

Petr Hájek

Partial conservativity revisited

Commentationes Mathematicae Universitatis Carolinae, Vol. 28 (1987), No. 4, 679--690

Persistent URL: <http://dml.cz/dmlcz/106582>

Terms of use:

© Charles University in Prague, Faculty of Mathematics and Physics, 1987

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://project.dml.cz>

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 Γ be a class of formulas of T . S is Γ -conservative over T if for each $\psi \in \Gamma$, provability of ψ in S implies provability of ψ in T . In particular, if φ is a T formula then φ is said to be Γ -conservative over T if the theory $(T+\varphi)$ is T -conservative over T . "Partial conservativity" means " Γ -conservativity for some Γ ". The first nontrivial example of partial conservativity was exhibited by Kressel [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 Π_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 Σ_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 properties 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 Σ_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 $S \geq T \geq PA$. Then the following are equivalent:

- (i) S is interpretable in T ,
- (ii) S is Π_1 -conservative over T ,
- (iii) for each k , $T \vdash \text{Con}_{S|k}$,
- (iv) there is a binumeration β of S in T such that $T \vdash \text{Con}_\beta$.

(Similarly for $S \geq T \geq ZF$ in the language of ZF and in general for $S \geq T$ where T contains PA and proves induction for all I -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 Σ_n -formula (whose free variables are x and possibly others). There is a Σ_n -formula $\psi(y)$ such that

- (1) $(\forall k) I\Sigma_1 \vdash \psi(k) \Leftrightarrow (\forall x \leq k) \varphi(x)$,
- (2) $I\Sigma_1 \vdash \psi(y) \rightarrow (\forall x \leq y) \varphi(x)$.

Proof. Trivial for $n=0$. For $n \geq 1$ and $\varphi(x) \equiv (\exists u)\alpha(x, u)$ let $\psi(y)$ be $(\exists s)(\text{Seq}(s) \& (\forall x \leq y)\alpha(x, (s)_x))$.

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

$$\begin{aligned}
(\forall k) I\Sigma_1 \vdash \Psi(k) &\equiv (\exists x \leq k) \varphi(x), \\
I\Sigma_1 \vdash (\exists x \leq y) \varphi(x) &\rightarrow \Psi(y).
\end{aligned}$$

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 Σ_n truth definition for all Σ_n sentences satisfying the usual Tarski's conditions and similarly for Π instead of Σ . If Γ is Σ_n or Π_n then $Tr_\Gamma(x)$ means the corresponding truth predicate; $I\Sigma_1 \vdash \varphi \equiv Tr(\overline{\varphi})$ for each $\varphi \in \Gamma$. We make the following definition (y is a proof of x from a true Γ -formula):

$$Prf_\Gamma(x, y) \equiv (\exists z \leq y)(z \in \Gamma \& Tr_\Gamma(z) \& Prf(z \rightarrow x, y)).$$

Clearly, if Γ is Σ_n then Prf_Γ is Σ_n in $I\Sigma_1$: a trivial transformation shows that if Γ is Π_n then Prf_Γ is Π_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 φ in T is a propositional proof of a disjunction of instances of the quantifier-free part of $He(\wedge T \rightarrow \varphi)$, where, for any Ψ , $He(\Psi)$ is the purely existential Herbrand form of Ψ .

Similarly we define $HPrf_\Gamma$ for Γ being Σ_n or Π_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 defines a cut in 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 sequential 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 $T \vdash HCon^J(S)$, 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

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

(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 ρ_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 Δ_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 Σ_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 φ a formula of T we say that φ is interpretable in T meaning that the theory $(T + \varphi)$ is interpretable in T .

2.3. Theorem. Let $T \supseteq I\Sigma_1$ 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 Π_1 conservative; Con is Σ_1 -conservative iff T is Σ_1 -sound (i.e. each provable Σ_1 -formula is true in N).

(3) Rosser's formula ρ is Π_1 -nonconservative; $\neg \rho$ is Σ_1 -conservative

iff Γ is Σ_1 -sound. (The same holds for the H-Rosser formula.)

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

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

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

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

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

(c) in full generality i.e. for any $\Gamma \supseteq I\Sigma_1$ 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 $\Gamma \supseteq I\Sigma_1$ 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 $(\Gamma+\phi)$ in $(\Gamma+\neg\phi)$ and vice versa.

