

# On Modal Systems with Rosser Modalities

Vítězslav Švejdar\*

Appeared in M. Bílková and O. Tomala eds., *The Logica Yearbook 2005: Proc. of the Logica 05 Int. Conference*, pp. 203–214, Philosophia Praha, 2006.

## 1 Introduction

Sufficiently strong axiomatic theories allow for the construction of self-referential sentences, i.e. sentences saying something about themselves. After the Gödel's paper on incompleteness (Gödel, 1931) the *self-reference* method found further applications—some are listed below—and became even more important. Around say 1970 it appeared that the self-reference method was not only a useful tool, but also an interesting field of study: it became clear that reasoning about self-referential sentences could be made more transparent and limitations of the self-reference method could be clarified by using *modal logic*. The connections of meta-mathematics to modal logic, especially after the Solovay's paper (Solovay, 1976), brought traditional modal logic to the attention of more mathematically oriented logicians and constituted a stimulus in modal logical studies.

In this paper we briefly discuss some of the existing modal systems, putting an emphasis on Rosser modalities, also known as *witness comparison* modalities, and present one new system of that kind. First we fix some symbolism and the way we speak about arithmetic and self-reference.

## 2 Some preliminaries

By *arithmetical language* we mean the language  $\{+, \cdot, 0, S, \leq, <\}$  with two binary function symbols, a constant, a unary function and two binary predicate symbols. Its intended realization, or the *standard model*, is the structure  $\mathbf{N} = \langle \mathbf{N}, +^{\mathbf{N}}, \cdot^{\mathbf{N}}, 0^{\mathbf{N}}, S, \leq^{\mathbf{N}}, <^{\mathbf{N}} \rangle$ , where  $\mathbf{N}$  is the set of all natural numbers (with zero),  $+^{\mathbf{N}}, \cdot^{\mathbf{N}}, \leq^{\mathbf{N}}, <^{\mathbf{N}}$  are the usual addition, multiplication, unstrict

---

\*This work is a part of the research plan MSM 0021620839 that is financed by the Ministry of Education of the Czech Republic.

and strict ordering on the set  $\mathbf{N}$ ,  $0^{\mathbf{N}}$  is the number zero, and  $s$  is the successor function  $a \mapsto a + 1$ . In the sequel we omit the superscript  $\mathbf{N}$ ; so e.g. “+” can be both a symbol and its standard realization, i.e. addition of natural numbers.

An example of an arithmetical formula is  $\exists v(v \cdot x = y)$ ; it can be read **the number  $x$  is a divisor of the number  $y$**  and denoted  $x \mid y$ . Another example, constructed using the formula  $x \mid y$ , is  $\forall z(z \mid x \rightarrow z = S(0) \vee S(S(0)) \mid z)$ . This formula says that each divisor of  $x$ , except the trivial divisor one, is even. Terms like  $S(S(0))$  are called *numerals*; the numeral  $S(S(\dots S(0)\dots))$  with  $n$  occurrences of the symbol  $S$  is denoted  $\bar{n}$ . In the arithmetical language, numerals make it possible to speak about particular numbers.

We sometimes use the **sans-serif** font for informal reading of syntactical objects. Whenever this font is used, the reader is expected to think a little bit about the formula (or, sometimes, the proof) it represents.

The sentence  $\bar{n} \mid \bar{m}$  resulting by substituting the numerals  $\bar{n}$  and  $\bar{m}$  for  $x$  and  $y$  in the formula  $x \mid y$  is valid in the structure  $\mathbf{N}$  if and only if, in reality,  $n$  is a divisor of  $m$ . This property of the formula  $x \mid y$  is expressed by saying that the formula  $x \mid y$  *defines* the divisibility relation in  $\mathbf{N}$ . One can check that our second example formula  $\forall z(z \mid x \rightarrow z = \bar{1} \vee \bar{2} \mid z)$  defines the set of all powers of two. So let us choose  $\text{Pow}(x)$  as a shorthand for this formula; we have e.g.  $\mathbf{N} \models \text{Pow}(\bar{8})$  and  $\mathbf{N} \not\models \text{Pow}(\bar{9})$ .

We identify formulas and other syntactical objects with their numerical codes under some fixed *coding of syntax*. So if, for example,  $\varphi$  is an arithmetical sentence then  $\text{Pow}(\bar{\varphi})$  is a sentence saying that (the numerical code of) the sentence  $\varphi$  is a power of two. We do not (need to) know whether this sentence is valid in  $\mathbf{N}$ , it may well depend on details of the chosen coding.

