortiori, we have the same theorem, when V is sequential.*

In Section 7, we provide an example to illustrate that one cannot drop the demand of identity preservation for any of the two interpretation. We provide two finitely axiomatized sequential theories that are bi-interpretable, but not synonymous. One of the two witnesses of the bi-interpretation is identity preserving.

6. APPLICATIONS

In this section we provide a number of applications of Corollary 5.5.

6.1. Natural Numbers and Rational Numbers. Julia Robinson, in her seminal paper [Rob49], shows that the natural numbers are definable in \mathbb{Q} , the field of the rationals. This gives us an identity preserving interpretation of $\text{Th}(\mathbb{N})$ in $\text{Th}(\mathbb{Q})$. Conversely, we can find an identity preserving interpretation of $\text{Th}(\mathbb{Q})$ in $\text{Th}(\mathbb{N})$ by using the Cantor pairing and by just considering pairs $\langle m, n \rangle$ where m and n have no common divisor except 1. Addition and multiplication are defined in the usual way. We can easily define internal isomorphisms witnessing that these interpretations form a bi-interpretation. Hence, $\text{Th}(\mathbb{N})$ in $\text{Th}(\mathbb{Q})$ are synonymous.

6.2. Finite Sets and Numbers. We consider the theory $ZF_{\text{fin}}^+ := (ZF - \text{INF}) + \neg\text{INF} + \text{TC}$. This is ZF in the usual formulation minus the axiom of infinity, plus the negation of the axiom of infinity and the axiom TC that tells us that every set has a transitive closure. Kaye and Wong in their paper [KW07] provide a careful verification that ZF_{fin}^+ and PA are synonymous. By Corollary 5.5, it is sufficient to show that the Ackermann interpretation of ZF_{fin}^+ in PA and the von Neumann interpretation of PA in ZF_{fin}^+ form a bi-interpretation. For further information about the related theory $ZF_{\text{fin}} = (ZF - \text{INF}) + \neg\text{INF}$, see [ESV10].

6.3. Sets with or without Urelements. Benedikt Löwe shows that a certain version of ZF with a countable set of urelements is synonymous with ZF. See [Löw06]. Again this result is easily obtained using Corollary 5.5.

6.4. Well-founded Sets versus Non-well-founded Sets. Let ZFC_{AFA} be ZFC minus *Foundation* plus the anti-foundation axiom AFA. See [Acz88] for an extensive treatment of ZFC_{AFA} .

We may interpret ZFC_{AFA} in ZFC, say via K , in the following way. The objects of the interpretation are rooted directed graphs (of set size) modulo bisimulation. The graph G is an ‘element’ of the graph H if there is an arrow from the root of H to a graph G' that is bisimilar to G .

We can eliminate the equivalence relation from the interpretation using Scott’s trick: we assign to an equivalence class the set of its elements of minimal rank. This does not give us canonical representatives, but just an injection of the equivalence classes to sets.⁴

We also have a backwards interpretation, say M , of ZFC in ZFC_{AFA} : we just relativize to the well-founded sets. The interpretation M is clearly identity preserving.

It is easy to see that the two interpretations form a bi-interpretation. Using Corollary 5.5, we see that ZFC_{AFA} and ZFC are synonymous.

The situation is even more interesting if we drop the axiom of choice. Let B be the axiom *every set is equinumerous to a well-founded set*. The two interpretations described above, transposed to the new context, form a bi-interpretation between $ZF_{\text{AFA}} + \text{B}$ and ZF. Thus, we find that $ZF_{\text{AFA}} + \text{B}$ is synonymous with ZF.

Thomas Forster shows, in his paper [For03], that the extra axiom B is independent of ZF_{AFA} . This means that ZF and ZF_{AFA} are not bi-interpretable via the given pair of interpretations. It is unknown whether ZF and ZF_{AFA} are bi-interpretable via another pair of interpretations. In an unpublished note Ali Enayat shows that ZF and $ZF - \text{Foundation}$ are not bi-interpretable.

7. FREGE MEETS CANTOR: AN EXAMPLE

In this section, we provide an example of two finitely axiomatized, sequential theories that are bi-interpretable but not synonymous. One of the two interpretations witnessing bi-interpretability is identity preserving. The example given is meaningful: it is the comparison of a Frege-style weak set theory and a Cantor style weak set theory. We show that these theories are not the same in the strictest sense, to wit *synonymy*, but that they are the same in a slightly weaker sense, to wit *bi-interpretability*.

