# PARTIALLY CONSERVATIVE EXTENSIONS OF ARITHMETIC

BY

## D. GUASPARI

ABSTRACT. Let $T$ be a consistent r.e. extension of Peano arithmetic; $\Sigma_n^0$, $\Pi_n^0$ the usual quantifier-block classification of formulas of the language of arithmetic (bounded quantifiers counting "for free"); and $\Gamma$, $\Gamma'$ variables through the set of all classes $\Sigma_n^0$, $\Pi_n^0$. The principal concern of this paper is the question: When can we find an independent sentence $\phi \in \Gamma$ which is $\Gamma'$-conservative in the following sense: Any sentence $\chi$ in $\Gamma'$ which is provable from $T + \phi$ is already provable from $T$? (Additional embellishments: Ensure that $\phi$ is not provably equivalent to a sentence in any class "simpler" than $\Gamma$; that $\phi$ is not conservative for classes "more complicated" than $\Gamma'$.) The answer, roughly, is that one can find such a $\phi$, embellishments and all, unless $\Gamma$ and $\Gamma'$ are so related that such a $\phi$ *obviously* cannot exist. This theorem has applications to the theory of interpretations, since "$\phi$ is $\Gamma$-conservative" is closely related to the property "$T + \phi$ is interpretable in $T$"–or to variants of it, depending on $\Gamma$. Finally, we provide simple model theoretic characterizations of $\Gamma$-conservativeness. Most results extend straightforwardly if extra symbols are added to the language of arithmetic, and most have analogs in the Levy hierarchy of set theoretic formulas ($T$ then being an extension of $ZF$).

**Introduction.** The main result of this paper is a theorem which is "obviously true" and of which the following is an instance: There is a $\Sigma_9^0$ sentence $\phi$ in the language of arithmetic such that: $\phi$ is essentially $\Sigma_9^0$ (i.e., not provably equivalent in arithmetic to a $\Pi_9^0$ sentence); any $\Sigma_5^0$ sentence provable in arithmetic $+\phi$ is already provable in arithmetic; some $\Pi_5^0$ sentence provable in arithmetic $+\phi$ is *not* provable in arithmetic. (According to Matijasevic's theorem it does not, up to an equivalence provable in arithmetic, matter whether or not we allow bounded quantifiers "for free" in the definition of $\Sigma_n^0$.) The same is true for many extensions of arithmetic, and an analogous result holds for the Levy hierarchy of formulas over set theory.

Consider, however, some obstacles to its proof: Let $T$ be an r.e. extension of Peano arithmetic. If $\Gamma$ is a set of formulas say that $\phi$ is $\Gamma$-con (over $T$) if every sentence from $\Gamma$ provable in $T + \phi$ is provable in $T$; that $\phi$ is $\Gamma$-non, otherwise. For the present $\Gamma$ will always be $\Pi_n^0$ or $\Sigma_n^0$. $\neg \mathrm{Con}(T)$ is a $\Sigma_1^0$ sentence which is $\Pi_1^0$-con over $T$, but might also be $T$ provable. A Rosser sentence for $T$ is guaranteed to be independent but will always be $\Pi_1^0$-non–as

---

will every independent sentence obtained by an "effectively inseparable sets" construction. So to produce an independent $\Sigma_1^0$ sentence which is $\Pi_1^0$-con requires, in some sense, another general construction for generating independent sentences.

Sentences which are $\Pi_1^0$-non make new *truths* provable. For example, one can produce a $\Sigma_1^0$ sentence $\sigma$ such that arithmetic $+ \sigma$ is consistent and proves $\Pi_1^0$ sentences (*truths*) not provable in *ZFC*, while *ZFC* $+ \sigma$ is inconsistent.

The existence of $\Sigma_1^0$ sentences which are $\Pi_1^0$-con can be equivalently stated: If $T$ is a consistent r.e. extension of arithmetic, there are $\Pi_1^0$ truths unprovable in any consistent extension of $T$ obtained by adding only $\Sigma_1^0$ sentences. This is a limitation on provability in false extensions of (false) theories.

After this paper was essentially completed Smorynski called attention to the papers [**Hajek, 1**], [**Hajek, 2**], [**Hajkova-Hajek**] dealing with the closely related notion of interpretability: Say that $\phi$ is interpretable in $T$ if $T + \phi$ is interpretable in $T$. For a broad class of theories "interpretable" is equivalent to "$\Pi_1^0$-con", and for such theories our methods immediately establish: There is a $\Sigma_1^0$ sentence $\phi$ which is interpretable in $T$ (so, of course, relatively consistent with $T$) but whose relative consistency with $T$ cannot be proven in, say, *PA*. (The point here is that $\phi$ is so simple.)

The paper is organized as follows: §1 lays out various examples from nature (Con($T$), $\neg$ Con($T$), Rosser sentences) and shows that we will have to look elsewhere for inspiration in proving the main (obviously true) existence theorem. The heart of the paper is §2. It contains the basic existence theorems and some embellishments thereof. §3 compares the notions "$\Pi_n^0$-con" with variants of the notion "interpretable." Straightforward generalizations to larger languages and to an analogous theory for the Levy hierarchy of formulas of set theory are noted in §4. The syntactical part of the paper is concluded in §5 by mentioning some open questions, in particular, that of classifying $\{\phi | \phi$ is $\Gamma$-con$\}$ in the arithmetical hierarchy. (Solovay has solved a particularly interesting special case originally proposed by Hajek.) §6 contains model theoretic characterizations of $\Gamma$-con. E.g., $\phi$ is $\Sigma_1^0$-con over *PA* (peano arithmetic) iff every countable model of *PA* contains an initial segment which models $PA + \phi$. ($\Pi_1^0$-con can be characterized dually.) The proofs of these theorems are independent of the previous sections. All follow easily from theorems of Friedman which characterize the end extensions of countable models of arithmetic. For models of set theory one and only one new twist occurs. Fact: $\phi$ is $\Pi_1$-con ("$\Pi_1$" refers to the Levy hierarchy) over *ZF* iff every countable *non-ω-model* of *ZF* can be end extended to model $ZF + \phi$. An example shows that the qualification "non-ω-model" cannot be dropped. Open question: If you *do* drop it, what do you get?

Thanks are due Craig Smorynski and Bob Solovay for helpful correspon-

dences and for generously allowing theorems of theirs to be included in this paper. Thanks are also due Tom McLaughlin who helped weed out some proofs of $0 = 1$ from those of $0 = 0$.

**0. Notation and preliminaries.** Every theory considered will be formulated in a language satisfying the following conditions: There is a special sort, the number sort, with $u$, $v$, ..., $z$ among the number variables. In addition to both unbounded quantifiers over each sort there are bounded quantifiers–$\forall x \leqslant y$, $\exists w \leqslant z$, etc.–over the number sort. Among the nonlogical signs are the signs of arithmetic–namely $+$, $\cdot$, $<$, $0$, $'$–applicable only to terms and variables of the number sort.

$LA$ is the language made up from the signs of arithmetic, variables of the number sort, and quantification over the number sort. The theory $N$ (see [**Shoenfield**, p. 22]) is a finitely axiomatizable theory of recursive arithmetic. That is, $N$ decides (correctly) every sentence of $LA$ containing no unbounded quantifiers. That fact is formalizable in $P$: (notation explained later in this section) for any sentence $\sigma$ in $\Sigma_1^0$, $PA \vdash \sigma \to \mathrm{Thm}_N(\ulcorner \sigma \urcorner)$. $PA$, Peano arithmetic, is $N$ + induction for all formulas of $LA$.

Notice that $ZF$ and $GB$, though not ordinarily expressed in languages that meet the requirements above, can easily be reformulated in languages that do so.

CONVENTION. "Theory" means "consistent theory containing $PA$".

So, e.g., $S$ is a subtheory of $T$ means that $PA \subseteq S \subseteq T$. This convention is too strong. In the sequel we can almost always get by on just the assumption that our theories contain $N$. Most of the time that part of the theory outside $LA$ will be quite hazy and unimportant.

Define the $\Sigma_n^0$, $\Pi_m^0$ formulas of $LA$ as usual: $\Sigma_0^0 = \Pi_0^0 = \Delta_0^0 =$ the class of formulas containing no unbounded quantifiers; $\Sigma_{n+1}^0$ is the class of formulas $\vec{Q}\vec{x}\chi$ with $\chi$ in $\Pi_n^0$, $\vec{x}$ a sequence of number variables, and $\vec{Q}$ a possibly empty sequence of quantifiers not containing $\forall$. $\Pi_{n+1}^0$ is defined dually. Notice that in this way each class of formulas is literally a subset of any class of "ostensibly more complicated" formulas. Every $\Sigma_n^0$ formula is equivalent (in $N$) to a 'prenex' $\Sigma_n^0$ formula: one with prefix a strictly alternating sequence of $n$ unbounded quantifiers followed by a 'matrix' containing no unbounded quantifiers. For any theory $T$, $T$-$\Sigma_n^0$ is the class of formulas $T$ provably equivalent to $\Sigma_n^0$ formulas; etc.; and $T$-$\Delta_n^0$ is $T$-$\Sigma_n^0 \cap T$-$\Pi_n^0$. $\Gamma$, $\Gamma'$, ... will vary over the classes $\Sigma_n^0$, $\Pi_m^0$; and $T$-$\Gamma$, ... over $T$-$\Sigma_n^0$, etc. $\check{\Gamma}$ is the dual class to $\Gamma$: i.e., $\check{\Sigma}_n^0 = \Pi_n^0$ and $\check{\Pi}_n^0 = \Sigma_n^0$.

DEFINITION 0.1. If $\phi$, $\theta$ are formulas of form $\exists x\chi(x)$, $\exists x\psi(x)$ respectively, then

$$\phi \preccurlyeq \theta =_{\mathrm{df}} \exists x(\chi(x) \wedge \forall y < x \neg \psi(y)),$$

$$\phi \prec \theta =_{\mathrm{df}} \exists x(\chi(x) \wedge \forall y \leqslant x \neg \psi(y)).$$

Notice that if $\phi$ and $\theta$ are 'prenex' $\Sigma_n^0$ formulas then so are $\phi \preccurlyeq \theta$ and $\phi \prec \theta$. More generally, if $\phi$ is $\Sigma_n^0$ and $\chi$ is $T\text{-}\Delta_n^0$, then $\phi \preccurlyeq \theta$ is $T\text{-}\Sigma_n^0$. Abuses of this notation will include writing $\phi \prec (\neg \theta)$ when, e.g., $\phi$ is $\Sigma_1^0$ and $\theta$ is $\Pi_1^0$. Smorynski has pointed out an unfortunate defect in this notation–namely, that it suggests that $\phi \preccurlyeq \phi$ is true, which is not always the case.

Notice also that if $\phi$ and $\theta$ are 'prenex' $\Sigma_1^0$ and either one of them is true, then $\phi \preccurlyeq \theta$ is decidable.

