Petr Hájek
Partial conservativity revisited

# PARTIAL CONSERVATIVITY REVISITED

## Petr HÁJEK

**Abstract:** We study the notions of partial conservativity and interpretability over fragments of Peano arithmetic.

**Key words:** Partial conservativity, interpretability, first order arithmetic.

**Classification:** 03F30, 03F25, 03D35

--------------------------------------------------------------------------

**Introduction.** Let S,T be theories, $S \subseteq T$, and let $\Gamma$ be a class of formulas of T. S is $\Gamma$ -conservative over T if for each $\psi \in \Gamma$ , provability of $\psi$ in S implies provability of $\psi$ in T. In particular, if $\varphi$ is a T formula then $\varphi$ is said to be $\Gamma$-conservative over T if the theory $(T+\varphi)$ is T-conservative over T. "Partial conservativity" means " $\Gamma$ -conservativity for some $\Gamma$ ". The first nontrivial example of partial conservativity was exhibited by Kreisel [62] ; the first systematic paper on partial conservativity is Guaspari [76] (containing also results by Solovay). Then various people contributed, among them Lindström, Smoryński and the present author (see references). A recent work is Bennet [86]. These papers typically discuss partial conservativity over theories in the language of arithmetic containing PA (Peano arithmetic); but equally typically, contain a remark saying that the assumptions on the underlying theory are in fact too strong and a weaker theory would suffice. One of the advantages of former theories is that for such a theory T, interpretability of $(T+\varphi)$ in T is equivalent to $\Pi_1$-conservativity of $(T+\varphi)$ over T and thus results on interpretability can be obtained as corollaries. Here we work systematically with theories containing $I\Sigma_1$ (i.e. arithmetic with induction restricted to $\Sigma_1$-formulas) but having possibly a richer language. One basic difference of fragments $I\Sigma_n$ from the whole PA is that they are finitely axiomatizable. As we shall see, this does not affect much proper ties of partial conservativity but does affect properties of interpretability. For positive results on interpretability we shall heavily use a result

due to Pudlák (formulated below) and its strengthening. A second obstacle is the fact that over a weak theory like $I\Sigma_1$, $\Sigma_n$ formulas are not closed under bounded universal quantification; but it turns out that this can be easily overcome. The result of our investigation is a systematic treatment of partial conservativity over theories containing $I\Sigma_1$ and of its relation to interpretability. We shall give only very sketchy hints on proofs; full proofs are contained in a Czech typescript not for publication; the proofs will also be incorporated in a forthcoming book. The paper is organized as follows: § 1 contains preliminaries, § 2 discusses prominent examples (Gödel's and Rosser's formulas), § 3 presents general theorems on partial conservativity, § 4 contains applications to interpretability and § 5 elaborates a classification of independent $\Sigma_1$ sentences.

**§ 1. Preliminaries.** We say "interpretation" meaning "relative interpretation with absolute equality" (cf. Tarski, Mostowski, Robinson [53]). First recall the old result on interpretability for PA and similar theories.

**1.1. Theorem** (Orey, Hájek, Guaspari). Let S, T be axiomatized theories in the language of PA and let $S \supseteq T \supseteq PA$. Then the following are equivalent:

(i)  S is interpretable in T,

(ii)  S is $\Pi_1$-conservative over T,

(iii)  for each k, $T \vdash Con_{S \restriction k}$,

(iv)  there is a binumeration $\beta$ of S in T such that $T \vdash Con_\beta$ .

(Similarly for $S \supseteq T \supseteq ZF$ in the language of ZF and in general for $S \supseteq T$ where T contains PA and proves induction for all T-formulas.)

See Orey [61], Hájek [72], Guaspari [76].

The following lemma is easy but basic for our considerations.

**1.2. Lemma.** Let $\varphi(x)$ be a $\Sigma_n$-formula (whose free variables are x and possibly others). There is a $\Sigma_n$-formula $\psi(y)$ such that

(1)  $(\forall k)\ I\Sigma_1 \vdash \psi(k) \equiv (\forall x \leqslant k)\varphi(x)$,

(2)  $I\Sigma_1 \vdash \psi(y) \rightarrow (\forall x \leqslant y)\varphi(x)$.

**Proof.** Trivial for $n=0$. For $n \geqslant 1$ and $\varphi(x) \rightleftharpoons (\exists u)\alpha(x,u)$ let $\psi(y)$ be $(\exists s)(Seq(s) \& (\forall x \leqslant y)\alpha(x,(s)_x)$.

