

INTERPRETATIONS, ACCORDING TO TARSKI

Harvey M. Friedman*

Nineteenth Annual Tarski Lectures

Interpretations of Set Theory in Discrete Mathematics and
Informal Thinking

Lecture 1

Delivered, April 9, 2007

Expanded May 24, 2007

1. Interpretations.
2. Tarski Degrees.
3. Adequate Sentences and Relative Consistency.
4. Predicative Extensions and Relative Consistency.
5. P degrees.
6. Infinite Theories.
7. Observed Linearity.

Special thanks to Albert Visser for considerable input.

1. Interpretations.

The notion of interpretation was first carefully defined and developed in the book [TMR53].

The notion of interpretation is absolutely fundamental to mathematical logic and the foundations of mathematics. It is also crucial for the foundations and philosophy of science - although here some crucial conditions generally need to be imposed; e.g., "the interpretation leaves the mathematical concepts unchanged".

The most obvious and direct use of interpretations is for relative consistency.

PROPOSITION 1.1. Suppose S is interpretable in T . If T is consistent then S is consistent.

But the notion was mostly used by Tarski to show that various formal systems are undecidable in the sense that there is no algorithm for determining whether a sentence in its language is provable.

PROPOSITION 1.2. Suppose S is interpretable in T . If T has a consistent decidable extension in its own language, then S has a consistent decidable extension in its own language. I.e., if S is essentially undecidable then T is essentially undecidable.

Rafael M. Robinson in [Ro52] set up a crucial very weak system of arithmetic, Q , and showed that Q is essentially undecidable.

Proposition 1.2 was successively applied, starting with Q , invoking transitivity: S interpretable in T , T interpretable in W , therefore S is interpretable in W .

We now define interpretations in the setting of ordinary first order predicate calculus with equality.

There are some choices along the road - the notion is still not entirely standard. We will work with the most liberal natural and convenient notion.

A theory T is a set of sentences together with a relational type, $L(T)$. Here we won't demand that we close T under logical consequence.

We begin semiformaly and semantically, with several details postponed.

\square is an interpretation of S in T iff

- i. \square specifies a domain of objects from viewpoint of T . Given by formula(s) in $L(T)$.
- ii. \square specifies an interpretation of all constant, relation, and function symbols of $L(S)$, on the domain in i, by means of formulas in $L(T)$.
- iii. This data induces a map \square from formulas of $L(S)$ into formulas of $L(T)$. We require that for all $\square \in S$, $T \vdash \square$.

We have hidden details. We take the most liberal natural view of what is allowed; i.e., "liberal" definability.

The specified domain D consists of tuples of any lengths ≥ 1 . Mixed lengths are allowed. There must be a bound on the lengths used. D is to be given by a finite set of formulas in $L(T)$, one for each length allowed. In every model of T , D is nonempty.

We have a binary relation E on D , given by a finite set of formulas in $L(T)$, one for each pair of lengths allowed.

In every model of T , E is an equivalence relation on D . We have formulas in $L(T)$ that interpret the constant, relation, function symbols of S , on D/E . We don't perform factoring, but instead treat the interpreted constants, relations, functions as relations on D that respect E , in all models of T .

We have lied a little: we will allow parameters. I.e., extra variables in the above, where we merely say that in any model of T , there exist choices for these variables (parameters) so that all of the above hold.

We are asserting that there is a way of defining a model of S inside any model of T , uniformly, multidimensionally, via equivalence, and with parameters. The definition is required to be uniform but the choice of parameters has no uniformity requirement. There may be many choices of parameters that work. We call this liberal definability.

This definition is particularly well behaved in the case where S is finite. To begin with, since S is finite, the data used and the requirements made in the definition are finite *before* we existentially quantify over choices of parameters. Therefore we are asserting that T logically implies a finite set of sentences. Hence we can replace these logical implications by provability in T . This gives us the *provability* version of (liberal) interpretability.

The following result shows a great deal of robustness. On the other hand, we can remove any requirements of uniformity.

THEOREM 1.3. A finite S is interpretable in T if and only if for all $M \models T$, some $M' \models S$ is liberally defined in M .

Proof: The forward direction is obvious. For the reverse direction, we use compactness and completeness.

We claim that M' can be liberally defined in M , uniformly, up to a choice of parameters. I.e., there is a fixed presentation of $L(S)$ data in $L(T)$ such that for all $M \models T$, there is some choice of parameters from $\text{dom}(M)$ such that the data constitutes a model $M' \models S$.

Suppose this is false. Then for any such presentation, there exists $M \models T$ such that for all choices of parameters from $\text{dom}(M)$, the presentation fails to define a model of S .

We then conclude that for any finite set of such presentations, there exists $M \models T$ such that for at least one of these presentations, for all choices of parameters from $\text{dom}(M)$, the presentation fails to define a model of S . For if we had a finite set of such presentations without this property, then we could get a single such presentation without this property through definition by finitely many cases.

Let X be the infinite set of sentences asserting T , and that for each of the countably many presentations, and all choices of parameters, the presentation fails to define a model of S . By hypothesis, we know that X is not satisfiable. Hence some finite subset of X is not satisfiable. This finite subset of X provides a counterexample to the previous paragraph. QED

The case of finite S, T (in finite relational types) is fundamental. We can create unimpeachable certificates of interpretability: the required definitions in $L(T)$ together with the required proofs in T .

We will identify these certificates with their Gödel numbers when convenient. Thus we speak of one certificate being smaller than another if and only if the Gödel number of the first certificate is less than the Gödel number of the second certificate.

Unimpeachable certification is important for relative consistency, where issues of impeachability are crucial.

After Tarski, many of the general investigations into interpretability, such as in Feferman, concentrated on the non finitely axiomatized cases (reflexive theories). We touch on this case in section 6, and infinite theories also are used in section 7.

In the infinite case (where S is infinite), it is natural to disallow the use of parameters, and consider only parameterless interpretations. Obviously the semantic version of parameterless interpretations is equivalent, by the completeness theorem, to the provability version of parameterless interpretations. However, we will lose Theorem 1.3.

We now focus on the finite case. By taking conjunctions, we can assume that we have single sentences A, B . We will have a limited discussion of the infinite case in section 6. Also infinite theories are present in section 7.

2. Tarski Degrees.

$S \sqsubseteq I T$ means that the theory S is interpretable in the theory T . Define

$$\begin{aligned} S =_I T &\iff S \sqsubseteq I T \text{ and } S \sqsupseteq I T. \\ S <_I T &\iff S \sqsubseteq I T \text{ and } \neg T =_I S. \end{aligned}$$

For the purposes of this definition, we view theories as sets of sentences without mentioning any underlying relational type. Underlying relational types play no significant role here.

Obviously $\sqsubseteq I$ is reflexive and transitive. The equivalence classes of theories, under the equivalence relation $=_I$, are called the Tarski degrees.

It is easy to show that there is a proper class of Tarski degrees, corresponding to the fact that there are a proper class of theories.

If we restrict to countable theories, then there are 2^{\aleph_0} countable theories, and 2^{\aleph_0} associated Tarski degrees. If we restrict to finite theories in finite relational types, there are \aleph_0 associated Tarski degrees. Obviously, restricting to finite theories is the same as restricting to single sentences.

The unary Tarski degrees are obtained by restricting $\sqsubseteq I$ to single sentences. I.e., they are the equivalence relations of single sentences under $=_I$. There are \aleph_0 unary Tarski degrees.

The unary Tarski degrees bear a rough resemblance to two very well studied countable structures from recursion theory - the Turing degrees (of subsets of \mathbb{N}) and the r.e. sets (subsets of \mathbb{N}). We know from experience that the latter two structures are very complicated in various ways. It is not clear how complicated the unary Tarski degrees are, and how these three countable structures are related.

THEOREM 2.1. The unary Tarski degrees form a (reflexive) partial ordering with a minimum and maximum element. The maximum element 1 is the equivalence class of all sentences with no models. The minimum element 0 is the equivalence class of all sentences with a finite model.

Proof: Obviously every A is interpretable in any B with no models. We claim that in any M , we can define all models of all nonzero finite cardinalities using a single parameter x and the multidimensional apparatus. I.e., we can use x , (x,x) , (x,x,x) , ..., (x,x,\dots,x) as the domain elements. Thus any A with a finite model is interpretable in any B . QED

THEOREM 2.2. The unary Tarski degrees form a distributive lattice with a minimum and maximum element.

Proof: Let A, B be sentences. Define $\inf([A], [B]) = [A \ \& \ B]$. To see that this is the inf, we have to verify that $C \sqsubseteq [A, B] \iff C \sqsubseteq [A \ \& \ B]$. The reverse direction is obvious. Now assume $C \sqsubseteq [A, B]$. Let $M \models A \ \& \ B$. If $M \models A$ then there is a model of C that is liberally defined in M . If $M \models B$ then there is a model of C that is liberally defined in M .

$\sup([A], [B])$ requires more care. Introduce a new monadic predicate symbol R . Take

$$\sup([A], [B]) = [(\exists x, y) (Rx \ \& \ \neg Ry) \ \& \ A^R \ \& \ B^{\exists R}].$$

To see that this is the sup, we have to verify that

$$A, B \sqsubseteq [C] \iff (\exists x, y) (Rx \ \& \ \neg Ry) \ \& \ A^R \ \& \ B^{\exists R} \sqsubseteq [C].$$

The reverse direction is obvious. Now assume $A, B \sqsubseteq [C]$. Let $M \models C$. Let M', M'' be liberally defined in M , where $M' \models A$, $M'' \models B$. Using multidimensionality, we can disjointify the domains of A, B in order to liberally define a model of $(\exists x, y) (Rx \ \& \ \neg Ry) \ \& \ A^R \ \& \ B^{\exists R}$ in M .

For distributivity, because we are in a lattice, it suffices to check

$$\sup(x, \inf(y, z)) = \inf(\sup(x, y), \sup(x, z)).$$

Choose representatives A, B, C of x, y, z . Choose unary R not in A, B, C . It suffices to show

$$\begin{aligned}
 & ((\exists x, y) (Rx \wedge \neg Ry) \wedge A^R \wedge (B \vee C))^{\exists R} = I \\
 & ((\exists x, y) (Rx \wedge \neg Ry) \wedge A^R \wedge B^{\exists R}) \quad ((\exists x, y) (Rx \wedge \neg Ry) \wedge A^R \wedge C^{\exists R}).
 \end{aligned}$$

In fact, we have logical equivalence. QED

THEOREM 2.3. $\text{sup}(x, y) = 1 \iff x = 1 \wedge y = 1$. $\text{inf}(x, y) = 0 \iff x = 0 \wedge y = 0$. I.e., in the Tarski degrees, 1 is not a sup, and 0 is not an inf.

Proof: Let $[A], [B] < I 1$. Then $(\exists x, y) (Rx \wedge \neg Ry) \wedge A^R \wedge B^{\exists R}$ has a model, and so $\text{sup}([A], [B]) < I 1$. Let $[A], [B] > 0$. Then $A \wedge B$ has no finite model, and so $\text{inf}([A], [B]) > 0$. QED

It will be convenient to write $A \wedge B$ for $(\exists x, y) (Rx \wedge \neg Ry) \wedge A^R \wedge B^{\exists R}$. Here R is not in A, B . We can live with the ambiguity of which R is chosen.

