

## Peano's Smart Children: A Provability Logical Study of Systems with Built-in Consistency

ALBERT VISSER\*

**Abstract** The systems studied in this article prove the same theorems (from the "extensional" point of view) as Peano Arithmetic, but are equipped with a self-correction procedure. These systems prove their own consistency and thus escape Gödel's second theorem. Here, the provability logics of these systems are studied. An application of the results obtained turns out to be the solution to a problem of Orey on relative interpretability.

*1 Introduction* Consistency can be built into a system in various ways. The two best known constructions are Rosser's and Feferman's, both of which take a given formal system in the usual sense as initial data. Consider, for example, Peano Arithmetic (PA). A proof in the Peano System will count as a proof in the Rosser System based on PA, if there is no shorter Peano proof of the negation of its conclusion. The Feferman System can be described in various interesting ways, modulo provable equivalence in PA of the formulas defining the set of theorems. One such way is this: A proof in the Peano System will count as a proof in the Feferman System based on PA, if the finite set of arithmetical Peano axioms smaller than or equal to the largest arithmetical Peano axiom used in the proof is consistent.

The reasons such constructions occur in the literature are various:

- (i) They serve as counterexamples in the study of the relations between Gödel's first and second Incompleteness Theorems (see [4]).
- (ii) They serve as didactic examples in philosophical discussions, like the

---

\*I would like to thank Johan van Benthem, Dirk van Dalen, Karst Koymans, Henryk Kotlarski, and Fer-Jan de Vries for stimulating discussions. I am also grateful to Erik Krabbe for carefully reading parts of an early draft of this paper. And I am especially thankful to George Kreisel without whose interest and questions the paper probably would never have seen the light of day.

debate on intensionality in mathematics (e.g. [1]) and the discussion on the possible bearing of the Incompleteness Theorems on the Minds & Machines problem (see [10], [21], [3]).

- (iii) Rosser's construction is used to sharpen Gödel's first Incompleteness Theorem.
- (iv) Feferman's construction is an important tool in the study of relative interpretability (see [4], [13]).

The main objects of study in the present article are certain variants of both Rosser's and Feferman's constructions. My motivations are closely related to (i)–(iv) above:

- (a) There is much interest in the study of bimodal systems in the current literature on provability logic (see e.g. [12] and [16]). There are two directions of research: The pure study of arithmetical self-reference and the study of arithmetical self-reference as a tool for unifying self-referential arguments in arithmetic (see [6], chapter 7). In the first line of study one aims at characterizing the modal logic for a certain 'given' class of interpretations. There is no objection here to having 'few' interpretations and strong modal systems. In the second line one looks primarily for a modal system which is sound for as many interpretations as possible, but is still rich enough to carry out the proofs of the arithmetical arguments under study. The distinction between the two lines described here is not precisely that between pure and applied. The first line also has its typical applications: Solovay-style completeness results yield a powerful machinery for producing arithmetical sentences with rich but controlled properties. These sentences can be used to prove various incompleteness and other kinds of results (for an example, see Section 7 of this paper).

The contribution of this study lies along the first line. I provide an example of a rich modal logic of not too standard a sort, valid for two different arithmetical interpretations. This example can be used to test conjectures concerning the conditions for uniqueness and explicit definability of fixed points (see [15] for a discussion of these matters). Questions of uniqueness and explicit definability generalize the problem of the precise connection between the first and second Incompleteness Theorems; in this sense (a) generalizes (i). The logic can also be used to illustrate the point that one can simulate the results of intensional self-reference (as with Rosser's Theorem) to quite an extent by applying provably extensional self-reference—the cost being an increase in the complexity (modulo provable equivalence) of the sentences involved.

- (b) The modal derivability conditions are an improvement in the presentation of systems with built-in consistency in the discussions mentioned under (ii) above.
- (c) The methods developed have as a spin-off an application to relative interpretability: I answer a question posed by Orey for the case of PA.

As far as prerequisites go, knowledge of [4] and [16] should bring the reader a long way.

**2 Contents of the article** In Section 3 the necessary notions and notations are introduced. Section 4 is a step-by-step introduction to the construction of the systems that are central in this article. This section also illustrates the powers of provably extensional self-reference and contains a discussion of the problem of uniqueness and explicitness of the Gödel and Henkin sentences of the various systems considered. Section 5 treats the bimodal principles valid for the two central systems; a Kripke model completeness theorem is proved. Section 6 has a partial result on embedding Kripke models for our modal system into arithmetic. In Section 7 this embedding result is applied to a problem about relative interpretability.

**3 Conventions, notions, and elementary facts**

**3.1 Point** All the arithmetical results in the next sections will be stated for Peano Arithmetic. PA, of course, is just a convenient peg to hang the discussion on; almost any RE theory into which PA, minus induction, plus  $\Sigma_2$ -induction, can be interpreted, would do. Where results on relative interpretability appear one must also demand that the theories considered be essentially reflexive.

**3.2  $\Box$  and  $\Delta$  (in different contexts)** Let  $\text{Proof}(x, y)$  be the  $\Delta_0$  arithmetical formula representing the relation:  $x$  is the Gödel number of a PA-proof of the formula with Gödel number  $y$ . We assume for convenience that  $\text{PA} \vdash \forall x \exists ! y \text{Proof}(x, y)$ . Let  $\text{Prov}(y) = \exists x \text{Proof}(x, y)$ . We write, *par abus de langage*, ‘ $\text{Proof}(x, A(x_1, \dots, x_n))$ ’ for  $\text{Proof}(x, \ulcorner A(\dot{x}_1, \dots, \dot{x}_n) \urcorner)$ , where:

- (i) all free variables of  $A$  are among those shown
- (ii)  $\ulcorner A(\dot{x}_1, \dots, \dot{x}_n) \urcorner$  is the “Gödel term” for  $A(x_1, \dots, x_n)$  as defined on p. 43 of [16].

The modal operators  $\Box$  and  $\Delta$  will appear both in the context of modal logic and in the context of arithmetic. ‘ $\Box A(x_1, \dots, x_n)$ ’ will stand for  $\text{Prov}(\ulcorner A(\dot{x}_1, \dots, \dot{x}_n) \urcorner)$ . In arithmetical contexts ‘ $\Delta A(x_1, \dots, x_n)$ ’ will stand for  $B(\ulcorner A(\dot{x}_1, \dots, \dot{x}_n) \urcorner)$ , where  $B(x)$  is the arithmetization of theoremhood in the particular system with built-in consistency that we are considering at the place of occurrence of ‘ $\Delta A(x_1, \dots, x_n)$ ’. To avoid confusion we will use  $\Delta^R, \Delta^K$ , etc. To differentiate arithmetical from modal contexts, we use  $A, B, \dots$ , for arithmetical formulas and  $\phi, \psi, \dots$ , for modal propositional formulas.

If  $t$  is a term for a provably recursive function we will have that (supposing that  $t$  is substitutable for  $x$  in  $A$ )  $\text{PA} \vdash (\Box A(x)) [t/x] \leftrightarrow \Box A(t)$ . We will employ terms for provably recursive functions only, so we may indeed treat  $x_1, \dots, x_n$  in  $\Box A(x_1, \dots, x_n)$  simply as free variables. Similarly for  $\Delta$ .

‘ $\Diamond$ ’ will be an abbreviation for  $\neg \Box \neg$ , and ‘ $\nabla$ ’ for  $\neg \Delta \neg$ .

When we want to consider systems with other axiom sets than PA, we will write  $\text{Proof}_\alpha, \text{Prov}_\alpha, \Box_\alpha$ , etc., where  $\alpha$  is a formula that represents the axiom set of the system under consideration in an intensionally correct way in PA. We fix a formula  $\pi$  correctly representing the axiom set of PA. Thus, our notation ‘ $\Box$ ’ is just short for  $\Box_\pi$ .

### 3.3 $\Box \uparrow x$ and $\Box^*$ Define:

$$\begin{aligned}\pi \uparrow x(y) &\Leftrightarrow \pi(y) \wedge y \leq x, \\ \Box \uparrow x A &\Leftrightarrow \Box_{\pi \uparrow x} A, \\ \Diamond \uparrow x A &\Leftrightarrow \neg \Box \uparrow x \neg A, \\ \Box^* A &\Leftrightarrow \exists x \Box \uparrow x A.\end{aligned}$$

It is clear that  $\text{PA} \vdash \Box A \leftrightarrow \Box^* A$ , but the difference in form will be of some importance when Rosser-orderings come into play. (The usefulness of  $\Box^*$  in this connection was discovered by Švejdar; see [18].)

**3.4 Witnessing and the Rosser-ordering** Let  $A$  be of the form  $\exists x A_0(x)$ . Define  $t$  wit  $A \Leftrightarrow A_0(t)$ . Here we assume that bound variables in  $A_0$  are renamed, if necessary, to make  $t$  substitutable for  $x$  in  $A_0$ .

Let  $A$  be of the form  $\exists x A_0(x)$  and  $B$  of the form  $\exists x B_0(x)$ . The Rosser-orderings between  $A$  and  $B$  are defined as follows:

$$\begin{aligned}A \leq B &\Leftrightarrow \exists x (A_0(x) \wedge \forall y < x \neg B_0(y)) \\ A < B &\Leftrightarrow \exists x (A_0(x) \wedge \forall y \leq x \neg B_0(y)).\end{aligned}$$

We will always apply witnessing and the Rosser-ordering to the precise forms in which the relevant arithmetical formulas are introduced.

In connection with the Feferman System we will consider formulas of the form  $\Box^* C < \Box^* D$ . These formulas are of the more general form  $A < B$ , where  $A$  is  $\exists x \exists y A_0(x, y)$  with  $A_0$  in  $\Delta_0$  and where  $B$  is  $\exists x \exists y B_0(x, y)$  with  $B_0$  in  $\Delta_0$ . It is of some interest to know the complexity of such formulas  $A < B$ ; *prima facie*  $A < B$  is  $\Sigma_2$ . We have the following theorem:

**Theorem**  $\text{PA} \vdash \exists x (\exists y A_0(x, y) \wedge \forall z \leq x \forall u \neg B_0(z, u)) \leftrightarrow \forall u \exists x (\exists y A_0(x, y) \wedge \forall z \leq x \neg B_0(z, u))$ .

*Proof:* The “ $\rightarrow$ ” side is trivial. For the “ $\leftarrow$ ” side reason in PA as follows: Suppose that  $\forall u \exists x (\exists y A_0(x, y) \wedge \forall z \leq x \neg B_0(z, u))$ . It follows that  $\exists x \exists y A_0(x, y)$ . Let  $x_0$  be the smallest such  $x$ . Consider any  $u$ . Pick an  $x$  such that  $\exists y A_0(x, y)$  and  $\forall z \leq x \neg B_0(z, u)$ . Clearly  $x_0 \leq x$  and hence  $\forall z \leq x_0 \neg B_0(z, u)$ . We then conclude that  $\exists y A_0(x_0, y) \wedge \forall z \leq x_0 \neg B_0(z, u)$ .

Both Švejdar [18] and Lindström [9] show that in every degree of relative interpretability over PA there is a sentence of the form  $A < B$  where  $A$  and  $B$  are as above. Thus, every degree of relative interpretability contains a  $\Delta_2$  sentence.

**3.5 Relative interpretability** ‘ $A \triangleleft B$ ’ stands for the arithmetization of:  $\text{PA} + A$  is relatively interpretable in  $\text{PA} + B$ .  $\text{PA} \vdash A \triangleleft B \leftrightarrow \forall x \Box (B \rightarrow \Diamond \uparrow x A)$  is a result due to Orey and Hájek. We list a number of principles valid for  $\Box$  and  $\triangleleft$ :

- (I1)  $\text{PA} \vdash \Box (B \rightarrow A) \rightarrow A \triangleleft B$
- (I2)  $\text{PA} \vdash (A \triangleleft B \wedge B \triangleleft C) \rightarrow A \triangleleft C$
- (I3)  $\text{PA} \vdash (A \triangleleft B \wedge A \triangleleft C) \rightarrow A \triangleleft (B \vee C)$
- (I4)  $\text{PA} \vdash A \triangleleft B \rightarrow (\Diamond B \rightarrow \Diamond A)$

- (I5)  $PA \vdash \diamond A \triangleleft B \rightarrow \Box (B \rightarrow \diamond A)$   
 (I6)  $PA \vdash A \triangleleft \diamond A$   
 (I7)  $PA \vdash A \triangleleft B \rightarrow (A \wedge \Box C) \triangleleft (B \wedge \Box C)$ .

The principle (I7) is new and is due to Franco Montagna. We will prove (I4), (I5), and (I7). First we treat of (I4) and (I5). Given (I1) and (I2) it is easily seen that (I4) is equivalent to

$$(I4') \quad PA \vdash \perp \triangleleft B \rightarrow \Box \neg B.$$

Thus it is sufficient to prove

$$(J1) \quad \text{For all } P \text{ in } \Pi_1, PA \vdash P \triangleleft B \rightarrow \Box (B \rightarrow P).$$

First note that for every  $n$ ,  $PA \vdash \forall x \Box (B \rightarrow \diamond \uparrow x A) \leftrightarrow \forall x > n \Box (B \rightarrow \diamond \uparrow x A)$ . Pick  $q$  so big that  $\Box \uparrow q$  contains Robinson's Arithmetic. We then have for  $S$  in  $\Sigma_1$  that  $PA \vdash \forall x > q \Box (\Box \uparrow x S \leftrightarrow S)$ , so for  $P$  in  $\Pi_1$   $PA \vdash \forall x > q \Box (\diamond \uparrow x P \leftrightarrow P)$ . Hence

$$\begin{aligned} PA \vdash P \triangleleft B &\leftrightarrow \forall x > q \Box (B \rightarrow \diamond \uparrow x P) \\ &\leftrightarrow \Box (B \rightarrow P). \end{aligned}$$

We now turn to (I7). We prove

$$(J2) \quad \text{For all } S \text{ in } \Sigma_1, PA \vdash A \triangleleft B \rightarrow (A \wedge S) \triangleleft (B \wedge S).$$

Suppose that  $S$  is  $\Sigma_1$ . Let  $q$  be as above. Note that

$$PA \vdash \forall x > q \Box (S \rightarrow \Box \uparrow x ((D \wedge S) \leftrightarrow D)).$$

It follows that

$$\begin{aligned} PA \vdash \forall x \Box (B \rightarrow \diamond \uparrow x A) &\rightarrow \forall x > q \Box ((B \wedge S) \rightarrow \diamond \uparrow x (A \wedge S)) \\ &\rightarrow \forall x \Box ((B \wedge S) \rightarrow \diamond \uparrow x (A \wedge S)). \end{aligned}$$

For further details see [19].

**3.6 On systems** Philosophically, I think it is best to make the whole apparatus for generating theorems part of the identity conditions of systems. For our purposes however, it is more convenient to confuse the systems considered with the arithmetical predicates that codify theoremhood in the system in an intensionally correct way. I will say that a system with associated arithmetical predicate  $A$  if a *variant* of a system with predicate  $B$  if  $PA \vdash \forall x (A(x) \leftrightarrow B(x))$ .

The notion of 'system' is kept more or less open in this paper. The usual formal systems are still paradigms of systemhood. The systems we consider here are in some sense derived from the usual systems: they use the proofs of formal systems as data. A second point is that the systems considered may be seen to be extensionally equal to the formal systems on which they are based, given the information that the original systems are consistent.

**4 Systems with built-in consistency, an introduction** This section serves several purposes. First, it exhibits various ways of 'loading' systems with desired 'modal' properties. Secondly, it contains brief discussions of the various systems with built-in consistency that can be found in the literature. Thirdly, the prob-

lems of uniqueness and of explicitness of Gödel and Henkin sentences of the systems introduced are considered. (The rationale behind the attention to these specific problems is that these problems were historically at the crib of provability logic for  $\Box$ , and also that these problems turn out to be a quite pleasant starting point when one wants to get acquainted with the systems studied here.) In the fourth place, I give examples of the powers and possibilities of provably extensional self-reference. Specifically, I show how to use provably extensional self-reference to construct four nonequivalent Orey sentences.

In this section ‘ $\vdash$ ’ stands for  $\text{PA} \vdash$ , and  $A, B, C$  stand for *formulas* of the language of  $\text{PA}$ . Note that by our conventions we have that  $\vdash A(x) \Rightarrow \vdash \forall x A(x)$ , but not  $\vdash \Box A(x) \rightarrow \Box \forall x A(x)$ .

For the record we state here the usual principles valid in  $\text{PA}$  for  $\Box$ .

**P (The Peano System)** The provability principles of  $\text{PA}$  are:

- (L1)  $\vdash A \Rightarrow \vdash \Box A$
- (L2)  $\vdash \Box (A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$
- (L3)  $\vdash \Box A \rightarrow \Box \Box A$
- (L4)  $\vdash \Box (\Box A \rightarrow A) \rightarrow \Box A$ .

We will use these principles without explicit mention.