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

2.4. Corollary. If Γ is sequential and finitely axiomatized then there is a Π_1 -formula ϕ such that both ϕ and $\neg\phi$ are interpretable in Γ . (For PA and similar theories there is no such ϕ .) In particular, if Γ is ACA_0 then $(ACA_0+\phi)$ is interpretable in ACA_0 but $(PA+\phi)$ is not interpretable in PA. Similarly for GB and ZF instead of ACA_0 and PA. First example of such a ϕ 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 $\Gamma \supseteq I\Sigma_1$ is a fixed consistent axiomatized theory (and we assume a Δ_1 binumeration of Γ to be fixed).

3.1. Notation. (1) Let $\alpha(u)$, $\beta(u)$ be two Γ -formulas, let Δ be $(\exists u)\alpha(u)$ and ∇ be $(\exists u)\beta(u)$. (α, β may contain parameters.) Following

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

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

(there is a witness for α less than each witness for β).

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

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

(assuming that n is clear from the context). Similarly, $\Delta \prec^* \nabla$ is

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

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

$$\Gamma \vdash \xi \equiv \Delta(\neg \bar{\xi}) \prec^* \nabla(\bar{\xi})$$

or, more generally, for each k ,

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

Observe that if α is Σ_n and β is Π_n then ξ is Σ_n in Γ .

(2) If Δ_i is $(\exists u)\alpha_i(u)$ and ∇_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) \mathcal{G} is hereditarily Γ -conservative over Γ if, for each Γ_0 such that $I\Sigma_1 \subseteq \Gamma_0 \subseteq \Gamma$, \mathcal{G} is Γ -conservative over Γ_0 .

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

We shall now formulate three general theorems on partial conservativity.

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

Examples (Γ is Σ_n , Λ is Π_n):

- (1) ξ such that $I\Sigma_1 \vdash \xi \equiv Pr_\Gamma(\neg \bar{\xi}) \prec^* Pr_\Gamma(\bar{\xi})$,
- (2) $(\neg \bar{\xi})$ such that $I\Sigma_1 \vdash \xi \vdash Pr(\neg \bar{\xi}) \prec^* Pr_\Lambda(\bar{\xi})$,
- (3) ξ such that $I\Sigma_1 \vdash \xi \subseteq Pr_\Gamma(\neg \bar{\xi}) \prec^* Pr_\Lambda(\bar{\xi})$.

If Γ is Σ_n -sound we may take in (1) a ξ such that $I\Sigma_1 \vdash \xi \equiv Pr_\Gamma(\neg \bar{\xi})$.

3.5. **Theorem** (on non-separability). Let Γ be Σ_n or Π_n ($n \geq 1$). Th be the set of all theorems of Γ , $Consv(\Gamma)$ and $hConsv(\Gamma)$ the set of all Γ -con-

servative and hereditarily Γ -conservative sentences respectively, NRef the set of all the sentences non-refutable in \mathcal{T} . Then obviously

$$\text{Th} \subseteq \text{hConsv}(\Gamma) \subseteq \text{Consv}(\Gamma) \subseteq \text{Consv}(\Sigma_1) \cap \text{Consv}(\Pi_1) \subseteq \text{NRef}$$

and there is no set X such that

- (1) X is Δ_1 and $\text{Th} \subseteq X \subseteq \text{NRef}$ (classical!),
- (2) X is Π_1 and $\text{Th} \subseteq X \subseteq \text{Consv}(\Gamma)$,
- (3) X is Σ_1 and $\text{hConsv}(\Gamma) \subseteq X \subseteq \text{NRef}$,
- (4) X is Σ_2 , $\Gamma \supseteq \Sigma_1$ and $\text{hConsv}(\Gamma) \subseteq X \subseteq \text{Consv}(\Sigma_1)$,
- (4') X is Σ_2 , $\Gamma \supseteq \Pi_1$ and $\text{hConsv}(\Gamma) \subseteq X \subseteq \text{Consv}(\Pi_1)$.

3.6. **Theorem** (Π_2 -completeness). For each $n \geq 1$ and $\Gamma = \Sigma_n$ or Π_n , both $\text{Consv}(\Gamma)$ and $\text{hConsv}(\Gamma)$ is Π_2 -complete.