THEOREM 2.4. [Fe60]. Every Tarski degree $< I 1$ lies strictly below some Tarski degree $< I 1$.

Proof: Let A be consistent. Consider $T \wedge \text{Con}(A)$, where T is the conjunction of a reasonably healthy finitely fragment of $PA =$ Peano Arithmetic. $A \wedge I T \wedge \text{Con}(A)$ by formalizing Gödel's completeness theorem in T . Since $T \wedge \text{Con}(A)$ is true, $T \wedge \text{Con}(A) < I 1$.

Suppose $T \wedge \text{Con}(A)$ is interpretable in A . Then

$T \vdash \text{Con}(A) \wedge \text{Con}(T \wedge \text{Con}(A))$.
 $T + \text{Con}(A) \vdash \text{Con}(T + \text{Con}(A))$.
 $T + \text{Con}(A)$ is inconsistent (by Gödel's second incompleteness theorem).

But $T + \text{Con}(A)$ is true. QED

A particularly appropriate choice of $T \wedge PA$ is EFA = exponential function arithmetic = $I\mathbb{N}_0(\text{exp})$. See [HP98], p. 37-44, p. 495.

EFA is in $0, S, +, \cdot, \text{exp}, <$ with the proper axioms

1. $Sx \neq 0$.
2. $Sx = Sy \iff x = y$.
3. $x+0 = x$.
4. $x+Sy = S(x+y)$.

5. $x \cdot 0 = 0$.
6. $x \cdot Sy = x \cdot y + x$.
7. $\exp(x, 0) = S0$.
8. $\exp(x, Sy) = \exp(x, y) \cdot x$.
9. $x < y \rightarrow (\exists z)(y = x + Sz)$.
10. Induction for all bounded formulas.

It is known that EFA is finitely axiomatizable. E.g., this can be obtained from [HP98], Theorem 5.8, p. 366.

A Σ_1^0 sentence is a sentence in $0, S, +, \cdot, <$, beginning with universal quantifiers, followed by a formula with only bounded quantifiers.

LEMMA 2.5. There is a polynomial time function σ from Σ_1^0 sentences to sentences in $0, S, +, \cdot, \exp, <$, such that

- if $EFA \vdash \sigma$ then $\sigma(\sigma)$ has a finite model.
 if $EFA \vdash \neg\sigma$ then $\sigma(\sigma)$ has no model (A is inconsistent).

Proof: Let σ be Σ_1^0 . Let $\sigma(\sigma)$ be the conjunction of:

1. $<$ is a linear ordering with minimum 0, and at least three elements.
2. x is greatest $\rightarrow Sx = x+y = y+x = x \cdot y = y \cdot x = \exp(x) = x$.
3. $x < y \rightarrow Sx$ is the immediate successor of x in $<$.
4. $x \neq 0 \rightarrow (\exists y)(x = Sy)$.
5. $(x, y < z \rightarrow Sx = Sy) \rightarrow x = y$.
6. $x < y \rightarrow x + 0 = x$.
7. $x < y \rightarrow x \cdot 0 = 0$.
8. $S(x+y) < z \rightarrow x+Sy = S(x+y)$.
9. $x \cdot y + x < z \rightarrow x \cdot Sy = x \cdot y + x$.
10. $x < y \rightarrow \exp(x, 0) = S0$.
11. $\exp(x, y) \cdot x < z \rightarrow \exp(x, Sy) = \exp(x, y) \cdot x$.
12. x is greatest \rightarrow there is a proof with Gödel number $< x$ of $EFA \rightarrow \sigma$, and there is no proof with Gödel number $< x$ of $EFA \rightarrow \neg\sigma$.
13. There is no greatest element $\rightarrow EFA \rightarrow \text{Con}(EFA + \sigma)$.

In the above, we assume a fixed finite axiomatization of EFA, and also that "x is the Gödel number of a proof of the formula with Gödel number y" is formalized by a bounded formula in the language of EFA, where y is formally required to bound all of the values of the relevant terms.

Suppose $EFA \vdash \sigma$. Let n be the Gödel number of a proof of $EFA \rightarrow \sigma$. We can build a model of $\sigma(\sigma)$ on $0, 1, \dots, n+1$ as

follows. Take $<$ to be as usual. Define $S, +, \cdot, \exp$ on $[0, n]$ as usual provided the result is $\leq n$; $n+1$ otherwise. Define $S, +, \cdot, \exp$ to be $n+1$ if at least one argument is $n+1$. Axioms 1-7, 10 are immediate.

For 8, let $S(x+y) < z$ in this model. Then $\text{val}(x+y) < n$. Hence $\text{val}(x+y) = \text{val}(x) + \text{val}(y) < n$. Hence $\text{val}(S(x+y)) = \text{val}(x) + \text{val}(y) + 1$. Also $\text{val}(x), \text{val}(y) < n$. Hence $\text{val}(Sy) = \text{val}(y) + 1 \leq n$. And $\text{val}(x) + \text{val}(Sy) = \text{val}(x) + \text{val}(y) + 1 \leq n$. Hence $\text{val}(x+Sy) = \text{val}(x) + \text{val}(y) + 1 = \text{val}(S(x+y))$.

For 9, let $\text{val}(x \cdot y + x) < z$ in this model. Then $x, y \leq n$. Also $\text{val}(x \cdot y + x) = \text{val}(x \cdot y) + x \leq n$. Hence $\text{val}(x \cdot y) = x \cdot y \leq n$. Hence $\text{val}(x \cdot y + x) = x \cdot y + x = x(y+1) \leq n$. Also $\text{val}(Sy) = y+1$. Hence $\text{val}(x \cdot Sy) = x(y+1)$ since $x(y+1) \leq n$.

For 11, assume $\exp(x, y) \cdot x < z$ in this model. Then $\text{val}(\exp(x, y)) \leq n$, and so $y \leq n$, $S(y) \leq n+1$, $\text{val}(\exp(x, y)) = x^y \leq n$, $x^y \cdot x \leq n$, $\text{val}(\exp(x, y) \cdot x) = x^y \cdot x = x^{y+1}$. Hence $\text{val}(\exp(x, Sy)) = x^{y+1}$ since $x^{y+1} \leq n$.

For 12, let x be greatest. Then $x = n+1$, and n is the Gödel number of a proof of $\text{EFA} \vdash \perp$. Hence this holds in the model. Also, there cannot be a proof with Gödel number $\leq n$ of $\text{EFA} \vdash \perp$, according to the model, since otherwise, there would actually be a proof of $\text{EFA} \vdash \perp$. But this contradicts that there is a proof of $\text{EFA} \vdash \perp$.

Obviously 13 is vacuously true in this model.

Now suppose $\text{EFA} \vdash \perp$. Let M be a model of $\perp(\perp)$. Suppose there is a greatest element x in M . In M , $<$ is a linear ordering with least element 0, and $<$ must be discrete by axioms 3, 4. Hence we can generate standard integers $0, 1, 2, \dots$ which either goes on forever, or stops when it reaches x .

If it goes on forever, then this standard part is a copy of the standard model with ordinary $+, \cdot, \exp$. Hence the second part of axiom 12 is violated.

If it does not go on forever, then x is standard. We again contradict the second part of axiom 12.

Now suppose there is no greatest element x in M . By axiom 13, $M \models \text{EFA} + \text{Con}(\text{EFA} + \perp)$. But this is again impossible by $\text{EFA} \vdash \perp$. QED

THEOREM 2.6. The unary Tarski degrees form a dense distributive lattice with $0 < 1$. I.e., $a <_I b \iff (\exists c)(a <_I c <_I b)$.

Proof: Let A, B be sentences, where $A <_I B$. By self reference, let φ be a Σ^0_1 sentence such that EFA proves φ if and only if

Every interpretation certificate of $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ in A has a smaller interpretation certificate of B in $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$.

Suppose $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B \not\leq_I A$, and let n be the least interpretation certificate.

case 1. There is an interpretation certificate of B in $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ that is smaller than n . Then $EFA \vdash \varphi$. By Lemma 2.5, $\varphi(\varphi)$ has a finite model. Now there is an interpretation of B in $A \leftrightarrow \varphi(\varphi)$. Since $A \leftrightarrow \varphi(\varphi) \leq_I A$, we have $B \leq_I A$, which is a contradiction.

case 2. There is no interpretation certificate of B in $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ that is smaller than n . Then $EFA \vdash \neg \varphi$. By Lemma 2.5, $\varphi(\varphi)$ has no model. Hence $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ is logically equivalent to B . Therefore $B \leq_I A$, which is a contradiction.

We have refuted $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B \not\leq_I A$.

Now suppose $B \leq_I (A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$. Let n be the least interpretation certificate.

case 3. There is no interpretation certificate of $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ in A smaller than or equaled to n . Then $EFA \vdash \neg \varphi$. By Lemma 2.5, $\varphi(\varphi)$ has a finite model. Clearly $B \leq_I A \leftrightarrow \varphi(\varphi)$. Hence $B \leq_I A$, which is a contradiction.

case 4. There is an interpretation certificate of $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ in A smaller than or equaled to n . Then $EFA \vdash \varphi$. By Lemma 2.5, $\varphi(\varphi)$ has no model. Hence $(A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$ is logically equivalent to B . Therefore $B \leq_I A$, which is a contradiction.

We have refuted $B \leq_I (A \leftrightarrow \varphi(\varphi)) \leftrightarrow B$.

Since $A \models (A \models \varphi(x)) \rightarrow B \models B$, we have $A <_I (A \models \varphi(x))$
 $B <_I B$, as required. QED

We can sharpen Theorem 2.6 as follows.

In any lattice, let a_1, a_2, \dots be a finite or infinite sequence of points. We say that this sequence is independent if and only if in the lattice, for all $n, m \geq 1$,

if $\inf(a_{i_1}, \dots, a_{i_n}) \leq \sup(a_{j_1}, \dots, a_{j_m})$ then
 $\{i_1, \dots, i_n\} \cap \{j_1, \dots, j_m\} = \emptyset$.

THEOREM 2.7. In the unary Tarski degrees, let $a <_I b$. There exists an independent infinite sequence c_1, c_2, \dots lying in (a, b) .

Proof: Let $A \models a$, $B \models b$. Then $A <_I B$. By self reference, let $\varphi(n)$ be a Σ_1^0 formula with only the free variable n , such that EFA proves $\varphi(n)$ if and only if

for every interpretation certificate of $(A \models \varphi(n^*)) \rightarrow B$ in A , or one which ruins the independence of the finite sequence

$$(A \models \varphi(1^*)) \rightarrow B, \dots, (A \models \varphi(n^*)) \rightarrow B$$

using $(A \models \varphi(n^*)) \rightarrow B$ on the left (inside the inf), there is a smaller interpretation certificate of B in $(A \models \varphi(n^*)) \rightarrow B$, or a smaller interpretation certificate which ruins the independence of the finite sequence

$$(A \models \varphi(1^*)) \rightarrow B, \dots, (A \models \varphi((n)^*)) \rightarrow B$$

using $(A \models \varphi(n^*)) \rightarrow B$ on the right (inside the sup).