**1.3. Notation.** The formula $\psi(y)$ from the previous lemma will be denoted $[(\forall x \leqslant y)\varphi(x)]^{*,\Sigma,n}$ or simply $[(\forall x \leqslant y)\varphi(x)]^*$ . Dually, for each $\Pi_n$-formula $\varphi(x)$ we have a $\Pi_n$-formula $\psi(y)$ (denoted by $[(\exists x \leqslant y)\varphi(x)]^{*,\Pi,n}$) such that

$$(\forall k)I\Sigma_1 \vdash \gamma(k) \equiv (\exists x \leq k)\,\varphi(x),$$
$$I\Sigma_1 \vdash (\exists x \leq y)\,\varphi(x) \rightarrow \gamma(y).$$

**1.4. Discussion and definition** (proofs from true formulas; herbrandian proofs). It is well known that in $I\Sigma_1$ we may define partial truth predicates: for each n we have a $\Sigma_n$ truth definition for all $\Sigma_n$ sentences satisfying the usual Tarski's conditions and similarly for $\pi$ instead of $\Sigma$ . If $\Gamma$ is $\Sigma_n$ or $\pi_n$ then $\mathrm{Tr}_\Gamma(x)$ means the corresponding truth predicate; $I\Sigma_1 \vdash \varphi \equiv \mathrm{Tr}(\overline{\varphi})$ for each $\varphi \in \Gamma$ . We make the following definition (y is a proof of x from a true $\Gamma$-formula):

$$\mathrm{Prf}_\Gamma(x,y) \equiv (\exists z \leq y)(z \in \Gamma \,\&\, \mathrm{Tr}_\Gamma(z) \,\&\, \mathrm{Prf}(z \rightarrow x, y)).$$

Clearly, if $\Gamma$ is $\Sigma_n$ then $\mathrm{Prf}_\Gamma$ is $\Sigma_n$ in $I\Sigma_1$: a trivial transformation shows that if $\Gamma$ is $\pi_n$ then $\mathrm{Prf}_\Gamma$ is $\pi_n$ in $I\Sigma_1$. Here Prf is the proof predicate for a given theory defined in "Hilbert style". We shall also use the "Herbrand style" proof predicate HPrf investigated by Pudlák [85], but only for finitely axiomatized theories T. In words, a herbrandian proof of $\varphi$ in T is a propositional proof of a disjunction of instances of the quantifier-free part of He($\Lambda T \rightarrow \varphi$ ), where, for any $\gamma$ , He($\gamma$) is the purely existential Herbrand form of $\gamma$ .
Similarly we define $\mathrm{HPrf}_\Gamma$ for $\Gamma$ being $\Sigma_n$ or $\pi_n$.

Pr(x) and HPr(x) are provability predicates (usual and herbrandian); clearly, $I\Sigma_1 \vdash (\forall x)(Pr(x) \equiv HPr(x))$. But this equivalence does not relativize to definable cuts - see next definition.

**1.5. Definition.** (1) Let $T \supseteq Q$ (Q is Robinson's arithmetic). A formula J(x) with one free variable <u>defines a cut in</u> T if T proves the following:
$$J(0) \,\&\, (\forall x)(J(x) \rightarrow J(x+1)) \,\&\, (\forall x,y)(y < x \,\&\, J(x) \rightarrow J(y)).$$
Note that if T is PA then $T \vdash (\forall x)J(x)$ for each such J; but this is not the case for any fragment $I\Sigma_n$ - there are cuts J such that $(\forall x)J(x)$ is unprovable.
(2) A theory $T \supseteq Q$ is <u>sequential</u> if it has coding of finite sequences of arbitrary objects, i.e. has a predicate Seq and functions $(s)_x$, lh(s) such that the following is provable in T: there is an empty sequence (of length 0), the length of each sequence is a number and for each sequence s of a length x and an arbitrary object z there is a sequence s´ of length x+1 prolonging s by z. (See Pudlák [85]).
(3) $ACA_0$ is the usual (fully) conservative second order extension of PA, i.e. $ACA_0$ has two sorts of variables (numbers and sets), axioms of PA for numbers with the induction scheme replaced by a single axiom

$0 \in X \& (\forall x)(x \in X \rightarrow x+1 \in X)$ and with comprehension for all formulas containing no set quantifiers. (Note that $ACA_0$ is sequential; $ACA_0$ is related to PA as GB to ZF.)

**1.6. Theorem** (Pudlák [85]). (1) Let T be a consistent sequential theory, S a finitely axiomatized theory. S is interpretable in T iff there is a definable cut J in T such that $T \vdash HCon_J^J(S)$, i.e. T proves that J does not contain any herbrandian proof of the contradiction from S.