3.7. **Remark.** Theorem 3.4 was obtained for $\mathcal{T} \supseteq \text{PA}$ 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 $\mathcal{T} = \text{ZF}$, $\Gamma = \Pi_1$, interpretability instead of partial conservativity in Hájek [71]; Lindström [84] generaliz for $\mathcal{T} \supseteq \text{PA}$, the present generalizati on is mine. (4) is contained (implicitly) in Lindström [84] for $\mathcal{T} \supseteq \text{PA}$. **Theorem 3.6:**

For $\mathcal{T} = \text{ZF}$ 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 $\mathcal{T} \supseteq \text{PA}$ see Lindström [84]; generalization to $\mathcal{T} \supseteq \text{IS}_1$ is mine.

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

3.8. **Sheperdson-Smorczyński's fixed point theorem.** Let Φ, Ψ be Σ_1 formulas.

- (1) Let $\text{IS}_1 \vdash \xi \equiv [(\text{Pr}(\neg \bar{\xi}) \vee \Phi) \wedge (\text{Pr}(\xi) \vee \Psi)]$. Then
 - (i) $\mathcal{T} \vdash \xi$ iff $N \models \Phi \wedge \Psi$ iff $N \models \xi$;
 - (ii) $\mathcal{T} \vdash \neg \xi$ iff $N \models \Psi \wedge \neg \Phi$.
- (2) More generally, let, for $i=1,2$, $\mathcal{T}_i \supseteq \text{IS}_1$, let Pr_i be the proof predicate based on a fixed Δ_1 binumeration of \mathcal{T}_i . Let

$$\text{IS}_1 \vdash \xi \equiv [(\text{Pr}_1(\neg \bar{\xi}) \vee \text{Pr}_2(\neg \bar{\xi}) \vee \Phi) \wedge (\text{Pr}_1(\xi) \vee \text{Pr}_2(\xi) \vee \Psi)]. \text{ Then}$$

- (i) $\mathcal{T}_1 \vdash \xi$ iff $\mathcal{T}_2 \vdash \xi$ iff $N \models \Phi \wedge \Psi$ iff $N \models \xi$;

(ii) $T_1 \vdash \neg \xi$ iff $T_2 \vdash \neg \xi$ iff $N \models \Psi \not\sim \Phi$.

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

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

(1) Let $\chi(y)$ be Σ_n and let

$I\Sigma_1 \vdash \xi \equiv \text{Pr}_{\Sigma_n}(\neg \xi) \leftrightarrow (\exists y) \neg \chi(y)$. Then

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

(ii) for each $I\Sigma_1 \subseteq T_0 \subseteq T$ and each Π_n -sentence σ , $(T_0 + \xi) \vdash \sigma$ implies $T_0 + \{ \chi(m) \mid m \in \mathbb{N} \} \vdash \sigma$.

(2) Let $\chi(y)$ be Π_n and let

$I\Sigma_1 \vdash \xi \equiv (\exists y) \neg \chi(y) \leftrightarrow \text{Pr}_{\Pi_n}(\xi)$. Then

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

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

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

3.10. Remark. Both fixed point theorems may be parametrized (by replacing $\xi, \Phi, \Psi, \chi(y)$ by $\xi(\bar{k}), \Phi(\bar{k}), \Psi(\bar{k}), \chi(y, \bar{k})$ 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 Σ_1 set and $I\Sigma_1 \subseteq T_0 \subseteq T_1$ then there is a Σ_1 formula $\sigma(x)$ and Π_1 formula $\pi(x)$ such that both σ and π 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 φ interpretable in T (i.e. such that the theory $(T + \varphi)$ is interpretable in T) is a complete Π_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 φ interpretable in T is Σ_1 .

4.1. Theorem. Let T be finitely axiomatized, $T \supseteq I\Sigma_1$ and let $\Gamma = \Sigma_n$ or $\Gamma = \Pi_n$ ($n \geq 1$). Then $\text{Consv}(\Gamma)\text{-Intp}_T$ is non-empty and contains a $\check{\Gamma}$ -sentence.

4.2. Remark. In particular, let T be ACA_0 and $\Gamma = \Pi_1$. We get a Σ_1 -formula φ which is Π_1 -conservative over ACA_0 but $(\text{ACA}_0 + \varphi)$ is not interpret-

able in ACA_0 . Since φ is Π_1 -conservative over ACA_0 it is Π_1 -conservative over PA and thus $(PA+\varphi)$ is interpretable in PA. Compare this with 2.4: we had there a Π_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 Σ_2 formula φ such that $ZF+\varphi$ is interpretable in ZF but $(GB+\varphi)$ is not in GB was constructed in Hájek [71] under the assumption of ω -consistency; this assumption was removed in Hájek and Hájková [72]. A Σ_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 Σ_1 sentences. First observe that each independent (nonprovable and nonrefutable) Σ_1 sentence is false and trivially it is Σ_1 -nonconservative. For each such Σ_1 sentence σ we may ask

- whether σ is Π_1 -conservative,
- whether σ is interpretable,
- whether $\neg\sigma$ is interpretable.