⁴There are two normal forms for graphs, to wit bisimulation-minimal graphs and canonical unravelings. However, these normal forms are only determined up to isomorphism.

The theory ACF_b is the one-sorted version of adjunctive class theory with Frege function.⁵ Our theory has unary predicates ob and cl , and binary predicates \in and F . Here F is the Frege relation. We will write $x : \text{ob}$ for $\text{ob}(x)$, $\forall x:\text{ob} \dots$ for $\forall x(\text{ob}(x) \rightarrow \dots)$, $\exists x:\text{ob} \dots$ for $\exists x(\text{ob}(x) \wedge \dots)$. Similarly for cl . We have the following axioms.

- $\text{ACF}_b1. \vdash \forall x (x : \text{ob} \vee x : \text{cl}),$
- $\text{ACF}_b2. \vdash \forall x \neg (x : \text{ob} \wedge x : \text{cl}),$
- $\text{ACF}_b3. \vdash \forall x, y (x \in y \rightarrow (x : \text{ob} \wedge y : \text{cl})),$
- $\text{ACF}_b4. \vdash \forall x, y (x F y \rightarrow (x : \text{ob} \wedge y : \text{cl})),$
- $\text{ACF}_b5. \vdash \exists x:\text{cl} \forall y:\text{ob} \ y \notin x,$
- $\text{ACF}_b6. \vdash \forall x:\text{cl} \forall y:\text{ob} \exists z:\text{cl} \forall w:\text{ob} (w \in z \leftrightarrow (w \in x \vee w = y)).$
- $\text{ACF}_b7. \vdash \forall x, y:\text{cl} (\forall z:\text{ob} (z \in x \leftrightarrow z \in y) \rightarrow x = y).$
- $\text{ACF}_b8. \vdash \forall x:\text{ob} \exists y:\text{cl} \ x F y,$
- $\text{ACF}_b9. \vdash \forall x:\text{ob} \forall y, y':\text{cl} ((x F y \wedge x F y') \rightarrow y = y'),$
- $\text{ACF}_b10. \vdash \forall x:\text{cl} \exists y:\text{ob} \ y F x.$

We provide 1-dimensional interpretations witnessing that AS and ACF_b are bi-interpretable. Note that by Theorem B.1, it follows that ACF_b is sequential.

In the context of AS , we write:

- $\text{pair}(x, y, z) :\leftrightarrow \exists u, v \forall w ((w \in u \leftrightarrow w = x) \wedge (w \in v \leftrightarrow (w = x \vee w = y)) \wedge (w \in z \leftrightarrow (w = u \vee w = v))),$
- $\text{Pair}(x) :\leftrightarrow \exists y, z \text{pair}(y, z, x),$
- $\pi_0(z, x) :\leftrightarrow \exists y \text{pair}(x, y, z), \pi_1(z, y) :\leftrightarrow \exists x \text{pair}(x, y, z),$
- $\text{empty}(x) :\leftrightarrow \forall y \ y \notin x, \text{inhab}(x) :\leftrightarrow \neg \text{empty}(x).$
- $x \approx y :\leftrightarrow \forall z (z \in x \leftrightarrow z \in y).$

We can verify the usual properties of pairing. The π_i are functional on Pair . We will write them using functional notation. We should remember that they are undefined outside Pair . We first define an interpretation $L : \text{ACF}_b \rightarrow \text{AS}$. We can

- $\delta_L(x) :\leftrightarrow \text{Pair}(x),$
- $\text{ob}_L(x) :\leftrightarrow \text{Pair}(x) \wedge \text{empty}(\pi_0(x)),$
- $\text{cl}_L(x) :\leftrightarrow \text{Pair}(x) \wedge \text{inhab}(\pi_0(x)),$
- $x =_L y :\leftrightarrow (x, y : \text{ob}_L \wedge \pi_1(x) = \pi_1(y)) \vee (x, y : \text{cl}_L \wedge \pi_1(x) \approx \pi_1(y)).$
- $x \in_L y :\leftrightarrow x : \text{ob}_L \wedge y : \text{cl}_L \wedge \pi_1(x) \in \pi_1(y),$
- $x F_L y :\leftrightarrow x : \text{ob}_L \wedge y : \text{cl}_L \wedge \pi_1(x) \approx \pi_1(y).$