*Formalizing syntax and semantics.* We will use $i, j, \ldots, n$ metamathematically for natural numbers, and $\mathbf{n}$ for the canonical term denoting $n$. If $k$ is the Gödel number of $\phi$ (according to some fixed standard numbering) then $\ulcorner \phi \urcorner$ is $\mathbf{k}$. (So the Gödel number (from now on, "g.n.") of $\phi$ is a natural number, while $\ulcorner \phi \urcorner$ is a term in $LA$.)

Procedures for obtaining partial truth definitions are well known. Fact 0.2 is stated primarily to establish notation.

*Fact* 0.2. For any $\Gamma$ there is a formula $\phi(x)$ in $\Gamma$ satisfying:

For any sentence $\chi \in \Gamma$, $PA \vdash \phi(\ulcorner x \urcorner) \leftrightarrow \chi$.

Choose one such $\phi$ and call it $\Gamma$-True.

If $\chi$ is $\exists x \phi(x)$ we will use "$y$ is a witness for $\chi$" to mean $\phi(y)$.

We want $\text{Proof}_T(y, x)$ to be a (simple) binumeration in $T$ expressing: $y$ codes a proof from $T$ of (the formula with g.n.) $x$. That may not be possible–e.g., if $T$ is r.e. but nonrecursive. However, if $T$ is r.e. there is a recursive $T'$ equivalent to $T$ and we can therefore take $\text{Proof}_T$ to be a $PA - \Delta_1^0$ binumeration (in $PA$) of: $y$ codes a proof from $T'$ of $x$. Let $\text{Thm}_T(x)$ be $\exists y \text{Proof}_T(y, x)$ and $\text{Con}(T)$ be $\neg \text{Thm}_T(\ulcorner \mathbf{0} = \mathbf{1} \urcorner)$.

All conventions for denoting formalized operations on syntax are unsatisfactory. Here is another. Suppose we want to formalize a function which, inputing $m$ and $n$–g.n.'s for $\phi$ and $\theta$ respectively–outputs a g.n. for $\phi \preccurlyeq \theta$. We add a defined function symbol, say $g$, which formalizes the function and agree to notate $gxy$ by: "$x \preccurlyeq y$". (The quotation marks are part of the notation.) So, e.g., $N \vdash$"$\ulcorner \phi \urcorner \preccurlyeq \ulcorner \theta \urcorner$" $= \ulcorner \phi \preccurlyeq \theta \urcorner$. Another example, of another function: $N \vdash$"$\ulcorner \theta \urcorner \to \ulcorner \chi \urcorner$" $= \ulcorner \theta \to \chi \urcorner$.

Finally, we will make extensive use of a well-known lemma of Gödel.

*Fact* 0.3 (Self-reference lemma). If $\phi(x) \in \Gamma$ has only $x$ free, a sentence $\chi \in \Gamma$ can be effectively found for which

$$PA \vdash \chi \leftrightarrow \phi(\ulcorner \chi \urcorner).$$

## 1. Basic definitions and some examples.

Definition 1.1. Let $T$ be a theory, $\phi$ a sentence, $X$ any set of formulas in $LA$. Then,

$\phi$ is $X$-con over $T$ iff every sentence in $X$ provable from $T + \phi$ is provable from $T$,

$\phi$ is $X$-non over $T$ iff $\phi$ is not $X$-con over $T$.

The next lemma has to go somewhere, so may as well go here. It characterizes $\Gamma$-con. Until referred to it is skippable.

LEMMA 1.2. *Let $T$ be an extension of $N$. Then*:

(1) *A 'prenex' $\Pi_n^0$ sentence $\phi$ is $\Sigma_n^0$-con over $T$ if and only if $\forall \sigma$ ($\sigma$ a 'prenex' $\Sigma_n^0$ sentence and $T \vdash \sigma \Rightarrow T \vdash \sigma \prec (\neg \phi)$).*

(2) *A 'prenex' $\Sigma_n^0$ sentence $\phi$ is $\Pi_n^0$-con over $T$ if and only if $\forall \sigma$ ($\sigma$ a 'prenex' $\Sigma_n^0$ sentence and $T \vdash \phi \to \phi \prec \sigma \Rightarrow T \vdash \sigma \to \phi$).*

PROOF. (1) Suppose first that $\phi$ is $\Sigma_n^0$-con over $T$ and that $T \vdash \sigma$. Then $T + \phi \vdash \sigma \prec (\neg \phi)$; and therefore $T \vdash \sigma \prec (\neg \phi)$. Suppose conversely that $T \vdash \sigma \Rightarrow T \vdash \sigma \prec (\neg \phi)$ for any 'prenex' $\Sigma_n^0$ sentence $\sigma$. We want to show $\phi$ $\Sigma_n^0$-con, so let $\sigma'$ be $\Sigma_1^0$ and $T + \phi \vdash \sigma'$. Then $T$ implies the $\Sigma_n^0$ sentence $\neg \phi \vee \sigma'$ and so, by assumption, $T \vdash (\neg \phi \vee \sigma') \prec (\neg \phi)$. But $((\neg \phi \vee \sigma') \prec \neg \phi) \to \sigma'$ and therefore $T \vdash \sigma'$.

(2) Suppose $\phi$ is a $\Sigma_n^0$ sentence. For reference, call the statement " for any 'prenex' sentence $\sigma$ in $\Sigma_n^0$, if $T \vdash \phi \to \phi \prec \sigma$ then $T \vdash \sigma \to \phi$" by the name of $(+)$. Suppose $(+)$, and that $\pi$ is 'prenex' $\Pi_n^0$ and $T + \phi \vdash \pi$. Then $(*)$ $T \vdash \phi \to \pi$; and so $T \vdash \neg \phi \to \phi \prec (\neg \pi)$. Applying $(+)$, $T \vdash \neg \pi \to \phi$; i.e. $(**)$ $T \vdash \neg \phi \to \pi$. Combining $(*)$ with $(**)$, $T \vdash \pi$. Now suppose that $(+)$ is false and produce a $\Sigma_n^0$ sentence $\sigma$ such that $T \vdash \phi \to \phi \prec \sigma$ but $T \nvdash \sigma \to \phi$. Consider the $\Pi_n^0$ sentence $\neg (\sigma \preccurlyeq \phi)$, which is a consequence of $\phi$ because $\phi \prec \sigma$ is a consequence of $\phi$. We will show that $\phi$ is $\Pi_n^0$-non by showing that $T \nvdash \neg (\sigma \preccurlyeq \phi)$: From $\neg (\sigma \preccurlyeq \phi)$ we can infer in $T$ that $\sigma \to \phi \prec \sigma$ and therefore $\sigma \to \phi$. But $\sigma$ was chosen so that $T \nvdash \sigma \to \phi$.  □

EXAMPLES OCCURRING IN NATURE. For the moment we are principally concerned with $\Sigma_1^0$ and $\Pi_1^0$.

DEFINITIONS 1.3. $T$ is $\Gamma$-*correct* iff whenever $\phi$ is a $\Gamma$-sentence and $T \vdash \phi$, then $\phi$ is true.

If $T$ is r.e., $\theta$ is a *Rosser sentence* for $T$ iff $PA \vdash \theta \leftrightarrow \text{Thm}_T(\ulcorner \neg \theta \urcorner) \prec \text{Thm}_T(\ulcorner \theta \urcorner)$. (The self-reference lemma guarantees the existence of Rosser sentences. A Rosser sentence is guaranteed to be independent of $T$.)

For the rest of this section $T$ is r.e.

EXAMPLE 1.4. $\neg \text{Con}(T)$ is $\Pi_1^0$-con over $T$.

This has been observed by several people, including the author. The first to do so was Kreisel, several years ago (see [**Kreisel**]). [**Macintyre-Simmons**] independently obtained an abstract version imposing only weak conditions on $T$ and $\text{Thm}_T$.

EXAMPLE 1.5. Any Rosser sentence for $T$ is $\Pi_1^0$-non over $T$.

PROOF. Let $\sigma$ be Rosser for $T$, and let $\phi$ be the sentence $\text{Thm}_T(\ulcorner \neg \sigma \urcorner) \prec \text{Thm}_T(\ulcorner \sigma \urcorner)$, and $\theta$ be $\text{Thm}_T(\ulcorner \sigma \urcorner)$. Then $T \vdash \phi \to \phi \prec \theta$, and so, to show that $\phi$–hence $\sigma$–is $\Pi_1^0$-non it will suffice (by Lemma 1.2) to

show that $T \nvdash \theta \to \phi$. But if $T \vdash \theta \to \phi$, then $T \vdash \mathrm{Thm}_T(\ulcorner \sigma \urcorner) \to \sigma$; and so by [Löb], $T \vdash \sigma$. That is impossible, because a Rosser sentence for $T$ cannot be $T$ provable.

EXAMPLE 1.6 (SMORYNSKI). $\mathrm{Con}(T)$ is $\Sigma_1^0$-con iff $T$ is $\Sigma_1^0$-correct.

PROOF. The difficult part is to show: $T$ $\Sigma_1^0$-incorrect implies $\mathrm{Con}(T)$ is $\Sigma_1^0$-non. (The easy part, besides being easy, is a consequence of Theorem 2.1; whose proof is also skipped.) Let $\psi$ be any $\Sigma_1^0$ sentence and $\phi$ be such that $T$ proves $\phi \leftrightarrow (\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner) \vee \psi) \preccurlyeq \mathrm{Thm}_T(\ulcorner \phi \urcorner)$. As usual,

(i) $T \nvdash \neg\phi$.

Here, for the first and last time, is a detailed elaboration of "as usual": Suppose $T \vdash \neg\phi$. Then there is a $k \in \omega$ such that $T$ proves: k witnesses $\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner)$. Since $T$ is consistent, $T \nvdash \phi$; and therefore (because $\mathrm{Proof}_T$ is a binumeration) for every $n \in \omega$, $T \vdash$ n does not witness $\mathrm{Thm}_T(\ulcorner \phi \urcorner)$. Reason now in $T$: $\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner) \preccurlyeq \mathrm{Thm}_T(\ulcorner \phi \urcorner)$; so $(\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner) \vee \psi) \preccurlyeq \mathrm{Thm}_T(\ulcorner \phi \urcorner)$; so $\phi$. That contradicts the consistency of $T$.

Notice that the proof of (i) depends only on the assumption that $T$ is consistent, hence:

(ii) $T + \mathrm{Con}(T) \vdash \neg\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner)$.

The next claim is the crucial one:

(iii) If $\psi$ is true, so is $\phi$.

Let $\sigma$ be the formula $(\mathrm{Thm}_T(\ulcorner \neg\phi \urcorner) \vee \psi) \preccurlyeq \mathrm{Thm}_T(\ulcorner \phi \urcorner)$. If $\psi$ is true then $\sigma$ is *decidable* (in $N$). What if $\sigma$ is false? Then it must be the case that $N \vdash \neg\sigma$, so $T \vdash \neg\sigma$, so $T \vdash \neg\phi$—a contradiction. So $\sigma$, hence $\phi$, must be true.