The formula $\neg\sigma$ is Π_1 and hence Π_1 -nonconservative; if T is Σ_1 -sound then $\neg\sigma$ is Σ_1 -conservative. For Σ_1 -ill theories the Σ_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 Σ_1 formulas.

4.4. **Theorem.** For each type, there is an independent Σ_1 formula of that type.

4.5. We shall describe examples in a rather uniform way. We shall use formulas HPr and HPr_{Σ_1} , i.e. herbrandian provability and herbrandian provability from a true Σ_1 sentence. Further we shall use a Σ_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 ξ such that

$$I\Sigma_1 \vdash \xi \equiv \Delta(\neg \bar{\xi}) \rightarrow^* \nabla(\xi)$$

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 possibilities 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 Σ_1 sentences.
 (2) ξ is interpretable iff ∇ does not contain Intp(x).
 (3) $\neg \xi$ is interpretable iff Δ does not contain Intp(x).
 (4) ξ is Π_1 -conservative iff Δ contains $\text{HPr}_{\Sigma_1}(x)$.

The proof 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 $T \geq I\Sigma_1$, $n, m \geq 0$, $\Gamma = \Sigma_n$, $\Lambda = \Pi_m$. Let $\varphi(x, y)$ be a Γ -formula, $\theta(x, y)$ a Λ -formula. Let

$$T \vdash \xi = [(\text{Pr}_\Gamma(\neg \bar{\xi}) \vee (\exists y) \varphi(\neg \bar{\xi}, y)) \leftarrow^* (\text{Pr}_\Lambda(\bar{\xi}) \vee (\exists u) \theta(\bar{\xi}, u))].$$

Then

- (1) for each m , $(T + \xi) \vdash (\exists y \leq m) \theta(\bar{\xi}, y) \rightarrow (\exists y \leq m) \varphi(\neg \bar{\xi}, y)$,
- (2) for each m , $(T + \neg \xi) \vdash (\exists y \leq m) \varphi(\neg \bar{\xi}, y) \rightarrow (\exists y \leq m) \theta(\bar{\xi}, y)$,
- (3) each \bar{m} -sentence provable in $(T + \xi)$ is provable in $T + \neg \theta(\bar{\xi}, \bar{m}) \upharpoonright \bar{m}$,
- (4) each \bar{m} -sentence provable in $(T + \neg \xi)$ is provable in $T + \neg \varphi(\bar{\xi}, \bar{m}) \upharpoonright \bar{m}$.
- (5) The above remains true if Pr is replaced by HPr in all occurrences.

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 ξ satisfies the following:

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

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

$$(\Delta \vee \nabla) \rightarrow ((\Delta \leftarrow \nabla) \vee (\nabla \leftarrow \Delta))$$

(provable in PA) does not immediately generalize to $I\Sigma_1$ using \leftarrow^* . Another example is in Lindström [79]: If $T = \text{PA}$ and $X \subseteq \text{NRef}$ then there is a Δ_1 binumeration γ of PA such that $(\text{PA} + \neg \text{Con}_\gamma)$ is interpretable in PA but $(\text{ACA}_0 + \neg \text{Con}_\gamma)$ is not interpretable in ACA_0 . It is not clear how to generalize this from ACA_0 to any finitely axiomatizable $T \geq I\Sigma_1$.

Lindström's papers contain (for $T \geq PA$) various stronger and more detailed theorems on partial conservativity. It has been our task to demonstrate the possibility and ways of generalization to $T \geq \Sigma_1$ 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 Σ_k -conservative sentences are Π_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 formal theories, Archiv f. math. Logik 5(1960), 119-127.
- C. SMORYŃSKI [81]: Fifty years of self-reference, Notre Dame Journal of Formal Logi 22(1981), 357-374.
- C. SMORYŃSKI [81a]: Calculating self-referential sentences: Guaspari sentences 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. Logic 48(1983), 986-999.
- A. TARSKI, A. MOSTOWSKI, R.M. ROBINSON [53]: Undecidable theories, North-Holland P.C. 1953.

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

(Oblatum 20.7. 1987)