Here we use n^* for the closed term $S \dots S_0$, with n S 's.

We now show that the infinite sequence

$$(A \models \varphi(1^*)) \rightarrow B, (A \models \varphi(2^*)) \rightarrow B, \dots$$

is independent and lies in (a, b) . Obviously, this infinite list lies in $[a, b]$.

Let n be least such that it is not the case that

$$(A \models \varphi(1^*)) \rightarrow B, (A \models \varphi(2^*)) \rightarrow B, \dots,$$

$$(A \sqcup \sqcup (\Phi(n^*))) \quad B$$

is independent. We will obtain a contradiction.

Consider the following two cases.

case 1. $(A \sqcup \sqcup (\Phi(n^*))) \quad B \sqcup I A$, or there is a ruining of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n^*))) \quad B$, using $(A \sqcup \sqcup (\Phi(n^*))) \quad B$ on the left.

Suppose $EFA \vdash \Phi(n^*)$. Then by Lemma 2.5, $\Phi(\Phi(n^*))$ has no model, and so $(A \sqcup \sqcup (\Phi(n^*))) \quad B$ is logically equivalent to B . Hence $B \sqcup I A$ or there is a ruining of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n-1^*))) \quad B, B$, using B on the left. Since $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n-1^*))) \quad B \sqcup I B$, the B , being on the left side, which is an inf, gets absorbed. Hence $B \sqcup I A$ or there is a violation of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n-1^*))) \quad B$. The former is impossible, and the latter contradicts the choice of n .

Since EFA does not refute $\Phi(n^*)$, we see that $EFA \vdash \Phi(n^*)$. Hence by Lemma 2.5, $\Phi(\Phi(n^*))$ has a finite model, and so $(A \sqcup \sqcup (\Phi(n^*))) \quad B = I A \quad B = I A$.

Also since EFA does not refute $\Phi(n^*)$, we see that $B \sqcup I (A \sqcup \sqcup (\Phi(n^*))) \quad B$, or there is a ruining of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n^*))) \quad B$ using $(A \sqcup \sqcup (\Phi(n^*))) \quad B$ on the right. Hence $B \sqcup I A$, or there is a ruining of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi((n-1)^*))) \quad B, A$, using A on the right. Since $A \sqcup I (A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi((n-1)^*))) \quad B$, the A , being on the right side, which is a sup, gets absorbed. Hence $B \sqcup I A$, or there is a violation of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi((n-1)^*))) \quad B$. The former is impossible, and the latter contradicts the choice of n .

case 2. $B \sqcup I (A \sqcup \sqcup (\Phi(n^*))) \quad B$, or there is a ruining of the independence of the finite sequence $(A \sqcup \sqcup (\Phi(1^*))) \quad B, \dots, (A \sqcup \sqcup (\Phi(n^*))) \quad B$, using $(A \sqcup \sqcup (\Phi(n^*))) \quad B$ on the right. Because case 1 is impossible, we have $EFA \vdash \Phi(n^*)$. Hence by Lemma 2.5, $\Phi(\Phi(n^*))$ has a finite model. Hence $(A \sqcup \sqcup (\Phi(n^*))) \quad B = I A \quad B = I A$.

Hence $B \not\models A$, or there is a ruining of the independence of the finite sequence $(A \models \varphi(1^*)) \quad B, \dots, (A \models \varphi((n-1)^*)) \quad B, A$, using A on the right. The former is impossible. For the latter, since $A \models (A \models \varphi(1^*)) \quad B, \dots, (A \models \varphi((n-1)^*)) \quad B$, the A , being on the right side, which is a sup, gets absorbed. Hence there is a violation of the independence of the finite sequence $(A \models \varphi(1^*)) \quad B, \dots, (A \models \varphi((n-1)^*)) \quad B$. This contradicts the choice of n .

Since both cases lead to a contradiction, we see that there is no violation of the independence of the finite sequence $(A \models \varphi(1^*)) \quad B, \dots, (A \models \varphi(n^*)) \quad B$ using $(A \models \varphi(n^*)) \quad B$. Hence there is a violation of the independence of the finite sequence $(A \models \varphi(1^*)) \quad B, \dots, (A \models \varphi((n-1)^*)) \quad B$. This contradicts the choice of n .

Thus we have established the independence of the infinite sequence

$$(A \models \varphi(1^*)) \quad B, (A \models \varphi(2^*)) \quad B, \dots$$

It remains to show that these sentences lie in (a, b) . Let m be least such that $(A \models \varphi(m^*)) \quad B \models A$. Since there is no interpretation of B in $(A \models \varphi(m^*)) \quad B$, we see that $EFA \vdash \neg \varphi(m^*)$. Hence $\varphi(m^*)$ has no model, and so $B \models A$, which is a contradiction. Hence for all m , not $(A \models \varphi(m^*)) \quad B \models A$.

Let m be least such that $B \models (A \models \varphi(m^*)) \quad B$. Since there is no interpretation of $(A \models \varphi(m^*)) \quad B$ in A , we see that $EFA \vdash \neg \varphi(m^*)$. Hence $\varphi(m^*)$ has a finite model, and so $B \models A \quad B$ and $B \models A$. This is impossible. Hence for all m , not $B \models (A \models \varphi(m^*)) \quad B$. QED

Do all unary Tarski degrees other than $0, 1$ look alike? The answer is strongly no.

We have already seen that 1 is not a sup. Therefore the sup of any two incomparable Tarski degrees is a sup $\neq 0, 1$. However, other unary Tarski degrees $\neq 0, 1$ are not sups.

Let M, M' be two structures with disjoint domains. We write $M + M'$ for the structure with domain $\text{dom}(M) \cup \text{dom}(M')$, where the relations and constants of M and M' remain

unchanged, the functions of M are extended to $\text{dom}(M')$ by returning the first argument, and a new unary relation is introduced caring out $\text{dom}(M)$.

LEMMA 2.8. Let M, M' be two structures with disjoint domains. Suppose $M + M'$ liberally defines a linear ordering with no greatest element. Then M or M' liberally defines a linear ordering with no greatest element.

LEMMA 2.9. Let M, M' be two structures with disjoint domains. Suppose $M + M'$ liberally defines a pairing function with at least two inequivalent elements. Then M or M' alone liberally defines a pairing function with at least two inequivalent elements.

THEOREM 2.10. In the unary Tarski degrees, "Linear ordering with no greatest element" is not a sup. "Pairing function with at least two elements" is not a sup.

Proof: Suppose [linear ordering with no greatest element] = $\text{sup}([A], [B])$. Let $M \models A \wedge \neg B$. Let M' be the part of M for A and M'' be the part of M for B . Since M liberally defines a linear ordering with no greatest element, M' or M'' liberally defines a linear ordering with no greatest element. Hence $[A]$ or $[B]$ is [linear ordering with no greatest element]. The analogous argument proves the second claim. QED

A sentence A is said to be complete if and only if for all sentences B in $L(A)$, $A \vdash B$ or $A \vdash \neg B$.

THEOREM 2.11. Let A be a complete sentence. Then $[A]$ is not an inf in the unary Tarski degrees.

Proof: Let A be a complete sentence. Let $[A] = \text{inf}([B], [C]) = [B \wedge C]$. Let $M \models A$. Then M has a liberally defined model of $B \wedge C$. Hence M has a liberally defined model of B or M has a liberally defined of C . We now use the completeness of A . In the former case, every $M \models A$ has a liberally defined model of B . In the latter case, every $M \models A$ has a liberally defined model of C . QED

3. Adequate Sentences and Relative Consistency.

A very important condition is adequacy of a theory T . It is convenient to first define adequacy for a model.

We say that M is adequate if and only if there is an M definable binary relation whose cross sections are closed under single point additions and single point deletions.

Another, more common, definition is that there is an M definable binary relation with an empty cross section, whose cross sections are closed under single point additions.

THEOREM 3.1. These two definitions of adequacy are equivalent.

Proof: Let M be a model with an M definable binary relation R whose cross sections are closed under single point additions and single point deletions. Fix $c \in \text{dom}(M)$. Define $S(x,y) \iff (R(x,y) \iff \neg R(c,y))$. Clearly the cross section of S at c is empty. Now let x,y be given. If $\neg R(c,y)$ then choose a cross section of R obtained by inserting y in the cross section of R at x . If $R(c,y)$ then choose a cross section of R obtained by deleting y from the cross section of R at x .

Let M be a model with an M definable binary relation R with an empty cross section, whose cross sections are closed under single point additions.

We first use the following construction for getting: an empty cross section, and closure under single point additions and pairwise intersection. Call a cross section good if and only if its intersection with every cross section is a cross section. Obviously, empty cross sections are good. Let X,Y be good cross sections and $x \in X$. To check that $X \setminus \{x\}$ is good, clearly $X \setminus \{x\} \cap \{y\}$ and $(X \setminus \{x\}) \cap Z$ are cross sections since $X \cap Z$ is a cross section. To check that $X \cap Y$ is good, clearly $(X \cap Y) \cap \{x\}$ and $(X \cap Y) \cap Z$ are cross sections, since $X \cap Y, Y \cap Z$ are cross sections.

Call a cross section X of R very good if and only if it is good, and any single point deletions of any good cross section that is a subset of X , is a good cross sections.

We claim that the very good cross sections of R are closed under single point additions and single point deletions. To see this, let X be a very good cross section and $x \in \text{dom}(M)$. We must show that $X \setminus \{x\}$ and $X \cup \{x\}$ are both very good cross sections. Clearly there are both good cross

sections. We can assume without loss of generality that $x \in X$.

Let $Y \cap X \cap \{x\}$ be a good cross section of R . Any $Y \setminus \{y\}$ is $Y \setminus \{x\} \setminus \{y\}$ or $(Y \setminus \{x\} \setminus \{y\}) \cap \{x\}$. Since X is good, $X \cap Y$ is good, and $X \cap Y = Y \setminus \{x\} \cap X$. Since X is very good, $Y \setminus \{x\}$ is a good subset of X , and so $Y \setminus \{x\} \setminus \{y\}$ and $(Y \setminus \{x\} \setminus \{y\}) \cap \{x\}$ are good. Therefore $X \cap \{x\}$ is very good.

Let $Y \cap X \setminus \{x\}$ be good. Since $Y \cap X$, single point deletions to Y remain good. QED

We give four equivalent forms of adequacy for a theory T . We begin with two syntactic definitions.