(2) In particular, if T is consistent, finitely axiomatizable and sequential then there is a definable cut J in T such that $T \vdash HCon^J(T)$.

(3) On the contrary, for such a T there is no cut J such that $T \vdash Con^J(T)$.

Point (2) can be strengthened as follows:

**1.7. Theorem.** Let T be finitely axiomatizable, sequential and let $T \vdash I\Sigma_1$. Then there is a cut J in T such that

$$T \vdash (\forall u)(Tr_{\Sigma_1}(u) \rightarrow HCon^J((T+u))).$$

(Proof by inspection of Pudlák [85] - tedious.)

## § 2. Some prominent examples

2.1. We shall investigate the properties of Gödel´s consistency formula $Con_T$ and Rosser´s formula $\varphi_T$, for T being either an extension of PA in the same language or a finitely axiomatized extension of $I\Sigma_1$. (In the former case we assume a fixed $\Delta_1$ binumeration of the axioms to be given; in the latter we work with the natural binumeration just listing the axioms.) Both formulas may be constructed using either Prf or HPrf; for $Con_T$ this is immaterial (see above), but for the Rosser´s formula (which we assume in $\Sigma_1$ form, i.e. saying "there is a proof y of my negation such that no $z < y$ is a proof of me") there may be differences: some results below hold only for the Rosser´s formula based on HPrf, say, the H-Rosser formula.

**2.2. Convention:** If T is a theory and $\varphi$ a formula of T we say that $\varphi$ is interpretable in T meaning that the theory $(T+ \varphi)$ is interpretable in T.

**2.3. Theorem.** Let $T \geq I\Sigma_1$ be consistent (and axiomatizable).
(1) Gödel´s formula Con is not interpretable in T and its negation $\neg Con$ is interpretable in T.
(2) $\neg Con$ is $\Pi_1$ conservative; Con is $\Sigma_1$-conservative iff T is $\Sigma_1$-sound (i.e. each provable $\Sigma_1$-formula is true in N).
(3) Rosser´s formula $\varphi$ is $\Pi_1$-nonconservative; $\neg \varphi$ is $\Sigma_1$-conservative

iff T is $\Sigma_1$-sound. (The same holds for the H-Rosser formula.)

(4) If T is sequential and has induction for all formulas then neither $\varphi$ nor $\neg\varphi$ is interpretable. (The same for H-Rosser formula.)

(5) But if T is sequential and finitely axiomatizable and $\varphi$ is the H-Rosser formula then both $\varphi$ and $\neg\varphi$ are interpretable.

**Comments on proofs.** (1) For the first claim see Feferman [60] and Švejdar [78]; for the second see Feferman [60] where $T \supseteq PA$ is assumed. The same result can be proved:

(a) for $T \supseteq B\Sigma_2$ using Low Basis Theorem (see Clote [83]) and the corresponding Low Arithmetized Completeness Theorem,

(b) for $T \supseteq I\Sigma_1$ finitely axiomatized using Second Gödel´s incompleteness theorem and Pudlák´s theorem 1.6 and

(c) in <u>full generality</u> i.e. for any $T \supseteq T\Sigma_1$ using a version of Low Basis Theorem in $I\Sigma_1$ (see Hájek and Kučera [∞] ).

(2) The first claim is due to Kreisel [68] and is the first example of non-trivial partial conservativity. Checking for $T \supseteq I\Sigma_1$ is immediate. The second claim is due to Smoryński [80].

(3) First claim is due to Kreisel [62], second to Švejdar (unpublished)

(4) Follows easily from (3) and from 1.1.

(5) Pudlák´s theorem 1.6 gives an interpretation of $(T+\varphi)$ in $(T+\neg\varphi)$ and vice versa.

**Problem:** does (5) also hold for the usual (non-herbrandian) Rosser formula?

**2.4. Corollary.** If T is sequential and finitely axiomatized then there is a $\Pi_1$ formula $\varphi$ such that both $\varphi$ and $\neg\varphi$ are interpretable in T.(For PA and similar theories there is no such $\varphi$ .) In particular, if T is $ACA_0$ then $(ACA_0+\varphi)$ is interpretable in $ACA_0$ but $(PA+\varphi)$ is not interpretable in PA. Similarly for GB and ZF instead of $ACA_0$ and PA. First example of such a $\varphi$ was constructed by Solovay (unpublished, cf. Hájek [81]).