That reasoning can also be formalized in $T + \mathrm{Con}(T)$, for it used only (i) and the knowledge that the truth of either $\Sigma_1^0$ sentence $\theta$ or $\Delta$ guarantees the decidability of $\theta \preccurlyeq \Delta$; which fact can be formalized in *PA*. So:

(iv) $T + \mathrm{Con}(T) \vdash \psi \to \phi$. (In fact, $T + \mathrm{Con}(T) \vdash \psi \leftrightarrow \phi$.)

Similar but easier reasoning shows that:

(v) $(T \vdash \phi)$ iff $\psi$ is true.

Now we are ready to go. Suppose that $T$ is $\Sigma_1^0$-incorrect, and let $\psi$ be a false $\Sigma_1^0$ sentence proved by $T$. Let $\phi$ be as above. By (v), $T \nvdash \phi$. But by (iv), $T + \mathrm{Con}(T) \vdash \psi \to \phi$ and therefore $T + \mathrm{Con}(T) \vdash \phi$. So $\mathrm{Con}(T)$ is $\Sigma_1^0$-non.

EXAMPLE 1.7 (EXPLOITING THE WEAKNESS OF *PA*). There is a $\Sigma_1^0$ sentence $\theta$ such that $PA + \theta$ is consistent and proves $\Pi_1^0$ sentences (truths) unprovable in *ZFC*; while $ZFC + \theta$ is inconsistent.

PROOF. Let $\sigma$ be Rosser over *ZFC*. Then $\sigma$ is independent of *ZFC* and a slight extension of Lemma 1.2 (proved by the same arguments) shows that $PA + \sigma$ proves $\Pi_1^0$ sentences not provable in *ZFC*. Let $\theta$ be such that *PA* proves: $\theta \leftrightarrow (\sigma \wedge \mathrm{Thm}_{PA}(\ulcorner \neg\theta \urcorner)) \preccurlyeq \mathrm{Thm}_{PA}(\ulcorner \theta \urcorner)$. Then $\theta$ is stronger than $\sigma$ (so $PA + \theta$ proves $\Pi_1^0$ sentences not provable in *ZFC*) and, since $\theta \to \neg\mathrm{Con}(PA)$, $\theta$ is disprovable in *ZFC*. (That this example is slightly phoney

can be seen by asking: Just how is it we know that $PA + \sigma$ is consistent? The $\Pi_1^0$ theory of $PA + \theta$ is contained in the $\Pi_1^0$ theory of $ZFC + \mathrm{Con}(ZFC)$.)

**2. Existence theorems.** Here is what the examples leave open: If $T$ is r.e. is there an *independent* $\Gamma$ sentence which is $\check{\Gamma}$-con over $T$? ($\neg \mathrm{Con}(T)$ need not be independent.) If $T$ is $\Sigma_1^0$-correct is there any independent $\Pi_1^0$ sentence which is $\Sigma_1^0$-non over $T$? The answers are, respectively, Yes and No. First, the No.

THEOREM 2.1. *If $T$ is r.e., the following are equivalent*:
(1) *$T$ is $\Sigma_1^0$-correct.*
(2) *Every $\Pi_1^0$ sentence consistent with $T$ is $\Sigma_1^0$-con over $T$.*
(3) *Every $T$-$\Delta_1^0$ sentence is decided by $T$.*

PROOF. $(1) \Rightarrow (2) \Rightarrow (3)$. Easy exercises. (The implications do not, in fact, depend on $T$ being r.e.)

$(3) \Rightarrow (1)$. Prove the contrapositive. Let $\sigma$ be a 'prenex' $\Sigma_1^0$ sentence which is false but $T$ provable. Let $\theta$ be a $\Sigma_1^0$ sentence such that $T$ proves: $\theta \leftrightarrow \mathrm{Thm}_T(\ulcorner \neg \theta \urcorner) \prec (\sigma \vee \mathrm{Thm}_T(\ulcorner \theta \urcorner))$. As usual, $\theta$ is independent of $T$. From the point of view of $T$, the truth of $\sigma$ entails the decidability–i.e., the $\Delta_1^0$-ness–of $\theta$. More exactly: $\theta$ literally *is* $\Sigma_1^0$; and $\sigma$ implies that $\neg \theta \leftrightarrow [(\sigma \vee \mathrm{Thm}_T(\ulcorner \theta \urcorner)) \preccurlyeq \mathrm{Thm}_T(\ulcorner \neg \theta \urcorner)]$. So $\theta$ is $T$-$\Delta_1^0$. $\square$

NOTES TO THEOREM 2.1. The implication $\neg(1) \Rightarrow \neg(3)$ is ineffective: For every $e \in \omega$, let $T_e = PA \cup$ all sentences of $LA$ with g.n.'s in the $e$th r.e. subset of $\omega$. There is no recursive function $f$ such that $T_e$ consistent and $\Sigma_1^0$-incorrect implies $f(e)$ codes an independent $T_e$-$\Delta_1^0$ pair (meaning a pair $(\sigma, \pi)$ such that $\sigma$ is $\Sigma_1^0$, $\pi$ is $\Pi_1^0$, $T_e \vdash \sigma \leftrightarrow \pi$, and $T \nvdash \sigma$). *Proof.* Suppose there *is* such an $f$. Notice that if $(\sigma, \pi)$ is an independent $T_e$-$\Delta_1^0$ pair then $\pi \to \sigma$ is false and $T_e$ provable. So from $f$ we get a map $e \mapsto \sigma_e \in \Sigma_1^0$ such that: $T_e$ consistent and $\Sigma_1^0$-incorrect $\Rightarrow \sigma_e$ is *false* and $T_e$ provable. Then "$T_e$ is consistent and $\Sigma_1^0$-correct" is, as a predicate of $e$, a Boolean combination of r.e. predicates, being equivalent to "$T_e$ is consistent and ($T \nvdash \sigma_e$ or $\sigma_e$ is true)". That is impossible, for $\{e | T_e$ is consistent and $\Sigma_1^0$-correct$\}$ can easily be seen to be a complete $\Pi_2^0$ subset of $\omega$.

It will be convenient to formulate two lemmas axiomatizing the proof of the existence theorems–at the cost of some annoying notation.

NOTATION AND SETTING FOR LEMMA 2.2. $T$ is r.e. and $\mathrm{Proof}_T$ is a binumeration. If $\psi$ is a 'prenex' $\Sigma_1^0$ and $\theta$ a 'prenex' $\Sigma_n^0$ sentence, then

$$\psi * \theta \underset{\mathrm{df}}{=} \exists u(u \text{ witnesses } \psi) \wedge \forall x,$$
$$y < u(\mathrm{Proof}_T(y, x) \to \Sigma_n^0\text{-true}(\text{``}x \prec \ulcorner \theta \urcorner \text{''})).$$

Note that $\psi * \theta$ is $T$-$\Sigma_n^0$. For an explanation of $\Sigma_n^0$-True, see Fact 0.2; for "$x \prec \ulcorner \theta \urcorner$" see the end of §0.

LEMMA 2.2. *Suppose that* $PA \vdash \theta \leftrightarrow \psi * \theta$. *Then,*

(1) $\psi$ *is true* $\Rightarrow T \vdash \theta$,

(2) $\psi$ *is false* $\Rightarrow \neg\theta$ *is* $\Sigma_n^0$-*con over every subtheory of* $T$.

PROOF. (1) Suppose that $\psi$ is true, and let $k$ be the least witness to $\psi$. Then $T$ (in fact, $PA$) proves:

$(+)$ $\psi * \theta \leftrightarrow \forall x, y < \mathbf{k}$ $(\mathrm{Proof}_T(y, x) \to \Sigma_n^0\text{-True}(\text{``}x < \ulcorner\theta\urcorner\text{''}))$. For any theory containing $N$, a bounded quantifier is the same as a disjunction. I.e., it is provable that $\forall u < \mathbf{n}\phi(u) \leftrightarrow (\phi(\mathbf{0}) \wedge \cdots \wedge\phi(\mathbf{n}))$. Therefore if $\chi_1, \ldots, \chi_m$ is the list of $\Sigma_n^0$ sentences $T$-provable by proofs with g.n.'s less than $k$, it follows from $(+)$ that $T$ proves:

$(*)$ $\bigwedge\!\!\bigwedge \chi_i$; and also $\bigwedge\!\!\bigwedge \Sigma_n^0\text{-True}(\ulcorner\chi_i < \theta\urcorner) \to \psi * \theta$, hence

$(**)$ $\bigwedge\!\!\bigwedge(\chi_i < \theta) \to \theta$.

Now reason in $T$: Suppose $\neg\theta$. By $(**)$, $\neg \bigwedge\!\!\bigwedge (\chi_i < \theta)$. Together with $(*)$ this implies $\bigvee\!\!\bigvee(\theta \leqslant \chi_i)$, hence $\theta$. We have shown $\neg\theta \to \theta$; so $\theta$.

(2) Now suppose that $\psi$ is false and $T' \subseteq T$. It will suffice, by Lemma 1.2, to show that for any $\sigma$ in $\Sigma_n^0$: If $T' \vdash \sigma$, then $T' \vdash \sigma < \theta$. So suppose $T' \vdash \sigma$. Then $T \vdash \sigma$. Let $k$ be a g.n. of a $T$-proof of $\sigma$. Since $\psi$ is false $T'$ proves:

$\mathrm{Proof}_T(\mathbf{k}, \ulcorner\sigma\urcorner) \wedge \mathbf{k}, \ulcorner\sigma\urcorner <$ any possible witness to $\psi$. Therefore $T'$ proves $\psi * \theta \to \Sigma_n\text{-True}(\ulcorner\sigma < \theta\urcorner)$, so proves $\theta \to \sigma < \theta$. From $\theta \to \sigma < \theta$ and $\sigma$ it follows that $\sigma < \theta$. So $T' \vdash \sigma < \theta$. $\quad\square$

NOTATION AND SETTING FOR LEMMA 2.3. With $T$, $\mathrm{Proof}_T$, $\psi$, and $\theta$ as before, define

$$\psi^+\theta \underset{\mathrm{df}}{=} \exists y, p\,((1) \text{ no witness for } \psi \text{ is } \leqslant y \text{ or } p;$$

$$(2)\, p \text{ is the g.n. of a } \Pi_n^0 \text{ sentence};$$

$$(3)\, \mathrm{Proof}_T(y, \text{``}\ulcorner\theta\urcorner \to p\text{''}); \text{ and}$$

$$(4)\, \Sigma_n^0\text{-True}(\text{``}\neg p\text{''})).$$

Again, $\psi^+\theta$ is $\Sigma_n^0$.

LEMMA 2.3. *Suppose that* $PA \vdash \theta \leftrightarrow \psi^+\theta$. *Then,*

(1) $\psi$ *is true* $\Rightarrow T \vdash \neg\theta$,

(2) $\psi$ *is false* $\Rightarrow \theta$ *is* $\Pi_n^0$-*con over every subtheory of* $T$.