It is easy to see that the specified translation does carry an interpretation of ACF_b in AS , as promised. Next, we define an interpretation $\text{ACF}_b. K : \text{AS} \rightarrow \text{ACF}_b.$

- $\delta_K(x) :\leftrightarrow x : \text{ob},$
- $x =_K y \leftrightarrow x, y : \text{ob} \wedge x = y,$
- $x \in_K y \leftrightarrow x, y : \text{ob} \wedge \exists z:\text{cl} (x \in z \wedge y F z).$

We note that K is identity preserving. The verification that our interpretations do indeed specify a bi-interpretation is entirely routine. For completeness' sake, we provide the computations involved in Appendix C.

⁵It would be more natural to give the example for the two-sorted theory ACF and to use the notion of piecewise interpretation. The definitions of the interpretations and the verification that they form a bi-interpretation would be simpler. In fact, we would avoid the use of coding in our definitions. However, we would need to develop more of the theory of many-sorted interpretations to handle the superior approach smoothly. This is beyond the scope of our present paper.

We show that AS and ACF_b are not synonymous —not even when we allow parameters. We build the following model \mathcal{M} of AS. The domain is inductively specified as the smallest set M such that if X is a finite subset of M , then $\langle 0, X \rangle$ and $\langle 1, X \rangle$ are in M . Let m, n, \dots range over M . We define: $m \in^* n$ iff $n = \langle i, X \rangle$ and $m \in X$. It is easily seen that \mathcal{M} is indeed a model of AS. Clearly, for any finite subset X_0 of M we can find an automorphism σ of \mathcal{M} of order 2 that fixes X_0 and fixes only finitely many elements of M .

Suppose AS and ACF_b were synonymous. Let \mathcal{N} be the internal model of ACF_b in \mathcal{M} given by the synonymy. Say the interpretation is P , involving a finite set of parameters X_0 . Let σ be an automorphism of order 2 that fixes X_0 and that fixes at most finitely many objects. Consider the classes $\{p, \sigma p\}^{\mathcal{N}}$, where p is in $\text{ob}^{\mathcal{N}}$. (Note that σ must send \mathcal{N} -objects to \mathcal{N} -objects.) Clearly there is an infinity of such classes. By extensionality these classes are fixed by σ . This contradicts the fact that σ has only finitely many fixed points.

It is well known that if two models are bi-interpretable (without parameters) then their automorphism groups are isomorphic. Our example shows that the *action* of these automorphism groups on the elements can be substantially different.