If  $\varphi$  is a formula then  $\forall v \leq z \varphi$  is a shorthand for  $\forall v(v \leq z \rightarrow \varphi)$ . Similarly, the formulas  $\forall v < z \varphi$ ,  $\exists v \leq z \varphi$ , and  $\exists v < z \varphi$  have an obvious meaning. The expressions  $\forall v \leq z$ ,  $\forall v < z$ ,  $\exists v \leq z$ , and  $\exists v < z$  are *bounded quantifiers*. A formula is *bounded*, or a  $\Delta_0$ -*formula*, if all its quantifiers are bounded. A formula is a  $\Sigma_1$ -*formula* if it has the form  $\exists x \varphi$  where  $\varphi \in \Delta_0$ ; a formula is a  $\Sigma$ -*formula* if it is obtained from bounded formulas using (any number of) conjunctions, disjunctions, bounded quantifiers and existential quantifiers. Since, for example, the divisibility relation is definable also by the formula  $\exists v \leq y(v \cdot x = y)$ , we may think that the formula  $x \mid y$  is bounded. Similarly, the quantifier  $\forall z$  in the formula  $\text{Pow}(x)$  can be written as  $\forall z \leq x$ ; so  $\text{Pow}(x)$  is another example of a bounded formula.

It is not straightforward to *characterize* the sets definable by  $\Delta_0$ -formulas. However, a characterization of sets definable by  $\Sigma$ -formulas is known. A deep theorem says that, although the set of all  $\Sigma_1$ -formulas is a proper subset of the set of all  $\Sigma$ -formulas, the classes of all sets definable by  $\Sigma_1$ -formulas and by  $\Sigma$ -formulas coincide: a set  $A \subseteq \mathbf{N}^k$  is definable by a  $\Sigma_1$ -formula iff it is definable by a  $\Sigma$ -formula iff it is *recursively enumerable* (r.e.).

Since the set  $\text{Thm}(T)$  of all sentences provable in a recursively axiomatizable theory  $T$  is recursively enumerable, one can choose an arithmetical formula  $\text{Pr}_T(x) \in \Sigma_1$  which defines the set  $\text{Thm}(T)$  in  $\mathbf{N}$ , i.e. satisfies the equivalence  $T \vdash \varphi \Leftrightarrow \mathbf{N} \models \text{Pr}_T(\overline{\varphi})$  for each sentence  $\varphi$  in the language of  $T$ . The formula  $\text{Pr}_T(x)$  can be read the sentence  $x$  is provable in  $T$ ; it is called a *provability predicate* of the theory  $T$ .

*Peano arithmetic* PA is a theory formulated in the arithmetical language; its axiom set consists of seven (eight, or nine, in some sources) single axioms and one schema, the *induction schema*. The  $\Sigma$ -*completeness theorem* says that each  $\Sigma$ -sentence valid in  $\mathbf{N}$  is provable in PA. A simple consequence of the  $\Sigma$ -completeness theorem is the following. If  $\text{Pr}_{\text{PA}}(x)$  is any  $\Sigma_1$ -formula which defines the set  $\text{Thm}(\text{PA})$  of all sentences provable in PA then, besides the equivalence  $\text{PA} \vdash \varphi \Leftrightarrow \mathbf{N} \models \text{Pr}_{\text{PA}}(\overline{\varphi})$  for any arithmetical sentence  $\varphi$ , it also satisfies the implication  $\text{PA} \vdash \varphi \Rightarrow \text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\varphi})$ . This implication is known as the *first derivability condition D1*; it is a property shared by all provability predicates. There are of course many (true, i.e. valid in  $\mathbf{N}$ ) sentences that are provable in PA but the provability of which is not obtained by an appeal to the  $\Sigma$ -completeness theorem. An example is the sentence for each  $x$  there is a prime number  $y$  greater than  $x$ ; this sentence is not  $\Sigma$ .

The sentence  $\neg \text{Pr}_{\text{PA}}(\overline{0 = S(0)})$ , saying that contradiction is not provable in PA, is denoted  $\text{Con}(\text{PA})$  and called a *consistency statement* of PA (based on the formula  $\text{Pr}_{\text{PA}}(x)$ ). Since  $\text{PA} \not\vdash 0 = S(0)$ , from the fact that the formula  $\text{Pr}_{\text{PA}}(x)$  defines the set  $\text{Thm}(\text{PA})$  we have  $\mathbf{N} \models \text{Con}(\text{PA})$ .