**R (The Rosser System)** The Rosser System is defined as follows:  $\Delta A \Leftrightarrow \Box A < \Box \neg A$ .

Some principles valid for the Rosser System in  $\text{PA}$  are:

- (1)  $\vdash A \Rightarrow \vdash \Delta A$
- (2)  $\vdash \neg \Delta \perp$
- (3)  $\vdash \Delta A \rightarrow \Box \Delta A$
- (4)  $\vdash \neg \Box \perp \rightarrow (\Delta A \leftrightarrow \Box A)$ .

Some direct consequences of (1)–(4) are:

- (5)  $\vdash \Delta A \rightarrow \Box A$  (4)
- (6)  $\vdash \Box A \rightarrow \Box \Delta A$ . (3),(4)

It is perhaps worth noting that the set of theorems of the Rosser System is *provably* infinite. Reason in  $\text{PA}$  as follows: in case  $\neg \Box \perp$  this is trivial. In case  $\Box \perp$  for any  $A$ , clearly one of  $A, \neg A, \neg \neg A, \neg \neg \neg A, \dots$ , will be Rosser-provable.

In the Rosser System we have two explicit but nonunique Henkin sentences:

- (7)  $\vdash \top \leftrightarrow \Delta \top$  (1)
- (8)  $\vdash \perp \leftrightarrow \Delta \perp$ . (2)

Consider a Gödel sentence of the Rosser System, i.e., a sentence  $G$  such that

- (9)  $\vdash G \leftrightarrow \neg \Delta G$ .

We have

- (10)  $\vdash \Box G \rightarrow (\Box \Delta G \wedge \Box \neg \Delta G) \rightarrow \Box \perp$ . (6),(9)

Of course we can also prove

$$(11) \quad \vdash \Box \neg G \rightarrow \Box \perp$$

but not from the modal principles collected up to now; we have to go back to the underlying Rosser-ordering. A slight change in the definition of  $\Delta$  removes this defect as we will see under K.

The uniqueness or nonuniqueness of Gödel sentences in the Rosser System is still an open problem. Guaspari and Solovay show that if one allows *variants* of Prov in the definition of  $\Delta$ , the answer may be yes and may be no (see [5]). I am not aware of any argument that there are no explicit Gödel sentences for  $\Delta$ .

**Open question** Are there explicit Gödel sentences for  $\Delta^R$ ?

**Open question** If one allows  $\Sigma_1$  variants (with one existential quantifier in front of the  $\Delta$ ) of Prov in the definition of  $\Delta^R$ , can there be explicit Gödel sentences for  $\Delta^R$ ?

**K (Kreisel's symmetrized Rosser System)** Kreisel's variation on the Rosser System is reported in [7], pp. 298-302.

Define  $\Delta A \Leftrightarrow \exists x[\text{Proof}(x, A) \wedge \forall u, v, b, c \leq x((\text{Proof}(u, b) \wedge \text{Proof}(v, c)) \rightarrow c \neq \text{neg}(b))]$ . Clearly,  $\Delta A$  is  $\Sigma_1$ . The Kreisel System satisfies principles (1)-(4) and the additional principle

$$(12) \quad \vdash \neg (\Delta A \wedge \Delta \neg A).$$

We can now prove (11) modally:

$$\begin{aligned} \vdash \Box \neg G &\rightarrow \Box \Delta G \wedge \Box \Delta \neg G && (6), (9) \\ &\rightarrow \Box (\Delta G \wedge \Delta \neg G) \\ &\rightarrow \Box \perp. && (12) \end{aligned}$$

Note also that  $\vdash \Box \perp \rightarrow \{A \mid \Delta A\}$  is finite".

**R' (A minor variation of Rosser's System)** Yet another defect of the Rosser System is that we have no appropriate bimodal counterpart for the underlying principle

$$\vdash \Box \perp \rightarrow (\Box A < \Box \neg A \vee \Box \neg A \leq \Box A).$$

This can be easily repaired. Define

$$\begin{aligned} R'_0 &= \emptyset \\ R'_{n+1} &= \begin{cases} R'_n \cup \{A\}, & \text{if } \text{Proof}(n, A) \text{ and } (\neg A) \notin R'_n \\ R'_n, & \text{otherwise.} \end{cases} \end{aligned}$$

Let  $\Delta A$  be the arithmetization of  $\exists n A \in R'_{n+1}$ . Clearly,  $\Delta A$  is  $\Sigma_1$ . We have that  $\vdash \Delta A \leftrightarrow \Box A < \Delta \neg A$ , and even  $\vdash \forall x(x \text{ wit } \Delta A \leftrightarrow x \text{ wit } \Box A < \Delta \neg A)$ . This last fact happens to characterize  $\Delta^{R'}$ .

Principles (1)-(4) hold for  $\Delta$ , plus the additional principle

$$(13) \quad \vdash \Box \perp \rightarrow (\Delta A \vee \Delta \neg A).$$

A direct consequence is

$$(14) \quad \vdash \Box A \rightarrow (\Delta A \vee \Delta \neg A). \quad (4),(13)$$

Finally, note that  $\vdash \Delta^R A \rightarrow \Delta^{R'} A$ .

**BM (The Bernardi-Montagna System)** We now jump up directly to a system which is richer, from the modal point of view, than both  $\Delta^K$  and  $\Delta^{R'}$ . This system was discovered by Claudio Bernardi and Franco Montagna (see [2]). Define

$$\begin{aligned} \text{BM}_0 &= \emptyset \\ \text{BM}_{n+1} &= \begin{cases} \text{BM}_n \cup \{A\}, & \text{if Proof}(n, A) \text{ and } \text{BM}_n \cup \{A\} \text{ is consistent} \\ & \text{in propositional logic} \\ \text{BM}_n, & \text{otherwise.} \end{cases} \end{aligned}$$

Let  $\Delta A$  be the arithmetization of  $\exists n A \in \text{BM}_{n+1}$ . Clearly,  $\Delta A$  is  $\Sigma_1$ . Also, principles (1)–(4) and (13) hold for  $\Delta$ . We also have, by elementary reasoning,

$$(15) \quad \vdash \Delta(A \rightarrow B) \rightarrow (\Delta A \rightarrow \Delta B).$$

(15) in combination with (2) entails (12), so the principles valid for  $\Delta^{\text{BM}}$  comprise those valid for  $\Delta^K$  and  $\Delta^{R'}$  (at least insofar as we have found such principles).

**mBM (The modified Bernardi-Montagna System)** For our purposes we want the following additional principle

$$(16) \quad \vdash \Delta A \rightarrow \Delta \Box A.$$

It is not plausible for one to prove (16) for the BM System without additional assumptions about the order of the proofs of PA; e.g., given  $\Box \perp$ , why would one have  $\Delta \Box \perp$  rather than  $\Delta \neg \Box \perp$ ? We can, however, modify the BM System in such a way that we get (16).

Let  $\vdash_{\text{Prop}}$  stand for derivability in propositional logic. Define

$$\begin{aligned} \text{mBM}_0 &= \emptyset \\ \text{mBM}_{n+1} &= \begin{cases} \text{mBM}_n \cup \{A\}, & \text{if Proof}(n, A) \text{ and } \text{mBM}_n \cup \{A\} \not\vdash_{\text{Prop}} \neg \Box \perp \\ \text{mBM}_n, & \text{otherwise.} \end{cases} \end{aligned}$$

Let  $\Delta A$  be the arithmetization of  $\exists n A \in \text{mBM}_{n+1}$ . Clearly,  $\Delta A$  is  $\Sigma_1$ .

It is easily seen that (1)–(3) and (15) are valid. We may verify (4) in PA as follows:

Trivially,  $\Delta A \rightarrow \Box A$ . Suppose that  $\neg \Box \perp$  and  $\Box A$ , say  $\text{Proof}(x, A)$ . The only reason  $A$  could be left out of  $\text{mBM}_{x+1}$  is that  $\text{mBM}_x \cup \{A\} \vdash_{\text{Prop}} \neg \Box \perp$ . But then  $\Box \neg \Box \perp$  and hence  $\Box \perp$ , which is a contradiction. So we conclude that  $A \in \text{mBM}_{x+1}$  and thus  $\Delta A$ .

(13) can be proved in PA as follows:

Suppose that  $\Box \perp$ . For certain  $x$  and  $y$  we will have that  $\text{Proof}(x, A)$  and  $\text{Proof}(y, \neg A)$ . Let  $z = \max(x, y)$ . If  $\neg \Delta A$  and  $\neg \Delta \neg A$  we will have that

$mBM_{z+1} \cup \{A\} \vdash_{\text{Prop}} \neg \Box \perp$  and  $mBM_{z+1} \cup \{\neg A\} \vdash_{\text{Prop}} \neg \Box \perp$ . Hence  $mBM_{z+1} \vdash_{\text{Prop}} \neg \Box \perp$ .

Finally, we turn to (16). Clearly,

$$(17) \quad \vdash \neg \Delta \neg \Box \perp$$

from which it follows that

$$\vdash \Box \perp \rightarrow \Delta \Box \perp \quad (13),(17)$$

moreover

$$\vdash \Delta \Box \perp \rightarrow \Delta \Box A \quad (1),(15)$$

hence

$$\vdash \Box \perp \rightarrow (\Box A \rightarrow \Delta \Box A)$$

also

$$\vdash \neg \Box \perp \rightarrow (\Delta \Box A \leftrightarrow \Box \Box A) \quad (4)$$

thus

$$\vdash \neg \Box \perp \rightarrow (\Box A \rightarrow \Delta \Box A).$$

So we may conclude that (16) holds. (Conversely, one can derive (17) from (4), (12), and (16).)

Let us list for the sake of convenience the principles valid for  $\Delta^{\text{mBM}}$  with brand new names:

- (B1)  $\vdash A \Rightarrow \vdash \Delta A$
- (B2)  $\vdash \Delta (A \rightarrow B) \rightarrow (\Delta A \rightarrow \Delta B)$
- (B3)  $\vdash \neg \Delta \perp$
- (B4)  $\vdash \Box A \rightarrow \Delta \Box A$
- (B5)  $\vdash \Delta A \rightarrow \Box \Delta A$
- (B6)  $\vdash \neg \Box \perp \rightarrow (\Delta A \leftrightarrow \Box A)$
- (B7)  $\vdash \Box \perp \rightarrow (\Delta A \vee \Delta \neg A)$ .

We note some important consequences of these principles. First, a strengthening of Löb's Axiom

$$(18) \quad \vdash \Delta (\Box A \rightarrow A) \rightarrow \Delta A.$$

$$\begin{aligned} \text{Proof: } \vdash \Delta (\Box A \rightarrow A) &\rightarrow \Box (\Box A \rightarrow A) & (5) \\ &\rightarrow \Box A \\ &\rightarrow \Delta \Box A & (B4) \\ &\rightarrow \Delta A. & (B2) \end{aligned}$$

The second consequence is the principle of provable extensionality

$$(19) \quad \vdash \Box (A \leftrightarrow B) \rightarrow \Box (\Delta A \leftrightarrow \Delta B).$$

$$\begin{aligned} \text{Proof: } \vdash \Box (A \leftrightarrow B) &\rightarrow \Box \Delta (A \leftrightarrow B) & (6) \\ &\rightarrow \Box (\Delta A \leftrightarrow \Delta B). & (B2) \end{aligned}$$

The next principle is an immediate consequence of (12) and (B7)

$$(20) \quad \vdash \Box \perp \rightarrow (\Delta A \leftrightarrow \nabla A).$$

Let us define  $\Box^0 \perp = \perp$ ,  $\Box^{n+1} \perp = \Box \Box^n \perp$ ,  $\Box^\omega \perp = \top$ . We will say that an arithmetical formula  $A$  is *modally closed* if  $A$  is built up from  $\top, \perp$  with the propositional connectives and  $\Box, \Delta$  (in other words, if  $A$  is an interpreted sentence of the closed fragment of the bimodal propositional logic with operators  $\Box$  and  $\Delta$ ).

(21) Suppose that  $A$  is modally closed. Then there is an  $\alpha \in \{0, \dots, \omega\}$  such that  $\vdash \Delta A \leftrightarrow \Box^\alpha \perp$ .

*Proof:* Consider  $B$  built from  $\top, \perp$  with the propositional connectives and  $\Box$ . First there is the familiar fact that

$$\vdash (B \wedge \Box B) \leftrightarrow \Box^\beta \perp, \text{ for some } \beta \in \{0, \dots, \omega\}.$$

Hence

$$\begin{aligned} \vdash \Delta B &\leftrightarrow \Delta (B \wedge \Box B) \\ &\leftrightarrow \Delta \Box^\beta \perp. \end{aligned}$$

Secondly we have that  $\vdash \Delta \Box^0 \perp \leftrightarrow \Box^0 \perp$ , and  $\vdash \Delta \Box^{1+\gamma} \perp \leftrightarrow \Box^{2+\gamma} \perp$  (as is easily seen by considering the cases  $\Box \perp$  and  $\neg \Box \perp$  separately). Combining these results we see that  $\vdash \Delta B \leftrightarrow \Box^\delta \perp$ , for some  $\delta \in \{0, \dots, \omega\}$ . (21) follows by a trivial induction on  $A$ . (Note that we didn't use (B7) in the argument.)

We now turn to the Henkin sentences of  $\Delta$ . We have already seen in (7) and (8) that  $\perp$  and  $\top$  are explicit Henkin sentences. For  $\Delta^{\text{mBM}}$  we can show that they are the only explicit Henkin sentences. Consider  $H$  satisfying:

$$(22) \quad \vdash H \leftrightarrow \Delta H.$$

If  $H$  is explicit, it follows that  $H$  is modally closed. Hence by (21)

$$\vdash H \leftrightarrow \Box^\alpha \perp, \text{ for some } \alpha \in \{0, \dots, \omega\}.$$

If  $\alpha \neq 0$ ,  $\alpha \neq \omega$ , it follows that for some  $n \in \{1, 2, \dots\}$

$$\begin{aligned} \vdash \Box^n \perp &\leftrightarrow \Delta \Box^n \perp \\ &\leftrightarrow \Box^{n+1} \perp. \end{aligned}$$

**Open question** Are there nonexplicit Henkin sentences of  $\Delta^{\text{mBM}}$ ?

Next we turn to the Gödel sentences of  $\Delta$ . Under R and K we have seen in (10) and (11) that these have the Rosser Property. We show that they are nonexplicit and nonunique.

Consider  $G$  satisfying (9). If  $G$  were explicit,  $G$  would be modally closed. Hence by (21)  $\vdash G \leftrightarrow \neg \Box^\alpha \perp$  ( $\alpha \in \{0, \dots, \omega\}$ ). If  $\alpha \neq 0$ , we get that

$$\begin{aligned} \vdash \Delta G &\leftrightarrow \Delta \neg \Box^\alpha \perp && \text{(B1), (B2)} \\ &\leftrightarrow \Delta \perp && \text{(18)} \\ &\leftrightarrow \perp. && \text{(B3)} \end{aligned}$$

So, by (9),  $\vdash G$ . Thus  $\alpha = 0$ , which is a contradiction. If  $\alpha = \omega$ , we have that  $\vdash G$  and thus  $\vdash \Delta G$ , i.e., by (9),  $\vdash \neg G$ , another contradiction. Hence  $G$  cannot be explicit.

To see that  $G$  is not unique we show that  $\nabla G$  is also a Gödel sentence and that  $\nabla G$  is not provably equivalent to  $G$ . First we show

$$(23) \quad \vdash \nabla G \leftrightarrow \neg \Delta \nabla G.$$

To prove (23) it is clearly sufficient to show

- (a)  $\vdash \neg \Box \perp \rightarrow \nabla G$
- (b)  $\vdash \neg \Box \perp \rightarrow \neg \Delta \nabla G$
- (c)  $\vdash \Box \perp \rightarrow (\nabla G \leftrightarrow \neg \Delta \nabla G)$ .