REFERENCES

- [Acz88] Peter Aczel. *Non-well-founded sets*, volume 14 of *CSLI Lecture Notes*. CSLI, Stanford, 1988.
- [AZ86] G. Ahlbrandt and M. Ziegler. Quasi finitely axiomatizable totally categorical theories. *Annals of Pure and Applied Logic*, 30(1):63–82, 1986.
- [CH70] G.E. Collins and J.D. Halpern. On the interpretability of Arithmetic in Set Theory. *Notre Dame Journal of Formal Logic*, 11(4):477–483, 1970.
- [dB65a] K. L. de Bouvère. Logical synonymy. *Indagationes Mathematicae*, 27:622–629, 1965.
- [dB65b] K. L. de Bouvère. Synonymous Theories. In J.W. Addison, L. Henkin, and A. Tarski, editors, *The Theory of Models, Proceedings of the 1963 International Symposium at Berkeley*, pages 402–406. North Holland, Amsterdam, 1965.
- [ESV10] A. Enayat, J. Schmerl, and A. Visser. ω -models of finite set theory. In J. Kennedy and R. Kossak, editors, *Set Theory, Arithmetic and Foundations of Mathematics: Theorems, Philosophies*, number 36 in ASL Lecture Notes in Logic, pages 43–65. ASL and Cambridge University Press, New York, 2010.
- [For03] T. Forster. ZF + “Every Set Is the Same Size as a Wellfounded Set”. *The Journal of Symbolic Logic*, 68(1):1–4, 2003.
- [Hod93] W. Hodges. *Model theory*. Encyclopedia of Mathematics and its Applications, vol. 42. Cambridge University Press, Cambridge, 1993.
- [HP93] P. Hájek and P. Pudlák. *Metamathematics of First-Order Arithmetic*. Perspectives in Mathematical Logic. Springer, Berlin, 1993.
- [JV00] J.J. Joosten and A. Visser. The interpretability logic of *all* reasonable arithmetical theories. *Erkenntnis*, 53(1–2):3–26, 2000.
- [KW07] Richard Kaye and Tin Lok Wong. On interpretations of arithmetic and set theory. *Notre Dame Journal of Formal Logic*, 48(4):497–510, 2007.
- [Löw06] Benedikt Löwe. Set Theory with and without urelements and categories of interpretations. *Notre Dame Journal of Formal Logic*, 47(1):83–91, 2006.
- [MM94] F. Montagna and A. Mancini. A minimal predicative set theory. *Notre Dame Journal of Formal Logic*, 35(2):186–203, 1994.
- [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.
- [Nel86] E. Nelson. *Predicative arithmetic*. Princeton University Press, Princeton, 1986.
- [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(2):423–441, 1985.
- [Rob49] Julia Robinson. Definability and Decision Problems in Arithmetic. *Journal of Symbolic Logic*, 14(2):98–114, 1949.
- [Smo85] 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.
- [ST52] W. Szmielew and A. Tarski. Mutual Interpretability of some essentially undecidable theories. In *Proceedings of the International Congress of Mathematicians (Cambridge, Massachusetts, 1950)*, volume 1, page 734. American Mathematical Society, Providence, 1952.
- [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.
- [Vis92] A. Visser. An inside view of EXP. *The Journal of Symbolic Logic*, 57(1):131–165, 1992.
- [Vis93] A. Visser. The unprovability of small inconsistency. *Archive for Mathematical Logic*, 32(4):275–298, 1993.
- [Vis98] A. Visser. An Overview of Interpretability Logic. In M. Kracht, M. de Rijke, H. Wansing, and M. Zakharyashev, editors, *Advances in Modal Logic*, volume 1, 87 of *CSLI Lecture Notes*, pages 307–359. Center for the Study of Language and Information, Stanford, 1998.
- [Vis05] A. Visser. Faith & Falsity: a study of faithful interpretations and false Σ_1^0 -sentences. *Annals of Pure and Applied Logic*, 131(1–3):103–131, 2005.
- [Vis06] A. Visser. Categories of Theories and Interpretations. In Ali Enayat, Iraj Kalantari, and Mojtaba Moniri, editors, *Logic in Tehran. Proceedings of the workshop and conference on Logic, Algebra and Arithmetic, held October 18–22, 2003*, volume 26 of *Lecture Notes in Logic*, pages 284–341. ASL, A.K. Peters, Ltd., Wellesley, Mass., 2006.
- [Vis08] A. Visser. Pairs, sets and sequences in first order theories. *Archive for Mathematical Logic*, 47(4):299–326, 2008.
- [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.
- [Vis13] A. Visser. What is the right notion of sequentiality? In P. Cégielski, C. Charampolas, and C. Dimitracopoulos, editors, *New Studies in Weak Arithmetics*, volume 211 of *CSLI Lecture Notes*, pages 229–272. CSLI Publications and Presses Universitaires du Pôle de Recherche et d’Enseignement Supérieur Paris-est, Stanford, 2013.

APPENDIX A. DEFINITIONS

In this appendix, we provide detailed definitions of translations, interpretations and morphisms between interpretations.

A.1. Translations. Translations are the heart of our interpretations. In fact, they are often confused with interpretations, but we will not do that officially. In practice it is often convenient to conflate an interpretation and its underlying translation.

We think of formulas *modulo* α -conversion. An n -ary proto-formula A is something like a formula with (at most) n free variables where we abstract away from the identity of the variables. Officially, an n -ary proto-formula A is a function from sequences of n variables to formulas, so that, for any substitution σ of variables with as domain the the variables in \vec{x} , we have $A(\sigma(\vec{x})) = \sigma(A(\vec{x}))$.

We define more-dimensional, one-sorted, one-piece relative translations without parameters. Let Σ and Θ be one-sorted signatures. A translation $\tau : \Sigma \rightarrow \Theta$ is given by a triple $\langle m, \delta, F \rangle$. Here δ will be a m -ary proto-formula of signature Θ . The mapping F associates to each relation symbol R of Σ with arity n an $m \times n$ -ary proto-formula of signature Θ .