We say that a theory T is adequate if and only if there exists a formula $\phi(x, y, z_1, \dots, z_n)$ in $L(T)$ such that the following is provable in T . There exists $1 \leq i \leq k$ (as a disjunction) and z_1, \dots, z_n (normal universal quantification) such that

- i. $(\forall x) (\forall y) (\exists \bigwedge_i (x, y, z_1, \dots, z_n))$.
- ii. $(\forall u, w) (\forall x) (\forall y) (\bigwedge_i (x, y, z_1, \dots, z_n) \rightarrow \bigwedge_i (u, y, z_1, \dots, z_n) \rightarrow y = w)$.
- i'. $(\forall u, w) (\forall x) (\forall y) (\bigwedge_i (x, y, z_1, \dots, z_n) \rightarrow \bigwedge_i (u, y, z_1, \dots, z_n) \rightarrow y = w)$.
- ii'. $(\forall u, w) (\forall x) (\forall y) (\bigwedge_i (x, y, z_1, \dots, z_n) \rightarrow \bigwedge_i (u, y, z_1, \dots, z_n) \rightarrow y \neq w)$.

In the case of finite T , it is easily seen that adequacy of T is equivalent to the model theoretic condition that every model of T is adequate (in any of the two senses shown equivalent by Theorem 3.1). This uses a compactness argument.

It will also be convenient to rephrase the adequacy of T in terms of interpretations.

We call an interpretation regular if the domain used for the interpretation is the full domain, and the equality relation used is equality.

It is clear that the adequacy of T is equivalent to the regular interpretability of the following weak system NW of set theory, formulated in the usual language $\in, =$.

N (Null axiom). $(\forall x)(\forall y)(y \subseteq x)$.

W (With axiom). $(\forall y)(\forall z)(z \subseteq y \subseteq (z \subseteq x \rightarrow z = x))$.

Or we can equivalently use WD.

W (With axiom). $(\forall y)(\forall z)(z \subseteq y \subseteq (z \subseteq x \rightarrow z = x))$.

D (Delete axiom). $(\forall y)(\forall z)(z \subseteq y \subseteq (z \subseteq x \rightarrow z \neq x))$.

The system NW is credited to Ed Nelson in [MM94], where it is shown that Q is interpretable in NW. (They also credit Jan Krajicek for an earlier and unpublished proof of this result, not known to them at the time of their publication).

The interpretability of NW in Q was already mentioned in [TMR53], p. 34, where the authors report that the interpretability of even NWE had been established by W. Szemielew and A. Tarski in 1950. Here E is

E. Extensionality axiom. $(\forall z)(z \subseteq x \subseteq z \subseteq y) \rightarrow x = y$.

The argument in Theorem 3.1 establishes the mutual regular interpretability of WD and NW. However, WD, NW, WDE, NWE are only mutually interpretable, not mutually regularly interpretable. E.g., they have the same Tarski degree.

In fact, WD represents a very important Tarski degree with many different kinds of representatives. A particularly important representative of this Tarski degree is the well known weak system of arithmetic called Q, or Robinson's arithmetic after R. Robinson [Ro52]. Also see [HP98], p. 28. The language is $0, S, +, \cdot, \subseteq, =$, and the nonlogical axioms are

1. $S(x) \neq 0$.
2. $S(x) = S(y) \rightarrow x = y$.
3. $x \neq 0 \rightarrow (\exists y)(S(y) = x)$.
4. $x + 0 = x$.
5. $x + S(y) = S(x + y)$.
6. $x \cdot 0 = 0$.
7. $x \cdot S(y) = (x \cdot y) + x$.
8. $x \subseteq y \rightarrow (\exists z)(z + x = y)$.

Obviously the explicit definitional axiom 8 is redundant. However, since for many purposes, Q is normally extended by additional axioms that do involve \subseteq , it is convenient to incorporate 8.

One such important extension of Q is the system $I\Delta_0$. This also goes under the name "bounded arithmetic". I prefer the name PFA = polynomial function arithmetic. In this system, we add the induction scheme for bounded formulas in $L(Q)$ to Q .

A bounded formula in $L(Q)$ is a formula all of whose quantifiers are bounded to terms that do not contain the variable; i.e., $(\exists x \leq t)$, $(\forall x \leq t)$, where t is a term of $L(Q)$ in which x does not appear. Thus the nonlogical axioms of PFA are 1-8 and

9. Bounded induction scheme. $(\exists [x/0] \rightarrow (\forall x) (\phi \rightarrow \phi[x/Sx])) \rightarrow \phi$, where ϕ is a bounded formula of $L(Q)$.

PFA is believed to be not finitely axiomatizable, although this is not known. See [HP98], p. 350.

It is well known that PFA is interpretable in Q . See [HP98], p. 366.

THEOREM 3.2. WD, NW, WDE, NWE, Q, PFA are mutually interpretable.

Theories of a tame character such as Presburger arithmetic, real closed fields, and algebraically closed fields, are interpretable in Q , but cannot interpret Q .

Let T be adequate. An adequacy mechanism for T consists of formulas of $L(T)$ that

- i. define the internal nonnegative integers by a unary predicate, together with an equivalence relation representing equality. T must prove that this is an equivalence relation E .
- ii. define $<, 0, S, +, \cdot$ on the natural numbers. T must prove everywhere definedness and univalence (with respect to E).
- iii. define the internal finite sequences by a unary predicate.
- iv. define the length function and value function for finite sequences (i.e., the i -th term).
- v. T proves that there is a finite sequence of length 0, and that the finite sequences are closed under appending any object on the right - and that appending raises the length by 1.

vi. T proves the axioms of PFA formulated using ii above, with bounded induction extended to incorporate all of the symbols of $L(T)$.

THEOREM 3.3. Every adequate theory has an adequacy mechanism.

We need a notion of restricted proof in first order predicate calculus with equality. For definiteness, we use the following logical axioms and rules in $\square, \square, \square, \square, \square, \square, \square, =$.

- i. All tautologies.
- ii. $(\square x)(\square) \square \square[x/t]$, t substitutable for x in \square .
- iii. $\square[x/t] \square (\square x)(\square)$, t substitutable for x in \square .
- iv. $x = x$.
- v. $x = y \square (\square \square \square')$, where \square is atomic and \square' is obtained from \square by replacing an occurrences of x in \square by y .
- vi. From $\square \square \square$ derive $\square \square (\square x)(\square)$, where x is not free in \square .
- vii. From $\square \square \square$ derive $(\square x)(\square) \square \square$, where x is not free in \square .
- viii. From $\square, \square \square \square$, derive \square .

A proof consists of a finite sequence of formulas, each of which is either an axiom above, or follows from previous entries in the sequence by one or more rules of inference above. A proof of \square is a proof whose last entry is \square . We write $\text{PROV}(\square)$.

For any formula \square , we define $\text{lth}(\square)$ to be the total number of symbols in \square , including parentheses, connectives, quantifiers, constants, relations, functions, and variables. Officially, the constants, relations, functions, and variables are given with a single symbol adorned with subscripts in binary. Each occurrence of these symbols is charged one plus the length of the subscript.

Note that there is a small fixed constant b such that the number of formulas \square with $\text{lth}(\square) \leq n$ is at most b^n . In fact, we take b to be the length of

$\square \square \square \square \square \square \square = R F c 0 1 () ,$

which is 16.

A restricted proof is a proof with the additional requirements that

- i. All nonlogical symbols used are present in the last entry, \square .
- ii. Every entry is a propositional combination of formulas, each of which have at most l th(\square) quantifiers.

We write $RPROV(\square)$ for " \square has a restricted proof". The size of a proof is the total number of symbols that it contains. This is the same as the sum of the lengths of its entries. We write $RPROV(\square, n)$ for " \square has a restricted proof of size n ".

We write $RCON(A)$ for $\square RPROV(\square A)$.

We write $2^{[p]}(n)$ for an exponential stack of p 2's with n on top, where if $p = 0$, this is n .

THEOREM 3.4. Let $(\square_n)(\square_m)(\square(n,m))$ be a \square_2^0 sentence, where \square is bounded in $L(Q)$. The following are equivalent.

- i. $EFA \vdash (\square_n)(\square_m)(\square(n,m))$.
- ii. There exists $p \leq N$ such that $PFA \vdash (\square_n, r)(r = 2^{[p]}(n) \square (\square_m \square r)(\square(n,m)))$.

Proof: In ii, we are using an adequate formalization of the exponentiation relation in PFA, which is well known to exist. Let $\square = (\square_n)(\square_m)(\square(n,m))$ be as given. Assume i. Suppose ii is false. Let T be the theory in $L(Q)$ with constants d, e , consisting of

- i. PFA.
- ii. The scheme $d = 2^{[p]}(e) \square (\square_r \square d)(\square\square(d,r)), p \geq 0$.

Since ii is false, every finite fragment of T is consistent. Hence by compactness, T is consistent. Let $M \models T$. Let M' be the initial segment of M determined by the iterated exponentials of d . Then $M' \models EFA \square (\square_n)(\square_m)(\square(n,m))$. Hence ii holds.

The converse ii \square i is trivial. QED

Let M be adequate (see Theorem 3.1), and let an adequacy mechanism be given (see Theorem 3.3). A cut in M is taken to be an initial segment of (the numerical part of) M , with no greatest element. A proper cut in M is the cut in M consisting of all nonnegative integers of M (under the given adequacy mechanism). We will work with M definable cuts in M , only.

The following result was essentially proved in [Fr76], [Fr80], and exposited in [Sm85], [Vi90].

THEOREM 3.5. (Interpretability = relative restricted consistency). Let A, B be sentences, where B is adequate. The following are equivalent.

i. $A \vdash B$.

ii. $EFA \vdash RCON(B) \rightarrow RCON(A)$.

$i \rightarrow ii$ does not require the adequacy of B . However, $ii \rightarrow i$ does require the adequacy of B .

Proof: $i \rightarrow ii$. Let σ be an interpretation of A in B . Now σ specifies a parameterized family of structures M^* . I.e., there are free variables for the unspecified parameters. We have a proof from B that parameters can be chosen so that σ actually defines a specific model M^* of A .

We now argue in EFA . Let σ be a restricted proof of σA . So σ converts σ to a proof from B of

there exists a choice of parameters such that the specific model M^* satisfies A as well as $\sigma_1, \dots, \sigma_k = \sigma A$.

Obviously, this gives us a proof of a contradiction from B , and hence a proof of σB .

But how complicated are the formulas involved in the proof of σB obtained in this way? There are a number of sources of complexity here, including the use of function symbols in A that give rise to terms of uncontrolled complexity that have to be interpreted in M^* .

So we only claim that the formulas involved in the refutation of B have a standard integer bound on their quantifier depth. Quantifier depth counts blocks of like quantifiers like a single quantifier.

Fortunately, cut elimination has been reworked in terms of quantifier depth. See, e.g., [Zh91], [Zh94], [Ge05]. Since we have only standard quantifier depth, we can reduce the cut complexity to zero in a finite number of steps within EFA . EFA will take care of the exponential blowup (no worse) at each stage.

The implication $i \Rightarrow i$ is considerably more involved. Here cuts in nonstandard models are used, combined with formalized completeness proofs.

Let A, B be given, where B is adequate. Let an adequacy mechanism for B be given (Theorem 3.3).