We prove (a) by contraposition

$$\begin{aligned} \vdash \Delta \neg G &\rightarrow \Box \neg G && (5) \\ &\rightarrow (\Box \Delta G \wedge \Box \Delta \neg G) && (6),(9) \\ &\rightarrow \Box \perp. && (12) \end{aligned}$$

To prove (b), we show first

$$\begin{aligned} \vdash \Box \Box \perp &\rightarrow (\Delta \nabla G \rightarrow \Box \nabla G) && (5) \\ &\rightarrow \Box \Delta G && (20) \\ &\rightarrow \Box \neg G && (9) \\ &\rightarrow \Box \perp). && \text{(as in the proof of (a))} \end{aligned}$$

Hence

$$\vdash \Delta \nabla G \rightarrow (\Box \Box \perp \rightarrow \Box \perp)$$

thus

$$\begin{aligned} \vdash \Delta \nabla G &\rightarrow \Box \Delta \nabla G && (B5) \\ &\rightarrow \Box (\Box \Box \perp \rightarrow \Box \perp) \\ &\rightarrow \Box \Box \perp \end{aligned}$$

and so, combining,

$$\vdash \Delta \nabla G \rightarrow \Box \perp.$$

The proof of (c) goes as follows:

$$\begin{aligned} \vdash \Box \perp &\rightarrow \Delta \Box \perp && (B4) \\ &\rightarrow \Delta (\nabla G \leftrightarrow \Delta G) && (B1),(B2),(20) \\ &\rightarrow (\Delta \nabla G \leftrightarrow \Delta \Delta G) && (B2) \\ &\rightarrow (\Delta \nabla G \leftrightarrow \Delta \neg G) && (9),(B1),(B2) \\ &\rightarrow (\neg \Delta \nabla G \leftrightarrow \nabla G). \end{aligned}$$

Next we show that  $\nabla G$  is not provably equivalent to  $G$ . It is clearly sufficient to prove

$$(24) \quad \vdash \Box (G \leftrightarrow \nabla G) \rightarrow \Box \perp.$$

Clearly,

$$\begin{aligned} \vdash (G \wedge \nabla G) &\rightarrow (\neg \Delta G \wedge \neg \Delta \neg G) \\ &\rightarrow \neg \Box \perp && (B7) \end{aligned}$$

and

$$\begin{aligned} \vdash (\neg G \wedge \neg \nabla G) &\rightarrow (\Delta G \wedge \Delta \neg G) \\ &\rightarrow \perp. \end{aligned} \quad (12)$$

Combining, we get that

$$\begin{aligned} \vdash \Box (G \leftrightarrow \nabla G) &\rightarrow \Box ((G \wedge \nabla G) \vee (\neg G \wedge \neg \nabla G)) \\ &\rightarrow \Box \neg \Box \perp \\ &\rightarrow \Box \perp. \end{aligned}$$

Another way to prove (24) is by noting that

$$\begin{aligned} \vdash \Box \Box \perp &\rightarrow (\Box (G \leftrightarrow \nabla G) \rightarrow \Box (G \leftrightarrow \Delta G) \\ &\rightarrow \Box \perp). \end{aligned} \quad (20)$$

We leave it to the reader to verify that our procedure yields no further independent Gödel sentences, i.e.,

$$(25) \quad \vdash G \leftrightarrow \nabla \nabla G.$$

In Section 5 we will see that as far as our modal principles are concerned we cannot show more than this: if there is a Gödel sentence of  $\Delta$  then there is a second, nonequivalent one.

**Open question** Are there three pairwise nonequivalent Gödel sentences of  $\Delta^{\text{mBM}}$ ?

The R, K, R', BM, and mBM systems are all  $\Sigma_1$ , yet they escape the second Incompleteness Theorem. By a well-known result of Feferman (see [4]) these systems cannot be *provably* closed under the axioms and rules of predicate logic; in other words, we do not have  $\vdash \Delta A \leftrightarrow \Box_{\Delta} A$ . If  $\Delta$  is one of  $\Delta^{\text{BM}}, \Delta^{\text{mBM}}$  we can say a bit more. Let  $\Delta$  be one of  $\Delta^{\text{BM}}, \Delta^{\text{mBM}}$ . Let  $Q$  be the conjunction of the axioms of Robinson's Arithmetic. We clearly have  $\vdash Q$ , and hence  $\vdash \Delta Q$ , and thus  $\vdash \Box_{\Delta} Q$ . It immediately follows from the provable  $\Sigma_1$ -completeness of Robinson's Arithmetic that

$$(26) \quad \vdash \Delta A \rightarrow \Box_{\Delta} \Delta A.$$

Let  $G$  be a Gödel sentence of  $\Delta$ . We have

$$(d) \quad \vdash \Box \perp \rightarrow (\Delta G \vee \Delta \neg G) \quad (B7)$$

$$(e) \quad \vdash \Delta G \rightarrow \Box_{\Delta} \Delta G \quad (26)$$

$$(f) \quad \begin{aligned} \vdash \Delta G &\rightarrow \Box_{\Delta} G \\ &\rightarrow \Box_{\Delta} \neg \Delta G \end{aligned} \quad (9)$$

$$(g) \quad \vdash \Delta G \rightarrow \Box_{\Delta} \perp \quad (e), (f)$$

$$(h) \quad \begin{aligned} \vdash \Delta \neg G &\rightarrow \Box_{\Delta} \Delta \neg G \\ &\rightarrow \Box_{\Delta} \neg \Delta G \end{aligned} \quad (26)$$

$$(i) \quad \begin{aligned} \vdash \Delta \neg G &\rightarrow \Box_{\Delta} \neg G \\ &\rightarrow \Box_{\Delta} \Delta G \end{aligned} \quad (B1), (12)$$

$$(j) \quad \vdash \Delta \neg G \rightarrow \Box_{\Delta} \perp \quad (h), (i)$$

$$(k) \quad \vdash \Box \perp \rightarrow \Box_{\Delta} \perp \quad (d), (g), (j)$$

$$(l) \quad \vdash \Box \perp \rightarrow (\Box_{\Delta} A \leftrightarrow \Box A) \quad (k)$$

$$(m) \quad \begin{aligned} \vdash \neg \Box \perp &\rightarrow (\Box_{\Delta} A \leftrightarrow \Box_{\square} A) \\ &\leftrightarrow \Box A \end{aligned}$$

$$(27) \quad \vdash \Box_{\Delta} A \leftrightarrow \Box A. \quad (l), (m)$$

So  $\Delta^{\text{BM}}$  and  $\Delta^{\text{mBM}}$  are provably axiom sets for the theorems of PA. The same thing can be proved for  $\Delta^{\text{R}'}$  by a slightly refined variant of the above argument. What about  $\Delta^{\text{R}}$  and  $\Delta^{\text{K}}$ ? I'm not certain, but a good guess is that it would be so for  $\Delta^{\text{R}}$  but not for  $\Delta^{\text{K}}$ .

We now turn to systems that are provably closed under the axioms and rules of predicate logic.

**F (The Feferman System)** The Feferman System was invented by Feferman (see [4]) as an illustration in the study of the conditions for Gödel's second Incompleteness Theorem. Orey discovered important applications of its provability predicate in the theory of relative interpretability (see [4], [13]). A modal study of this provability predicate was made by Montagna [11].

Let us start by giving two rather different intuitive descriptions of the Feferman System (or, to be faithful to the conventions of Section 3.6, I should say: let us describe two variants of the Feferman System).

Suppose the arithmetical axioms of PA are enumerated by  $A_1, A_2, A_3, \dots$ , in the order of their Gödel numbers (i.e.,  $i < j \Rightarrow \ulcorner A_i \urcorner < \ulcorner A_j \urcorner$ ). We call a set  $X$  of arithmetical axioms of PA *initial* if  $A_i \in X$  and  $j < i \Rightarrow A_j \in X$ .

The Feferman System is simply the first-order system in the language of PA axiomatized by  $F = \bigcup \{X \mid X \text{ is a finite, initial, consistent set of arithmetical axioms of PA}\}$ .

Clearly, from the extensional point of view F coincides with the usual axiom set of PA. The Feferman System can be viewed as a system, where to be licensed to use axiom  $A_i$  one needs the *external information* that  $\{A_j \mid j \leq i\}$  is consistent.

The second way to introduce the Feferman System is as follows. Suppose we enumerate the *proofs* in the system PA by  $\pi_1, \pi_2, \pi_3, \dots$ . As soon as we hit upon a proof  $\pi_i$  of  $\perp$ , we extract the axiom  $A_j$  with largest Gödel number from  $\pi_i$ . We backtrack and scratch out all the proofs employing axioms  $A_k$  with  $k \geq j$ . Then we go on enumerating proofs, skipping those employing axioms  $A_k$  with  $k \geq j$ . As soon as we meet another proof of  $\perp$  we repeat the procedure. We call a proof *stable* if it occurs in our enumeration and is never scratched out. The stable proofs are the proofs of the Feferman System.

Under this last description the Feferman System can be seen as a fully effective procedure that will eventually yield all stable proofs. The catch here is, of course, that someone who does not know the consistency of PA will not be able to predict, at least not *prima facie*, when a proof is stable. In fact, the situation is even subtler: someone knowing PA, but not its consistency, will *ipso facto* not know that *all proofs* are stable, but he will know *of every proof* that it is stable.

We now turn to the formal definition of the Feferman System. Define

$$\begin{aligned} \pi^*(x) &\Leftrightarrow \pi(x) \wedge \diamond \ulcorner x \urcorner \\ \Delta A &\Leftrightarrow \Box_{\pi^*} A. \end{aligned}$$

The following are equivalents of  $\Delta A$ :

$$(28) \quad \vdash \Delta A \leftrightarrow \exists x (\Box \ulcorner x \urcorner \wedge \diamond \ulcorner x \urcorner).$$

where  $f$  is a primitive recursive function with

$$f(n) = \begin{cases} \text{the largest of the Gödel numbers of the arithmetical axioms} \\ \text{occurring in } \pi, \text{ if } n = \ulcorner \pi \urcorner \text{ for some proof } \pi \\ 0, \text{ otherwise} \end{cases}$$

we have

$$(29) \quad \vdash \Delta A \leftrightarrow \exists x(\text{Proof}(x, A) \wedge \diamond \uparrow f(x) \top).$$

Remembering that  $\Box^* A \Leftrightarrow \exists x \Box \uparrow x A$ , we have

$$(30) \quad \vdash \Delta A \leftrightarrow \Delta^* A < \Box^* \perp$$

$$(31) \quad \vdash \Delta A \leftrightarrow \Box^* A < \Box^* \neg A.$$

(31) brings out the similarity between the Feferman System and the Rosser System. By 3.4 and (30) or (31) we see that  $\Delta$  is  $\Delta_2$ .

(B1)–(B6) are valid for  $\Delta$ . In [4] all of these except (B4) are mentioned. In [11] a modal study is made of (B1), (B2), (B3), (B5), and (B6). The validity of (B2), (B3), and (B6) is immediate. To prove the other principles we will use the well-known fact that PA is provably essentially reflexive, and hence

$$\vdash \forall x \Box (\Box \uparrow x A \rightarrow A)$$

from which it follows that

$$\vdash \forall x \Box \diamond \uparrow x \top.$$

*Proof of B1:*

$$\begin{aligned} \vdash A &\Rightarrow \text{for some } n \vdash \Box \uparrow n A \\ &\Rightarrow \text{for some } n \vdash \Box \uparrow n A \wedge \diamond \uparrow n \top \\ &\Rightarrow \vdash \Delta A. \end{aligned}$$

*Proof of B4:* We will prove the stronger principle

$$(32) \quad \text{Let } S \text{ be } \Sigma_1, \text{ then } \vdash S \rightarrow \Delta S.$$

Let  $\Box_Q$  stand for provability in Robinson's Arithmetic. For some  $q$ ,  $\text{PA} \vdash \Box_Q A \rightarrow \Box \uparrow q A$ . Let  $S$  be  $\Sigma_1$ . We have that

$$\begin{aligned} \vdash S &\rightarrow \Box_Q S \\ &\rightarrow \Box \uparrow q S \\ &\rightarrow (\Box \uparrow q S \wedge \diamond \uparrow q \top) \\ &\rightarrow \Delta S. \end{aligned}$$

Note that we *do* have (B4) for  $\Delta^{\text{mBM}}$ , but not (32). (In fact, *assuming* (32) for  $\Delta^{\text{mBM}}$  quickly leads to the inconsistency of PA.)

*Proof of B5:* It is clearly sufficient to prove (6); i.e.,  $\vdash \Box A \rightarrow \Delta \Box A$ . To do this we formalize the reasoning for (B1) as follows:

$$\begin{aligned} \vdash \Box A &\rightarrow \exists x \Box \Box \uparrow x A \\ &\rightarrow \exists x \Box (\Box \uparrow x A \wedge \diamond \uparrow x \top) \\ &\rightarrow \Box \Delta A. \end{aligned}$$

Just as for  $\Delta^{\text{mBM}}$ ,  $\Delta^{\text{F}}$  has precisely two nonequivalent *explicit* Henkin sentences. I will now show that  $\Delta^{\text{F}}$  has in fact infinitely many pairwise nonequiva-

lent Henkin sentences. First we need to know a bit about  $\Sigma$ -minded sentences. A sentence  $A$  is  $\Sigma$ -minded if both  $A \wedge \Box A$  and  $\neg A \wedge \Box \neg A$  are provably equivalent (in PA) to a  $\Sigma_1$  formula. A good example of a  $\Sigma$ -minded sentence is the ordinary  $\Sigma_1$  Rosser sentence. We have that

$$(33) \text{ If } A \text{ is } \Sigma\text{-minded then } \vdash \Delta A \leftrightarrow (\Box A \wedge (\Box \perp \rightarrow A)).$$

*Proof:* Suppose that  $A$  is  $\Sigma$ -minded. To prove left to right, it is clear that  $\vdash \Delta A \rightarrow \Box A$ . Moreover,

$$\begin{aligned} \vdash (\Delta A \wedge \Box \perp) &\rightarrow (\neg A \rightarrow (\neg A \wedge \Box \neg A)) \\ &\rightarrow \Delta (\neg A \wedge \Box \neg A) && \text{(B1),(B2),(32)} \\ &\rightarrow \Delta \perp && \text{(B1),(B2)} \\ &\rightarrow \perp. && \text{(B3)} \end{aligned}$$

To prove right to left, we have that

$$\begin{aligned} \vdash (\Box A \wedge (\Box \perp \rightarrow A)) &\rightarrow (\Box \perp \rightarrow (A \wedge \Box A)) \\ &\rightarrow \Delta (A \wedge \Box A) && \text{(B1),(B2),(32)} \\ &\rightarrow \Delta A. && \text{(B1),(B2)} \end{aligned}$$

Moreover,

$$\vdash (\Box A \wedge (\Box \perp \rightarrow A)) \rightarrow (\neg \Box \perp \rightarrow \Delta A). \quad \text{(B6)}$$

Note that our proof uses only (B1), (B2), (B3), (B6), and (32). (33) is an example of the phenomenon of *reduction*: an arithmetical predicate takes a simple, uncharacteristic form on some restricted set of formulas. A further, more involved example of reduction will be given in Section 7.1.

To prove that there are infinitely many pairwise nonequivalent Henkin sentences I have to borrow some material and definitions from [20]. The reader not familiar with this paper can at least get the essential idea of the argument by considering the ordinary  $\Sigma_1$  Rosser sentence  $R$  (i.e., any sentence satisfying  $\vdash R \leftrightarrow \Box \neg R < \Box R$ ) and  $S = \Box R \leq \Box \neg R$ , and by proving for himself that  $R$  and  $S$  are  $\Sigma$ -minded and satisfy  $\vdash R \leftrightarrow (\Box R \wedge (\Box \perp \rightarrow R))$ ,  $\vdash S \leftrightarrow (\Box S \wedge (\Box \perp \rightarrow S))$ .

Consider a tail model  $\mathbf{K}$ . We write  $\llbracket \phi \rrbracket$  for the set of nodes that force  $\phi$ ,  $[\phi]$  for  $[\phi](\mathbf{K}, \text{PA})$ , and  $\langle \phi \rangle$  for  $\langle \phi \rangle(\mathbf{K}, \text{PA})$ . Note that  $\llbracket \phi \wedge \Box \phi \rrbracket$  is upwards closed and that  $\vdash ([\phi] \wedge \Box [\phi]) \leftrightarrow [\phi \wedge \Box \phi]$ . It follows that  $\vdash ([\phi] \wedge \Box [\phi]) \leftrightarrow \exists x h(x) \in \llbracket \phi \wedge \Box \phi \rrbracket$ , and hence that  $[\phi] \wedge \Box [\phi]$  is provably equivalent to a  $\Sigma_1$  sentence. Combining this with the fact that  $\vdash \neg [\phi] \leftrightarrow [\neg \phi]$ , we find that  $[\phi]$  is  $\Sigma$ -minded.

Now consider the tail model shown in Figure 1, where the  $p_i$  are only forced as shown. Clearly,  $\Vdash p_i \leftrightarrow (\Box p_i \wedge (\Box \perp \rightarrow p_i))$ , hence, by the Embedding Lemma and the fact that  $[p_i]$  is  $\Sigma$ -minded  $\vdash [p_i] \leftrightarrow \Delta [p_i]$ . On the other hand, for  $i \neq j \Vdash (p_i \leftrightarrow p_j) \rightarrow \neg \Box \perp$ , hence  $\vdash [(p_i \leftrightarrow p_j) \rightarrow \neg \Box \perp]$ , and so  $\vdash ([p_i] \leftrightarrow [p_j]) \rightarrow \neg \Box \perp$ . Thus the  $[p_i]$  are pairwise nonequivalent Henkin sentences of  $\Delta$ . Since  $n$  can be freely chosen it follows that there are infinitely many pairwise nonequivalent Henkin sentences of  $\Delta$ .

We state two open problems:

**Open problem** Are there Henkin sentences of  $\Delta^F$  that are not provably equivalent to  $\Sigma_1$  sentences?

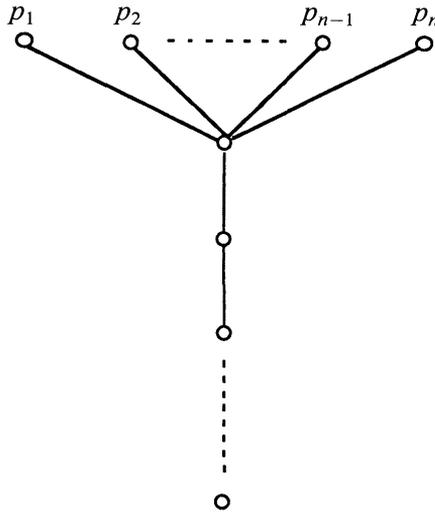


Figure 1.