**§ 3. General theorems on partial conservativity.** We shall present several theorems on partial conservativity. Their proofs use various generalizations of Rosser´s formula. In the whole section $T \supseteq I\Sigma_1$ is a fixed consistent axiomatized theory (and we assume a $\Delta_1$ binumeration of T to be fixed).

**3.1. Notation.** (1) Let $\alpha(u)$, $\beta(u)$ be two T-formulas, let $\Delta$ be $(\exists u)\alpha(u)$ and $\nabla$ be $(\exists u)\beta(u)$. ($\alpha, \beta$ may contain parameters.) Following

Guaspari we denote by $\Delta \prec \nabla$ the formula

$$(\exists u)(\alpha(u) \& (\forall v \le u)\neg \beta(v))$$

(there is a witness for $\alpha$ less than each witness for $\beta$ ).

(2) If $\beta$ is $\pi_n$ and $\beta'$ is the $\Sigma_n$ formula naturally equivalent to $\neg\beta$ then $\Delta \prec^* \nabla$ will denote the formula

$$(\exists u)(\alpha(u) \& [(\forall v \le u)\ \beta'(v)]^{*,\Sigma,n})$$

(assuming that n is clear from the context). Similarly, $\Delta \preceq^* \nabla$ is $(\exists u)(\alpha(u) \& [(\forall v < u)\ \beta'(v)]^{*,\Sigma,n})$.

3.2. **Remark.** We shall investigate selfreferential formulas satisfying

$$T \vdash \xi \equiv \Delta(\neg \overline{\xi}) \prec^* \nabla(\overline{\xi})$$

or, more generally, for each k,

$$T \vdash \xi(\overline{k}) \equiv \Delta(\neg \overline{\xi},\overline{k}) \prec^* \nabla(\xi,\overline{k}).$$

Observe that if $\alpha$ is $\Sigma_n$ and $\beta$ is $\pi_n$ then $\xi$ is $\Sigma_n$ in T.

(2) If $\Delta_i$ is $(\exists u)\alpha_i(u)$ and $\nabla_i$ is $(\exists v)\beta_i(v)$ (i=1,2) then $(\Delta_1 \vee \Delta_2) \prec (\nabla_1 \vee \nabla_2)$ means the formula saying "there is a witness for $\alpha_1 \vee \alpha_2$ less than each witness for $\beta_1 \vee \beta_2$"; similarly for $\prec^*$ instead of $\prec$ .

3.3. **Definition.** (1) $\varphi$ is hereditarily $\Gamma$-conservative over T if, for each $T_0$ such that $I\Sigma_1 \subseteq T_0 \subseteq T$, $\varphi$ is $\Gamma$-conservative over $T_0$.

(2) $\varphi$ is doubly $\Gamma$-conservative over T if $\varphi$ is $\Gamma$-conservative over T and $\neg\varphi$ is $\check{\Gamma}$-conservative over T (where $\check{\Gamma}$ is the dual class of $\Gamma$ ).

We shall now formulate three general theorems on partial conservativity.

3.4. **Theorem.** For each $n \ge 1$ there is (1) a hereditarily $\pi_n$-conservative $\Sigma_n$-sentence, (2) a hereditarily $\Sigma_n$-conservative $\pi_n$-sentence, (3) a doubly $\pi_n$-conservative $\Sigma_n$-sentence (its negation is thus a doubly $\Sigma_n$conservative $\pi_n$-sentence).

Examples ( $\Gamma$ is $\Sigma_n$, $\Lambda$ is $\pi_n$):

(1) $\xi$ such that $I\Sigma_1 \vdash \xi \equiv Pr_\Gamma(\neg \overline{\xi}) \prec^* Pr(\overline{\xi})$,

(2) $(\neg \xi)$ such that $I\Sigma_1 \vdash \xi \vdash Pr(\neg \overline{\xi}) \prec^* Pr_\Lambda(\overline{\xi})$,

(3) $\xi$ such that $I\Sigma_1 \vdash \xi \subseteq Pr_\Gamma(\neg \overline{\xi}) \prec^* Pr_\Lambda(\xi)$.

If T is $\Sigma_n$-sound we may take in (1)a $\xi$ such that $I\Sigma_1 \vdash \xi \equiv Pr_\Gamma(\neg \overline{\xi})$.