PROOF. (1) Suppose that $\psi$ is true and let $k$ witness $\psi$. Then $T$ proves: $\psi^+\theta \leftrightarrow \exists y, p < \mathbf{k}$ $(p \text{ is } \Pi_n^0 \wedge \mathrm{Proof}_T(y, \text{``}\ulcorner\theta\urcorner \to p\text{''}) \wedge \Sigma_n\text{-True}(\text{``}\neg p\text{''}))$. As before there is a list $\chi_1, \ldots, \chi_m$ of sentences (this time each is $\Pi_n^0$) such that $T$ proves

$(*)$ $\bigwedge\!\!\bigwedge(\theta \to \chi_i)$; and $\psi^+\theta \leftrightarrow \bigvee\!\!\bigvee \Sigma_n\text{-True}(\ulcorner\neg\chi_i\urcorner)$, hence

$(**)$ $\theta \to \bigvee\!\!\bigvee \neg\chi_i$.

But $(*)$ and $(**)$ together imply $\neg\theta$.

(2) Suppose $\psi$ is false. Let $T' \subseteq T$ and $\pi$ be a $\Pi_n^0$ sentence such that

$T' \vdash \theta \to \pi$. Then $T \vdash \theta \to \pi$, so let $k$ be a g.n. of a $T$-proof of $\theta \to \pi$. $T'$ proves: $\mathbf{k}$, $\ulcorner \pi \urcorner <$ any possible witness to $\psi$. Now reason in $T'$. Suppose $\neg \pi$. Then $\Sigma_n^0$-True($\ulcorner \neg \pi \urcorner$); therefore $\psi^+\theta$, and therefore $\theta$. But we know that $\theta \to \pi$. So $\pi$. We have shown that $\neg \pi \to \pi$. Therefore $\pi$. $\square$

NOTE TO LEMMA 2.3. The definition of $\psi^+\theta$ used above is due to Solovay and is considerably simpler and leads to a considerably simpler proof of the lemma than the author's original definition.

THEOREM 2.4. *Let $T$ be r.e. Then we can find, effectively from any r.e. index for $T$, an independent $\Gamma$ sentence which is $\check{\Gamma}$-con over every subtheory of $T$.*

PROOF. We will prove this for the case $\Gamma = \Sigma_n^0$. The proof for $\Pi_n^0$ is similar. By the self-reference lemma we can effectively find a $\Sigma_n^0$ sentence $\theta$ such that $PA$ proves:

$$\theta \leftrightarrow \mathrm{Thm}_T(\ulcorner \neg\theta \urcorner) * \theta.$$

The sentence we want is $\neg\theta$. It will suffice to show that $T \nvdash \neg\theta$; for that guarantees that $\mathrm{Thm}_T(\ulcorner \neg\theta \urcorner)$ is false, hence by Lemma 2.2(2) that $\neg\theta$ is $\Sigma_n^0$-con over any suitable $T'$ (and "nonprovable and $\Sigma_n^0$-con" certainly guarantees "independent"). So suppose instead that $T \vdash \neg\theta$. Then $\mathrm{Thm}_T(\ulcorner \neg\theta \urcorner)$ is true; so by Lemma 2.2(1) $T \vdash \theta$, contradicting the consistency of $T$.

If $\Gamma = \Pi_n^0$, find a $\Sigma_n^0 \phi$ such that $\phi \leftrightarrow \mathrm{Thm}_T(\ulcorner \phi \urcorner)^+ \phi$. $\square$

We will call $\phi$ essentially $\Gamma$ over $T$ if $\Gamma$ is its simplest classification in $T$, and exactly $\Gamma$-con if $\Gamma$ is the most complicated class for which $\phi$ is conservative. More exactly:

DEFINITION 2.5. $\phi$ is *essentially* $\Gamma$ over $T$ iff $\phi \in \Gamma$ but $\phi \notin T - \check{\Gamma}$.

$\phi$ is *exactly $\Gamma$-con* over $T$ iff $\phi$ is $\Gamma$-con but $\check{\Gamma}$-non over $T$.

The next theorem says just what you might expect–that we can obtain any not-obviously-absurd combination of essentially $\Gamma$ and exactly $\Gamma'$-con.

THEOREM 2.6. *Let $T$ be r.e. and $\Gamma' \subseteq \Gamma$. Then, effectively from an r.e. index for $T$ we can find a sentence which is essentially $\Gamma$ and exactly $\check{\Gamma}'$-con over every subtheory of $T$.*

PROOF. To keep the notation in hand take a special case: suppose $\Gamma' = \Pi_5^0$ and $\Gamma = \Sigma_8^0$. Let $\pi$ be an independent $\Pi_5^0$ sentence which is $\Sigma_5^0$-con over every subtheory of $T$; and $\sigma$ be an independent $\Sigma_8^0$ sentence which is $\Pi_8^0$-con over every subtheory of $T + \pi$. Let $\phi$ be $\pi \wedge \sigma$, and $T'$ be any subtheory of $T$. Then $\phi$ is $\Pi_5^0$-non over $T'$ because $\phi$ implies $\pi$. Further, $\phi$ is $\Sigma_5^0$-con, and therefore exactly $\Sigma_5^0$-con: For if $\theta$ is $\Sigma_5^0$ and $T + \pi + \sigma \vdash \theta$, then (because $\sigma$ is $\Pi_8^0$-con over $T + \pi$) $T + \pi \vdash \theta$; and so (because $\pi$ is $\Sigma_5^0$-con over $T$) $T \vdash \theta$. To complete the proof we need to show that $\phi$ is not equivalent in $T'$ to any $\Pi_8^0$ sentence. So suppose that $\theta$ is $\Pi_8^0$ and that $T' \vdash \phi \leftrightarrow \theta$. Then infer succes-

sively: $T' \vdash (\pi \wedge \sigma) \leftrightarrow \theta$; $T' + \pi \vdash \theta$–because $\sigma$ is $\Pi_8^0$-con over $T + \pi$; $T' + \pi \vdash \phi$; $T' + \pi \vdash \sigma$. The last statement contradicts the $T + \pi$ independence of $\sigma$.   $\square$

NOTE TO THEOREMS 2.4 AND 2.6. That the sentences provided in Theorems 2.4 and 2.6 are conservative over so many theories is somewhat surprising, since "$\phi$ is $\Gamma$-con over $T$" is, at least by its looks, a truly global property of $T$. It is of course trivial to produce $\phi$, $T'$, and $T$ with $T' \subseteq T$ and $\phi$ $\Gamma$-con over $T$ but not over $T'$.

THEOREM 2.7 (SOLOVAY). *Let $T$ be r.e. Then for any $\Gamma$ there is a $\Gamma$ sentence $\phi$ such that:*

   $\phi$ is $\check{\Gamma}$-con over $T$ and $\neg\phi$ is $\Gamma$-con over $T$.

PROOF. Notice that the "subtheory property" is not claimed for $\phi$. The parts of the argument through which that claim cannot pass will be starred (and remarked upon at the end).

Let $\sigma$ be a $\Sigma_n^0$ sentence such that $PA$ proves:

$\sigma \leftrightarrow \exists y, p((1)\, y, p <$ any $T$ proof of $\ulcorner\sigma\urcorner \wedge p$ is the g.n. of a $\Sigma_n^0$ sentence;

$\qquad$ (2) $\text{Proof}_T(y, \text{``}\ulcorner\sigma\urcorner \to p\text{''})$;

$\qquad$ (3) $\Sigma_n^0\text{-True}(\text{``}\neg p\text{''})$;

$\qquad$ (4) $\forall s, t < y(\text{Proof}_T(t, s) \to \Sigma_n^0\text{-True}(\text{``}s < \ulcorner\sigma\urcorner\text{''})))$.

Denote the formula to the right of the $\leftrightarrow$ sign by '$\Delta$'. By ignoring (4) we see that, in the notation of Lemma 2.3, $PA \vdash \sigma \to \text{Thm}_T(\ulcorner\sigma\urcorner)^+\sigma$. That fact suffices to carry out the argument of part (1) of Lemma 2.3. So, if $\text{Thm}_T(\ulcorner\sigma\urcorner)$ is true then $T \vdash \neg\sigma$–a contradiction. We have

*Fact* 1. $T \not\vdash \sigma$.

*Fact* 2. $\sigma$ is $\Pi_n^0$-con over $T$.

PROOF. Suppose $T \vdash \sigma \to \pi$ and let $k$ be the g.n. of a $T$-proof of $\sigma \to \pi$. Reason in $T$: Suppose $\neg\pi$. Then $\neg\sigma$. If we substitute $\mathbf{k}$ for "$y$" and $\ulcorner\pi\urcorner$ for "$p$" then clauses (1)–(3) of $\Delta$ are true. So the only way that $\sigma$ can be false is for (4) to fail, i.e.:

   (i) $\exists s, t < \mathbf{k}(\text{Proof}_T(t, s) \wedge \neg\Sigma_n^0\text{-True} (\text{``}s < \ulcorner\sigma\urcorner\text{''}))$.

(Now step outside of $T$ for a moment and produce $\chi_1, \ldots, \chi_m$–the $\Sigma_n^0$ sentences provable in $T$ by proofs with g.n.'s $< k$. Go back to $T$.) Therefore,

   (∗) $\bigwedge\chi_i$; and, because we are assuming $\neg\pi$, $\bigvee\neg(\chi_i < \sigma)$.

(That last remark is a consequence of the tail end of (i).) Since $\neg\sigma$, the only way to have $\neg(\chi_i < \sigma)$ is to have $\neg\chi_i$. We have shown: $\bigwedge\chi_i \wedge (\neg\pi \to \bigvee \neg\chi_i)$. I.e., we have shown $\pi$.

*Fact* 3. $\neg\sigma$ is $\Sigma_n^0$-con over $T$.

PROOF. Suppose $\pi$ is $\Pi_n^0$ and $T \vdash \pi \to \sigma$. It suffices to show that $T \vdash \neg\pi$. Let $\chi$ be $\Sigma_n^0$ such that $\chi \leftrightarrow (\neg\pi \vee \sigma)$ *and*

(ii) $T \vdash (\neg \pi \lor \sigma) \preccurlyeq \chi$;

this last is easily arranged. Let $k$ be a g.n. of a $T$-proof of $\chi$. We know that $\sigma$ is $\Pi_n^0$-con over $T$ and hence that for any $\psi \in \Pi_n^0$, if $T \vdash \sigma \to \psi$ then $T \vdash \psi$. Therefore

$(**) \; \forall y, p < k(p$ is $\Pi_n^0 \land \text{Proof}_T(y, \text{``}\ulcorner \sigma \urcorner \to p\text{''}) \to \Pi_n^0\text{-True}(p))$

is equivalent to a finite conjunction of $T$ provable statements, so is itself $T$ provable. What $(**)$ entails, and this entailment is formalizable in $T$, is that *if $\sigma$ is to be true, any candidate to play the role of the "$y$" in $\Delta$ must be $> k$.* Hence: *if $\sigma$ is true, $\ulcorner \chi \urcorner$ and $k$ must fall within the range of the bounded quantifier in clause (4) of $\Delta$*, and therefore $\sigma \to \Sigma_n^0\text{-True}(\ulcorner \chi \prec \sigma \urcorner)$. Formalizing in $T$ gives $T \vdash \sigma \to \chi \prec \sigma$. Now reason in $T$: We know that $\sigma \to \chi \prec \sigma$ and that $\chi$. Therefore $\chi \prec \sigma$. But that, together with $(\neg \pi \lor \sigma) \preccurlyeq \chi$–see (ii)–implies $\neg \pi$. We have proven $\neg \pi$. $\square$

NOTES TO THEOREM 2.7. In the proof of Fact 2, the subtheory property escapes at $(*)$. Suppose, at that point, that we are trying to reason in $T' \subseteq T$. We know that $\neg \pi \to \bigvee \neg(\chi_i \prec \sigma)$, but not that $\bigwedge \chi_i$–for the $\chi$'s are consequences of $T$. There is no way, in $T'$, to bring those sentences into collision. One trouble spot in the proof of Fact 3 is $(**)$. We need all of $T$ to prove $(**)$–and would still need all of $T$ even if, by magic, $\sigma$ were $\Pi_n^0$-con over $T'$.