Assume $EFA \vdash RCON(B) \Rightarrow RCON(A)$. Obviously

$EFA \vdash (\exists n)(\exists m)(RPROV(\ulcorner A, n \urcorner) \Rightarrow RPROV(\ulcorner B, m \urcorner))$.

Hence by Theorem 3.4, let p be such that

$I\ulcorner_0 \vdash (\exists n, r)(r = 2^{\lceil p \rceil}(n) \Rightarrow RPROV(\ulcorner A, n \urcorner) \Rightarrow RPROV(\ulcorner B, r \urcorner))$.

Let $M \models B$. We need to define a model $M^* \models A$ in M .

We use an attempted Henkin construction of a model of A in M . There are a number of problems that we encounter in order to make this work. All of these problems stem from the fact that we don't have much induction in M . However, we recover by making things work in M definable cuts.

Let $L(A)'$ be $L(A)$ together with the Henkin constants d_0, d_1, \dots . Let W be the set of all sentences \ulcorner of $L(A)'$ which are propositional combinations of sentences whose number of quantifiers is at most $\#(\ulcorner A)$, and whose nonlogical symbols appear in A .

We begin with listing certain elements of W . If EFA holds in (the numerical part of) M , then there is no problem listing all elements of W . Specifically, for each k , there is a unique (coded) finite sequence consisting of the first k elements of W , ordered first by length and second lexicographically. These finite sequences are initial segments of each other. Thus we can form the M definable $\ulcorner_0, \ulcorner_1, \dots$ which enumerates all elements of W in order.

If EFA fails in M , then there is an M definable cut. By cut shortening techniques, we can construct an M definable cut C which is closed under multiplication, and where there exists $t \in C$ such that 2^{2^t} exists. Then we can easily enumerate the first t elements of W , lexicographically, in order, as

$\ulcorner_0, \ulcorner_1, \dots, \ulcorner_t$.

We can restrict the indices to elements of C , obtaining $\{\varphi_i : i \in C\}$ as a sequence.

It is easy to see that $\{\varphi_i : i \in C\}$ enumerates exactly the elements of W whose length n has $2^n \in C$. (Recall the discussion above that $\{\varphi : \text{length}(\varphi) \in n\}$ has at most 16^n elements, and that C is closed under addition).

Note that $\{\varphi_i : i \in C\}$ is closed under propositional combinations since C is closed under propositional combinations roughly add lengths, and C is closed under multiplication. Also note that the terms appearing in the elements of $\{\varphi_i : i \in C\}$ are closed under the term building operations, since the term building operations roughly add lengths, and C is closed under multiplication.

So far, we can consolidate the cases where M satisfies EFA and does not satisfy EFA, by allowing C to be the improper cut of all $n \geq 0$, or C to be a proper cut. In either case, C is an M definable cut closed under multiplication.

We write $\log(C)$ for the cut $\{2^i : i \in C\}$. Then $\log(C)$ is an M definable cut closed under addition, and $\{\varphi_i : i \in C\} = \{\varphi \in W : \text{length}(\varphi) \in \log(C)\} = W[\log(C)]$.

Recall that we are working entirely within M .

case 1. $\text{RCON}(A)$. A finite sequence of sentences $\varphi_0, \dots, \varphi_{2k+1}$ from $W[\log(C)]$ is called good for A if and only if it satisfies the following conditions for all $0 \leq i \leq k$.

- i. $k \in \log(C)$.
- ii. φ_{2i} is φ_i or $\neg\varphi_i$.
- iii. If φ_{2i} is $\neg(\exists x)\varphi(x)$ and c is the first constant not appearing in $\varphi_0, \dots, \varphi_{2i}$, then φ_{2i+1} is $\varphi(x/c)$.
- iv. If φ_{2i} is not of the form $\neg(\exists x)\varphi(x)$, then $\varphi_{2i+1} = \varphi_{2i}$.
- v. There is no proof of $\varphi(A \wedge \varphi_0 \wedge \dots \wedge \varphi_{2k+1})$ of any size, whose entries lie in $W[\log(C)]$.
- vi. For all $0 \leq i \leq k$, if $\varphi_{2i} = \varphi_i$, then there is a proof of $\varphi(A \wedge \varphi_1 \wedge \dots \wedge \varphi_{2i-2} \wedge \varphi_i)$ of any size, whose entries lie in $W[\log(C)]$.

It is easy to see that there exists φ_0, φ_1 that is good for A . Also, suppose $\varphi_0, \dots, \varphi_{2k+1}$ is good for A , $k \geq 0$. Then there is an extension $\varphi_0, \dots, \varphi_{2k+3}$ which is still good for A .

We now assume that there are good sequences for A of every even length $\leq \log(C)$, and also for any two sequences good for A , one is an extension of the other. Of course, these two suppositions are obvious if we had enough induction in M - which we don't.

The previous paragraph defines a particular "path" if we think of this set of M sequences as forming an internal tree.

We can define a structure on the closed terms in the obvious way from the sequences good for A . Note that the atomic sentences in $L(A)'$ among the \square_i are exactly the atomic sentences whose length is in $\log(C)$, which are the atomic sentences whose closed terms have length in $\log(C)$. These closed terms are closed under the term building operations since $\log(C)$ is closed under addition.

The domain of the prospective model of A is the set of closed terms of length in $\log(C)$. We make this into a structure interpreting $L(A)'$ in the obvious way. The interpretation of equality is: $s \equiv t \iff s = t$ lies on the path.

We have defined an interpretation M^* of $L(A)'$. Obviously M^* is M definable. We have to show that as a structure external to M , $M^* \models A$.

Let q be the number of quantifiers in A . We show by external induction on j from 0 through q , that for all standard formulas \square in $L(A)'$ with at most j quantifiers,

\square holds in M at any assignment comprised of internal closed terms of M of length $\leq \log(C)$, in the external sense, if and only if the result of making the substitution, creating a sentence in $W[\log(C)]$ internally in M , lies on the internal path.

This follows from the usual basic facts about what is sitting on the path. Since A is a sentence on the path, we see that M^* is in fact a model of A . Both M and the outside world agree on this.

case 2. $\square \text{RCON}(A)$. Since $M \models \text{PFA}$, we have

$$M \models (\square n, r) (r = 2^{[p]}(n) \iff \text{RPROV}(\square A, n) \iff \text{RPROV}(\square B, 2^{[p]}(n)))$$

where p was chosen above in advance of the choice of M .

We now produce an M definable cut C such that $(\exists n \square C) (RCON(A, n))$.

If there is no least n such that $RPROV(\square A, n)$, then we immediately have our desired cut. Namely the cut of all n such that $RCON(A, n)$.

Now let n be least such that $RPROV(\square A, n)$.

case 2a. $2^{[p]}(n)$ does not exist. By cut shortening techniques, we can explicitly construct an M definable cut C below n .

case 2b. $2^{[p]}(n)$ exists. Then $RPROV(\square B, 2^{[p]}(n))$. Let \square be a witness to this statement. We now construct an M definable cut below $2^{[p]}(n)$ as follows.

We would like to give an M definable valuation of all terms in $L(B)$ at all assignments of objects. However, we may not have enough induction in M . In M , let C be the set of all m such that for every finite list t_1, \dots, t_n of terms in $L(B)$, with a total of at most m symbols, and every assignment of objects to the variables in t_1, \dots, t_n , there is a unique assignment of objects to all subterms of t_1, \dots, t_n obeying the obvious valuation clauses. It is easy to see that C is an M definable cut (not necessarily proper).

Let C' be the set of all $m \square C$ such that there is a satisfaction relation for all propositional combinations of formulas in $L(B)$ with at most $\text{lth}(\square B)$ quantifiers each, whose total number of symbols is $\square m$, based on the above definition of valuation. Again, $C' \square C$ is a cut.

Finally, let C'' be the set of all $m \square C'$ such that there is a satisfaction relation for all finite sequences of propositional combinations of formulas in $L(B)$ with at most $\text{lth}(\square B)$ quantifiers each, whose total number of symbols (in the whole finite sequence) is $\square m$, based on the above definition of satisfaction. Once again, $C'' \square C' \square C$ is a cut.

We cannot have $2^{[p]}(n) \square C''$ because the satisfaction relation provides a sequence of truth values for the universal closures of the formulas in this Hilbert style proof, \square . We can then do an induction on the resulting bit

sequence, taking the first entry which is false, and obtaining a contradiction.

In particular, C'' is a proper cut below $2^{[p]}(n)$. We can then do cut shortening to obtain an M definable cut below n .

The cut shortening procedure can be made to yield an M definable cut C^* below n in which PFA holds. We can also obtain an M definable cut $C^{**} \sqsubset C^*$ several exponentials lower, where C^{**} is closed under multiplication.

We are now in the same position that we were at the beginning of case 1. Here the full numerical part of M corresponds to the cut C^* here, and the cut C there corresponds to the cut C^{**} here. The identical argument will produce an M definable model of A. QED

The following result is used crucially in section 5.

THEOREM 3.6. Let A be a Σ^0_1 sentence. There exists an adequate sentence B in $L(Q)$ such that $EFA \vdash A \sqsubset RCON(B)$. B can be taken to be the conjunction of any sufficiently large finite fragment of PFA with some Σ^0_1 sentence in $L(Q)$.

Proof: Firstly, we may assume A is in $L(Q)$, since every Σ^0_1 sentence is provably equivalent, in EFA, to a Σ^0_1 sentence in $L(Q)$.

Choose a sufficiently large finite fragment PFA' of PFA. Then PFA' is adequate. By self reference, let C be a Σ^0_1 sentence in $L(Q)$ provably equivalent over PFA' to

C^* . there exists a restricted proof \square of $\square(PFA' \sqsubset C)$ such that A is true (with outermost universal quantifiers restricted to natural numbers) \sqsubset lth(\square).

The adequate sentence B that we want is the sentence $B = PFA' \sqsubset C^*$. Note that C^* is Σ^0_1 . To see that this works, we argue in EFA.

Suppose A. If $\square RCON(B)$ then C^* . Hence C. Since C is a true Σ^0_1 sentence, we can prove $PFA' \sqsubset C$ using formulas in $L(Q)$ which are propositional combinations of formulas each of which have at most a number of quantifiers that depends only on the choice of PFA' and C. Also, we have $RCON(PFA')$. Hence we conclude that $RCON(PFA' \sqsubset C)$, because we can

perform cut elimination several times (EFA is our background theory). I.e., $\text{RCON}(B)$.

Now suppose $\text{RCON}(B)$. Let n be a counterexample to A . Note that PFA' sees that there is no restricted proof \square of $\square B$ with $\text{lth}(\square) \leq n$, by $\text{RCON}(B)$. Hence $\text{PFA}' \vdash \square C^*$ using formulas in $L(Q)$ which are propositional combinations of formulas each of which have at most a number of quantifiers that depends only on the choice of PFA' and C . Hence $\text{PFA}' \vdash \square C$ using formulas in $L(Q)$ which are propositional combinations of formulas each of which have at most a number of quantifiers that depends only on the choice of PFA' and C . The same must be true of a proof of $\square B = \text{PFA}' \vdash C$. By successive cut elimination, we obtain $\square \text{RCON}(B)$. Hence the counterexample n to A must not exist. Therefore A . QED

Here is a sharpening of Theorem 3.6.