We demand that predicate logic proves $F(R)(\vec{x}_0, \dots, \vec{x}_{n-1}) \rightarrow (\delta(\vec{x}_0) \wedge \dots \wedge \delta(\vec{x}_{n-1}))$. Of course, given any candidate proto-formula $F(R)$ not satisfying the restriction, we can obviously modify it to satisfy the restriction.

We translate Σ -formulas to Θ -formulas as follows.

- $(R(x_0, \dots, x_{n-1}))^\tau := F(R)(\vec{x}_0, \dots, \vec{x}_{n-1})$.
We use sloppy notation here. The single variable x_i of the source language needs to have no obvious connection with the sequence of variables \vec{x}_i of the target language that represents it. We need some conventions to properly handle the association $x_i \mapsto \vec{x}_i$. We do not treat these details here. We demand that the \vec{x}_i are fully disjoint when the x_i are different.
- $(\cdot)^\tau$ commutes with the propositional connectives;
- $(\forall x A)^\tau := \forall \vec{x} (\delta(\vec{x}) \rightarrow A^\tau)$;
- $(\exists x A)^\tau := \exists \vec{x} (\delta(\vec{x}) \wedge A^\tau)$.

Here are some convenient conventions and notations.

- We write δ_τ for ‘the δ of τ ’ and F_τ for ‘the F of τ ’.
- We write R_τ for $F_\tau(R)$.
- We write $\vec{x} \in \delta$ for: $\delta(\vec{x})$.

There are some natural operations on translations. The identity translation $\text{id} := \text{id}_\Theta$ is one-dimensional and it is defined by:

- $\delta_{\text{id}}(x) := (x = x)$,
- $R_{\text{id}}(\vec{x}) := R\vec{x}$.

We can compose relative translations as follows. Suppose τ is an m -dimensional translation from Σ to Θ , and ν is a k -dimensional translation from Θ to Ξ . We define:

- We suppose that with the variable x we associate under τ the sequence x_0, \dots, x_{m-1} and under ν we send x_i to \vec{x}_i .
 $\delta_{\tau\nu}(\vec{x}_0, \dots, \vec{x}_{m-1}) := (\delta_\nu(\vec{x}_0) \wedge \dots \wedge \delta_\nu(\vec{x}_{m-1}) \wedge (\delta_\tau(x))^\nu)$,
- Let R be n -ary. Suppose that under τ we associate with x_i the sequence $x_{i,0}, \dots, x_{i,m-1}$ and that under ν we associate with $x_{i,j}$ the sequence $\vec{x}_{i,j}$. We take:
 $R_{\tau\nu}(\vec{x}_{0,0}, \dots, \vec{x}_{n-1,m-1}) = \delta_\tau(\vec{x}_{0,0}) \wedge \dots \wedge \delta_\tau(\vec{x}_{n-1,m-1}) \wedge (R_\tau(x_0, \dots, x_{n-1}))^\nu$.

A one-dimensional translation τ *preserves identity* if $(x =_\tau y) = (x = y)$. A one-dimensional translation τ *is unrelativized* if $\delta_\tau(x) = (x = x)$. An one-dimensional translation τ is *direct* if it is unrelativized and preserves identity. Note that all these properties are preserved by composition.

Consider a model \mathcal{M} with domain M of signature Σ and k -dimensional translation $\tau : \Sigma \rightarrow \Theta$. Suppose that $N := \{\vec{m} \in M^k \mid \mathcal{M} \models \delta_\tau \vec{m}\}$. Then τ specifies an internal model \mathcal{N} of \mathcal{M} with domain N and with $\mathcal{N} \models R(\vec{m}_0, \dots, \vec{m}_{n-1})$ iff $\mathcal{M} \models R_\tau(\vec{m}_0, \dots, \vec{m}_{n-1})$. We will write $\tilde{\tau}(\mathcal{M})$ for the internal model of \mathcal{M} given by τ . Treating the mapping $\tau, \mathcal{M} \mapsto \tilde{\tau}\mathcal{M}$ as a partial function that is defined precisely if $\delta_\tau^{\mathcal{M}}$ is non-empty. Let Mod or (\cdot) be the function that maps τ to $\tilde{\tau}$. We have:

$$\text{Mod}(\tau \circ \rho)(\mathcal{M}) = (\text{Mod}(\rho) \circ \text{Mod}(\tau))(\mathcal{M}).$$

So, Mod behaves contravariantly.