**Open problem** What are the possible truth values of the *literal* Henkin sentences of  $\Delta^F$ ?

We turn next to Gödel sentences. Let  $G$  satisfy (9), so as in the case of  $\Delta^{\text{mBM}}$   $G$  is nonexplicit. (We can, using the observations about tail models above, also see this “vom höheren Standpunkt,” for consider the ‘minimal’ tail model, i.e., the linear one. The ‘propositions’ of this model correspond precisely to the closed fragment of Löb’s logic. Clearly, the interpretations of this closed fragment are going to be closed under  $\Delta$  (modulo provable equivalence). It follows that the modally closed sentences are provably equivalent to arithmetical interpretations of elements of the closed fragment of Löb’s logic. In this model the ‘equation’  $(\phi \leftrightarrow \neg(\Box\phi \wedge (\Box\perp \rightarrow \phi)))$  has no solution, hence no modally closed sentence solves the equation in PA!) The argument for the nonuniqueness of the Gödel sentences of  $\Delta^{\text{mBM}}$  depended upon (B7), so we can’t use it here. The problem of the uniqueness of  $G$  thus remains open. This problem was first posed by Montagna [11].

**Montagna’s problem** Is  $G$  unique?

The Gödel sentence of  $\Delta^F$  is an Orey sentence. Before defining what an Orey sentence is, I want to note that the fact that  $G$  is such a sentence only depends upon (B1), (B2), (B3), (6), and (32). Let us call a  $\Delta$  that satisfies these principles *precocious*.

An *Orey sentence* is a sentence  $A$  that has the property  $A \triangleleft \top$  and  $\neg A \triangleleft \top$ . (Strictly speaking, my usage is at variance with the tradition: e.g., Gödel sentences and Rosser sentences are sentences that solve certain fixed equations; on the other hand, we say of a sentence  $A$  satisfying  $\vdash(\Box A \vee \Box \neg A) \rightarrow \Box \perp$  that *it has the Rosser Property*. Rosser sentences have the Rosser Property, but other sentences do as well. So the more correct usage would be: sentence with

the Orey Property.) Trivially, the negation of an Orey sentence is again an Orey sentence.

We now show that the Gödel sentence of any precocious  $\Delta$  is an Orey sentence. Suppose that  $\Delta$  is precocious. First we prove

$$(34) \quad \vdash A \triangleleft \nabla A.$$

*Proof:* Let  $B = \neg A$ ; we have

$$\begin{aligned} \vdash \forall x \square (\square \uparrow x B \rightarrow B) &\Rightarrow \\ \vdash \forall x \square \Delta (\square \uparrow x B \rightarrow B) &\Rightarrow (6) \\ \vdash \forall x \square (\Delta \square \uparrow x B \rightarrow \Delta B) &\Rightarrow (B2) \\ \vdash \forall x \square (\square \uparrow x B \rightarrow \Delta B) &\Rightarrow (32) \\ \vdash \forall x \square (\nabla A \rightarrow \diamond \uparrow x A) &\Rightarrow \\ \vdash A \triangleleft \nabla A. & \end{aligned}$$

Secondly, one easily proves, using (I1), (I2), (I3)

$$(35) \quad \vdash A \triangleleft \neg A \rightarrow A \triangleleft \top.$$

We have

$$\begin{aligned} \vdash G \triangleleft \nabla G & (34) \\ \triangleleft \Delta G & (B1),(B2),(B3),(I1),(I2) \\ \triangleleft \neg G & (9),(I1),(I2) \end{aligned}$$

hence by (35)

$$\vdash G \triangleleft \top$$

moreover

$$\begin{aligned} \vdash \neg G \triangleleft \nabla \neg G & (34) \\ \triangleleft \neg \Delta G & (B1),(B2),(I1),(I2) \\ \triangleleft G & (9),(I1),(I2) \end{aligned}$$

hence by (I1), (I2), and (35)

$$\vdash \neg G \triangleleft \top.$$

A curious fact is that the Gödel sentences of  $\Delta^F$  is *precisely* the Orey sentence discovered independently by Lindström and Švejdar (see [9] and [18]); by 3.4 *this* Orey sentence is  $\Delta_2$ . In Section 7 we will see that there are infinitely many nonequivalent Orey sentences.

Before leaving the subject of Orey sentences, I want to note that Orey sentences are  $\Sigma_1$ - and  $\Pi_1$ -flexible and that they are Kent sentences. Let  $\Gamma$  be a set of formulas. A formula  $A$  is  $\Gamma$ -flexible if, for all  $B$  in  $\Gamma$ ,  $\vdash \square \neg (A \leftrightarrow B) \rightarrow \square \perp$ . A sentence  $A$  is a *Kent sentence* if  $(A \wedge \square A)$  is not provably equivalent to a  $\Sigma_1$  sentence. I will show that an Orey sentence is a Kent sentence and leave the proof that Orey sentences are  $\Sigma_1$ - and  $\Pi_1$ -flexible to the reader. Suppose that  $A$  is an Orey sentence and suppose, for a reductio, that  $A$  is a Kent sentence. Then clearly  $(\neg A \vee \diamond \neg A)$  is provably equivalent to a  $\Pi_1$  sentence and hence

$$\begin{aligned}
\vdash (A \triangleleft \top \wedge \neg A \triangleleft \top) &\rightarrow (\neg A \vee \diamond \neg A) \triangleleft \top && \text{(I1),(I2)} \\
&\rightarrow \Box (\neg A \vee \diamond \neg A) && \text{(I1),(I2),(J1)} \\
&\rightarrow \Box (A \rightarrow \diamond \neg A) \\
&\rightarrow \diamond \neg A \triangleleft \top && \text{(I1),(I2)} \\
&\rightarrow \Box \diamond \neg A && \text{(I5)} \\
&\rightarrow \Box \perp.
\end{aligned}$$

**mF (The modified Feferman System)** The modified Feferman System is a modification of both the Feferman System and of the BM System. Define

$$\begin{aligned}
\text{mF}_0 &= \emptyset \\
\text{mF}_{n+1} &= \begin{cases} \text{mF}_n \cup \{A\}, & \text{if Proof}(n, A) \text{ and } \text{mF}_n \cup \{A\} \text{ is consistent} \\ \text{mF}_n, & \text{otherwise.} \end{cases}
\end{aligned}$$

Let  $\Delta_x A$  be the arithmetization of  $A \in \text{mF}_{x+1}$ . Define further

$$\begin{aligned}
\Delta A &\Leftrightarrow \exists x \Delta_x A \\
\Delta^* A &\Leftrightarrow \exists x \Box_{\Delta_x} A.
\end{aligned}$$

It is easily seen that  $\vdash \Delta A \leftrightarrow \Box_{\Delta} A$ , and hence that  $\vdash \Delta A \leftrightarrow \Delta^* A$ . Moreover we have that  $\vdash \Delta A \leftrightarrow \Box A < \Delta^* \neg A$ , and even

$$\vdash \forall x (x \text{ wit } \Delta A \leftrightarrow x \text{ wit } \Box A < \Delta^* A).$$

This last observation brings out the Rosserlike character of  $\Delta$ .

We claim that  $\Delta$  satisfies (B1)–(B7). The argument for the validity of (B1)–(B6) is similar to the one for the case of  $\Delta^F$ . We treat of example (B5).

Define  $\Box_x A \Leftrightarrow \exists y \leq x \text{ Proof}(y, A)$ . Clearly  $\vdash \forall x \Box \neg \Box_{\Box_x} \perp$ , and hence, by induction on  $x$  in PA,  $\vdash \forall x \Box (\Delta_x A \leftrightarrow \Box_x A)$ . It follows that

$$\begin{aligned}
\vdash \Delta A &\rightarrow \Box A \\
&\rightarrow \exists x \Box \Box_x A \\
&\rightarrow \Box \Delta A.
\end{aligned}$$

The argument for (B7) is similar to the one for the case of  $\Delta^{\text{mBM}}$ . Just like  $\Delta^F, \Delta^{\text{mF}}$  satisfies (32), i.e.,  $\Delta^{\text{mF}}$  is provably  $\Sigma_1$ -complete. *Prima facie*,  $\Delta$  is  $\Sigma_2$ . It is seen to be  $\Delta_2$  by the following observation

$$(36) \quad \vdash \neg \Delta A \leftrightarrow (\neg \Box A \vee (\Box \perp \wedge \Delta \neg A)). \quad \text{(B6),(B7)}$$

Concerning the Henkin sentences of  $\Delta^{\text{mF}}$  the same remarks can be made as for  $\Delta^F$ . Just like  $\Delta^{\text{mBM}}$ ,  $\Delta^{\text{mF}}$  has at least two nonequivalent Gödel sentences. Clearly,  $\Delta^{\text{mF}}$  is precocious. It is now easy to see that the two nonequivalent Gödel sentences and their negations give us four pairwise nonequivalent Orey sentences. In Section 6 we will show that  $\Delta^{\text{mF}}$  has in fact infinitely many pairwise nonequivalent Gödel sentences; thus, there are infinitely many pairwise nonequivalent Orey sentences.

$\Delta^{\text{mF}}$  is our final system and the main object of study of this paper. In Section 5 we will study the principles (B1)–(B7) from the modal point of view. In Section 6 we give a partial result on embedding Kripke models for our modal system into arithmetic. In Section 7 we will apply the result of Section 6 to relative interpretability.

**5 The system BMF**

**5.1 Description of the system** BMF is the smallest system, containing the tautologies of propositional logic, closed under modus ponens and the following axioms and rules:

- (L1)  $\vdash \phi \Rightarrow \vdash \Box \phi$
- (L2)  $\vdash \Box(\phi \rightarrow \psi) \rightarrow (\Box \phi \rightarrow \Box \psi)$
- (L3)  $\vdash \Box \phi \rightarrow \Box \Box \phi$
- (L4)  $\vdash \Box(\Box \phi \rightarrow \phi) \rightarrow \Box \phi$
- (B1)  $\vdash \phi \Rightarrow \vdash \Delta \phi$
- (B2)  $\vdash \Delta(\phi \rightarrow \psi) \rightarrow (\Delta \phi \rightarrow \Delta \psi)$
- (B3)  $\vdash \neg \Delta \perp$
- (B4)  $\vdash \Box \phi \rightarrow \Delta \Box \phi$
- (B5)  $\vdash \Delta \phi \rightarrow \Box \Delta \phi$
- (B6)  $\vdash \neg \Box \perp \rightarrow (\Delta \phi \leftrightarrow \Box \phi)$
- (B7)  $\vdash \Box \perp \rightarrow (\Delta \phi \vee \Delta \neg \phi)$ .

This list is very long and rather redundant. A more economical list would consist of (B1), (B2), (B3), (B5), (B7), and the principles

- (B8)  $\vdash \Box \phi \leftrightarrow (\Delta \phi \vee \Box \perp)$
- (B9)  $\vdash \Delta(\Box \phi \rightarrow \phi) \rightarrow \Delta \phi$ .

(B8) is easily derived from (L1), (L2), and (B6). (B9) follows from (L1), (L2), (L4), (B2), (B4), and (B6).

Let me briefly indicate how to derive the long list from the short one: (L1) follows from (B1) and (B8); (L2) from (B2) and (B8); (L4) from (B8) and (B9); and (B6) from (B8). We show how to derive (B4) by a familiar trick

- (a)  $\vdash \phi \rightarrow (\Box(\phi \wedge \Box \phi) \rightarrow (\phi \wedge \Box \phi))$  (L1)
- (b)  $\vdash \Delta \phi \rightarrow \Delta(\Box(\phi \wedge \Box \phi) \rightarrow (\phi \wedge \Box \phi))$  (a),(B1),(B2)
- (c)  $\vdash \Delta \phi \rightarrow \Delta(\phi \wedge \Box \phi)$  (b),(B9)
- (d)  $\vdash \Delta \phi \rightarrow \Delta \Box \phi$  (c),(B1),(B2)
- (e)  $\vdash \neg \Box \perp \rightarrow (\Box \phi \rightarrow \Delta \Box \phi)$  (d),(B6)
- (f)  $\vdash \Box \perp \rightarrow (\Delta \Box \perp \vee \Delta \neg \Box \perp)$  (B7)
- (g)  $\vdash \neg \Delta \neg \Box \perp$  (B3),(B9)
- (h)  $\vdash \Box \perp \rightarrow \Delta \Box \perp$  (f),(g)
- (i)  $\vdash \Box \perp \rightarrow \Box \phi$  (L1),(L2)
- (j)  $\vdash \Delta \Box \perp \rightarrow \Delta \Box \phi$  (i),(B1),(B2)
- (k)  $\vdash \Box \perp \rightarrow \Delta \Box \phi$  (h),(j)
- (l)  $\vdash \Box \phi \rightarrow \Delta \Box \phi$  (e),(h)

Finally, (L3) follows from (B4) and (B8).

We list a few further convenient consequences of BMF:

- (B10)  $\vdash \Delta \phi \rightarrow \Box \phi$
- (B11)  $\vdash \Box \phi \rightarrow \Box \Delta \phi$
- (B12)  $\vdash \Box \perp \rightarrow (\Delta(\phi \vee \psi) \leftrightarrow (\Delta \phi \vee \Delta \psi))$
- (B13)  $\vdash \Box(\phi \leftrightarrow \psi) \rightarrow \Box(\Delta \phi \leftrightarrow \Delta \psi)$  (Provable Extensionality)
- (S)  $\vdash \Box \phi \rightarrow (\psi \leftrightarrow \chi) \Rightarrow$   
 $\vdash \Box \phi \rightarrow (\nu[\psi/p] \leftrightarrow \nu[\chi/p])$ . (Substitution Rule)

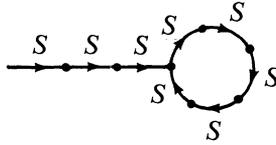
**5.2 Nonuniqueness and nonexplicitness in BMF** Clearly, the discussion in Section 4 on nonuniqueness and nonexplicitness of Henkin and Gödel sentences under mBM can be carried out in BMF; e.g., one can show

$$\begin{aligned} \vdash \Box(p \leftrightarrow \neg \Delta p) &\rightarrow \Box(\nabla p \leftrightarrow \neg \Delta \nabla p) \\ \vdash \Box(p \leftrightarrow \neg \Delta p) &\rightarrow (\Box(p \leftrightarrow \nabla p) \rightarrow \Box \perp). \end{aligned}$$

### 5.3 Kripke Semantics for BMF

#### 5.3.1 Definitions

- (a) Let  $K$  be a finite, nonempty set. Let  $S$  be a binary relation on  $K$ . The structure  $\langle K, S \rangle$  is called a *lolly* if
- (i) for each  $k, k'$  in  $K$ ,  $kS^T k'$ . Here  $S^T$  is the transitive, symmetric, and reflexive closure of  $S$
  - (ii) for each  $k$  in  $K$ , there is *precisely one*  $k'$  in  $K$  such that  $kSk'$ . We will call  $k'$  with  $kSk'$  the *S-successor* of  $k$
  - (iii) there is *at most one*  $k$  in  $K$  such that, for no  $k'$  in  $K$ ,  $k'Sk$ .  
It should be clear that a lolly looks like this:



- (b) A lolly such that for every  $k$  in  $K$  there is a  $k'$  in  $K$  such that  $k'Sk$  is called a *circle*.
- (c) A structure  $\langle K, R, S \rangle$  is called a *lolly-frame* if  $K$  is nonempty,  $R$  and  $S$  are binary relations in  $K$ , and
  - (i)  $R$  is transitive
  - (ii)  $R$  is upwards wellfounded
 Let  $K_0 = \{k \text{ in } K \mid \text{for no } k' \text{ in } K, kRk'\}$   
 $K_1 = K \setminus K_0$   
 $S_0 = S \upharpoonright K_0$   
 $S_0^T =$  the transitive, symmetric, and reflexive closure of  $S_0$ .
  - (iii)  $k \in K_1 \Rightarrow (kSk' \Leftrightarrow kRk')$
  - (iv) Suppose that  $k \in K_0$ . Let  $[k] = \{k' \mid k'S_0^T k\}$ . Then  $\langle [k], S_0 \upharpoonright [k] \rangle$  is a lolly. Moreover, if  $k'Rk$ , then  $k'Rk''$  for all  $k''$  in  $[k]$
  - (v)  $k \in K_0$  and  $kSk' \Rightarrow k' \in K_0$ .
- (d) A *lolly-model* is a structure  $\langle K, R, S, \Vdash \rangle$ , where  $\langle K, R, S \rangle$  is a lolly-frame and  $\Vdash$  is a relation between elements of  $K$  and formulas of the language of BMF, satisfying:
  - (i)  $k \Vdash \top$
  - (ii)  $k \not\Vdash \perp$
  - (iii)  $k \Vdash (\phi \wedge \psi) \Leftrightarrow (k \Vdash \phi \text{ and } k \Vdash \psi)$
  - (iv)  $k \Vdash (\phi \vee \psi) \Leftrightarrow (k \Vdash \phi \text{ or } k \Vdash \psi)$

- (v)  $k \Vdash (\phi \rightarrow \psi) \Leftrightarrow (k \Vdash \phi \Rightarrow k \Vdash \psi)$
- (vi)  $k \Vdash \neg\phi \Leftrightarrow k \not\Vdash \phi$
- (vii)  $k \Vdash \Box\phi \Leftrightarrow$  for all  $k'$  such that  $kRk'$ ,  $k' \Vdash \phi$
- (viii)  $k \Vdash \Delta\phi \Leftrightarrow$  for all  $k'$  such that  $kSk'$ ,  $k' \Vdash \phi$ .