DEFINITION 2.8. $\phi$ is *essentially* $\Delta_n^0$ over $T$ iff $\phi$ is $T$-$\Delta_n^0$ but not $T$-$\Sigma_{n-1}^0$ or $T$-$\Pi_{n-1}^0$.

$\phi$ is *exactly* $\Delta_n^0$-con over $T$ iff $\phi$ is $\Delta_n^0$-con, but $\Sigma_n^0$-non and $\Pi_n^0$-non, over $T$.

COROLLARY 2.9. *Let $T$ be r.e. and $n > m \geqslant 1$. Then there is a sentence which is essentially $\Delta_n^0$ and exactly $\Delta_m^0$-con over $T$.*

PROOF. It will suffice to prove the theorem for the case $n = m + 1$, the general result following by iterating in the manner of Theorem 2.6. Say $m = 4$, $n = 5$. Let $\sigma$ be a $\Sigma_4^0$ sentence which is $\Pi_4^0$-con over $T$ and for which $\neg \sigma$ is $\Sigma_4^0$-con over $T$. Let $\pi$ be a $\Pi_4^0$ sentence which is $\Sigma_4^0$-con and independent over $T + \sigma$. We want $\phi =_{\text{df}} \pi \land \sigma$. Clearly $\phi$ is $\Delta_5^0$, $\Pi_4^0$-non, and $\Sigma_4^0$-non over $T$. To see that $\phi$ is $\Delta_4^0$-con: Suppose $P$ is $\Pi_4^0$, $S$ is $\Sigma_4^0$, $T \vdash P \leftrightarrow S$, and $T + \phi \vdash P(\land S)$. Then $T + \sigma \vdash \pi \leftrightarrow S$ and so by choice of $\pi$, $T + \sigma \vdash S$. So $T + \sigma \vdash P$, and by choice of $\sigma$, $T \vdash P$. To see that $\phi$ is neither $T$-$\Sigma_4^0$ nor $T$-$\Pi_4^0$: Suppose $S$ is $\Sigma_4^0$ and $T \vdash \phi \leftrightarrow S$. Then $T + \sigma \vdash \pi \leftrightarrow S$ contradicting the essential $\Pi_4^0$-ness of $\pi$ over $T + \sigma$. Suppose $P$ is $\Pi_4^0$ and $T \vdash \phi \leftrightarrow P$. Then $T \vdash \neg \sigma \to \neg P$, and by choice of $\sigma$, $T \vdash \neg P$. So $T \vdash \neg \phi$, which is impossible. $\square$

NOTES TO COROLLARY 2.9. Only the last step required use of Theorem 2.7. That is where the subtheory property is lost. Mixing Theorem 2.4 and Corollary 2.9 does *not* yield, e.g., an essentially $\Sigma_4^0$ sentence which is exactly

$\Delta_4^0$-con. Whether all reasonable combinations of $\Pi$, $\Sigma$, and $\Delta$ can be obtained I do not know.

Our last existence theorem simply codifies the effectiveness available in the proofs of Lemmas 2.2 and 2.3. It helps partially to classify $\{\phi | \phi$ is $\Gamma$-con over $T\}$ for r.e. $T$ by showing that it cannot be an r.e. set of formulas.

LEMMA 2.10. *Let $X$ be an r.e. subset of $\omega$ and $T$ an r.e. theory. Then, for any $\Gamma$, there is a recursive map $k \mapsto \theta_k \in \Gamma$, such that*:
$$k \in X \Rightarrow T \vdash \neg \theta_k.$$
$k \notin X \Rightarrow \theta_k$ *is independent and $\check{\Gamma}$-con over every subtheory of $T$.*

PROOF. Consider the case $\Gamma = \Pi_n^0$. Let $\psi(x)$ be a $\Sigma_1^0$ predicate which, in the real world, defines $X$. For each $k$ produce $\phi_k$ such that (in the notation of Lemma 2.2), $PA$ proves:
$$\phi_k \leftrightarrow \big(\psi(\mathbf{k}) \vee \mathrm{Thm}_T(\ulcorner \neg \phi_k \urcorner)\big) * \phi_k.$$

We are going to apply Lemma 2.2 repeatedly. Since the truth of either $\mathrm{Thm}_T(\ulcorner \neg \phi_k \urcorner)$ or $\psi(\mathbf{k})$ implies that $T \vdash \phi_k$, we have for any $k$:

(i) $T \nvdash \phi_k$,

(ii) $k \in X \Rightarrow \psi(\mathbf{k})$ is true $\Rightarrow T \vdash \phi_k$;

and by the other half of Lemma 2.2,

(iii) $k \notin X \Rightarrow \psi(\mathbf{k}) \vee \mathrm{Thm}_T(\ulcorner \neg \phi_k \urcorner)$ is false $\Rightarrow \neg \phi_k$ is $\Sigma_n^0$-con and independent over any subtheory of $T$.

Let $\theta_k$ be (a $\Pi_n^0$ sentence equivalent to) $\neg \phi_k$. The other case is similar.    $\square$

By iterating the steps in the lemma as in the proof of Theorem 2.6, we get:

THEOREM 2.11. *Let $T$ be r.e. and $\Gamma' \subseteq \Gamma$. Then every $\Pi_1^0$ subset of $\omega$ is reducible to either of the sets*:

$\{\phi | \phi$ *is essentially $\Gamma$ and exactly $\check{\Gamma}'$-con over $T\}$.*

$\{\phi | \phi$ *is essentially $\Gamma$ and exactly $\check{\Gamma}'$-con over every subtheory of $T\}$.*

It seems worth noting one more restatement of the existence lemmas which is cute and sometimes useful. Say that $\psi(x)$ is $\Gamma$-disjunctive (over $T$) iff for every sentence $\chi \in \Gamma$ and any $n \in \omega$, $T \vdash \psi(\mathbf{n}) \vee \chi$ implies $T \vdash \psi(\mathbf{n})$ or $T \vdash \chi$.

THEOREM 2.12. *Let $X$ be an r.e. subset of $\omega$ and $T$ an r.e. theory. Then there is a $\psi(x) \in \Gamma$ such that $\psi$ numerates $X$ in $T$ and $\psi$ is $\Gamma$-disjunctive.*

PROOF. To be consistent with the notation of 2.10 we will make $\psi \in \check{\Gamma}$ and $\check{\Gamma}$-disjunctive. Let the map $k \mapsto \theta_k \in \Gamma$ be as in 2.10, and let $x \mapsto \dot{\theta}_x$ be the formalization of that map. Set
$$\psi(x) \underset{\mathrm{df}}{=} \check{\Gamma}\text{-True}(\text{``}\neg \dot{\theta}_x\text{''}).$$
Since $T \vdash \psi(\mathbf{n}) \leftrightarrow \neg \theta_n$ for each $n \in \omega$, $\psi$ numerates $X$. Suppose now that $\chi \in \Gamma$ and $T \vdash \psi(\mathbf{n}) \vee \chi$, but $T \nvdash \psi(\mathbf{n})$. Then $n \notin X$ and therefore $\theta_n$ is $\check{\Gamma}$-con over $T$; and, since $T \vdash \neg \psi(\mathbf{n}) \to \chi$, $T \vdash \theta_n \to \chi$. So $T \vdash \chi$.    $\square$

### 3. Interpretability.

DEFINITIONS 3.1. $\phi$ is *interpretable* in $T$ iff $T + \phi$ is interpretable in $T$ in the sense of [**Shoenfield**, pp. 61–62].

$T$ is *essentially reflexive* iff for every sentence $\phi$ in the language of $T$, $T \vdash \phi \rightarrow \mathrm{Con}(\{\phi\})$.

Interpretability is discussed in [**Hajek, 1**], [**Hajek, 2**], [**Hajkova-Hajek**] and [**Solovay**].

Whether $T$ is essentially reflexive depends solely on the nonarithmetical part of $T$. All extensions of $PA$ in the same language as $PA$ are essentially reflexive. $ZF$ (and any extension in the same language) is essentially reflexive. $GB$ is not (for an essentially reflexive theory cannot be finitely axiomatized). (See [**Montague**].)

The connections between interpretability and $\Pi_1^0$-con are easily stated:

THEOREM 3.2 (HÁJEK, LARGELY). *Let $T$ be r.e. and essentially reflexive. Then, $\phi$ is interpretable in $T$ iff $\phi$ is $\Pi_1^0$-con over $T$.*

PROOF. Say that $\phi$ is strongly consistent with $T$ if for every finite $F \subseteq T$, $T \vdash \mathrm{Con}(F + \phi)$. [**Hájek,1**] shows that for r.e. essentially reflexive $T$, $\phi$ is strongly consistent with $T$ iff $\phi$ is interpretable in $T$. Thus it will more than suffice to show:

*Claim.* If $T$ is essentially reflexive (not necessarily r.e.), then $\phi$ is $\Pi_1^0$-con over $T$ iff $\phi$ is strongly consistent with $T$.

PROOF. Suppose first that $\phi$ is $\Pi_1^0$-con over $T$ and let $F \subseteq T$ be finite. Since $T$ is essentially reflexive, $T \vdash \bigwedge F \wedge \phi \rightarrow \mathrm{Con}(\bigwedge F \wedge \phi)$–i.e., $T \vdash \phi \rightarrow \mathrm{Con}(F + \phi)$. Since $\phi$ is $\Pi_1^0$-con over $T$, $T \vdash \mathrm{Con}(F + \phi)$.

Suppose, conversely, that for every finite $F \subseteq T$, $T \mathrm{Con} \vdash (F + \phi)$. Suppose that $\pi$ is $\Pi_1^0$ and $T + \phi \vdash \pi$. Produce a finite $F \subseteq T$ such that $F + \phi \vdash \pi$. We may as well assume $N \subseteq F$. Then $T \vdash \mathrm{Thm}_{F+\phi}(\ulcorner \pi \urcorner)$. Now reason in $T$: If $\neg \pi$, then $\mathrm{Thm}_N(\ulcorner \neg \pi \urcorner)$, so $\mathrm{Thm}_{F+\phi}(\ulcorner \pi \wedge \neg \pi \urcorner)$; contradicting $\mathrm{Con}(F + \phi)$. Since $\neg \pi$ implies a contradiction, $\pi$. □

NOTES TO THEOREM 3.2. [**Hájek, 2**] shows that if $T$ contains induction for *all* formulas of $T$ (no assumption about reflexiveness or recursive enumerability) then every sentence interpretable in $T$ is $\Pi_1^0$-conservative over $T$. A slight elaboration of that argument is used in 6.5.