A.2. Relative Interpretations. A translation τ supports a *relative interpretation* of a theory U in a theory V , if, for all U -sentences A , $U \vdash A \Rightarrow V \vdash A^\tau$. Note that this automatically takes care of the theory of identity and assures us that δ_τ is inhabited. We will write $K = \langle U, \tau, V \rangle$ for the interpretation supported by τ . We write $K : U \rightarrow V$ for: K is an interpretation of the form $\langle U, \tau, V \rangle$. If M is an interpretation, τ_M will be its second component, so $M = \langle U, \tau_M, V \rangle$, for some U and V .

Par abus de langage, we write ‘ δ_K ’ for: δ_{τ_K} ; ‘ R_K ’ for: R_{τ_K} ; ‘ A^K ’ for: A^{τ_K} , etc. Here are the definitions of three central operations on interpretations.

- Suppose T has signature Σ . We define:
 $\text{id}_T : T \rightarrow T$ is $\langle T, \text{id}_\Sigma, T \rangle$.
- Suppose $K : U \rightarrow V$ and $M : V \rightarrow W$. We define:
 $M \circ K : U \rightarrow W$ is $\langle U, \tau_M \circ \tau_K, W \rangle$.

It is easy to see that we indeed correctly defined interpretations between the theories specified.

A.3. Equality of Interpretations. Two interpretations are *equal* when the *target theory thinks they are*. Specifically, we count two interpretations $K, K' : U \rightarrow V$ as equal if they have the same dimension, say m , and:

- $V \vdash \forall \vec{x} (\delta_K(\vec{x}) \leftrightarrow \delta_{K'}(\vec{x}))$,
- $V \vdash \forall \vec{x}_0, \dots, \vec{x}_{n-1} \in \delta_K (R_K(\vec{x}_0, \dots, \vec{x}_{n-1}) \leftrightarrow R_{K'}(\vec{x}_0, \dots, \vec{x}_{n-1}))$.

Modulo this identification, the operations identity and composition give rise to a category INT_0 , where the theories are objects and the interpretations arrows.⁶

Let MOD be the category with as objects classes of models and as morphisms all functions between these classes. We define $\text{Mod}(U)$ as the class of all models of U . Suppose $K : U \rightarrow V$. Then, $\text{Mod}(K)$ is the function from $\text{Mod}(V)$ to $\text{Mod}(U)$ given by: $\mathcal{M} \mapsto \tilde{K}(\mathcal{M}) := \tilde{\tau}_K(\mathcal{M})$. It is clear that Mod is a *contravariant functor* from INT_0 to MOD .

A.4. Maps between Interpretations. Consider $K, M : U \rightarrow V$. Suppose K is m -dimensional and M is k -dimensional. A V -definable, V -provable morphism from K to M is a triple $\langle K, F, M \rangle$, where F is a $m+k$ -ary proto-formula.⁷ We write $\vec{x} F \vec{y}$ for $F(\vec{x}, \vec{y})$. We demand that F has the following properties.

- $V \vdash \vec{x} F \vec{y} \rightarrow (\vec{x} \in \delta_K \wedge \vec{y} \in \delta_M)$.
- $V \vdash \vec{x} =_K \vec{u} F \vec{v} =_M \vec{y} \rightarrow \vec{x} F \vec{y}$.
- $V \vdash \forall \vec{x} \in \delta_K \exists \vec{y} \in \delta_M x F y$.
- $V \vdash (\vec{x} F \vec{y} \wedge \vec{x} F \vec{z}) \rightarrow \vec{y} =_M \vec{z}$.
- $V \vdash (\vec{x}_0 F \vec{y}_0 \wedge \dots \wedge \vec{x}_{n-1} F \vec{y}_{n-1} \wedge R_K(\vec{x}_0, \dots, \vec{x}_{n-1})) \rightarrow R_M(\vec{y}_0, \dots, \vec{y}_{n-1})$.

We will call the arrows between interpretations: *i-maps* or *i-morphisms*. We write $F : K \rightarrow M$ for: $\langle K, F, M \rangle$ is a V -provable, V -definable morphism from K to M . Remember that the theories U and V are part of the data for K and M . We consider $F, G : K \Rightarrow M$ as *equal* when they are V -provably the same.

⁶For many reasons, the choice for the reverse direction of the arrows would be more natural. However, our present choice coheres with the extensive tradition in degrees of interpretability. So, we opted to adhere to the present choice.