**5.3.2 Remark** It is easy to verify that lolly-frames also satisfy

$$kRk'Sk'' \Rightarrow kRk'', \text{ and } kSk'Rk'' \Rightarrow kRk''.$$

A lolly-frame is best visualized as a conventional frame for provability logic where the top nodes are blown up to lollies, as shown in Figure 2. Here, e.g., Figure 3 means Figure 4. Note that we don't draw the arrows to exhibit the transitivity of  $R$ . Also, since  $R \subseteq S$ , we don't write 'S' next to  $R$ -arrows.

**5.3.3 Soundness** Consider any lolly-model  $\mathbf{K} = \langle K, R, S, \Vdash \rangle$ . We write  $\mathbf{K} \Vdash \phi$  for: for all  $k \in K$ ,  $k \Vdash \phi$ . We then have that  $\text{BMF} \vdash \phi \Rightarrow \mathbf{K} \Vdash \phi$ .

*Proof:* The proof is entirely routine.

**5.3.4 Completeness** Suppose that  $\text{BMF} \not\vdash \phi$ ; then there is a finite lolly-model  $\mathbf{K}$  such that  $\mathbf{K} \not\Vdash \phi$ .

We proceed with some preliminaries for the proof of 5.3.4.

**5.3.5 Definition** Let  $\Gamma$  and  $\Delta$  be sets of formulas of the language of BMF.

- (a)  $\Gamma \vdash \Delta \Leftrightarrow$  there are finite  $\Gamma_0 \subseteq \Gamma$ ,  $\Delta_0 \subseteq \Delta$  such that  $\text{BMF} \vdash \bigwedge \Gamma_0 \rightarrow \bigvee \Delta_0$  (the empty conjunction is  $\top$ , the empty disjunction  $\perp$ )
- (b) Let  $X$  be a set of formulas.  $\Gamma$  is  $X$ -saturated if  $\Gamma$  is consistent and, for each  $\Delta \subseteq X$ ,  $\Gamma \vdash \Delta \Rightarrow$  there is a  $\phi \in \Delta$  such that  $\phi \in \Gamma$ .

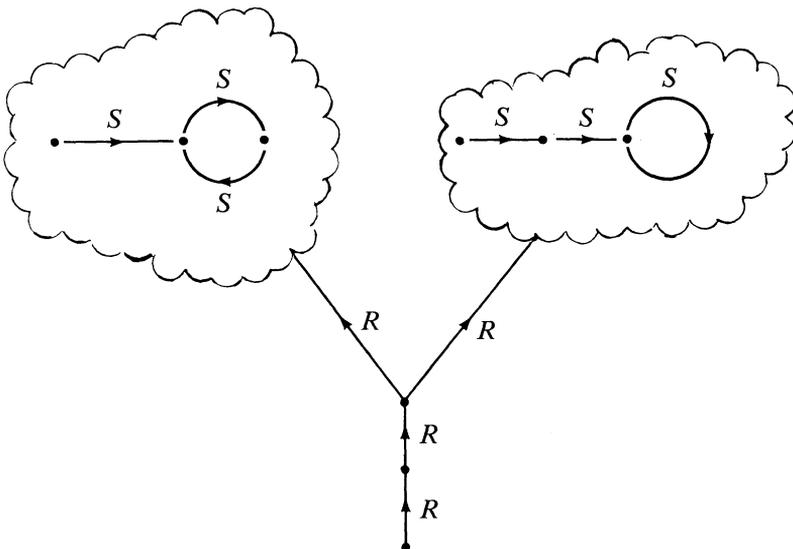


Figure 2.

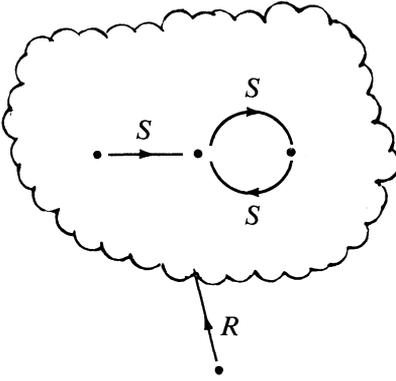


Figure 3.

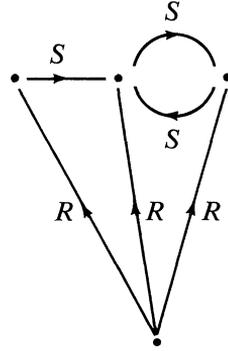


Figure 4.

**5.3.6 Lemma** *Suppose that  $\Gamma \Vdash \Delta$ , and let  $X$  be a set of formulas. There is a set  $Y \subseteq X$  such that  $Y \cup \Gamma$  is  $X$ -saturated and  $Y, \Gamma \Vdash \Delta$ .*

*Proof:* The proof is entirely routine.

**5.3.7 The Henkin construction** Let  $X$  be a finite set of formulas that is closed under subformulas, such that  $\neg \Box \perp \in X$  and  $(\Delta \phi \in X \Leftrightarrow \Box \phi \in X)$ . We now construct a Kripke Model.  $K$ , the set of nodes of the Kripke model we are constructing, consists of those sets of formulas  $y$  such that

- (i)  $y$  is  $X$ -saturated
- (ii) if  $\phi$  is in  $y$  and not in  $X$ , then  $\phi$  is of the form  $\Delta \psi$  and both  $\psi$  and  $\Delta \Delta \psi$  are in  $y$ .

Clearly,  $y$  consists of elements of  $X$  plus, for certain  $\chi$  in  $X \cap y$ ,  $\Delta \chi, \Delta \Delta \chi, \Delta \Delta \Delta \chi, \dots$ . As is easily seen,  $K$  is finite and nonempty (by 5.3.6). For  $x, y \in k$  we define:

$$xRy \Leftrightarrow (\Box \phi \in x \Rightarrow \phi, \Delta \phi, \Delta^2 \phi, \dots \in y) \\ \text{and (there is a } \Box \psi \in y \text{ with } \Box \psi \notin x)$$

$$xSy \Leftrightarrow ((\neg \Box \perp) \in x \text{ and } xRy) \\ \text{or } (\Box \perp \in x, \Box \perp \in y \text{ and } (\Delta \phi \in x \Rightarrow \phi \in y) \\ \text{and } ((\Delta \phi \notin x \text{ and } \Delta \phi \in X) \Rightarrow \phi \notin y)).$$

Finally, we define  $x \Vdash p_i \Leftrightarrow p_i \in x$ .

**Claim 1**  $R$  is transitive and irreflexive (and hence, upwards wellfounded).

**Claim 2**  $xRySz \Rightarrow xRz$ .

The simple proofs of Claims 1 and 2 are left to the reader.

**Claim 3** For  $\phi \in X$ ,  $x \Vdash \phi \Leftrightarrow \phi \in x$ .

*Proof:* We prove this claim by induction on  $\phi$  in  $X$ . The only nontrivial cases are when  $\phi$  is of the form  $\Box \psi$  or  $\Delta \psi$ .

*Case 1:* Suppose that  $\phi = \Box\psi$ .

*Subcase 1.1:* Suppose that  $\Box\psi \in x$  and  $xRy$ . Then  $\psi \in y$ , hence by IH  $y \Vdash \psi$ . Therefore,  $x \Vdash \Box\psi$ .

*Subcase 1.2:* Suppose that  $\Box\psi \notin x$ . Let  $x^R = \{\chi, \Delta\chi, \Delta^2\chi, \dots \mid \Box\chi \in x\}$ . We claim that  $x^R, \Box\psi \not\vdash \psi$ , for otherwise  $\{\Box\chi, \Box\Delta\chi, \dots \mid \Box\chi \in x\} \vdash \Box(\Box\psi \rightarrow \psi)$ ; hence  $\{\Box\chi \mid \Box\chi \in x\} \vdash \Box\psi$  and thus  $\Box\psi \in x$ . *Quod non.* By 5.3.6 there is a set  $x_0 \subseteq X$  such that  $x_0 \cup x^R \cup \{\Box\psi\}$  is  $X$ -saturated and  $x_0 \cup x^R \cup \{\Box\psi\} \not\vdash \psi$ . As is easily seen,  $x_0 \cup x^R \cup \{\Box\psi\} \in K$ . Define  $y = x_0 \cup x^R \cup \{\Box\psi\}$ . Clearly,  $xRy$  and  $\psi \notin y$  and so by IH  $y \not\vdash \psi$ . Hence  $x \not\vdash \Box\psi$ .

*Case 2:* Suppose that  $\phi = \Delta\psi$ . In case  $(\neg\Box\perp) \in x$ , this reduces to the previous case. So we assume that  $\Box\perp \in x$ .

*Subcase 2.1:* Suppose that  $\Delta\psi \in x$  and  $xSy$ . Then  $\psi \in y$ , hence by IH  $y \Vdash \psi$ . So we conclude that  $x \Vdash \Delta\psi$ .

*Subcase 2.2:* Suppose that  $\Delta\psi \notin x$ . Let  $x^S = \{\chi \mid \Delta\chi \in x\} \cup \{\Box\perp\}$  and  $x_S = \{\chi \mid \Delta\chi \notin x, \Delta\chi \in X\}$ . We claim that  $x^S \not\vdash x_S$ , for otherwise  $x \vdash \Delta\mathbb{W}x_S$ , ergo (by (B12), using the fact that  $\Box\perp \in x$ )  $x \vdash \mathbb{W}\{\Delta\chi \mid \chi \in x_S\}$ . Hence  $x \vdash \Delta\chi$  for *some*  $\chi$  in  $x_S$ , which is a contradiction. By 5.3.6 there is an  $x_0 \subseteq X$  such that  $x_0 \cup x^S$  is  $X$ -saturated and  $x_0 \cup x^S \not\vdash x_S$ . Let  $y = x_0 \cup x^S$ . We now show that  $y \in K$ . Suppose that  $\nu \in y$  and  $\nu \notin X$ . Clearly  $\nu \in x^S$ , hence  $\Delta\nu \in x$ . Since  $\Delta\nu \notin X$ ,  $\nu$  and  $\Delta\Delta\nu$  are in  $x$ , and since  $\nu \in x$  and  $\nu \notin X$ ,  $\nu$  must be of the form  $\Delta\rho$ . We conclude that  $\rho$  and  $\Delta\Delta\rho$  are in  $x^S$ . Next we show that  $xSy$ . We have  $\Box\perp \in x$ ,  $\Box\perp \in y$ , and  $\Delta\chi \in x \Rightarrow \chi \in y$ . If  $\Delta\chi \notin x$  and  $\Delta\chi \in X$ , then  $\chi \in x_S$ , so  $\chi \notin y$ . So we conclude that  $xSy$ .

Since  $\psi \in x_S$ , we have that  $\psi \notin y$ . So by IH  $y \not\vdash \psi$  and hence  $x \not\vdash \Delta\psi$ .

**Claim 4** *There is a  $y$  such that  $xRy \Leftrightarrow (\neg\Box\perp) \in x$ .*

**Claim 5** *For every  $x$  there is a  $y$  such that  $xSy$ .*

We leave the simple proofs to the reader. (For the proof of Claim 5, note that  $\Delta\perp \in X$ .)

The model we constructed is not quite a lolly-model yet, so a small transmutation is needed. Consider any  $x$  such that  $\Box\perp \in x$ . Clearly we can produce a sequence  $x = x_0 S x_1 S \dots S x_{n+1}$ , where  $x_i = x_{n+1}$  for some  $i < n + 1$  and where if  $k < j$  and  $x_k = x_j$  then  $k = i$  and  $j = n + 1$ . We define a small lolly model  $L_x$  as follows:  $\langle \{x_0, \dots, x_n\}, R', S', \Vdash' \rangle$ , where

- (i)  $R'$  is empty
- (ii)  $yS'z \Leftrightarrow y = x_j$  and  $z = x_{j+1}$  for some  $j \in \{0, \dots, n\}$
- (iii)  $y \Vdash' p_i \Leftrightarrow p_i \in y$ .

**Claim 6** *For  $y \in \{x_0, \dots, x_n\}$  and  $\phi \in X$ ,  $y \Vdash' \phi \Leftrightarrow y \Vdash \phi$ .*

*Proof:* By induction on  $\phi$  in  $X$  for all  $x_j$  simultaneously. The atomic case and the cases of  $\wedge$ ,  $\vee$ ,  $\neg$ , and  $\rightarrow$  are trivial. If  $\phi$  is  $\Box\psi$  it is sufficient to note that, since  $R'$  is empty,  $x_j \Vdash' \Box\psi$  and that, on the other hand,  $\Box\perp \in x_j$  for each  $x_j$ . Hence by Claim 3,  $x_j \Vdash \Box\perp$ , so  $x_j \Vdash \Box\psi$ .

Suppose that  $\phi = \Delta\psi$ . Note that  $x_i \Vdash \Delta\psi \Rightarrow x_{i+1} \Vdash \psi$ , and  $x_i \not\vdash \Delta\psi \Rightarrow x_{i+1} \not\vdash \psi$ . Hence  $x_i \Vdash' \Delta\psi \Leftrightarrow x_{i+1} \Vdash' \psi \stackrel{\text{IH}}{\Leftrightarrow} x_{i+1} \Vdash \psi \Leftrightarrow x_i \Vdash \Delta\psi$ .

With each  $x$  such that  $\Box \perp \in x$  we associate a small lolly-model  $L_x$  as above. Define

$$\begin{aligned}
 K^* &= \{x \in K \mid \neg \Box \perp \in x\} \cup \{\langle y, x \rangle \mid \Box \perp \in x \in K, \text{ where } y \text{ is in the} \\
 &\quad \text{domain of } L_x\} \\
 K_0^* &= \{\langle y, x \rangle \mid \Box \perp \in x \in K \text{ and } y \text{ is in the domain of } L_x\} \\
 K_1^* &= \{x \in K \mid \neg \Box \perp \in x\} \\
 uR^*v &\Leftrightarrow (u, v \text{ are in } K_1^* \text{ and } uRv) \text{ or} \\
 &\quad (u \text{ is in } K_1^*, v \text{ is in } K_0^*, \text{ where } v = \langle y, x \rangle \text{ and } uRx) \\
 uS^*v &\Leftrightarrow uR^*v \text{ or } (u \text{ is in } K_0^*, v \text{ is in } K_0^*, u = \langle y, x \rangle, v = \langle z, x \rangle \text{ and} \\
 &\quad yS'z, \text{ where } S' \text{ is the relevant relation of } L_x) \\
 u \Vdash^* p_i &\Leftrightarrow (u \in K_1^* \text{ and } p_i \in u) \text{ or } (u \in K_0^*, u = \langle y, x \rangle \text{ and } p_i \in y) \\
 \mathbf{K}^* &= \langle K^*, R^*, S^*, \Vdash^* \rangle.
 \end{aligned}$$

We claim that

(A)  $\mathbf{K}^*$  is a lolly-model

(B) If  $u \in K_1^*$ , then  $(u \Vdash^* \phi \Leftrightarrow u \Vdash \phi)$  for  $\phi \in X$ . If  $u \in K_0^*$  and  $u = \langle y, x \rangle$ , then  $(u \Vdash^* \phi \Leftrightarrow y \Vdash \phi)$  for  $\phi \in X$ .

*Proof of (A):* One easily verifies that  $K^*$  is finite and that  $R^*$  is transitive, irreflexive, and hence upwards wellfounded. Moreover,  $u \in K_1^* \Leftrightarrow$  there is a  $v$  in  $K^*$  such that  $uR^*v$ . The definition of  $S^*$  implies that  $u \in K_1^* \Rightarrow (uS^*v \Leftrightarrow uR^*v)$ . We leave to the reader the easy verification that  $\langle [\langle y, x \rangle], S_0^* \uparrow [\langle x, y \rangle] \rangle$  is isomorphic to the lolly-frame part of  $L_x$  (for  $\langle y, x \rangle$  in  $K_0^*$ ). Clearly, if  $uR^*\langle y, x \rangle$  then  $uR^*\langle z, x \rangle$  for all  $z$  in the domain of  $L_x$ . Also, if  $u \in K_0^*$  and  $uS^*v$ , then  $v \in K_0^*$ .