The hypothesis "$T$ is essentially reflexive" cannot merely be omitted from the hypothesis of 3.2. Consider $GB$, which is not essentially reflexive (nor does it contain induction for all formulas). Then $I$, the set of sentences interpretable in $GB$, is r.e. (because $GB$ is finitely axiomatizable), while $C$, the set of sentences $\Pi_1^0$-con over $GB$, is not (by 2.11). So $I \neq C$. More is known: Solovay has shown that $I \setminus C \neq 0$; and Hájek has observed that $C \setminus I \neq 0$ follows easily from Lemma 2.3.

I owe thanks to Hájek and Švejdar for clearing up my confusions about these things.

An interpretation of $\phi$ in $T$ provides a proof of the consistency of $T + \phi$ relative to $T$. Whether this consistency proof is "elementary" depends on the presentation of that interpretation. Suppose our measure of elementary is this: $\text{Con}(T) \to \text{Con}(T + \phi)$ can be proven in $PA$. Our reductive attitude about the formulation of $\text{Con}(T)$ will come in handy (and may make the next result appear to say more than it does say). If, as in [**Feferman**], we let $\text{Con}_\psi$ be the consistency statement naturally based on $\psi$ (and let $\psi'(x)$ abbreviate $\psi(x) \lor x = \ulcorner \phi \urcorner$), then the meaning–and presumably the provability–of $\text{Con}_\psi \to \text{Con}_{\psi'}$ depends strongly on $\psi$. By choosing particular representations of $\text{Con}(T)$ and $\text{Con}(T + \phi)$–which are related in the natural way– such difficulties are dodged. If $\phi$ is a Rosser sentence for $T$, then there is an elementary proof of the relative consistency of $T + \phi$, even though $T + \phi$ cannot be interpreted in $T$. The next theorem says that even if $\phi$ is $\Sigma_1^0$, the existence of an interpretation of $T + \phi$ in $T$ need not guarantee the existence of an elementary relative consistency proof for $\phi$. (The point, again, is that $\phi$ is so simple.)

**THEOREM 3.3.** *Let $T$ be r.e. and essentially reflexive. Then there is a $\Sigma_1^0$ sentence $\phi$ such that $\phi$ is interpretable in $T$ but $PA \nvdash \text{Con}(T) \to \text{Con}(T + \phi)$.*

PROOF. Choose $\phi$ so that, in the notation of Lemma 2.3,

$$\phi \leftrightarrow \text{Thm}_{PA}(\text{Con}(T) \to \text{Con}(T + \phi))^+ \phi.$$

We will apply Lemma 2.3 repeatedly. If $PA \vdash \text{Con}(T) \to \text{Con}(T + \phi)$ then by part (1) of the lemma, $T \vdash \neg \phi$. So $PA \vdash \neg \text{Con}(T + \phi)$ and therefore $PA \vdash \neg \text{Con}(T)$–which is impossible because $PA$ is $\Sigma_1^0$-correct. Since $PA \nvdash \text{Con}(T) \to \text{Con}(T + \phi)$, part (2) of the lemma guarantees that $\phi$ is $\Pi_1^0$-con and therefore interpretable in $T$.   □

Here is how $\Pi_n^0$-con can be understood in terms of interpretability. Say that an interpretation $I$ of $T'$ in $T (\subseteq T')$ is *provably* $\Gamma$-*faithful* if for every sentence $\chi$ in $\Gamma$, $T \vdash \chi_I \to \chi$. (Here $\chi_I$ is the interpretation of $\chi$. See [**Shoenfield**].)

**THEOREM 3.4.** *Let $T$ be r.e. and essentially reflexive and $\phi$ any sentence in the language of $T$. Then, $\phi$ is $\Pi_n^0$-con over $T$ iff there is a provably $\Pi_n^0$-faithful interpretation of $T + \phi$ in $T$.*

PROOF. The implication from right to left is immediate. Suppose conversely, that $\phi$ is $\Pi_n^0$-con over $T$. For each finite $F \subseteq T$ consider the sentence $\phi_F$: $\forall x(\Sigma_n^0 - \text{True}(x) \to$ the theory whose axioms are $\ulcorner F \urcorner$, $\ulcorner \phi \urcorner$, and $x$ is consistent). Then $T + \phi \vdash \phi_F$; and since $\phi_F$ is $\Pi_n^0$, $T \vdash \phi_F$. Using that fact and the tricks in [**Feferman**] we can find a formula $\psi_1(x)$ binumerating $T$ in $T$

such that $T \vdash \mathrm{Con}\ \psi_2$, where $\psi_2(x)$ is ($\Sigma_n^0$-True$(x) \vee \psi_1(x) \vee x = \ulcorner\phi\urcorner$). From the point of view of $T$, $\psi_2$ is the theory consisting of $T$ (or at least that fragment of $T$ described by $\psi_1$, which of course is a numeration of $T$, but is not $T$ *provably* equal to $T$), $\phi$, and all the true $\Sigma_n^0$ sentences. By 5.9 of [**Feferman**], slightly modified, there is a formula $\chi(x)$ such that $T$ proves: "$\{x|\chi(x)\}$ is a complete Henkin extension of $\{x|\psi_2(x)\}$." (A Henkin theory is one in which every provable existential statement is witnessed by some constant.) In particular, if $\sigma$ is $\Sigma_n^0$ then, since $T \vdash \sigma \to \Sigma_n^0$-True$(\ulcorner\sigma\urcorner)$, we have $T \vdash \sigma \to \chi(\ulcorner\sigma\urcorner)$. By imitating the usual construction of a model from a complete Henkin theory we can extract from $\chi$ an interpretation $I$ of $T + \phi$ in $T$ such that for any sentence $\theta$ of $T$, $T \vdash \chi(\ulcorner\theta\urcorner) \leftrightarrow \theta_I$. Now suppose that $\pi$ is $\Pi_n^0$ and reason in $T$: $\neg\pi \to \chi(\ulcorner\neg\pi\urcorner), \to (\neg\pi)_I, \to \neg(\pi_I)$. So $\pi_I \to \pi$. $\square$

The argument in effect constructs inside $T$ a model of $T + \phi$ which extends the "standard" one (standard from the point of view of $T$) and is elementary for $\Sigma_{n-1}^0$ formulas. An attempt to produce some corresponding theorem for $\Sigma_n^0$ would seem to require us to build not extensions but submodels; yet $T$ thinks that its number structure is the minimal one.

**4. Generalizations.** *Extending the language of arithmetic.* If $L$ extends $LA$ define $PA_L$ to be $PA$ + induction for all formulas in $L$. Define $\Sigma_n^0(L)$ and $\Pi_n^0(L)$ in the obvious way. The self-reference lemma holds for $PA_L$; and if $L' \subseteq L$ has only *finitely many* symbols other than constant symbols (in particular, if $L$ is finite) we can define $\Sigma_n^0(L')$ truth by a $\Sigma_n^0(L')$ formula. The only important difference between $L$ and $LA$ is this: $T$ might not decide all the $\Delta_0^0(L)$ sentences. If one wants to check the extendability of an argument from $LA$ to $L$ that is the first point to check. In particular, if $L$ is finite the existence theorems 2.4, 2.6, 2.7, 2.9, 2.10 all hold with "$\Gamma(L)$" replacing "$\Gamma$" everywhere. (Note: the $\psi$ of Lemmas 2.2 and 2.3 must be $\Sigma_1^0$, not $\Sigma_1^0(L)$–indeed, the statement "$\psi$ is true" is probably nonsensical otherwise.) It is easy to construct an r.e. theory $T$ in an infinite language $L$ for which *every* formula is $T$-$\Delta_0^0(L)$: let $\langle S_n | n \in \omega \rangle$ be a sequence of new symbols such that each $S_n$ defines truth for the formulas in $\Sigma_n^0(PA \cup \{S_0, \ldots, S_{n-1}\})$. To extend the theorems mentioning reflexiveness it is necessary to formulate a stronger notion of reflexive. Roughly, $T$ is essentially-$L$-reflexive iff for every $\phi$, $T \vdash \phi \to \mathrm{Con}(\phi + \text{all the true } \Delta_0^0(L) \text{ sentences})$. If $L$ is finite all extensions of $PA_L$ having language $L$ are essentially $L$-reflexive. All the results of §3 go through if "essentially reflexive" is everywhere strengthened to "essentially-$L$-reflexive."

*Set theory.* Use $\Sigma_n$, $\Pi_n$ for the classes in the Levy hierarchy. The existence theorems will follow once we have a workable definition of $\phi \leqslant \chi$. Define $x <^* y$ to mean: $x \in$ the transitive closure of $\{y\}$; and $x <^* y$ iff $x \leqslant^* y$ and $x \neq y$. Define $\phi \leqslant \chi$ as before, using $<^*$ in place of $<$. If $\phi$ and $\chi$ are

'prenex' $\Sigma_n$, so is $\phi \preccurlyeq \chi$. Our only worry: $\preccurlyeq^*$ is a partial, but not total, order. One of the arrows in Lemma 1.2 fails, but it is not one that hurts. The proofs of Lemmas 2.2 and 2.3 go through if we simply replace $\Sigma_n^0$-True by $\Sigma_n$-True. So Theorems 2.6, 2.7, 2.10 generalize. We could beef up the language and talk about $\Sigma_n(L)$–there are no problems.

The theorems on interpretations go through as before, but there is a slight difference in spirit. The models corresponding to the interpretations are not even *extensions* of the originals, let alone true outer models (i.e., *end extensions*). To "correct" this we could try reworking some of the theorems to apply to $L_{\infty\omega}$. Here is how one looks–one which will be useful in §6. (Elementary knowledge of $L_{\infty\omega}$ is assumed. What we need is contained in the first few chapters of [**Barwise, 1**] and [**Keisler**].) Our language contains $\in$ and, for each set $x$, a constant $\mathbf{x}$ which will turn out to be a name for $x$. Let $\lambda_x$ be the usual infinitary sentence saying that $\mathbf{x}$ denotes $x$. (A transitive set satisfies $\lambda_x$ iff it contains $x$.) Define $\mathrm{Proof}_\infty(y; z, x)$ to mean: using the usual axioms and rules for $L_{\infty\omega}$, $x$ is a deduction of $z$ from assumptions $y \cup \{\lambda_x | x \in V\}$. $\mathrm{Proof}_\infty$ is $ZF\text{-}\Delta_1$. We obtain from $\mathrm{Proof}_\infty$ the $\Sigma_1$ formulas $\mathrm{Thm}_\infty(y; z)$–"$z$ is a theorem of $y \cup \{\lambda_x | x \in V\}$"; and $\mathrm{Con}_\infty(y)$–"$y \cup \{\lambda_x | x \in V\}$ is consistent." The proof of the claim in 3.2 is a roadmap for the proof of

THEOREM 4.1. *A sentence $\phi$ is $\Pi_1$-con over ZF iff for every finite $F \subseteq ZF$,* $ZF \vdash \mathrm{Con}_\infty(F + \phi)$.