⁷Since, in this stage, we are looking at definitions without parameters we could, perhaps, better speak of *V-0-definable*. Parameters may be added but in the context where we consider theories rather than models some extra details are needed to make everything work smoothly.

An isomorphism of interpretations is easily seen to be a morphism with the following extra properties.

- $V \vdash \forall \vec{y} \in \delta_M \exists \vec{x} \in \delta_K x F y$,
- $V \vdash (\vec{x} F \vec{y} \wedge \vec{z} F \vec{y}) \rightarrow \vec{x} =_K \vec{z}$,
- $V \vdash (\vec{x}_0 F \vec{y}_0 \wedge \dots \wedge \vec{x}_{n-1} F \vec{y}_{n-1} \wedge R_M(\vec{x}_0, \dots, \vec{x}_{n-1})) \rightarrow R_K(\vec{y}_0, \dots, \vec{y}_{n-1})$.

We call such isomorphisms: i-isomorphisms. By a simple compactness argument one may prove:

Theorem A.1. *Suppose the signature of U is finite. Consider $K, M : U \rightarrow V$. Suppose that, for every model \mathcal{N} of V , there is an \mathcal{N} -definable isomorphism between $\tilde{K}(\mathcal{N})$ and $\tilde{M}(\mathcal{N})$. Then, K and M are i-isomorphic.*

A.5. Adding Parameters. We can add parameters in the obvious way. An interpretation $K : U \rightarrow V$ with parameters will have a k -dimensional parameter domain α (officially a proto-formula), where $V \vdash \exists \vec{x} \alpha \vec{x}$. We allow the extra variables \vec{x} to occur in the translations of the U formulas. We have to take the appropriate measures to avoid variable-clashes. The condition for K to be an interpretation changes into: $\vdash \forall \vec{x} (\alpha \vec{x} \rightarrow A^{K, \vec{x}})$, where A is an axiom of U .

We note that, in the presence of parameters, the function \tilde{K} associates a class of models of U to a model of V .

Similar adaptations are needed to define i-isomorphisms with parameters.

APPENDIX B. BACKGROUND FOR SEQUENTIALITY AND CONCEPTUALITY

The notion of sequential theory was introduced by Pavel Pudlák in his paper [Pud83]. Pudlák uses his notion for the study of the degrees of local multi-dimensional parametric interpretability. He proves that sequential theories are prime in this degree structure. In [Pud85], sequential theories provide the right level of generality for theorems about consistency statements.

The notion of sequential theory was independently invented by Friedman who called it *adequate theory*. See Smoryński's survey [Smo85]. Friedman uses the notion to provide the Friedman characterization of interpretability among finitely axiomatized sequential theories. (See also [Vis90] and [Vis92].) Moreover, he shows that ordinary interpretability and faithful interpretability among finitely axiomatized sequential theories coincide. (See also [Vis93] and [Vis05].)

The story of the weak set theory AS can be traced in a sequence of papers the following papers: [ST52], [CH70], [Pud85], [Nel86], [MM94], [MPS90] (appendix III), [Vis08], [Vis09]. The connection between AS and sequentiality is made in [Pud85] and [MPS90].

For further work concerning sequential theories, see, e.g., [HP93], [Vis93], [Vis98], [JV00], [Vis05].

A theorem that is relevant in this paper is Theorem 10.7 of [Vis06]:

Theorem B.1. *Sequentiality is preserved to INT_1 -retracts for one-dimensional interpretations. In other words: if V is sequential and if U is a one-dimensional retract in INT_1 of V , then U is sequential.*

Proof. Suppose $K : U \rightarrow V$ and $M : V \rightarrow U$ are one-dimensional and $M \circ K$ is i-isomorphic to id_U via F . Let \in^* be the V -formula witnessing the sequentiality of V . We define the U -formula \in^* witnessing the sequentiality of U by: $x \in^* y$ iff y is in δ_M and, for some z with $z F x$, we have $(z \in^* y)^M$. \square

This result holds only for one-dimensional interpretability. There are examples of non-sequential theories that are bi-interpretable with a sequential theory. Since bi-interpretability is such a good notion of sameness of theories, one could argue that the failure of closure of sequential theories under bi-interpretability is a defect and that we need a slightly more general notion to fully reflect the intuitions that sequentiality is intended to capture. For an elaboration of this point, see [Vis13].