*Proof of (B):* By induction on  $\phi$  in  $X$ , simultaneously for all  $u$  in  $K^*$ .

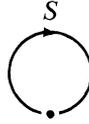
- If  $u \in K_0^*$  and  $u = \langle y, x \rangle$  then  $u \Vdash^* \phi \Leftrightarrow y \Vdash' \phi$   
 $\Leftrightarrow y \Vdash \phi$ .

(The first equivalence is by a completely trivial induction.)

- Suppose that  $u \in K_1^*$ . The case where  $\phi$  is atomic is trivial and so are the cases of  $\wedge$ ,  $\vee$ ,  $\rightarrow$ , and  $\neg$ .
- Suppose that  $\phi = \Box \psi$ 
  - Suppose that  $u \Vdash \Box \psi$  and  $uR^*v$ . If  $v$  is in  $K_1^*$ , we have that  $uRv$ , hence  $v \Vdash \psi$ , so by IH  $v \Vdash^* \psi$ . If  $v$  is in  $K_0^*$ , say  $v = \langle y, x \rangle$ , we have that  $uRx$ . Using Claim 2, we can show that  $uRy$ . It follows that  $y \Vdash \psi$ . Hence by IH  $v \Vdash^* \psi$ , so we conclude that  $u \Vdash^* \Box \psi$ .
  - Suppose that  $u \Vdash^* \Box \psi$  and  $uRy$ . If  $(\neg \Box \perp) \in y$ , we have that  $uR^*y$ ; hence  $y \Vdash^* \psi$ , so by IH  $y \Vdash \psi$ . If  $\Box \perp \in y$ , then  $uR^*\langle y, y \rangle$  and  $\langle y, y \rangle \Vdash^* \psi$ . Hence by IH  $y \Vdash \psi$ , so we conclude that  $u \Vdash \Box \psi$ .
- The case where  $\phi = \Delta \psi$  is similar.

*Proof of 5.3.4:* Suppose that  $\text{BMF} \not\vdash \phi$ . Let  $X_0$  be the smallest set that is closed under subformulas and contains  $\phi$  and  $\neg \Box \perp$ , and let  $X = X_0 \cup \{\Delta \psi \mid \Box \psi \in X_0\} \cup \{\Box \psi \mid \Delta \psi \in X_0\}$ . Construct a finite lolly-model  $\mathbf{K}^*$  as in 5.3.7 for  $X$ . By 5.3.6 there is an  $X$ -saturated  $x_0 \subseteq X$  such that  $x_0 \not\vdash \phi$ .  $x_0$  will correspond to a node of  $\mathbf{K}^*$ , say  $u$ , and  $u \not\vdash^* \phi$ .

**5.3.8 Application** In Section 4 under mBM we showed that neither in BMF nor in PA is there an explicit Gödel sentence for  $\Delta$ , where  $\Delta$  is interpreted in PA as  $\Delta^{\text{mBM}}$ , or  $\Delta^{\text{F}}$ , or  $\Delta^{\text{mF}}$ . In the case of BMF this fact can be easily shown by considering the following Kripke Model:



**6 Embedding circle-tail models in arithmetic** We would like to generalize the result of Solovay [17] to the logic BMF, interpreting  $\Delta$  as  $\Delta^{\text{mF}}$ . To do this we must embed lolly-models in arithmetic. This program, however, meets with a difficulty I could not solve: in a nutshell, the problem is how to handle the sticks of the lollies. It turns out that *if the sticks are absent* a straightforward embedding is possible. For the record I state the obvious open problem:

**Stick problem** Can lolly-models be embedded in arithmetic?

Even if we do not achieve arithmetical completeness for BMF, it seems to me that the partial result proved here is of interest: the Embedding Theorem gives us the powerful machinery to construct arithmetical sentences (see also Section 7). Moreover, the methods employed add to our experience with Solovay-style arguments: we have the first example here of an embedding of structures that are not (completely) upwards well-founded. (In this section I follow the presentation of [20].)

To get a true embedding of circle-models in arithmetic we must add a tail to the circle-models. Consider a finite circle-model; we hang a down-going  $\omega + 1$ -tail (in R) under it, as in Figure 5. We can arrange it so that the nodes of the finite model at the top are numbered  $1, \dots, N$ , and the nodes of the down-going tail (except the bottom)  $N + 1, N + 2, \dots$ , and the bottom is numbered by 0. The nodes numbered  $N + 1, N + 2, \dots, 0$  will be called *tail elements*. We stipulate that at each of the nodes only finitely many atoms are forced and that on all elements of the tail, including 0, the same atoms are forced. We call the resulting models circle-tail models. Clearly, a circle-tail model *is* a circle-model.

An immediate consequence of our definition is

**6.1 Tail Lemma**

- $0 \Vdash \phi \Leftrightarrow$  there is a  $k$  such that, for all  $m > k$ ,  $m \Vdash \phi$
- $0 \nVdash \phi \Leftrightarrow$  there is a  $k$  such that, for all  $m > k$ ,  $m \nVdash \phi$ .

*Proof:* By a simple induction on  $\phi$ .

Let  $\llbracket \phi \rrbracket = \{k \mid k \Vdash \phi\}$ . Then, by the Tail Lemma,  $\llbracket \phi \rrbracket$  is either finite or cofinite.

Note that circle-tail models satisfy the principle:

$$(C) \quad \vdash \Box(\Box \perp \rightarrow \Delta \phi) \rightarrow \Box(\Box \perp \rightarrow \phi).$$

I would be very surprised if (C) were arithmetically valid. A lolly-model to refute (C) can be easily found.

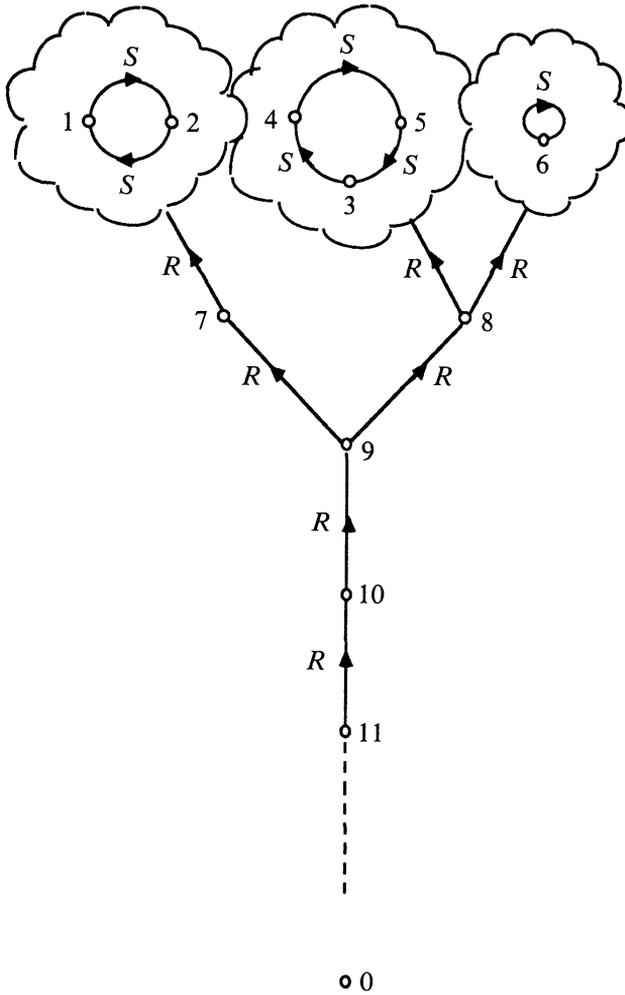


Figure 5

**Open problem** Is (C) arithmetically valid when  $\Delta$  is interpreted as  $\Delta^{\text{mF}}$ ?

For the rest of this section the following are to be kept in mind:

- (i) We fix a circle-tail model  $\mathbf{K}$
- (ii) We assume that  $\mathbf{K}$  is suitably described in arithmetic. Specifically, we assume that  $R$  and  $S$  are given  $\Delta_0$  definitions in such a way that all their simple properties are verifiable
- (iii) We interpret  $\Delta$  as  $\Delta^{\text{mF}}$  in arithmetical contexts
- (iv) ‘ $\vdash$ ’ stands for  $\text{PA}\vdash$
- (v) We assume that ‘Proof’ satisfies the following plausible assumptions:

$$\begin{aligned} &\vdash \Box A \rightarrow \forall x \exists y > x \text{ Proof}(y, A) \\ &\vdash \forall u \forall v (\text{Proof}(u, v) \rightarrow v < u). \end{aligned}$$

We now need to define a variant of Solovay's reluctant function: the function that dares not go anywhere for fear of having to stay. Our variant will not dare to go anywhere for fear of coming there too often. To this end, define by the Recursion Theorem

$$\text{COF}a \Leftrightarrow \forall x \exists y > x \text{ } hy = a$$

$$h0 = 0$$

$$h(k+1) = \begin{cases} a, & \text{if for some } a \text{ such that } hkRa \text{ Proof}(k, \neg \text{COF}a) \\ a, & \text{if for some } a \text{ such that } hkSa \text{ Proof}(k, \neg \text{COF}a) \\ & \text{and } (\neg \text{COF}a) \in \text{mF}_{k+1} \\ hk, & \text{otherwise.} \end{cases}$$

It is easy to see that the arithmetization of ' $(\neg \text{COF}a) \in \text{mF}_{k+1}$ ' is  $\Delta_2$  and hence that  $h$  is  $\Delta_2$ . An important difference with Solovay's original construction is that we use 'COF $a$ ' instead of ' $l = a$ '. Later we will see that  $\vdash \text{COF}a \leftrightarrow l = a$ ; but to show this we need  $h$  to be defined using 'COF' rather than ' $l$ '.

We now prove a sequence of lemmas about  $h$ .

**6.2 Lemma**     *Let  $S^*$  be the transitive reflexive closure of  $S$ . Then*

$$\vdash \forall x \forall y (x \leq y \rightarrow hxS^*hy).$$

*Proof:* By a trivial induction on  $z$  with  $x + z = y$ .

**6.3 Lemma**

- (i)  $\vdash ((\forall z < x \text{ } hz \in K_1) \wedge hx = y) \rightarrow \Delta hx = y$
- (ii)  $\vdash \forall x \forall y (hx = y \rightarrow \square hx = y)$ .

*Proof of (i):* The proof is conducted in PA as follows. The proof is by induction on  $x$ . The case where  $x = \underline{0}$  is trivial. Suppose that  $x = u + \underline{1}$ ,  $\forall z < x \text{ } hz \in K_1$ ,  $hx = y$ , and  $hu = v$ . There are three possibilities:

- (a)  $hx$  was computed by the first clause of the definition of  $h$ . So we then have that  $vRy$  and  $\text{Proof}(u, \neg \text{COF}y)$ , hence,  $\Delta vRy$  and  $\Delta \text{Proof}(u, \neg \text{COF}y)$ . By the induction hypothesis,  $\Delta hu = v$ , so we conclude that  $\Delta hx = y$ .
- (b)  $hx$  was computed by the second clause of the definition of  $h$ . So we then have that  $vSy$  and  $\text{Proof}(u, \neg \text{COF}y)$ . Since  $v \in K_1$ , we also have  $vRy$ . So  $hx$  was also computed by the first clause.
- (c)  $hx$  was computed by the third clause. Clearly,  $hu = v = y = hx$  and  $\forall z < u \text{ } hz \in K_1$ . By the induction hypothesis  $\Delta hu = v$ . Moreover, either for no  $w \text{ } \text{Proof}(u, \neg \text{COF}w)$  or for some  $w \text{ } \text{Proof}(u, \neg \text{COF}w)$  and not  $vRw$  (whence not  $vSw$ , since by hypothesis  $u < x$  implies that  $v \in K_1$ ).

Hence we have  $\Delta \forall w < u \neg \text{Proof}(u, \neg \text{COF}w)$  or  $(\Delta \text{Proof}(u, \neg \text{COF}w)$  and  $\Delta (\neg vRw \wedge \neg vSw))$ , so we conclude that  $\Delta hx = y$ .

*Proof of (ii):* By the argument used in the proof that  $\vdash \Delta A \rightarrow \Box \Delta A$  one shows that  $\vdash \forall u \forall v (u \in \text{mF}_v \rightarrow \Box (u \in \text{mF}_v))$ . The claim now follows by an easy induction in PA.

We define  $\text{LIM}a \Leftrightarrow \exists x \, hx = a \wedge \forall x \forall y ((hx = a \wedge x \leq y) \rightarrow hy = a)$ .

**6.4 Lemma**  $\vdash \forall a (\text{COF}a \rightarrow \text{LIM}a)$ .

*Proof:* By 6.2 it is clear that  $\vdash \forall a ((\text{COF}a \wedge a \in K_1) \rightarrow \text{LIM}a)$ , so it is sufficient to show that  $\vdash \forall a ((\text{COF}a \wedge a \in K_0) \rightarrow \text{LIM}a)$ . We reason in PA as follows:

Suppose that  $\text{COF}a$  and  $a \in K_0$ . Assume that  $a$  is on the circle  $C$  with, say,  $a = \underline{a}_1 S \underline{a}_2 S \dots S \underline{a}_n S \underline{a}_{n+1} = \underline{a}_1$ . Let  $x_0$  be the unique number such that  $hx_0 \in K_1$  and  $h(x_0 + 1) \in K_0$ . Clearly,  $h(x_0 + 1) = \underline{a}_j$  for some  $j$ . By 6.3(i)  $\Delta h(x_0 + 1) = \underline{a}_j$ , so we conclude that  $\Delta \mathbb{W}\{\text{COF}\underline{a}_i \mid i = 1, \dots, n\}$ .

Now suppose for a reductio that  $\neg \text{LIM}a$ . Clearly, by 6.2  $\text{COF}\underline{a}_1, \text{COF}\underline{a}_2, \dots, \text{COF}\underline{a}_n$ . It follows from the definition of  $h$  that  $\Delta \neg \text{COF}\underline{a}_1, \Delta \neg \text{COF}\underline{a}_2, \dots, \Delta \neg \text{COF}\underline{a}_n$  (or how else could  $h$  move on and on?). Hence,  $\Delta \mathbb{W}\{\text{COF}\underline{a}_i \mid i = 1, \dots, n\}$  and  $\Delta \mathbb{M}\{\neg \text{COF}\underline{a}_i \mid i = 1, \dots, n\}$ , therefore  $\Delta \perp$  and thus  $\perp$ . So we conclude that  $\text{LIM}a$ .

**6.5 Lemma**  $\vdash \exists a \text{LIM}a$ .

*Proof:* It is easily seen that  $\vdash \exists a \text{COF}a$ .

**6.6 Lemma**  $\vdash \exists x \, hx \in K_0 \leftrightarrow \Box \perp$ .

*Proof:* Reason in PA as follows:

Right to left is trivial.

From left to right, suppose that  $hx = \underline{a}_1$ , where  $\underline{a}_1$  is on the circle  $C$ , given by  $\underline{a}_1 S \underline{a}_2 S \dots S \underline{a}_n S \underline{a}_{n+1} = \underline{a}_1$ . We have that  $\Box hx = \underline{a}_1$  by 6.3(ii), hence  $\Box \mathbb{W}\{\text{COF}\underline{a}_i \mid i = 1, \dots, n\}$ .  $h$  moved up to  $\underline{a}_1$  by the first or by the second clause. In either case we have  $\Box \neg \text{COF}\underline{a}_1$ .

We now show for  $k = 0, \dots, n - 1$  that  $\mathbb{M}\{\Box \Delta^k \neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}$ , by (external) induction on  $k$ . The case where  $k = 0$  is simply  $\Box \neg \text{COF}\underline{a}_1$ . Suppose that  $\mathbb{M}\{\Box \Delta^k \neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}$ . By (B11)  $\mathbb{M}\{\Box \Delta^{k+1} \neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}$ . We now need to show  $\Box \Delta^{k+1} \neg \text{COF}\underline{a}_{k+2}$ . Clearly

$$\Box ((hx = \underline{a}_1 \wedge \mathbb{M}\{\neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}) \rightarrow \exists y \geq x \, hy = \underline{a}_{k+2})$$

hence

$$\Box ((hx = \underline{a}_1 \wedge \mathbb{M}\{\neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}) \rightarrow \Delta \neg \text{COF}\underline{a}_{k+2}).$$

We conclude using (L1), (L2), (B1), (B2), and (B11) that

$$(\Box \Delta^k hx = \underline{a}_1 \wedge \mathbb{M}\{\Box \Delta^k \neg \text{COF}\underline{a}_j \mid j = 1, \dots, k + 1\}) \rightarrow \Box \Delta^{k+1} \neg \text{COF}\underline{a}_{k+2}.$$

Moreover, by (B11) we have from  $\Box hx = \underline{a}_1$  that  $\Box \Delta^k hx = \underline{a}_1$ . So finally

$$\Box \Delta^{k+1} \neg \text{COF}\underline{a}_{k+2}.$$

We have found that  $\mathbb{M}\{\Box \Delta^{n-1} \neg \text{COF}\underline{a}_j \mid j = 1, \dots, n\}$ . On the other hand, we have  $\Box \mathbb{W}\{\text{COF}\underline{a}_j \mid j = 1, \dots, n\}$ , hence  $\Box \Delta^{n-1} \mathbb{W}\{\text{COF}\underline{a}_j \mid j = 1, \dots, n\}$ . Combining we find  $\Box \Delta^{n-1} \perp$  and hence  $\Box \perp$ .