THEOREM 3.7. Let A be a \square^0_1 sentence and B be adequate, with an adequacy mechanism. There exists a \square^0_1 sentence C such that $\text{EFA} \vdash (A \square \text{RCON}(B)) \square \text{RCON}(B \square C)$.

Proof: We again assume that A is in $L(Q)$, and choose a sufficiently large finite fragment PFA' of PFA . Then PFA' is adequate. By self reference, let C be a \square^0_1 sentence in $L(Q)$ provably equivalent over PFA' to

C^* . there exists a restricted proof \square of $\square (B \square C)$ such that A is true (with outermost universal quantifiers restricted to natural numbers) $\square \text{lth}(\square)$.

We argue in EFA. Suppose $A \square \text{RCON}(B)$. If $\square \text{RCON}(B \square C)$ then C^* . Therefore C . Hence B proves C with a restricted proof. Therefore $\text{RCON}(B \square C)$, which is a contradiction. (For these last two statements, we use cut elimination for a standard number of steps in the background theory EFA). Hence $\text{RCON}(B \square C)$.

Now suppose $\text{RCON}(B \square C)$. Let n be a counterexample to A . Note that B sees that there is no restricted proof \square of $\square (B \square C)$ with $\text{lth}(\square) \leq n$, by $\text{RCON}(B \square C)$. Hence $B \vdash \square C^*$ with restricted proof. Therefore $B \vdash \square C$ with restricted proof. (We have again used cut elimination for a standard number of steps in the background theory EFA). This contradicts $\text{RCON}(B \square C)$. QED

COROLLARY 3.8. Let A, B be consistent sentences, where B is adequate. There exists consistent adequate $B \sqsupseteq_I A$. In fact, under any adequacy mechanism for B , C can be taken to be \sqsupseteq_1^0 .

Proof: Let A, B be as given. By Theorem 3.7, let C be such that $EFA \vdash (RCON(A) \sqsupseteq RCON(B)) \sqsupseteq RCON(B \sqsupseteq C)$. In particular, $EFA \vdash RCON(B \sqsupseteq C) \sqsupseteq RCON(A)$, and $B \sqsupseteq C$ is adequate. By Theorem 3.5, $A \sqsupseteq_I B \sqsupseteq C$. QED

An interpretation \sqsupseteq from S to T is said to be faithful iff it does not use any parameters, and for all sentences $\phi \in L(S)$, $S \vdash \phi \iff T \vdash \phi$.

THEOREM 3.9. (Interpretability = faithful interpretability). Let A, B be sentences, where B is adequate without parameters. The following are equivalent.

- i. A is interpretable in B .
- ii. A is faithfully interpretable in B .

For a proof of Theorem 3.9, see [Vi05].

Theorem 3.9 is easier if we assume that there is an adequacy mechanism for B , without parameters, which satisfies all true \sqsupseteq_1^0 sentences. For this result, we need a Lemma.

LEMMA 3.10. Let B be consistent, with an adequacy mechanism. There is a formula $\phi(n)$ in $L(B)$ such that for all $n \geq 0$, $B \sqsupseteq \phi(0) \sqsupseteq \dots \sqsupseteq \phi(n-1) \sqsupseteq \phi(n)$ is conservative over B for all \sqsupseteq_1^0 sentences. In fact, we can replace \sqsupseteq_1^0 by \sqsupseteq_k^0 , for any k given in advance of ϕ .

Proof: By a well known Rosser construction. QED

Here is the weakened form of Theorem 3.9 that was promised.

THEOREM 3.11. (Interpretability = faithful interpretability). Let A, B be sentences, where B is given an adequate mechanism without parameters. Assume that B is consistent with the true \sqsupseteq_1^0 sentences. The following are equivalent.

- i. A is interpretable in B .
- ii. A is faithfully interpretable in B .

Proof: Let σ be an interpretation from A into B , where B is adequate without parameters. Let $\sigma(n)$ be as given by Lemma 3.10. We now define a faithful interpretation of A to B .

Let $M \models B$. We define a model of A within M .

Let n be least such that $\sigma(n)$. If this does not exist, use the default value $n = 0$.

Let E be the n -th sentence in $L(A)$. If M satisfies the consistency of $A \cup E$ then return the model of $A \cup E$ obtained using the formalized completeness proof. Otherwise, return the model of A given by σ .

We claim that for any $E \in L(A)$ such that $A \cup E$ is consistent, some $M \models B$ returns a model of $A \cup E$.

Let E be as given, and let n be its Gödel number. Let $B^* = B \cup \sigma(0) \cup \dots \cup \sigma(n-1) \cup \sigma(n)$. Then B^* remains consistent with the true Σ_1^0 sentences, according to Lemma 3.10.

Let $M \models B^*$. Then the calculation in M of n , and therefore E , is correct. Hence $M \models \text{Con}(A \cup E)$ since $\text{Con}(A \cup E)$ is a true Σ_1^0 sentence.

Now under this interpretation σ' of A into B , which sentences ϕ in $L(A)$ are satisfied in all models arising in the interpretation? I.e., in all $\sigma'M$, $M \models B$?

Such a sentence ϕ must be consistent with every consistent sentence $A \cup E$. So if A does not prove ϕ , then $A + \neg\phi$ is consistent, and ϕ must be consistent with $A + \neg\phi$. This is patently absurd, and so ϕ must be provable from A , as required. QED

4. Predicative Set Extensions and Relative Consistency.

The predicative set extension of a theory T is defined as follows. First we introduce a second sort for "sets", with the binary relation \in between objects of the original sort and objects of the new sort. Then we add the predicative comprehension axiom

$$(\exists A) (\forall x) (x \in A \leftrightarrow \phi)$$

where φ is a formula in the extended language without any quantifiers over the second sort, in which A is not free.

Since we are officially working in single sorted logic, convert this two sorted system to a single sorted system in the obvious way by introducing a unary predicate.

We write this extension as $PSE(T)$. It is well known, by a simple model theoretic argument, that $PSE(T)$ is a conservative extension of T .

THEOREM 4.1. Let A be an adequate sentence. $PSE(A)$ is adequate, and finitely axiomatized. EFA proves $RCON(PSE(A)) \equiv CON(A)$. EFA does not prove $RCON(A) \equiv CON(A)$. EFA does not prove $CON(A) \equiv CON(PSE(A))$. $PSE(A)$ is not interpretable in A .

Proof: Well known. See Theorems 5.5 and 7.1 of [Vi06]. QED

THEOREM 4.2. Let A, B be adequate sentences. The following are equivalent.

- i. $PSE(A)$ is interpretable in $PSE(B)$.
- ii. EFA proves $CON(B) \equiv CON(A)$.

Proof: Apply Theorem 3.5 to the adequate theories $PSE(A)$, $PSE(B)$. We obtain that i \equiv (EFA proves $RCON(PSE(B)) \equiv RCON(PSE(A))$) \equiv ii, using Theorem 4.1. QED

5. Adequate Tarski Degrees and P Degrees.

The adequate Tarski degrees are the Tarski degrees of adequate sentences. There is another kind of degree which is intimately connected with Tarski degrees.

It is most convenient to work in the language of EFA, $0, S, +, \cdot, \exp, <$. For Σ_1^0 sentences A, B , we write

$A \equiv_P B$ \iff B logically implies A .

\equiv_P is reflexive and transitive. Define $=_P, \geq_P$ in the obvious way. The P degrees are the equivalence classes under $=_P$.

EFA represents a particular P degree.

Recall Theorem 3.6.

THEOREM 3.6. Let A be a Σ_1^0 sentence. There exists an adequate sentence B in $L(Q)$ such that $EFA \vdash A \leftrightarrow RCON(B)$. B can be taken to be the conjunction of any sufficiently large finite fragment of PFA with some Σ_1^0 sentence in $L(Q)$.

THEOREM 5.1. The adequate Tarski degrees are isomorphic to the P degrees $\geq EFA$.

Proof: Let A be adequate. Map it to $A^* = EFA \leftrightarrow RCON(A)$. By Theorem 3.5, for adequate A, B , we have

$A \leq B \leftrightarrow EFA \vdash RCON(B) \leftrightarrow RCON(A)$.
 $A \leq B \leftrightarrow EFA \leftrightarrow RCON(B)$ logically implies $EFA \leftrightarrow RCON(A)$.
 $A \leq B \leftrightarrow EFA \leftrightarrow RCON(A) \leftrightarrow_P EFA \leftrightarrow RCON(B)$.
 $A \leq B \leftrightarrow A^* \leftrightarrow_P B^*$.

This defines an embedding from the adequate Tarski degrees into the P degrees $\geq_P EFA$.

To see that the embedding is onto, let $EFA \leftrightarrow A$ be given, where A is a Σ_1^0 sentence. By Theorem 3.6, Let B be an adequate sentence such that $EFA \vdash A \leftrightarrow RCON(B)$. Then

$B^* = EFA \leftrightarrow RCON(B) \equiv_P EFA \leftrightarrow A$.

QED

THEOREM 5.2. The adequate Tarski degrees, and the P degrees $\geq EFA$, are isomorphic distributive lattices with $0, 1$. The P degrees are a distributive lattice with $0, 1$.

Proof: By Theorem 5.1, we need only deal with the P degrees. The P degrees $\geq EFA$ have EFA as 0 and $1 = 0$ as 1 . Obviously \wedge and \vee are the inf and sup. The P degrees are also a distributive lattice $0 = 0$ as the 0 , and $1 = 0$ as the 1 . QED

THEOREM 5.3. In the adequate Tarski degrees, every degree < 1 is an inf. 1 is not an inf in the adequate Tarski degrees (obvious).

Proof: It suffices to work within the P degrees $\geq EFA$, where it is basically well known. Let A be a consistent Σ_1^0 sentence that logically implies EFA . By self reference, let φ be provably equivalent, over EFA , to

*) every proof of $A \leftrightarrow \varphi$ has a smaller proof of $A \leftrightarrow \varphi\varphi$.

Let σ be

every proof of $A \wedge \sigma$ has a smaller proof of $A \wedge \sigma$.

Note that σ is the usual Rosser sentence construction over A , and $\bar{\sigma}$ is its dual.

We claim that $EFA \vdash \sigma \wedge \bar{\sigma}$, so that A is the inf of $A \wedge \sigma$ and $A \wedge \bar{\sigma}$. To see this, argue in EFA . Suppose $\sigma \wedge \bar{\sigma}$. Then $\sigma^* \wedge \bar{\sigma}$. Then we get two distinct Gödel numbers, neither of which is $>$ the other. This is impossible.

It remains to refute $A \vdash \sigma$ and $A \vdash \bar{\sigma}$.