3.5. **Theorem** (on non-separability). Let $\Gamma$ be $\Sigma_n$ or $\pi_n$ ($n \ge 1$), Th be the set of all theorems of T, Consv($\Gamma$) and hConsv($\Gamma$) the set of all $\Gamma$-con-

- 684 -

servative and hereditarily $\Gamma$-conservative sentences respectively , NRef the set of all the sentences non-refutable in T. Then obviously

$$Th \subseteq hConsv(\Gamma) \subseteq Consv(\Gamma) \subseteq Consv(\Sigma_1) \cap Consv(\Pi_1) \subseteq Nref$$

and there is no set X such that

(1)  X is $\Delta_1$ and $Th \subseteq X \subseteq NRef$ (classical!),

(2)  X is $\Pi_1$ and $Th \subseteq X \subseteq Consv(\Gamma)$,

(3)  X is $\Sigma_1$ and $hConsv(\Gamma) \subseteq X \subseteq NRef$,

(4)  X is $\Sigma_2$, $\Gamma \supseteq \Sigma_1$ and $hConsv(\Gamma) \subseteq X \subseteq Consv(\Sigma_1)$,

(4´)  X is $\Sigma_2$, $\Gamma \supseteq \Pi_1$ and $hConsv(\Gamma) \subseteq X \subseteq Consv(\Pi_1)$.


**3.6. Theorem** ($\Pi_2$-completeness). For each $n \geq 1$ and $\Gamma = \Sigma_n$ or $\Pi_n$, both $Consv(\Gamma)$ and $hConsv(\Gamma)$ is $\Pi_2$-complete.


**3.7. Remark.** Theorem 3.4 was obtained for $T \supseteq PA$ by Guaspari and Solovay, see Guaspari [67]; their examples are more complicated than ours. Theorem 3.5:

(1) is very classical, (2) seems to be new. (3) was first proved in the particular case T=ZF, $\Gamma = \Pi_1$, interpretability instead of partial conservativity in Hájek [71]; Lindström [84] generaliz  for $T \supseteq PA$, the present generalization is mine. (4) is contained (implicitly) in Lindström [84] for $T \supseteq PA$.
Theorem 3.6:

For T=ZF and $\Gamma = \Pi_1$ (and interpretability) Solovay; his proof works for PA. For PA, $\Gamma = \Pi_n$ and Cons see Hájek [79], for PA, $\Gamma = \Sigma_n$ and Consv see Quinsey [81]. In full generality but for $T \supseteq PA$ see Lindström [84]; generalization to $T \supseteq I\Sigma_1$ is mine.

We shall present two general fixed point theorems that form the main means of proofs of the preceding theorems.


**3.8. Shepherdson-Smoryński's fixed point theorem.** Let $\Phi$ , $\Psi$ be $\Sigma_1$ formulas.

(1)  Let $I\Sigma_1 \vdash \xi \equiv [(Pr(\neg \overline{\xi}) \vee \Phi) \prec (Pr(\xi) \vee \Psi)]$. Then

(i)  $T \vdash \xi$ iff $N \models \Phi \prec \Psi$ iff $N \not\models \xi$;

(ii)  $T \vdash \neg \xi$ iff $N \models \Psi \preceq \Phi$.

(2)  More generally, let, for i=1,2, $T_i \supseteq I\Sigma_1$, let $Pr_i$ be the proof predicate based on a fixed $\Delta_1$ binumeration of $T_i$. Let

$$I\Sigma_1 \vdash \xi \equiv [(Pr_1(\neg \overline{\xi}) \vee Pr_2(\neg \xi) \vee \Phi) \prec (Pr_1(\xi) \vee Pr_2(\xi) \vee \Psi)]. \text{ Then}$$

(i)  $T_1 \vdash \xi$ iff $T_2 \vdash \xi$ iff $N \models \Phi \prec \Psi$ iff $N \models \xi$;

(ii)  $T_1 \vdash \neg \xi$  iff  $T_2 \vdash \neg \xi$  iff  $N \vDash \Psi \preccurlyeq \Phi$.

For proofs see Shepherdson [60], Smoryński [80].

3.9. **Lindström's fixed point theorem.** (Let T be as above.)

(1)  Let  $\chi(y)$  be  $\Sigma_n$  and let

$$I\Sigma_1 \vdash \xi \equiv Pr_{\Sigma_n}(\neg \widehat{\xi}) \preccurlyeq^* (\exists y) \neg \chi(y). \text{ Then}$$

(i)  for each m,  $(T+\xi) \vdash \chi(m)$,

(ii)  for each  $I\Sigma_1 \subseteq T_0 \subseteq T$  and each  $\Pi_n$-sentence  $\pi$,  $(T_0+\xi) \vdash \pi$  implies  $T_0 + \{\chi(m) \mid m\} \vdash \pi$ .