Consider  $i$  in  $K_0$ . We call the  $S$ -successor of  $i$   $si$ , and the  $S$ -predecessor  $\pi i$ .

### 6.7 Lemma

- (i)  $\vdash(\text{COF}u \wedge uSv) \rightarrow \nabla \text{COF}v$
- (ii)  $\vdash(y \in K_0 \wedge \text{COF}y) \rightarrow \Delta \text{COF}\sigma y$
- (iii)  $\vdash(y \in K_0 \wedge \Box \perp) \rightarrow (\Delta \text{COF}\sigma y \rightarrow \text{COF}y)$
- (iv)  $\vdash(y \in K_0 \wedge \Box \perp) \rightarrow (\text{COF}y \leftrightarrow \Delta \text{COF}\sigma y)$ .

*Proof:*

(i) Reason in PA as follows. Suppose that  $\text{COF}u$ ,  $uSv$ , and  $\Delta \neg \text{COF}v$ . By 6.4 LIM $u$ . Suppose that  $hx = u$  and for all  $y \geq x$   $hy = u$ . For some  $z$   $(\neg \text{COF}v) \in \text{mF}_z$ . Consider  $w$  such that  $w > x$ ,  $w > z$ , and  $\text{Proof}(w, \neg \text{COF}v)$ . Clearly,  $(\neg \text{COF}v) \in \text{mF}_{w+1}$ . Hence  $h$  would move up to  $v$  at  $w + 1$ . *Quod non*. So we conclude that  $\neg \Delta \neg \text{COF}v$ .

(ii) Immediate from (i) using  $\vdash(y \in K_0 \wedge \text{COF}y) \rightarrow \Box \perp$ , which follows directly from 6.6, and  $\vdash \Box \perp \rightarrow (\nabla A \leftrightarrow \Delta A)$ .

(iii) Reason in PA as follows. Suppose that  $y \in K_0$ ,  $\Box \perp$ , and  $\Delta \text{COF}\sigma y$ . From  $\Box \perp$  we have by 6.6 that for some  $z \in K_0$   $\text{COF}z$ . By (ii)  $\Delta \text{COF}\sigma z$ . Hence by 6.4, (B1), (B2), and  $\Delta \text{COF}\sigma y$ ,  $\Delta y = z$  and thus  $y = z$ . (We have  $\Pi_1$ -Reflection for  $\Delta$ !)

(iv) By (ii) and (iii).

### 6.8 Definitions

- (i) Let  $f$  be a function from the propositional variables of the language of BMF to the sentences of PA. We define  $( )^f$  from the formulas of the language of BMF as follows:

- $(p_i)^f = f(p_i)$
- $( )^f$  commutes with the propositional connectives (including  $\top, \perp$ )
- $(\Box \phi)^f = \Box (\phi)^f$  (note that ' $\Box$ ' shifts its meaning!)
- $(\Delta \phi)^f = \Delta (\phi)^f$

- (ii) Consider  $\phi$  in the language of BMF. If  $\llbracket \phi \rrbracket$  is finite, we set

$$[\phi] = \mathbb{W}\{\text{COF}i \mid i \Vdash \phi\} \quad (\text{we take } \mathbb{W}\emptyset = (\underline{0} = \underline{1}))$$

If  $\llbracket \phi \rrbracket$  is cofinite, we set

$$[\phi] = \mathbb{M}\{\neg \text{COF}i \mid i \nVdash \phi\} \quad (\text{we take } \mathbb{M}\emptyset = (\underline{0} = \underline{0}))$$

Note that  $[\phi]$  is simply an arithmetization of  $\exists x \in \llbracket \phi \rrbracket \text{COF}x$

- (iii) Define  $Fp_i = [p_i]$ , and  $\langle \phi \rangle = (\phi)^F$ .

### 6.9 Embedding Theorem $\vdash \langle \phi \rangle \leftrightarrow [\phi]$ .

*Proof:* It is clearly sufficient to show in PA that  $[ ]$  'commutes' with the logical constants, including  $\Box$  and  $\Delta$ . The cases of the propositional constants are trivial (using 6.4). We show that (i)  $\vdash [\Box \psi] \leftrightarrow \Box [\psi]$  and (ii)  $\vdash [\Delta \psi] \leftrightarrow \Delta [\psi]$ .

*Proof of (i):* In case  $\{i \mid i \Vdash \Box \psi\}$  is infinite, we have that  $\llbracket \Box \psi \rrbracket = \llbracket \psi \rrbracket = \omega$ , hence  $[\Box \psi] = [\psi] = (\underline{0} = \underline{0})$ . It follows that  $\vdash [\Box \psi] \leftrightarrow \Box [\psi]$ . Suppose that  $\{i \mid i \Vdash \psi\}$  is finite. Reason in PA as follows:

From right to left, let  $j_1, \dots, j_s$  be the complete set of nodes such that

$j_k \Vdash \Box\psi$  and  $j_k \nVdash \psi$ . Suppose that  $\Box[\psi]$ . Clearly,  $\Box\neg\text{COF}j_k$ . Suppose that  $\text{Proof}(p, \neg\text{COF}j_k)$  and  $hp = y$ . There are two possibilities:

*Case 1:*  $yRj_k$ . If  $yRj_k$ , clearly  $h(p+1) = j_k$ .

*Case 2:*  $\neg yRj_k$ . It follows that if  $\text{COF}x$ , then  $\neg xRj_k$ , for if we had  $yS^*xRj_k$ , it would follow that  $yRj_k$ .

In both cases,  $\text{COF}x \rightarrow \neg xRj_k$ .

On the other hand, it is easily seen that if  $x \nVdash \Box\psi$  then  $xRj_k$  for some  $k$ . Hence,  $\text{COF}x \rightarrow x \Vdash \Box\psi$ , so we conclude that  $\Box[\psi]$ .

From left to right, suppose that  $\text{COF}i$  for some  $i$  such that  $i \Vdash \Box\psi$ . Because  $i \neq 0$ ,  $h$  must have moved up to  $i$  at a certain point by clause 1 or clause 2 of the definition of  $h$ . In either case we have that  $\Box\neg\text{COF}i$ . Suppose that  $hx = i$ . By 6.3(ii),  $\Box hx = i$ . If  $i \in K_0$  we have by 6.6  $\Box\perp$ , and hence  $\Box[\psi]$ . If  $i \in K_1$  we see that  $\Box\neg\text{COF}i$  and  $\Box hx = i$  imply  $\Box\forall y(\text{COF}y \rightarrow iRy)$ , thus we conclude that  $\Box[\psi]$ .

*Proof of (ii):* Clearly,  $\vdash[\neg\Box\perp \rightarrow (\Delta\psi \leftrightarrow \Box\psi)]$ , hence  $\vdash\neg\Box\perp \rightarrow ([\Delta\psi] \leftrightarrow \Box[\psi])$  by the fact that  $[ ]$  ‘commutes’ with the propositional connectives and  $\Box$ . Also,  $\vdash\neg\Box\perp \rightarrow (\Delta[\psi] \leftrightarrow \Box[\psi])$ . So we may conclude that  $\vdash\neg\Box\perp \rightarrow ([\Delta\psi] \leftrightarrow \Delta[\psi])$ .

To complete the argument we need to show that  $\vdash\Box\perp \rightarrow ([\Delta\psi] \leftrightarrow \Delta[\psi])$ .

$$\begin{aligned}
\vdash\Box\perp \rightarrow ([\Delta\psi] \leftrightarrow ([\Delta\psi] \wedge \Box\perp)) & \\
& \leftrightarrow [\Delta\psi \wedge \Box\perp] \\
& \leftrightarrow \mathbb{W}\{\text{COF}j \mid j \Vdash \Delta\psi \wedge \Box\perp\} \\
& \leftrightarrow \mathbb{W}\{\text{COF}\pi i \mid i \Vdash \psi \wedge \Box\perp\} \\
& \leftrightarrow \mathbb{W}\{\Delta\text{COF}i \mid i \Vdash \psi \wedge \Box\perp\} & (6.7(\text{iv})) \\
& \leftrightarrow \Delta\mathbb{W}\{\text{COF}i \mid i \Vdash \psi \wedge \Box\perp\} & (\text{B12}) \\
& \leftrightarrow \Delta[\psi \wedge \Box\perp] \\
& \leftrightarrow \Delta[\psi]. & (\text{B1}), (\text{B2}), (\text{B4})
\end{aligned}$$

**6.10 Remark** The reduction result proved as (33) in Section 4 clearly applies to  $\Delta^{\text{mF}}$ . It implies that for the arithmetical embedding of *traditional* tail models we have that  $\vdash\Delta[\phi] \leftrightarrow (\Box[\phi] \wedge (\Box\perp \rightarrow [\phi]))$ . We can now understand this result in a new way: the arithmetical embedding of traditional tail models is similar to the arithmetical embedding of circle-tail models *which have just singleton circles!* (This point will become even clearer in the light of Lemma 7.3.)

**6.11 Application** There are infinitely many nonequivalent Gödel sentences for  $\Delta^{\text{mF}}$ .

*Proof:* It is clearly sufficient to prove that for any  $n$  there are  $n$  nonequivalent Gödel sentences for  $\Delta$ . Consider the circle-tail model shown in Figure 6.

Let  $s$  be a sequence  $c_1c_2 \dots c_n$  of 0’s and 1’s. Consider an atom  $p_s$ . Let

$$\begin{aligned}
a_{0i} \Vdash p_s & \Leftrightarrow c_i = 0 \\
a_{1i} \Vdash p_s & \Leftrightarrow c_i = 1 \\
b_j \Vdash p_s & \text{ for all } j \\
0 \Vdash p_s & .
\end{aligned}$$

Define  $G_s = [p_s]$ . It follows immediately from the Embedding Theorem that  $\vdash G_s \leftrightarrow \neg\Delta G_s$ . Moreover, if  $s \neq s'$  then  $\vdash\Box(G_s \leftrightarrow G_{s'}) \rightarrow \Box\perp$ .

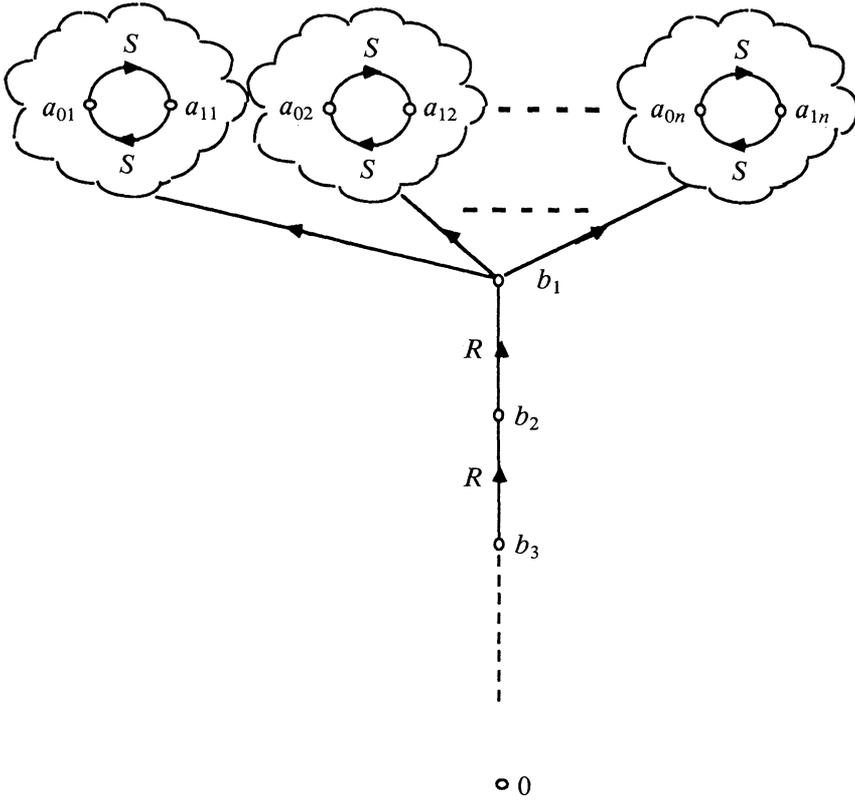


Figure 6.

Because Gödel sentences of  $\Delta$  are Orey sentences it follows that there are infinitely many nonequivalent Orey sentences.

**7  $\Delta^{\text{mF}}$  meets relative interpretability** In this section ‘ $\Delta$ ’ will stand for ‘ $\Delta^{\text{mF}}$ ’ in arithmetical contexts, ‘ $\vdash$ ’ will stand for ‘ $\text{PA}\vdash$ ’. We fix a circle-tail model  $\mathbf{K}$ .

For convenience we now repeat the derivability conditions we collected for relative interpretability in 3.5:

- (I1)  $\vdash \Box (B \rightarrow A) \rightarrow A \triangleleft B$
- (I2)  $\vdash (A \triangleleft B \wedge B \triangleleft C) \rightarrow A \triangleleft C$
- (I3)  $\vdash (A \triangleleft B \wedge A \triangleleft C) \rightarrow A \triangleleft (B \vee C)$
- (I4)  $\vdash A \triangleleft B \rightarrow (\Diamond B \rightarrow \Diamond A)$
- (I5)  $\vdash \Diamond A \triangleleft B \rightarrow \Box (B \rightarrow \Diamond A)$
- (I6)  $\vdash A \triangleleft \Diamond A$
- (I7)  $\vdash A \triangleleft B \rightarrow (A \wedge \Box C) \triangleleft (B \wedge \Box C)$
- (J1) for all  $P$  in  $\Pi_1$ ,  $\text{PA} \vdash P \triangleleft B \rightarrow \Box (B \rightarrow P)$
- (J2) for all  $S$  in  $\Sigma_1$ ,  $\text{PA} \vdash A \triangleleft B \rightarrow (A \wedge S) \triangleleft (B \wedge S)$ .

We add the (for our purposes) essential (34) of Section 4

$$(J3) \quad \vdash A \triangleleft \nabla A.$$

Note that (I5), (I6), and (I7) are redundant in our present list.

We list some immediate consequences of our list:

$$(J4) \quad \vdash A \triangleleft \Delta A. \quad (B1),(B2),(B3),(I1),(J3),(I2)$$

Define  $A \equiv B \Leftrightarrow A \triangleleft B \wedge B \triangleleft A$ .

$$(J5) \quad \vdash (A \equiv A' \wedge B \equiv B') \rightarrow (A \triangleleft B \leftrightarrow A' \triangleleft B') \quad (I2)$$

$$(J6) \quad \vdash \bigwedge \{A_i \triangleleft B_i \mid i = 1, \dots, n\} \\ \rightarrow \mathbb{W} \{A_i \mid i = 1, \dots, n\} \triangleleft \mathbb{W} \{B_i \mid i = 1, \dots, n\} \quad (I1),(I2),(I3)$$

$$(J7) \quad \vdash B \equiv (B \vee \diamond B) \quad (I1),(I6),(I3)$$

$$(J8) \quad \text{If } P \in \Pi_1, \text{ then } \vdash \square((B \vee \diamond B) \leftrightarrow P) \\ \rightarrow (B \triangleleft C \leftrightarrow \square(C \rightarrow (B \vee \diamond B))) \quad (I1),(J5),(J7),(J1)$$

We now wish to take a closer look at the interaction between  $\triangleleft$  and the sentences  $[\phi]$  constructed in Section 6. The classes of sentences  $[\phi]$ , constructed for different circle-tail models, are too poor to refute all modal principles *not* valid in PA in a language with  $\square$  and  $\triangleleft$ . For example, Per Lindström has shown that there is a  $\Sigma_1$  sentence  $A$  such that

$$\not\vdash A \triangleleft \top \rightarrow \square(A \triangleleft \top).$$

On the other hand we will see that

$$\vdash [\phi] \triangleleft \top \rightarrow \square([\phi] \triangleleft \top).$$

This weakness, however, turns out to be a strength:  $[\phi] \triangleleft \top$  *reduces* to a simpler formula. (We encountered the phenomenon of reduction before in connection with Feferman's Predicate.)

We define an *ad hoc* modal operator  $(\ )^*$  as follows:  $\llbracket \phi^* \rrbracket$  is the smallest set  $X$  such that  $\llbracket \phi \rrbracket \subseteq X$ , and if  $j \in X \cap K_0$  then  $\sigma j \in X$ . In other words,  $\llbracket \phi^* \rrbracket$  is obtained by adding to  $\llbracket \phi \rrbracket$  all circles  $C$  such that  $C \cap \llbracket \phi \rrbracket \neq \emptyset$ .

$$\mathbf{7.1 Reduction Theorem} \quad \vdash [\phi] \triangleleft A \leftrightarrow \square(A \rightarrow (\llbracket \phi^* \rrbracket \vee \diamond \llbracket \phi^* \rrbracket)).$$

To prove this we need a few lemmas.