Assume $A \vdash \sigma$. Let n be the least proof. If there is a smaller proof of $A \wedge \sigma$ then $EFA \vdash \bar{\sigma}A$. This is impossible. Hence there is no smaller proof of $A \wedge \sigma$. So σ is true. Since σ is σ_1^0 , we have $EFA \vdash \bar{\sigma}A$. This is a contradiction.

Assume $A \vdash \bar{\sigma}$. Let n be the least proof. By the previous paragraph, there is no proof of $A \wedge \sigma$. Hence $EFA \vdash \sigma$, and so $EFA \vdash \bar{\sigma}A$. This is a contradiction.

Note that the previous two paragraphs were conducted within $EFA + \text{Con}(A)$. Thus we obtain

$$EFA + \text{Con}(A) \vdash \text{Con}(A + \sigma) \wedge \text{Con}(A + \bar{\sigma}).$$

Note that $EFA + \sigma \wedge \bar{\sigma} \vdash \text{Con}(A)$.

Assume $A \vdash \sigma$. Then $A + \bar{\sigma} \vdash \text{Con}(A)$. Hence $A + \bar{\sigma} \vdash \text{Con}(A + \bar{\sigma})$. By Gödel's second incompleteness theorem, $A + \bar{\sigma}$ is inconsistent. I.e., $A \vdash \bar{\sigma}$. This is impossible. QED

THEOREM 5.4. In the adequate Tarski degrees, every degree other than $0, 1$, is a sup. $0, 1$ are not sups in the adequate Tarski degrees (obvious for 0).

Proof: (repaired by Albert Visser, private communication). We work in the P degrees $\geq EFA$. Let A be a consistent σ_1^0 sentence that proves EFA , but is not provable in EFA . We show that A is a sup in the P degrees $\geq EFA$.

By self reference, let σ be provably equivalent, over EFA , to

\square^* . Every proof of $EFA \vdash \square \square \square \vdash A$ has a smaller proof of $EFA \vdash \square A$ or of $EFA \vdash \square \square \square \vdash A$.

Let \square be

every proof of $EFA \vdash \square A$ or of $EFA \vdash \square \square \square \vdash A$ has a smaller proof of $EFA \vdash \square \square \square \vdash A$.

We claim that A is the sup of $EFA \vdash (A \rightarrow \square)$ and $EFA \vdash (A \rightarrow \square \square \square)$. Obviously A logically implies these two sentences.

Note that $EFA \vdash \square \square \square \square \vdash \text{Con}(A) \vdash A$. To see this, argue in EFA , and assume \square, \square . There cannot be a proof of $EFA \vdash \square \square \square \square \vdash A$, and so $\text{Con}(A)$.

Assume $EFA \vdash \square \square \square \vdash A$. There is no proof of $EFA \vdash \square \square A$ and no proof of $EFA \vdash \square \square \square \vdash A$. Hence $EFA \vdash \square \square$, $EFA \vdash A$. This is a contradiction.

Assume $EFA \vdash \square \square \vdash A$. Let n be the least proof. If there is a smaller proof of $EFA \vdash \square \square \square \square \vdash A$ then, since there is no proof of $EFA \vdash \square \square A$, we have $EFA \vdash \square \square$, and hence $EFA \vdash A$. This is impossible. If there is no smaller proof of $EFA \vdash \square \square \square \square \vdash A$ then $EFA \vdash \square$. Hence $EFA \vdash A$, which is again a contradiction.

Assume $EFA \vdash \square \square \vdash A$. Since $EFA \vdash \square \square \square \square$, we have $EFA \vdash \square \square \square \square \vdash A$. But this is impossible. QED

Theorems 5.3 and 5.4 may suggest that in the adequate Tarski degrees, or alternatively, in the P degrees $\geq EFA$, the degrees $\neq 0,1$ look alike. This is not the case.

THEOREM 5.5. In the P degrees $\geq EFA$, $(\square d < 1) (\square e > 0) (\inf(d,e) > 0)$. The d 's with this property are closed upwards. This property holds of $EFA \vdash \text{RCON}(EFA)$, but not of $EFA \vdash \square$, where \square is the Rosser sentence over EFA (see the proof of Theorem 5.3.).

Proof: Obviously the property is closed upwards. Let B be a $\square 01$ sentence. Assume $EFA \vdash \text{RCON}(EFA) \vdash B$. Then

$EFA + \square B \vdash \text{RCON}(EFA)$.

$EFA + \square B \vdash \text{RCON}(EFA + \square B)$.

$EFA + \square B$ is inconsistent (Gödel's second incompleteness theorem for restricted consistency).

EFA \vdash B.

By the proof of Theorem 5.3, $\inf(\delta, \epsilon) = 0$, and $\delta > 0$. QED

COROLLARY 5.6. In the adequate Tarski degrees, $(\delta < 1) (\epsilon > 0) (\inf(d, e) > 0)$. The d 's with this property are closed upwards. This property holds of EFA, but not of all adequate Tarski degrees.

Proof: Immediate from Theorem 5.5, using the isomorphism from the adequate Tarski degrees onto the P degrees \geq EFA, defined in the proof of Theorem 5.1. The isomorphism maps EFA to EFA \square RCON(A). QED

6. Infinite Theories.

We scratch the surface concerning interpretability between infinite theories. The nature of this subject is quite different than the finite case, as is clear from Theorems 6.2 - 6.4. Here we will assume that we are dealing with parameterless interpretations only.

We begin with a well known sufficient condition for interpretability of infinite theories.

THEOREM 6.1. Every r.e. set of sentences S is interpretable in $Q + \{\text{Con}(S') : S' \text{ is a finite subset of } S\}$. As a consequence, if S, T are theories, S is r.e., T is adequate with an adequacy mechanism, and T proves the consistency of each finite fragment of S , then S is interpretable in T .

Proof: Let $M \models Q + \{\text{Con}(S') : S' \text{ is a finite subset of } S\}$. Pass to an M definable cut C in M satisfying PFA. Then $C \models \{\text{Con}(S') : S' \text{ is a finite subset of } S\}$.

If $C \models \text{Con}(S)$ then we can use formalized completeness to return a model of S . Suppose $C \models \neg \text{Con}(S)$. In C , let n be largest such that in C , the first n axioms of S are consistent. Since in C , for each standard k , the first k axioms of S are consistent, this n must be a nonstandard integer. Again use formalized completeness to return a model M^* which, according to C , is a model of the first n axioms of S . Hence M^* really is a model of S . The stated consequence is immediate. QED

Theorem 6.1 applies to such pairs of theories as ZFC and ZFC + \square CH.

THEOREM 6.2. (No certificates). There is a Σ_2^0 sentence A such that

- i. $PA + A$ is interpretable in PA .
- ii. ZFC doesn't prove that $PA + A$ is interpretable in PA , provided ZFC is consistent.
- iii. Moreover, " $PA + A$ interpretable in PA " is provably equivalent to $\text{Con}(\text{ZFC})$ over PA .
- iv. ZFC does not prove $(R)\text{Con}(PA) \leftrightarrow (R)\text{Con}(PA+A)$, provided ZFC is consistent.

Proof: Let $A =$ "the greatest n such that I_{\Box_n} is consistent is less than any inconsistency in ZFC".

We claim that PA proves the consistency of every $I_{\Box_n} + A$. To see this, we argue in PA , and fix n . We know that $I_{\Box_n} + \text{Con}(I_{\Box_n}) + \Box\text{Con}(I_{\Box_{n+1}})$ is consistent. By exhaustive search, we see that n is less than the Gödel number of any inconsistency in ZFC. Hence $I_{\Box_n} + A$ is consistent.

By Theorem 6.1, $PA + A$ is interpretable in PA . Note that this argument is verifiable in $PA + \text{Con}(ZF)$.

Now suppose $PA + A$ is interpretable in PA . Suppose ZFC is inconsistent, and let n be the least Gödel number of an inconsistency in ZFC. Since $PA \vdash \text{Con}(I_{\Box_{n+1}})$, we see that $PA + A$ is inconsistent, and therefore not interpretable in PA . Hence ZFC is consistent. This establishes ii,iii.

In particular, we have shown that if $PA + A$ is consistent then ZFC is consistent. Hence iv. QED

For additional results along these lines, see [Sh97].

THEOREM 6.3. Let A be a Σ_1^0 sentence. $PA + A$ is interpretable in PA if and only if PA proves A .

THEOREM 6.4. There is a Σ_1^0 sentence A such that

- i. $PA + A$ is not interpretable in PA .
- ii. $EFA \vdash \text{Con}(PA) \leftrightarrow \text{Con}(PA + A) \leftrightarrow \text{RCON}(PA)$.

Proof: Set A to be a usual Rosser sentence over PA . Since PA does not prove A , we have i by Theorem 6.3. Also ii is clear as in the proof of Theorem 5.2. QED

On the other hand, interpretations from r.e. theories into single sentences behave like interpretations from single

sentences into single sentences, as the following result indicates.

THEOREM 6.5. Let T be r.e. and B be adequate. If $EFA \vdash RCON(B) \square RCON(T)$ then T is interpretable in B . T is interpretable in B if and only if T is faithfully interpretable in B .

Let T be r.e. and adequate, with finite $L(T)$. In $PSE(T)$ we can construct a truth definition for $L(T)$ in the standard way. We can replace T by "all axioms of T are true" in $PSE(T)$, resulting in the finitely axiomatized system $PSE^*(T)$, which is a conservative extension of T .

Note that $PSE^*(T)$ depends on the r.e. presentation of T . Depending on some details about how to formulate $PSE^*(T)$, one can prove the following result.

THEOREM 6.6. Let S, T be adequate r.e. theories with adequacy mechanisms. If EFA proves $CON(T) \square CON(S)$ then $PSE^*(S)$ is interpretable in $PSE^*(T)$.

[Sh97] tells us, for example, that the set of \square_1^0 sentences \square such that $ZF + \square$ is interpretable in NBG is complete \square_3^0 .

Section 7. Observed Linearity.

We restrict attention to the adequate theories from the literature which constitute a coherent system of axioms for mathematical reasoning. Here we don't require philosophically or foundationally coherent, as this would arguably be quite restrictive. We mean coherent only in that there is a modicum amount of naturalness. But we do require that there be no metamathematical ingredients present beyond a very standard base theory.

In particular, we include any theory of the form

$EFA + \{A_1, \dots, A_k\}$
 $PA + \{B_1, \dots, B_n\}$
 $RCA_0 + \{C_1, \dots, C_m\}$
 $Z + \{D_1, \dots, D_r\}$
 $ZF + \{E_1, \dots, E_s\}$

and much more, where A_1, \dots, A_k are robust formulations of published mathematical theorems in the language of EFA , B_1, \dots, B_n are such in the language of PA , C_1, \dots, C_m are such

in the language of RCA_0 , and $D_1, \dots, D_r, E_1, \dots, E_s$ are such in the language of set theory.

The striking observation is that one finds a remarkable linearity. This linearity is found not only with finitely axiomatized systems - which we have emphasized here - but with the non finitely axiomatized systems such as PA and ZFC.

This is perhaps the most intriguing, thought provoking, fundamental, and deep phenomenon in the whole of the foundations of mathematics.