We can easily adapt Theorem B.1, to obtain:

Theorem B.2. *Conceptuality is preserved to INT_1 -retracts.*

APPENDIX C. VERIFICATION OF BI-INTERPRETABILITY

We verify that the interpretations K and L of Section 7 do indeed form a bi-interpretation. We first compute $M := (L \circ K) : \text{AS} \rightarrow \text{AS}$. We find:

- $\delta_M(x) \leftrightarrow (\delta_L(x) \wedge (x \in \delta_K)^L) \leftrightarrow (x : \text{Pair} \wedge \text{empty}(\pi_0(x)))$,
- We have (using the contextual information that x and y are in δ_M):

$$\begin{aligned} x =_M y &\leftrightarrow (x =_K y)^L \\ &\leftrightarrow \pi_1(x) = \pi_1(y) \end{aligned}$$

- We have (using the contextual information that x and y are in δ_M):

$$\begin{aligned} x \in_M y &\leftrightarrow (\exists z : \text{cl}(x \in z \wedge y \text{ F } z))^L \\ &\leftrightarrow \exists z : \text{pair}(\text{inhab}(\pi_0(z)) \wedge \pi_1(x) \in \pi_1(z) \wedge \pi_1(y) \approx \pi_1(z)) \\ &\leftrightarrow \pi_1(x) \in \pi_1(y) \end{aligned}$$

Clearly π_1 is the desired isomorphism from L to id_{AS} .

In the other direction, let $N := (K \circ L) : \text{ACF}_b \rightarrow \text{ACF}_b$. We first note a simple fact about K . Since F is functional on classes we will use functional notation for it. We have, for $u, v : \text{ob}$,

$$\begin{aligned} u \approx_K v &\leftrightarrow \forall w (w \in_K u \leftrightarrow w \in_K v) \\ &\leftrightarrow \forall w (w \in \text{F}(u) \leftrightarrow w \in \text{F}(v)) \\ &\leftrightarrow \text{F}(u) = \text{F}(v). \end{aligned}$$

We have:

- $\delta_N(x) \leftrightarrow (x \in \delta_K(x) \wedge (\delta_L(x))^K) \leftrightarrow (x : \text{ob} \wedge \text{Pair}^K(x))$,
- $\text{ob}_N(x) \leftrightarrow (\text{empty}(\pi_0(x)))^K$,
- $\text{cl}_N(x) \leftrightarrow (\text{inhab}(\pi_0(x)))^K$,
- $x =_N y \leftrightarrow ((\text{ob}_N(x) \wedge \text{ob}_N(y) \wedge \pi_1^K(x) = \pi_1^K(y)) \vee (\text{cl}_N(x) \wedge \text{cl}_N(y) \wedge \text{F}(\pi_1^K(x)) = \text{F}(\pi_1^K(y))))$,
- $x \in_N y \leftrightarrow (\text{ob}_N(x) \wedge \text{cl}_N(y) \wedge \pi_1^K(x) \in_K \pi_1^K(y))$.
- $x \text{ F}_N y \leftrightarrow (\text{ob}_N(x) \wedge \text{cl}_N(y) \wedge \text{F}(\pi_1^K(x)) = \text{F}(\pi_1^K(y)))$.

We define $G : N \rightarrow \text{id}_{\text{ACF}_b}$ as follows:

- $x \text{ G } y \leftrightarrow (\text{ob}_N(x) \wedge \text{ob}(y) \wedge \pi_1^K(x) = y) \vee (\text{cl}_N(x) \wedge \text{cl}(y) \wedge \text{F}(\pi_1^K(x)) = y)$.

We note that K is identity preserving. Thus, e.g., we find that pair^K is a true pairing on ob . It is easy to see that, in ACF_b , the virtual classes ob_N and cl_N for a partition of δ_N . It is now trivial to check that G is indeed an isomorphism.

DEPARTMENT OF MATHEMATICS, MATHEMATICS BUILDING, 231 WEST 18TH AVENUE, COLUMBUS, OH 43210, USA

E-mail address: `friedman@math.ohio-state.edu`

PHILOSOPHY, FACULTY OF HUMANITIES, UTRECHT UNIVERSITY, JANSKERKHOF 13, 3512 BL UTRECHT, THE NETHERLANDS

E-mail address: `a.visser@uu.nl`