**7.2 Definition** We define a *recursive* function  $h_0$  as follows:

$$h_0 0 = 0$$

$$h_0(k+1) = \begin{cases} a, & \text{if for some } a \text{ such that } h_0 k R a, \text{ Proof}(k, \neg \text{COF} a) \\ hk, & \text{otherwise.} \end{cases}$$

Here 'COF' is as in Section 6. Note that COF is based on  $h$  and not on  $h_0$ .

### 7.3 Lemma

$$(i) \quad \vdash (\forall z < x \ hz \in K_1) \rightarrow hx = h_0 x$$

$$(ii) \quad \text{Let } S^* \text{ be the transitive, reflexive closure of } S; \text{ then } \vdash \forall x \ h_0 x S^* h x.$$

*Proof:* The proof in both cases is by a simple induction on  $x$  in PA. These inductions are much like the proof of 6.3(i).

**7.4 Corollary**  $[\phi]$  is  $\Delta_2$ .

*Proof:* It is clearly sufficient to show that sentences of the form  $\text{COF}i$  are  $\Delta_2$ . In case  $i \in K_0$  we have by 6.9 that  $\vdash \text{COF}i \leftrightarrow \Delta \text{COF}\sigma_i$ , hence  $\text{COF}i$  is in  $\Delta_2$ . In case  $i \in K_1$  we have by 6.4 and 7.3 that

$$\vdash \text{COF}i \leftrightarrow (\exists x h_0x = i \wedge \forall x \forall y ((h_0x = i \wedge x \leq y) \rightarrow h_0 = i)).$$

**7.5 Definition** Consider  $X \subseteq K$ . We call  $X$  *upwards persistent* if  $(i \in X$  and  $iSj) \Rightarrow j \in X$ .

**7.6 Lemma** Suppose that  $\llbracket \phi \rrbracket$  is upwards persistent. Then  $[\phi]$  is provably equivalent to a  $\Sigma_1$  sentence.

*Proof:* In case  $\llbracket \phi \rrbracket$  is infinite this is trivial. So, supposing that  $\llbracket \phi \rrbracket$  is finite, we show that

$$\vdash [\phi] \leftrightarrow \mathbb{W}\{\exists x h_0x = i \mid i \Vdash \phi\}.$$

We reason in PA as follows:

From right to left, suppose that  $h_0x = i$  for  $i \in \llbracket \phi \rrbracket$ .  $iS^*hx$  by 7.3(ii), hence by the upwards persistence of  $\llbracket \phi \rrbracket$ ,  $hx \in \llbracket \phi \rrbracket$ . Thus  $\forall z > x hz \in \llbracket \phi \rrbracket$ , so we conclude that  $[\phi]$ .

From left to right, suppose that  $\text{COF}i$  for  $i \in \llbracket \phi \rrbracket$ . In case  $i \in K_1$  we have by 7.3(i)  $\exists x h_0x = i$ . Suppose that  $i \in K_0$ , say  $i$  is on circle  $C$ . Clearly there is a  $u$  on  $C$  and a  $y$  such that  $hy = u$  and for all  $z < y$ ,  $hz \in K_1$ . By 7.3(i)  $h_0y = u$ .

Then  $\llbracket \phi \rrbracket$  is upwards persistent,  $i$  is in  $\llbracket \phi \rrbracket$ ,  $i$  is on  $C$ , hence  $C \subseteq \llbracket \phi \rrbracket$ . We conclude that  $u \in \llbracket \phi \rrbracket$ , and so  $\exists y h_0y \in \llbracket \phi \rrbracket$ .

**7.7 Lemma** Suppose that  $i$  is on circle  $C$ . Then  $\vdash \text{COF}i \triangleleft \mathbb{W}\{\text{COF}j \mid j \in C\}$ .

*Proof:* Reason in PA as follows: By 6.7(ii) we have that  $\square(\text{COF}\pi i \rightarrow \Delta \text{COF}i)$ ; hence by (I1)  $(\Delta \text{COF}i) \triangleleft \text{COF}\pi i$ . By (J4) and (I2)

$$\text{COF}i \triangleleft \text{COF}\pi i$$

and similarly we have that

$$\begin{aligned} \text{COF}\pi i &\triangleleft \text{COF}\pi^2 i \\ &\vdots \\ \text{COF}\pi^{n-2} i &\triangleleft \text{COF}\pi^{n-1} i. \end{aligned}$$

Here we suppose that  $n$  is the number of elements of  $C$ . By (I1), (I2) and the above we have that

$$\begin{aligned} \text{COF}i &\triangleleft \text{COF}i \\ \text{COF}i &\triangleleft \text{COF}\pi i \\ &\vdots \\ \text{COF}i &\triangleleft \text{COF}\pi^{n-1} i; \end{aligned}$$

hence by (I3)

$$\text{COF}i \triangleleft \mathbb{W}\{\text{COF}\pi^k i \mid 0 \leq k < n\}.$$

In other words

$$\text{COF}i \triangleleft \mathbb{W}\{\text{COF}j \mid j \in C\}.$$

**7.8 Lemma**  $\vdash[\phi] \equiv [\phi^*]$ .

*Proof:* It is immediate that  $\vdash[\phi^*] \triangleleft [\phi]$ , so we need to show that  $\vdash[\phi] \triangleleft [\phi^*]$ . Reason in PA as follows: First, note that by 7.9  $\Box([\phi] \leftrightarrow ([\phi \wedge \Box\perp] \vee [\phi \wedge \neg\Box\perp]))$ . Hence by (J6) and (I1)  $[\phi \wedge \Box\perp] \triangleleft [\phi^* \wedge \Box\perp] \rightarrow [\phi] \triangleleft [\phi^*]$ . It follows that we may restrict ourselves to  $\phi$  with  $\llbracket\phi\rrbracket \subseteq K_0$ . So suppose that  $\llbracket\phi\rrbracket \subseteq K_0$ . Clearly  $\llbracket\phi^*\rrbracket$  consists precisely of those circles  $C$  such that  $\llbracket\phi\rrbracket \cap C$  is not empty. We have by 7.7 and (J6)

$$\mathbb{W}\{\text{COF}_i | i \Vdash \phi\} \triangleleft \mathbb{W}\{\mathbb{W}\{\text{COF}_j | j \in C\} | C \cap \llbracket\phi\rrbracket \neq \emptyset\}.$$

In other words,  $[\phi] \triangleleft [\phi^*]$ .

**7.9 Lemma**  $([\phi^*] \vee \Diamond[\phi^*])$  is provably equivalent to a  $\Pi_1$  sentence.

*Proof:* Note that  $\vdash\neg([\phi^*] \vee \Diamond[\phi^*]) \leftrightarrow [\neg\phi^* \wedge \Box\neg\phi^*]$ . Moreover, as is easily seen,  $\llbracket\neg\phi^* \wedge \Box\neg\phi^*\rrbracket$  is upwards persistent. Apply 7.6, and we are done.

*Proof of 7.1:* We have

$$\begin{aligned} \vdash[\phi] \triangleleft A &\leftrightarrow [\phi^*] \triangleleft A && (7.8),(J5) \\ &\leftrightarrow \Box(A \rightarrow ([\phi^*] \vee \Diamond[\phi^*])). && (7.9),(J8) \end{aligned}$$

**7.10 Corollary**  $\vdash[\phi] \triangleleft \top \leftrightarrow \Box(\Box\perp \rightarrow [\phi^*])$ .

*Proof:* We leave it as an exercise to the reader to show that

$$\vdash\Box(A \vee \Diamond A) \leftrightarrow \Box(\Box\perp \rightarrow A).$$

**7.11 On a question of Orey** Orey asks: for which sets  $\Gamma$  of propositional formulas in the variables  $p_1, \dots, p_n$  are there arithmetical sentences  $B_1, \dots, B_n$  such that  $\Gamma = \{\phi | \phi(B_1, \dots, B_n) \triangleleft \top\}$ ? (I learned this formulation of Orey's problem from Per Lindström. Actually, the question is asked for arbitrary essentially reflexive theories  $T$ . I think that inspection of the argument of this paper shows that the answer given here applies to consistent essentially reflexive RE theories  $T$  into which PA restricted to  $\Sigma_2$ -induction can be translated.)

Let us say that  $\{\phi | \phi(B_1, \dots, B_n) \triangleleft \top\}$  is *the interpretability class* of  $B_1, \dots, B_n$ . A moment's reflection shows that interpretability classes  $\Gamma$  should satisfy

- (i)  $\top \in \Gamma$
- (ii)  $\perp \notin \Gamma$
- (iii)  $\phi \vdash_{\text{Prop}} \psi$  and  $\phi \in \Gamma \Rightarrow \psi \in \Gamma$ .

We will show that, conversely, every set  $\Gamma$  of propositional formulas in  $p_1, \dots, p_n$  satisfying (i), (ii), and (iii) is an interpretability class.

*Proof:* Let  $\Gamma$  be a class of propositional formulas in  $p_1, \dots, p_n$  satisfying (i), (ii), and (iii). The plan of the proof is to construct a circle-tail model  $\mathbf{K}$  and to take  $B_i = [p_i]$ . 7.10 tells us that what happens below the circles is really irrelevant, so we start by stipulating an arbitrary tail, say  $b_0 \dots b_3 R b_2 R b_1$ , where no atom is true at the nodes  $b_j$ . We then proceed to construct the circles.

$\Gamma^C = \{\phi | \phi$  is a propositional formula in the variables  $p_1, \dots, p_n$  and  $\phi \notin \Gamma\}$ . Note that  $\Gamma$  and  $\Gamma^C$  are both closed under provable equivalence in propositional logic (in the language based on  $p_1, \dots, p_n$ ). Let  $\phi_1, \dots, \phi_k$  be represen-

tatives of the equivalence classes of  $\Gamma$  and let  $\psi_1, \dots, \psi_m$  be representatives of the equivalence classes of  $\Gamma^C$ . Define:

$$K_0 = \{\langle i, j \rangle \mid 1 \leq i \leq k, 1 \leq j \leq m\}$$

$$\langle i, j \rangle S \langle i', j' \rangle \Leftrightarrow j = j' \text{ and } ((1 \leq i < k \text{ and } i' = i + 1) \\ \text{or } (i = k \text{ and } i' = 1)).$$

Let us say that the nodes  $\langle i, j \rangle$  for fixed  $j$  form a circle  $C_j$ .

Consider a node  $\langle i, j \rangle$ . Clearly  $\phi_i \not\vdash_{\text{PROP}} \psi_j$ , so there is an assignment  $f$  of truth values to  $p_1, \dots, p_n$  under which  $\phi_i$  is true and  $\psi_j$  is false. Pick such an assignment  $f$  and put  $\langle i, j \rangle \Vdash p_s \Leftrightarrow fp_s = \top$ .

Clearly, on every circle  $C_j$  there is a node  $\langle i, j \rangle$  such that  $\langle i, j \rangle \Vdash \phi_i$ . Hence  $(\Box \perp \rightarrow \phi_i^*)$  and  $\Box(\Box \perp \rightarrow \phi_i^*)$  are forced everywhere in the model.

On the other hand, no node  $\langle i, j \rangle$  on  $C_j$  forces  $\psi_j$ , hence  $\langle i, j \rangle \not\Vdash \Box \perp \rightarrow \psi_j^*$ . It follows that  $(\Box(\Box \perp \rightarrow \psi_j^*) \rightarrow \Box \perp)$  is forced everywhere in the model.

Put  $B_s = [p_s]$ . Note that for any propositional formula  $\chi$  in  $p_1, \dots, p_n$  we have that  $\chi(B_1, \dots, B_n) = \langle \chi \rangle$ . We have by 6.9 that

$$\begin{aligned} \vdash \Box(\Box \perp \rightarrow [\phi_i^*]) &\Rightarrow & (7.10) \\ \vdash [\phi_i] \triangleleft \top &\Rightarrow & (7.9), (I1), (J5) \\ \vdash \langle \phi_i \rangle \triangleleft \top &\Rightarrow & \\ \vdash \phi_i(B_1, \dots, B_n) \triangleleft \top &\Rightarrow & \text{(Reflection Principle)} \\ \phi_i(B_1, \dots, B_n) \triangleleft \top. & & \end{aligned}$$

Moreover, by 6.9

$$\begin{aligned} \vdash \Box(\Box \perp \rightarrow [\psi_j^*]) \rightarrow \Box \perp &\Rightarrow & (7.10) \\ \vdash [\psi_j] \triangleleft \top \rightarrow \Box \perp &\Rightarrow & (6.9), (I1), (J5) \\ \vdash \langle \psi_j \rangle \triangleleft \top \rightarrow \Box \perp &\Rightarrow & \\ \vdash \psi_j(B_1, \dots, B_n) \triangleleft \top \rightarrow \Box \perp &\Rightarrow & \\ \neg(\psi_j(B_1, \dots, B_n) \triangleleft \top). & & \text{(Reflection Principle)} \end{aligned}$$

Note that the uses of the Reflection Principle are eliminable here: we could just have proved the necessary lemmas externally, i.e., in nonformalized form. (In case the theory under consideration is not PA it may even be necessary to reason externally.)

It follows immediately that  $\Gamma = \{\phi \mid \phi(B_1, \dots, B_n) \triangleleft \top\}$ .

**7.12 Remark** Note that in the proof of 7.11 it would have sufficed to consider representatives  $\phi_i$  of the equivalence classes in  $\Gamma$  that are *minimal* in the implication ordering. Similarly, we need only consider representatives  $\psi_j$  of the equivalence classes of  $\Gamma^C$  that are *maximal* in the implication ordering.

## REFERENCES

- [1] Auerbach, D. D., "Intentionality and the Gödel theorems," *Philosophical Studies*, vol. 48 (1985), pp. 337-351.
- [2] Bernardi, C. and F. Montagna, "Equivalence relations induced by extensional formulae: classification by means of a new fixed point property," *Rapporto Matematico* 63, Dipartimento di Matematica, Via del Capitano 15, 53100 Siena, Italia.

- [3] Bowie, G. L., "Lucas' number is finally up," *Journal of Philosophical Logic*, vol. 11 (1982), pp. 279–285.
- [4] Feferman, S., "Arithmetization of metamathematics in a general setting," *Fundamenta mathematica*, vol. 49 (1960), pp. 35–92.
- [5] Guaspari, D. and R. M. Solovay, "Rosser sentences," *Annals of Mathematical Logic*, vol. 16 (1979), pp. 81–99.
- [6] Hájek, P., "Experimental logics and  $\Pi_3$ -theories," *The Journal of Symbolic Logic*, vol. 42 (1977), pp. 515–522.
- [7] Hilbert, D. and P. Bernays, *Grundlagen der Mathematik*, 2nd edition, Springer-Verlag, Berlin, 1970.
- [8] Jeroslow, R. G., "Experimental logics and  $\Delta_2$ -theories," *Journal of Philosophical Logic*, vol. 4 (1975), pp. 253–267.
- [9] Lindström, P., "Some results on interpretability," pp. 329–361 in *Proceedings of the 5th Scandinavian Logic Symposium*, Aalborg, 1979.
- [10] Lucas, J. R., "Minds, machines and Gödel," *Philosophy*, vol. 36 (1961), pp. 120–124.
- [11] Montagna, F., "On the algebraization of Feferman's predicate," *Studia Logica*, vol. 37 (1978), pp. 221–263.
- [12] Montagna, F., "Provability in finite subtheories of PA and relative interpretability: A modal investigation," *The Journal of Symbolic Logic*, vol. 52 (1987), pp. 494–511.
- [13] Orey, S., "Relative interpretations," *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik*, vol. 7 (1961), pp. 146–153.
- [14] Putnam, H., "Trial and error predicates and the solution to a problem of Mostowski," *The Journal of Symbolic Logic*, vol. 30 (1965), pp. 49–57.
- [15] Smoryński, C., "Modal logic and self-reference," pp. 441–496 in *Handbook of Philosophical Logic*, edited by D. Gabbay and F. Guenther, Reidel, Boston, 1984.
- [16] Smoryński, C., *Self-reference and Modal Logic*, Springer-Verlag, New York, 1985.
- [17] Solovay, R. M., "Provability interpretations of modal logic," *Israel Journal of Mathematics*, vol. 25 (1976), pp. 287–304.
- [18] Švejdar, "Degrees of interpretability," *Commentationes Mathematicae Universitatis Carolinae* 19, pp. 789–813.
- [19] Švejdar, "Modal analysis of generalized Rosser sentences," *The Journal of Symbolic Logic*, vol. 48 (1983), pp. 986–999.
- [20] Visser, A., "The provability logics of recursively enumerable theories extending Peano arithmetic at arbitrary theories extending Peano arithmetic," *Journal of Philosophical Logic*, vol. 13 (1984), pp. 79–113.
- [21] Webb, J., *Mechanism, Mentalism and Metamathematics*, Reidel, Dordrecht, 1980.

Central Interfaculty Department of Philosophy  
 State University of Utrecht  
 3508 TC Utrecht  
 The Netherlands