(2)  Let  $\chi(y)$  be  $\Pi_n$  and let

$$I\Sigma_1 \vdash \xi \equiv (\exists y) \neg \chi(y) \preccurlyeq^* Pr_{\Pi_n}(\xi). \text{ Then}$$

(i)  for each m,  $(T+\neg \xi) \vdash \chi(m)$,

(ii)  for each  $I\Sigma_1 \subseteq T_0 \subseteq T$  and each  $\Sigma_n$-sentence  $\epsilon$,  $(T_0+\neg \xi) \vdash \sigma$  implies  $T_0 + \{\chi(m) \mid m\} \vdash \sigma$.

For  $T \supseteq PA$  see Lindström [84].

3.10. **Remark.** Both fixed point theorems may be <u>parametrized</u> (by replacing  $\xi$,  $\Phi$,  $\Psi$,  $\chi(y)$  by  $\xi(\overline{k})$,  $\Phi(\overline{k})$,  $\Psi(\overline{k})$,  $\chi(y,\overline{k})$  respectively); details are evident. Shepherdson-Smoryński theorem is used to prove 3.5 (1),(2) (and is very useful at many other occasions); Lindström's theorem is used for the rest of 3.5 and for 3.4. The proof of 3.6 uses Lindström's theorem and the following consequence of Shepherdson-Smoryński theorem: if X is a  $\Sigma_1$  set and  $I\Sigma_1 \subseteq T_0 \subseteq T_1$  then there is a  $\Sigma_1$  formula  $\sigma(x)$  and  $\Pi_1$  formula  $\pi(x)$  such that both  $\sigma$  and  $\pi$  numerate X both in  $T_0$  and in  $T_1$ .

**§ 4. Applications to interpretability.** It follows from 1.1 and 3.6 that if T is consistent sequential and has full induction then the set of all  $\varphi$  interpretable in T (i.e. such that the theory  $(T+\varphi)$  is interpretable in T) is a complete  $\Pi_2$  set. We focus our attention to finitely axiomatized theories containing  $I\Sigma_1$. Note that if. T is finitely axiomatizable then the set  $Intp_T$  of all  $\varphi$  interpretable in T is  $\Sigma_1$.

4.1. **Theorem.** Let T be finitely axiomatized,  $T \supseteq I\Sigma_1$  and let  $\Gamma = \Sigma_n$ or  $\Gamma = \Pi_n$  (n ≥1). Then  Consv( $\Gamma$ )-Int $p_T$  is non-empty and contains a  $\breve{\Gamma}$ -sentence.

4.2. **Remark.** In particular, let T be  $ACA_0$  and  $\Gamma = \Pi_1$. We get a  $\Sigma_1$-formula  $\varphi$  which is  $\Pi_1$-conservative over  $ACA_0$  but  $(ACA_0+\varphi)$  is not interpret-

able in $ACA_0$. Since $\varphi$ is $\Pi_1$-conservative over $ACA_0$ it is $\Pi_1$-conservative over PA and thus $(PA+\varphi)$ is interpretable in PA. Compare this with 2.4: we had there a $\Pi_1$ formula such that $(ACA_0+\varphi)$ is interpretable in $ACA_0$ but $(PA+\varphi)$ is not in PA. Similarly for GB and ZF.

The first example of a $\Sigma_2$ formula $\varphi$ such that $ZF+\varphi$ is interpretable in ZF but $(GB+\varphi)$ is not in GB was constructed in Hájek [71] under the assumption of $\omega$-consistency; this assumption was removed in Hájek and Hájková [72]. A $\Sigma_1$ formula of the desired properties was first constructed by Solovay.

4.3.  We shall investigate the situation more closely. In the sequel let T be a finitely axiomatized consistent theory containing $I\Sigma_1$. We focus our attention to independent $\Sigma_1$ sentences. First observe that each independent (nonprovable and nonrefutable) $\Sigma_1$ sentence is false and trivially it is $\Sigma_1$-nonconservative. For each such $\Sigma_1$ sentence $\sigma$ we may ask
- whether $\sigma$ is $\Pi_1$-conservative,
- whether $\sigma$ is interpretable,
- whether $\neg\sigma$ is interpretable.

The formula $\neg\sigma$ is $\Pi_1$ and hence $\Pi_1$-nonconservative; if T is $\Sigma_1$-sound then $\neg\sigma$ is $\Sigma_1$-conservative. For $\Sigma_1$-ill theories the $\Sigma_1$-conservativity of $\neg\sigma$ is a reasonable question but we shall not discuss it.

Our three questions admit eight combinations of answers, say, eight types of independent $\Sigma_1$ formulas.