**5. Questions.** *Existence theorems.* (1) Can Theorem 2.4 be made uniform? I.e., if $\langle T_i | i \in \omega \rangle$ is an r.e. sequence of r.e. theories is there a $\Gamma$ sentence which is independent and $\check{\Gamma}$-con over each $T_i$? The question is open even for sequences of length 2.

(2) Can Theorem 2.6 be extended to allow all not-obviously-absurd combinations of $\Pi$, $\Sigma$, and $\Delta$?

*Classifying $\Gamma$-con.* Let $T$ be r.e. Here is what is known:

(a) $\{\phi | \phi$ is $\Gamma$-con over $T\}$ is, by inspection, a $\Pi_2^0$ subset of $\omega$; and, by 2.11, cannot be r.e.

(b) Solovay has shown that if $T$ is essentially reflexive, then $\{\phi | \phi$ is $\Pi_1^0$-con over $T\}$ is a complete $\Pi_2^0$ subset of $\omega$. So:

(3) Is $\{\phi | \phi$ is $\Gamma$-con over $T\}$ a complete $\Pi_2^0$ subset of $\omega$?[1]

**6. Partially conservative extensions: Semantics.** This section contains model theoretic characterizations of partially conservative sentences–characterizations which have been forshadowed by the syntactical arguments of §3 which mimicked model constructions. Some of the proofs can be carried out very smoothly using the machinery of admissible covers [**Barwise, 1**] but others

---

[1] Hajek has shown that $\{\phi | \phi$ is $\Gamma$-con over $PA\}$ is a $\Pi_2^0$ complete subset of $\omega$.

cannot–not, at least, in any obvious way. The problem seems to be that compactness arguments, useful as they are for extending models, are not so good at producing submodels. Instead we will use and/or modify a series of theorems and definitions from [**Friedman**], few of which will be stated in their most general form.

The results about models of arithmetic can be read independently of those about set theory. We will assume the reader knows some basic facts about models of set theory–in particular, that he knows the definition of "standard part". (See e.g., [**Friedman**] or [**Barwise, 1**].) We will always identify the standard part of a model of set theory with the transitive set to which it is isomorphic.

From now on say that $T$ is a number theory if for some not necessarily finite language $L$, $T$ is an extension of $PA_L$ with language $L$ (i.e., the only sort is the number sort). $T$ is a set theory if for some language $L$, $T$ is an extension of $ZF_L$ with language $L$. ($ZF_L$ is $ZF$ + comprehension, collection, and foundations for all formulas of $L$.) We will perpetrate the small abuse of crediting every model of set theory with being a model of $PA$.

DEFINITION 6.1. $X$ is $c$-closed (for "completion closed") iff $X$ is a collection of subsets of $\omega$ such that: (i) $X$ is closed under Turing join and "recursive in"; (ii) whenever $y \in X$ codes an infinite binary tree, some path through that tree is in $X$.

The terminology "completion closed" is suggested by the fact that if $X$ is $c$-closed and $T \in X$ is a consistent theory, then some complete extension of $T$ is an element of $X$.

DEFINITION 6.2. If $\mathfrak{M}$ is an $\omega$-nonstandard model of arithmetic then $x \in ss(\mathfrak{M})$ if $x \subseteq \omega$ and for some $\phi(x, \vec{y})$ and some $\vec{b} \in |\mathfrak{M}|$, $x = \{n\varepsilon\omega | \mathfrak{M} \vDash \phi(\mathbf{n}, \mathbf{b})\}$. The sets in $ss(\mathfrak{M})$–read "the standard system of $\mathfrak{M}$"–are all in fact initial segments of $\Delta_0^0$-definable classes. Define $u\varepsilon v$ to mean that the $u$th prime divides $v$. Then standard sets all have the form $\{n\varepsilon\omega | \mathfrak{M} \vDash \mathbf{n} \in \mathbf{a}\}$ for some suitable (nonstandard) $a \in |\mathfrak{M}|$. The corresponding syntactical notion is:

$$bi(T) = \{x \subseteq \omega | x \text{ is binumerated in the theory } T\}.$$

THEOREM 6.3 ([**Friedman**], MODIFIED). (i) $\forall\mathfrak{M}(\mathfrak{M}$ models $PA \Rightarrow ss(\mathfrak{M})$ is $c$-closed).

(ii) *If $X$ is countable and $c$-closed, $T$ a consistent set or number theory, and* $bi(T) \subseteq X$, *then* $\exists\mathfrak{M}(\mathfrak{M} \vDash T$ and $ss(\mathfrak{M}) = X)$.

These results are more or less contained in Theorems 2.4 and 2.5 of [**Friedman**]–see especially the discussion of $Z$ on pp. 541–542. Part (ii) above does not quite follow from Friedman's 2.5 but is proved by the same argument. Our primary interest in 6.3 lies in the following embedding

theorem– which shows that an obviously necessary condition for "$\mathfrak{N}$ is an end extension of $\mathfrak{M}$" is also sufficient.

THEOREM 6.4 ([**Friedman**], MODIFIED). *Let L be a countable language and T a number theory (resp. set theory) with language L. Suppose that $\mathfrak{M}$ and $\mathfrak{N}$ are countable $\omega$-nonstandard models of T. Then, $\mathfrak{N}$ is (isomorphic to) an end extension of $\mathfrak{M}$ if and only if* ss($\mathfrak{M}$) = ss($\mathfrak{N}$) *and every* $\Sigma_1^0(L)$ *(resp., $\Sigma_1(L)$) sentence true in $\mathfrak{M}$ is true in $\mathfrak{N}$.*

Theorem 4.2 of [**Friedman**] is not about end extensions of models for set theory, but rather what are sometimes called "rank extensions"–all the new elements in the larger model having ordinal rank greater than any ordinal of the smaller model. However, a proof of 6.4 (above) is embedded in Friedman's proof.

If $\mathfrak{M}$ models $PA$, say that $\mathfrak{M}$ $T$-models $T$ if $\mathfrak{M} \vDash T$ and $T \in$ ss($\mathfrak{M}$). This stands in for the property that $T$ be numerable in $T$. It is trivial to check that if $\mathfrak{M} \vDash T$ there is some elementary extension $\mathfrak{N}$ of $\mathfrak{M}$ which $T$-models $T$.

THEOREM 6.5. *Let T be a number theory in a finite language L and $\phi$ be any sentence of L. Then,*

(i) $\phi$ *is $\Pi_1^0(L)$-con over $T \Leftrightarrow$ every $T$-model of T can be end extended to a model of $T + \phi$.*

(ii) $\phi$ *is $\Sigma_1^0(L)$-con over $T \Leftrightarrow$ every countable $T$-model of T is an end extension of a model of $T + \phi$.*

Matijasevic's theorem, hereinafter referred to as [M], says that every $\Sigma_1^0$ formula is equivalent in $PA$ to a purely existential formula. (This does not generalize to $\Sigma_1^0(L)$!) From [M], 6.5, and the Lowenheim-Skolem theorem it easily follows that

THEOREM 6.6 [M]. *If $L = LA$, then 6.5 remains true if "end extension" is replaced by "extension"; and even if "countable" is dropped from the statement of* (ii).

PROOF OF 6.5. (i) The implication from right to left is trivial: Suppose that every $T$-model of T can be end extended to model $T + \phi$, and that $\pi$ is $\Pi_1^0(L)$ and $T + \phi \vdash \pi$. We will show that every model of T models $\phi$. Let $\mathfrak{M} \vDash T$. There exists a $T$-model of $T$, $\mathfrak{M}'$, which is an elementary extension of $\mathfrak{M}$. By assumption we can produce an end extension $\mathfrak{N}$ of $\mathfrak{M}'$ modelling $T + \phi$. Then $\mathfrak{N} \vDash \pi$; and since $\Pi_1^0(L)$ sentences persist downward to initial segments, $\mathfrak{M}' \vDash \pi$. Therefore $\mathfrak{M} \vDash \pi$.

Suppose conversely that $\phi$ is $\Pi_1^0(L)$-con over $T$. Examine the proof of 3.4. The recursive enumerability of $T$ was used only to guarantee the existence of a numeration of $T$ in $T$. So expand $T$ to $T'$ by adding a one place predicate

$\mathcal{T}$, induction for all new formulas, and $\{\mathcal{T}(\ulcorner\chi\urcorner)|\chi \in T\} \cup \{\neg\mathcal{T}(\ulcorner\chi\urcorner)|\chi \notin T\}$. $T'$ is a conservative extension of $T$, so $\phi$ is $\Pi_1^0(L)$-con over $T'$ and the proof of 3.4 proceeds as before. It provides an interpretation $I$ with the following properties: (a) There is a $T'$ definable function $u \mapsto u_I$ with meaning "$u_I$ is the interpretation of the $u$th numeral (hence $u$ is the denotation of $u_I$)". (b) $T'$ proves that function embeds $<$ as an initial segment of $<_I$ and, for any $R \in L$, that $R(u, \ldots, v) \leftrightarrow R_I(u_I, \ldots, v_I)$. If $\mathfrak{M}$ $T$-models $T$ we can expand $\mathfrak{M}$ to a model $\mathfrak{M}'$ of $T'$ and then pull out the interpretation $I$ to a structure $\mathfrak{N}$ with domain $= \{b|\mathfrak{M}' \vDash \mathbf{b}$ is in the universe of $I\}$ and relations of form $R^{\mathfrak{N}} = \{(a, \ldots, b)|\mathfrak{M}' \vDash R_I(\mathbf{a}, \ldots, \mathbf{b})\}$. The function $f = \{(a, b)|\mathfrak{M}' \vDash \mathbf{b} = (\mathbf{a})_I\}$ embeds $\mathfrak{M}$ as an initial segment of $\mathfrak{N}$. And since $T'$ proves $\chi_I$ for each $\chi \in T + \phi$, $\mathfrak{N} \vDash T + \phi$.

(ii) The proof of the right to left implication is dual to the proof in (i). So suppose that $\phi$ is $\Sigma_1^0(L)$-con over $T$ and that $\mathfrak{M}$ is countable and $\mathfrak{M}$ $T$-models $T$. Let $S = \{\pi \in \Pi_1^0(L)|\mathfrak{M} \vDash \pi\}$. Because $\phi$ is $\Sigma_1^0(L)$-con over $T$, $T' =_{df} T \cup S \cup \{\phi\}$ is consistent. Since $T$ admits a $\Pi_1^0(L)$ truth definition $S \in ss(\mathfrak{M})$–and by the closure properties of $ss(\mathfrak{M})$, $T' \in ss(\mathfrak{M})$ and $bi(T) \subseteq ss(\mathfrak{M})$. Apply 6.3 to get an $\mathfrak{N} \vDash T'$ with $ss(\mathfrak{N}) = ss(\mathfrak{M})$. Because $\mathfrak{N} \vDash S$, every $\Sigma_1^0(L)$ sentence true in $\mathfrak{N}$ is true in $\mathfrak{M}$. So $\mathfrak{N}$ is a model of $T + \phi$ which (by 6.4) can be embedded as an initial segment of $\mathfrak{M}$. $\square$

NOTES TO 6.5 AND 6.6. (1) The method of 6.5(i) for constructing end extensions seems to be well known.