In practice, the non finitely axiomatized systems encountered are almost finitely axiomatized, in the sense that the axiomatizations consist of finitely many axioms and finitely many axiom schemes. E.g., PA and ZFC meet this criteria.

Of course, for any adequate system T based on finitely many axioms and axiom schemes, there is a canonical version of $PSE(T)$, which is finitely axiomatizable. Under interpretability, $PSE(T)$ is a little bit higher than T . However, in the other orderings that we consider, with one exception, $\Box S$, $PSE(T)$ and T are equivalent. Examples: PA and ACA_0 . ZF and NBG.

The formal systems from the literature range from very weak axioms of arithmetic, to the much stronger axioms of infinite set theory such as ZFC - and beyond, with the extensions of ZFC by the so called large cardinal axioms.

We consider the following relations between theories. For each of these notions, we allow systems with finitely many axioms and finitely many axiom schemes.

$\Box I$. $T \Box I T' \Box T$ is interpretable in T' .

$\Box S$. Logical strength. $T \Box S T' \Box EFA \vdash Con(T') \Box Con(T)$.

$\Box S^*$. Liberal logical strength. $T \Box S^* T' \Box PA \vdash Con(T') \Box Con(T)$.

$\Box R$. Provably recursive functions. We say that $f: N \Box N$ is a provably recursive function of T if and only if there is a Turing machine code e computing f such that $T \vdash "e$ computes a function from N into $N"$. We define $T \Box R T' \Box$

every provably recursive function of T is a provably recursive function of T' .

\square_1 . $T \square_1 T'$ \square every \square_1^0 sentence provable in T is provable in T' .

\square_2 . $T \square_2 T'$ \square every \square_2^0 sentence provable in T is provable in T' .

\square . $T \square T'$ \square every arithmetic sentence provable in T is provable in T' .

\square_0 . Provable ordinals. For all systems below discussing subsets of ω (or general sets), we say that α is a provable ordinal of T if and only if there is a Turing machine code e computing a well ordering of \mathbb{N} of order type α , such that $T \vdash$ "e computes a well ordering of \mathbb{N} ".

Here are some some observed linearity phenomena.

Note that comparability under \square_0 is entirely automatic. This kicks in when we are extending RCA_0 .

A. Any two natural systems interpreting Q are comparable under \square_1 .

B. Any two natural systems extending EFA are comparable under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2$. They agree under $\square_{S^*}, \square_R, \square_1, \square_2$. If they are finitely axiomatized, then they agree under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2$.

C. Any two natural systems extending PA are comparable under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2, \square$. They agree under $\square_{S^*}, \square_R, \square_1, \square_2, \square$. If they are finitely axiomatized, then they agree under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2, \square$.

D. Any two natural systems extending RCA_0 are comparable under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2, \square, \square_0$. They agree under $\square_{S^*}, \square_R, \square_1, \square_2, \square, \square_0$. If they are finitely axiomatized, then they agree under $\square_1, \square_S, \square_{S^*}, \square_R, \square_1, \square_2, \square, \square_0$.

E. If T and T' are natural systems extending PA , and $T < T'$ under any of $<_{S^*}, <_R, <_1, <_2, <, <_0$, then T' proves the consistency of T .

Here is a table that lists the finite number of levels that have figured prominently in f.o.m. by means of preferred

representatives. We have linearity under ΠI (with axiom schemes allowed here). We also exemplify B - E above.

The relevant interpretations here can all be taken to be parameterless.

PFA.
 EFA.
 SEFA.
 PRA.
 RCA₀.
 I Π_2 .
 I Π_3 .
 PA.
 ACA₀.
 ACA₀ + $(\exists n, x) (TJ(n, x) \Pi)$.
 ACA.
 RCA₀ + TJ(Π) Π .
 ACA₀ + TJ(Π) Π .
 ACA + TJ(Π) Π .
 ACA₀ + $(\exists x) (TJ(\Pi, x) \Pi)$.
 ACA₀ + $\{(\exists \Pi, x) (TJ(\Pi, x) \Pi) : \Pi < \Pi^0\}$.
 ACA₀ + $\{(\exists \Pi < \Pi^0) (\exists x) (TJ(\Pi, x) \Pi)\}$.
 RCA₀ + TJ(Π^0) Π .
 ACA₀ + TJ(Π^0) Π .
 ACA₀ + $\{(\exists x) (TJ(\Pi^0, x) \Pi)\}$.
 ACA₀ + $\{(\exists x) (TJ(\Pi, x) \Pi) : \Pi < \Pi_0\}$.
 Π_1^1 -CA.
 RCA₀ + TJ(Π_0) .
 ACA₀ + TJ(Π_0) .
 ACA + TJ(Π_0) .
 ACA₀ + $(\exists x) (TJ(\Pi_0, x) \Pi)$.
 $\{ATI(\Pi) : \Pi < \Pi_0\}$.
 ATR₀.
 ATI($<\Pi_0$) .
 ATR.
 Π_2^1 -TI₀.
 Π_2^1 -TI.
 TI.
 ID₂.
 ID $<\Pi$.
 Π_1^1 -CA₀.
 Π_1^1 -CA.
 Π_1^1 -CA + TI.
 Π_1^1 -TR₀.
 Π_1^1 -TR.
 Π_2^1 -CA₀.

\aleph_2^1 -CA.
 \aleph_2^1 -CA + TI.
 Z_2 .
 Z_3 .
 Type Theory.
 Weak Zermelo.
 ZC.
 ZC + $(\aleph_1 < \aleph_2)$ (V(\aleph_1)).
 KP(\emptyset).
 ZFC.
 ZFC + strongly inaccessible \aleph_1 .
 ZFC + strongly Mahlo \aleph_1 .
 ZFC + {strongly n-Mahlo \aleph_1 : $n < \aleph_1$ }.
 ZFC + $(\aleph_1 < \aleph_2)$ (strongly n-Mahlo \aleph_1).
 ZFC + (weakly compact \aleph_1).
 ZFC + (indescribable \aleph_1).
 ZFC + (subtle \aleph_1).
 ZFC + (almost ineffable \aleph_1).
 ZFC + (ineffable \aleph_1).
 ZFC + {n-subtle \aleph_1 : $n < \aleph_1$ }.
 ZFC + $(\aleph_1 < \aleph_2)$ (n-subtle \aleph_1).
 ZFC + $\aleph_1 \aleph_2 \aleph_3$.
 ZFC + $(\aleph_1 < \aleph_2)$ ($\aleph_1 \aleph_2 \aleph_3$).
 ZFC + $0\# \aleph_1$.
 ZFC + $(\aleph_1 \times \aleph_2)$ ($\aleph_1 \# \aleph_2$).
 ZFC + $\aleph_1 \aleph_2 \aleph_3 \aleph_4$.
 ZFC + Ramsey \aleph_1 .
 ZFC + Measurable \aleph_1 .
 ZFC + Concentrating Measurable \aleph_1 .
 ZFC + Strong \aleph_1 .
 ZFC + Woodin \aleph_1 .
 ZFC + Superstrong \aleph_1 .
 ZFC + Supercompact \aleph_1 .
 ZFC + Extendible \aleph_1 .
 ZFC + Vopenka \aleph_1 .
 ZFC + Almost Huge \aleph_1 .
 ZFC + Huge \aleph_1 .
 ZFC + Superhuge \aleph_1 .
 ZFC + $(\aleph_1 < \aleph_2)$ (n-huge \aleph_1).
 ZFC + Rank into Itself \aleph_1 .
 ZFC + Rank + 1 into Itself \aleph_1 .
 VB + V into V \aleph_1 .

There are extremely natural theories for which linearity fails, in which Q is not interpretable. For example, let

T = axioms for discrete linear orderings without endpoints (every point has an immediate predecessor and an immediate successor).

T' = axioms for dense linear orderings without endpoints (between any two elements there is a third, and no endpoints).

Then $T \not\equiv T'$ and $T' \not\equiv T$ both fail. Also note that T, T' are finitely axiomatized and complete.

Interpretability is very interesting in this kind of tame world, below Q - where it takes on an entirely different character than the theory $\geq Q$ that we have focused on.

In particular, it would be interesting to understand interpretability among subsystems of RCF (real closed fields) and ACF (algebraically closed fields) and Presburger arithmetic (the semigroup of nonnegative integers), finitely axiomatized or otherwise.

REFERENCES

[Fe60] S. Feferman, Arithmetization of metamathematics in a general setting, *Fundamenta Mathematicae* 49 (1960), 35-92.

[Fr76] Translatability and relative consistency, November, 1976, 7 p.

[Fr80] Translatability and relative consistency II, Ohio State University, unpublished, September, 1980, 6 p.

[Ge05] P. Gerhardy, The Role of Quantifier Alternations in Cut Elimination, *Notre Dame Journal of Formal Logic*, vol. 46, no. 2, pp. 165-171 (2005).

[HP98] P. Hajek, P. Pudlak, *Metamathematics of First-Order Arithmetic*, *Perspectives in Mathematical Logic*, Springer, 1998, 460 p.

[MM94] A. Mancinci, F. Montagna, A minimal predicative set theory, *Notre Dame J. of Formal Logic* 35 (1994) 186-203.

[Ro52] R.M. Robinson, An essentially undecidable axiom system, *Proceedings of the 1950 International Congress of Mathematicians*, Cambridge MA, 1952, pp. 729-730.

[Sh97] V.Y. Shavrukov, *INterpreting Reflexive Theories in Finitely Many Axioms*, *Fundamenta Mathematicae*, vol. 152, p. 99-116, 1997.

[Sm85] C. Smorynski, *Nonstandard models and related developments*, in: *Harvey Friedman's Research on the Foundations of Mathematics*, North Holland: Amsterdam, 1985 pp. 179-229.

[Ta56] A. Tarski, *Logic, Semantics, Metamathematics: Papers from 1923 to 1938* (translated by J.H. Woodger), Clarendon Press, Oxford 1956, 471 p.

[TMR53] A. Tarski, A. mostowski, R.M. Robinson, *Undecidable Theories*, North Holland: Amsterdam 1953, 98 p.

[Vi90] A. Visser, *Interpretability logic*, in: *Mathematical logic*, (Proceedings of the Heyting 1988 summer school in Vjarna, Bulgaria, Plenum Press, Boston, 1990, p. 175-209.

[Vi05] A. Visser, *Faith and Falsity: a study of Faithful Interpretations and flase Π_1 -sentences*, *Annals of Pure and APplied Logic*, volume 131, p. 103-131, 2005.

[Vi06] A. Visser, *The Predicative Frege Hierarchy*, October, 2006, #246, <http://www.phil.uu.nl/preprints/lgps/list.html>.

[Zh91] W. Zhang, *Cut elimination and automatic proof procedures*, *Theoretical Computer Science* 91 (1991), 265-284.

[Zh94] W. Zhang, *Depth of proofs, depth of cut-formulas and complexity of cut formulas*, *Theoretical Computer Science* 129 (1994), 193-206.

* This research was partially supported by NSF DMS 0245349.