4.4. **Theorem.**  For each type, there is an independent $\Sigma_1$ formula of that type.


4.5.  We shall describe examples in a rather uniform way. We shall use formulas HPr and $HPr_{\Sigma_1}$, i.e. herbrandian provability and herbrandian provability from a true $\Sigma_1$ sentence. Further we shall use a $\Sigma_1$ formula $Intp(x)$ formally expressing (i.e. numerating-in-$I\Sigma_1$) interpretability of a formula in T. All examples have the form of a self-referential formula $\xi$ such that

$$I\Sigma_1 \vdash \xi \equiv \Delta(\neg\overline{\xi})\prec^{\#}\nabla(\xi)$$

where $\Delta(x)$ has one of the following forms:

$$HPr(x),\ HPr_{\Sigma_1}(x),\ HPr(x)\vee Intp(x),\ HPr_{\Sigma_1}(x)\vee Intp(x).$$

$\nabla(x)$ has one of the following forms:

$$HPr(x),\ HPr(x)\vee Intp(x).$$

This gives eight possibilites that exactly give our eight examples. For several claims it is immaterial whether we use Pr or Hpr; but for some claims (using 1.7) it is not.

**4.6. Theorem.** (1) All eight examples are independent $\Sigma_1$ sentences.

(2) $\xi$ is interpretable iff $\nabla$ does not contain Intp(x).

(3) $\neg\xi$ is interpretable iff $\Delta$ does not contain Intp(x).

(4) $\xi$ is $\Pi_1$-conservative iff $\Delta$ contains $\text{HPr}_{\Sigma_1}(x)$.

The <u>proof</u> uses 1.6, 1.7, provability of Herbrand's theorem in $I\Sigma_1$ and the following very important theorem (more precisely, its corollary 4.8):

**4.7. Second Lindström's fixed point theorem.** Let $T\supseteq I\Sigma_1$, $n,m\geq 0$, $\Gamma = \Sigma_n$, $\Lambda=\Pi_m$. Let $\varphi(x,y)$ be a $\Gamma$-formula, $\Theta(x,y)$ a $\Lambda$-formula. Let

$$T \vdash \xi \equiv [(\text{Pr}_n(\neg\overline{\xi})\vee(\exists y)\varphi(\neg\overline{\xi},y))\preceq^{*}(\text{Pr}_\Lambda(\overline{\xi})\vee(\exists u)\Theta(\overline{\xi},u))].$$

Then

(1) for each m, $(T+\xi)\vdash(\exists y\leq\overline{m})\Theta(\overline{\xi},y)\rightarrow(\exists y\leq\overline{m})\varphi(\neg\xi,y)$,

(2) for each m, $(T+\neg\xi)\vdash(\exists y\leq\overline{m})\varphi(\neg\overline{\xi},y)\rightarrow(\exists y\leq\overline{m})\Theta(\xi,y)$,

(3) each $\check{\Gamma}$-sentence provable in $(T+\xi)$ is provable in $T+\{\neg\Theta(\overline{\xi},\overline{m})|m\}$,

(4) each $\check{\Gamma}$-sentence provable in $(T+\neg\xi)$ is provable in $T+\{\neg\varphi(\overline{\xi},\overline{m})|m\}$.

(5) The above remains true if Pr is replaced by HPr in all occurences.

**4.8. Corollary.** If $(\exists y)\Theta(x,y)$ defines a set $X\subseteq\text{NRef}$ (of non-refutable sentences) and $(\exists y)\varphi(x,y)$ defines a set Y of non-refutable sentences then $\xi$ satisfies the following:

$\xi$ is $\Pi_n$-conservative, $\neg\xi$ is $\Sigma_n$-conservative, $\xi\notin X$, $(\neg\xi)\notin Y$.

**4.9. Remark.** Note that Lindström has 4.7 and 4.8 for $T\supseteq PA$. The generalization of constructions of partially conservative sentences presented up to now (from $T\supseteq PA$ to $T\supseteq I\Sigma_1$) could make the reader think that everything generalizes smoothly. It is indeed remarkable that in most constructions the replacement of $\preceq$ by $\preceq^{*}$ works. But there are some things that do not generalize: for example, the implication

$$(\Delta\vee\nabla)\rightarrow((\Delta\prec\nabla)\vee(\nabla\preceq\Delta))$$

(provable in PA) does not immediately generalize to $I\Sigma_1$ using $\preceq^{*}$. Another example is in Lindström [79]: If T=PA and $X\subseteq\text{NRef}$ then there is a $\Delta_1$ binumeration $\gamma$ of PA such that $(PA+\neg\text{Con}_\gamma)$ is interpretable in PA but $(ACA_0+\neg\text{Con}_\gamma)$ is not interpretable in $ACA_0$. It is not clear how to generalize this from $ACA_0$ to any finitely axiomatizable $T\supseteq I\Sigma_1$.