(2) For any theory $S$ in any language it is trivial to see that a sentence $\phi$ is (universal sentences)-con over $S$ iff every model of $S$ can be extended to model $S + \phi$. Accordingly, the characterization of $\Pi_1^0$-con provided in 6.6 is an immediate consequence of [M]. The corresponding fact about $\Sigma_1^0$ does *not* follow from [M] by "trivial model theory".

(3) Model theoretic equivalents of $\Pi_n^0(L)$-con and $\Sigma_n^0(L)$-con are immediate: every $T$-model of $T$ can be end extended to (is an end extension of) a model of $T + \phi$ which is elementary for $\Sigma_{n-1}^0(L)$ sentences.

The proofs of corresponding theorems for set theories are identical with (or dual to) the proof of 6.5(ii). (Notice that every model of set theory has an elementary extension which is a non-$\omega$-model.) They are stated separately in order to emphasize: (i) that we consider only $\omega$-nonstandard models of $ZF$; (ii) that 6.5(i) alone has no cardinality restriction.

THEOREM 6.7. *Let $T$ be a set theory in a countable language $L$ and $\phi$ any sentence of $L$. Then,*

(i) *$\phi$ is $\Pi_1(L)$-con over $T \Leftrightarrow$ every countable $\omega$-nonstandard $T$-model of $T$ can be end extended to a model of $T + \phi$.*

(ii) *$\phi$ is $\Sigma_1(L)$-con over $T \Leftrightarrow$ every countable $\omega$-nonstandard $T$-model of $T$ is an end extension of a model of $T + \phi$.*

The necessity of restricting attention ot $\omega$-nonstandard models of set theory is shown by the following example.

THEOREM 6.8. *There is a $\Sigma_1$ sentence $\phi$ such that $\phi$ is $\Pi_1$-con over ZF but ZF proves that the theory ZF + $\phi$ has no $\omega$-model.*

PROOF. Let $\phi$ be a $\Sigma_1$ sentence such that $ZF$ proves $\phi \leftrightarrow \exists$ finite $x \subseteq ZF(\neg \operatorname{Con}_\infty(x \cup \{\ulcorner \phi \urcorner\}))$. §4 contains the definition of $\operatorname{Con}_\infty$. We will show first that $ZF + \phi$ cannot have an $\omega$-model.

Let $LST$ be the (finitary) language of set theory. Consider the following metamathematical fact (proof deferred):

(i) For any sentence $\phi$ of $LST$, $ZF \vdash \phi \to \operatorname{Con}_\infty(\{\phi\})$.

The proof of (i) can be arithmetized, to give:

(ii) $PA \vdash$ For any $x$ in $LST$, $\operatorname{Thm}_{ZF}("x \to \operatorname{Con}_\infty(\{x\})")$.

Applying (ii) inside $ZF$, we get a theorem of $ZF$ about models of $ZF$:

(iii) $ZF \vdash$ For any $x$ in $LST$, if $\mathfrak{M} \vDash ZF$, then $\mathfrak{M} \vDash "x \to \operatorname{Con}_\infty(\{x\})"$.

Now reason in $ZF$: Suppose that $\mathfrak{M}$ is an $\omega$-model of $ZF + \phi$. Produce some $b \in |\mathfrak{M}|$ such that $\mathfrak{M}$ satisfies:

$$\mathbf{b} \text{ is finite } \bigwedge \mathbf{b} \subseteq ZF \wedge (*) \neg \operatorname{Con}_\infty(\mathbf{b} \cup \{\ulcorner \phi \urcorner\}).$$

Because $\mathfrak{M}$ is an $\omega$-model, $b$ is in the standard part of $\mathfrak{M}$ and really a finite subset of $ZF$. So $\bigwedge\!\!\bigwedge b \wedge \ulcorner \phi \urcorner$ is a sentence of $LST$ satisfied by $\mathfrak{M}$. And by (iii), $\mathfrak{M} \vDash \operatorname{Con}_\infty(\{\bigwedge\!\!\bigwedge \mathbf{b} \wedge \ulcorner \phi \urcorner\})$–contradicting (*).

Proof of (i). We will prove the contrapositive. Suppose $\neg \operatorname{Con}_\infty(\{\phi\})$. Produce a set $p$ such that $\operatorname{Proof}_\infty(0; \ulcorner \neg\phi \urcorner, p)$. By the reflection principle there is a transitive set $M$ such that $p \in M$ and $(*)$ $\phi \leftrightarrow (\phi)^M$. The axioms from which the proof $p$ proceeds are all of form $\lambda_x$ for $x \in M$. (That each such $x$ is in $M$ is a consequence of the usual procedure for coding proofs; alternately, the reflection principle guarantees that we can choose $M$ large enough to contain them all.) So $M$ satisfies each axiom of $p$. Since $L_\infty\omega$ is sound, $M$ satisfies the conclusion of $p$: i.e., $(\neg\phi)^M$. (Actually, what we know is $\langle M, \in \rangle \vDash \ulcorner \neg\phi \urcorner$, which is provably equivalent to $(\neg\phi)^M$.) By (*), then, $\neg\phi$.

What is left is to show that $\phi$ is $\Pi_1$-con over $ZF$. We will use 4.1. Let $F_0$ be a fairly large finite chunk of $ZF$ (how much we will need will become evident soon). Let $F_1$ be any finite subset of $ZF$. We want to show that $ZF$ proves $\operatorname{Con}_\infty(F_1 + \phi)$. Let $F = F_0 \cup F_1$. Then $F + \neg\operatorname{Con}_\infty(F)$ is a theory at least as strong as $F_1 + \phi$, so it will suffice to prove $\operatorname{Con}_\infty(F + \neg\operatorname{Con}_\infty(F))$. To do that imitate the proof of Gödel's 2nd Incompleteness Theorem. The key step is justified by the following fact:

There is a finitely axiomatizable subtheory $F$ of $ZF$ such that for any $\Sigma_1$ sentence $\sigma$,

$F \vdash \sigma \to \mathrm{Thm}_\infty(\ulcorner F \urcorner;\ \ulcorner \sigma \urcorner)$; (cf. analogous theorem for $N$). (A similar argument is carried out in detail in [**Krivine-McAloon**].) $\square$

FINAL QUESTIONS AND COMMENTS. (1) In Theorems 6.5(ii) and 6.7 can the restriction to countable models be dropped?[2]

(2) Let $EE = \{\phi|$ every countable model of $ZF$ can be end extended to a model of $ZF + \phi\}$. Every sentence whose consistency is provable by forcing is in $EE$–because generic extensions (even of nonstandard models) are end extensions. It is shown in [**Barwise, 1**] or [**Barwise, 2**] that $V = L$ is an element of $EE$. (Note: That $V = L$ is both $\Pi_1$ and $\Sigma_1$-con over $ZF$ follows from the Shoenfield Absoluteness Lemma, which can be formulated as follows: For every $\Sigma_1$ *sentence* $\phi$, $ZF \vdash \phi \leftrightarrow (\phi)^L$.) It is mildly interesting to note that for Barwise's argument the $\omega$-nonstandard models were the hard ones to deal with–while from the present point of view the $\omega$-nonstandard models are the understandable ones. How understandable is $EE$? Is it a complete $\Pi_2^1$ set?

### REFERENCES

[**Barwise, 1**] *Admissible sets and structures*, Springer-Verlag, Berlin, 1975.

[**Barwise, 2**] *Infinitary methods in the model theory of set theory*, Logic Colloquium, 1969 (R. O. Gandy and M. Yates, Editors), North Holland, Amsterdam, 1971, pp. 53–66.

[**Hajkova-Hajek**] *On interpretability in theories containing arithmetic*, Fund. Math. **76** (1972), 131–137.

[**Hajek, 1**] *On interpretability in set theories*, Comment. Math. Univ. Carolinae **12** (1971), 73–79.

[**Hajek, 2**] *On interpretability in set theories*. II, Comment. Math. Univ. Carolinae **13** (1972), 445–455.

[**Feferman**] *Arithmetization of metamathematics in a general setting*, Fund. Math. **49** (1960), 35–92.

[**Friedman**] *Countable models of set theories*, Proc. Cambridge Summer School in Math Logic, 1971 (A. R. D. Mathis and H. Rogers, Editors), Lectures Notes in Math., vol. 337, Springer, Berlin and New York, 1973, pp. 539–573.

[**Keisler**] *Model theory for infinitary logic*, North Holland, Amsterdam, 1971.

[**Kreisel**] *On weak completeness of intuitionistic predicate logic*, J. Symbolic Logic **27** (1962), 139–158.

[**Krivine-McAloon**] *Some true unprovable formulas for set theory*, Bertrand Russell Memorial Logic Conference, Leeds, 1973, pp. 332–341.

[**Levy**] *A hierarchy of formulas in set theory*, Mem. Amer. Math. Soc. No. 57 (1965).

[**Löb**] *Solution of a problem of Leon Henkin*, J. Symbolic Logic **20** (1955), 115–118.

[**Macintyre-Simmons**] *Gödel's diagonalization technique and related properties of theories*, Colloq. Math. **28** (1973), 165–180.

[**Matijasevic**] *Diophantine representation of recursively enumerable predicates*, Proc. Second Scandinavian Logic Symposium (J. E. Fenstad, Editor), North Holland, Amsterdam, 1971.

---

[2]J. Quinsey has pointed out that the following elementary argument establishes 6.7(ii) directly without the restriction to countable models: If $\phi$ is $\Sigma_1(L)$-con over $T$ then for each finite subset $F$ of $T$, $T$ proves that there is a transitive model of $F + \phi$. Inside an $\omega$-nonstandard $T$-model $\mathfrak{M}$ we can now by overspill produce a transitive (in the sense of $\mathfrak{M}$) model of $T + \phi$. He has also pointed out that a similar argument (using Kripke's notion of fulfillability) removes the cardinality restriction from 6.5(ii).

D. GUASPARI

**[Montague]** *Semantical closure and nonfinite axiomatizability*, in Infinitistic Methods, Pergamon Press, New York, 1959, pp. 45–69.

**[Shoenfield]** *Mathematical logic*, Addison-Wesley, Reading, Mass., 1967.

**[Solovay]** *On interpretability in set theories* (to appear).

DEPARTMENT OF MATHEMATICS, TEXAS TECHNICAL UNIVERSITY, LUBBOCK, TEXAS 79409

*Current address*: St. John's College, Annapolis, Maryland 21404