Lindström's papers contain (for $T \supseteq PA$) various stronger and more detailed theorems on partial conservativity. It has been our task to demonstrate the possibility and ways of generalization to $T \supseteq I\Sigma_1$ on the most important theorems rather than to cover everything.

## References

C. BENNET [86]: On some orderings of extensions of arithmetic (thesis), University of Göteborg 1986.

P. CLOTE [83]: Partition relations in arithmetic, Proc. Sixth Latin-American Symp. on Math. Logic, Venezuela 1983.

P. CLOTE [85]: Applications of the low-basis theorem in arithmetic, Proceedings of Recursion-theory week at Oberwolfach, Springer 1985.

S. FEFERMAN [60]: Arithmetization of metamathematics in a general setting, Fund. Math. 49(1960), 35-92.

D. GUASPARI [79]: Partially conservative extensions of arithmetic, Trans. Amer. Math. Soc. 254(1979), 47-68.

P. HÁJEK [71]: On interpretability in set theories, Comment. Math. Univ. Carolinae 12(1971), 73-79.

P. HÁJEK [72]: On interpretability in set theories II, Comment. Math. Univ. Carolinae 13(1972), 445-455.

P. HÁJEK [79]: On partially conservative extensions of arithmetic, in: Logic Colloquium 78, North-Holland Pub. (1979), 225-234.

P. HÁJEK [81]: On interpretability in theories containing arithmetic II, Comment. Math. Univ. Carolinae 22(1981), 617-688.

P. HÁJEK [84]: On a new notion of partial conservativity, in: Logic Colloquium 83, Lect. Notes in Math. vol. 1104.

P. HÁJEK, M. HÁJKOVÁ [72]: On interpretability in theories containing arithmetic, Fund.Math. 76(1972), 131-137.

P. HÁJEK, A. KUČERA [∞] : On recursion theory in fragments of arithmetic (to appear).

G. KREISEL [62]: On weak completeness of intuitionistic predicate logic, J. Symb. Logic 27(1962), 139-158.

G. KREISEL [68]: A survey on proof theory, J. Symb. Logic 33(1968), 321-388.

P. LINDSTRÖM [79]: Some results on interpretability (Jensen, Mayoh, Motler, ed.), Proc. 5th Scandinavian Logic Symp., Aalborg Univ. Press 1979.

P. LINDSTRÖM [84]: On partially conservative sentences and interpretability, Proceedings Amer. Math. Soc. 91(1984), 436-443.

P. LINDSTRÖM [84a]: On faithful interpretability, in: Logic Colloquium 83, Lect. Notes in Math. vol. 1104, Springer-Verlag 1984.

S. OREY [61]: Relative interpretations, Z. Math. Logik und Grund. Math. 7 (1961), 146-153.

J. QUINSEY [81]: Sets of $\Sigma_k$-conservative sentences are $\Pi_2$-complete, J. Symb Logic 46(1981), 442 (abstract.

P. PUDLÁK [85]: Cuts, consistency statements and interpretations, J. Symb.
    Logic 50(1985), 423-441.

J.S. SHEPERDSON [60]: Representability of recursively enumerable sets in for-
    mal theories, Archiv f. math. Logik 5(1960), 119-127.

C. SMORYŃSKI [81]: Fifty years of self-reference, Notre Dame Journal of Form-
    al Logi 22(1981), 357-374.

C. SMORYŃSKI [81a]: Calculating self-referential sentences: Guaspari senten-
    ces of the first kind, J. Symb. Logic 46(1981), 329-344.

C. SMORYŃSKI [85]: Self-reference and modal logic, Springer-Verlag 1985.

V. ŠVEJDAR [81]: A sentence that is difficult to interpret, Comment. Math.
    Univ. Carolinae 22(1981), 661-666.

V. ŠVEJDAR [83]: Modal analysis of generalized Rosser sentences, J. Symb. Lo-
    gic 48(1983), 986-999.

A. TARSKI, A. MOSTOWSKI, R.M. ROBINSON [53]: Undecidable theories, North -Hol-
    land P.C. 1953.

Mathematical Institute of the Czechoslovak Academy of Sciences, Žitná 25,
115 67 Praha 1, Czechoslovakia