Given a notion like being a power of two or being a (numerical code of a) sentence provable in PA, it is natural to ask the following question: if  $\psi$  is a formula that defines that notion in  $\mathbf{N}$ , can we take for granted that “canonical facts” about that notion, if expressed in the arithmetical language using the formula  $\psi$ , are provable in PA? The answer is no in general, but yes in many cases and if the formula  $\psi$  is chosen properly. A canonical fact about powers of two is that a number  $x$  is a power of two if and only if  $x = \overline{1}$  or  $x = \overline{2} \cdot y$  where  $y$  is a power of two. If powers of two are represented by the formula  $\text{Pow}(x)$  defined above then indeed, this fact is provable in PA. It is however not so difficult to find a formula  $\psi(x)$  which also defines the set of all powers of two but does not have this additional property.

As to formalized provability, the canonical facts, i.e. the desired additional properties of the formula  $\text{Pr}_{\text{PA}}(x)$ , are the *second and third derivability conditions D2 and D3*, where D2 is  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\varphi \rightarrow \psi}) \rightarrow (\text{Pr}_{\text{PA}}(\overline{\varphi}) \rightarrow \text{Pr}_{\text{PA}}(\overline{\psi}))$  and D3 is  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\varphi}) \rightarrow \text{Pr}_{\text{PA}}(\overline{\text{Pr}_{\text{PA}}(\overline{\varphi})})$ . With some effort one can show that the “normal” choice of the formula  $\text{Pr}_{\text{PA}}(x)$  defining the set  $\text{Thm}(\text{PA})$  satisfies the conditions D2 and D3 (i.e. satisfies all derivability conditions D1–D3) for all sentences  $\varphi$  and  $\psi$ . The condition D3 is in fact a consequence of a more general fact, namely the *formalized  $\Sigma$ -completeness theorem* saying that  $\text{PA} \vdash \sigma \rightarrow \text{Pr}_{\text{PA}}(\overline{\sigma})$  for each  $\Sigma$ -sentence  $\sigma$ . The last thing we need to

mention is the proof predicate. One can assume that the provability predicate  $\text{Pr}_{\text{PA}}(x)$  has the form  $\exists y \text{Proof}_{\text{PA}}(x, y)$  where  $\text{Proof}_{\text{PA}}(x, y) \in \Delta_0$ . The formula  $\text{Proof}_{\text{PA}}(x, y)$  is a *proof predicate* of PA; it can be read the number  $y$  is a proof of the sentence  $x$ .

### 3 Meta-mathematics and self-reference

The *self-reference theorem* says that for any arithmetical formula  $\psi(x)$  there exists a sentence  $\varphi$  such that  $\text{PA} \vdash \varphi \equiv \psi(\overline{\varphi})$ . The sentence  $\varphi$  satisfying  $\text{PA} \vdash \varphi \equiv \psi(\overline{\varphi})$ , i.e. provably equivalent to the sentence **the sentence  $\varphi$  has the property  $\psi$** , can naturally be viewed as a first-person sentence, saying **I have the property  $\psi$** . It is often convenient to view the expression  $\text{PA} \vdash \varphi \equiv \psi(\overline{\varphi})$  as an *equation*, determined by the formula  $\psi$ , with an unknown sentence  $\varphi$ . So the self-reference theorem ensures that every self-referential equation has a solution.

A prominent example on the use of the self-reference theorem is *Gödel sentence* saying **I am not provable** in PA, i.e. a sentence  $\nu$  satisfying the condition  $\text{PA} \vdash \nu \equiv \neg \text{Pr}_{\text{PA}}(\overline{\nu})$ . A second important example is *Rosser sentence*  $\rho$  satisfying

$$\text{PA} \vdash \rho \equiv \exists y (\text{Proof}_{\text{PA}}(\overline{\neg \rho}, y) \ \& \ \forall v \leq y \neg \text{Proof}_{\text{PA}}(\overline{\rho}, v)), \quad (1)$$

i.e. saying there is a proof of my negation such that beneath it there is no proof of myself. Both Gödel and Rosser sentences demonstrate incompleteness of Peano arithmetic. An important distinction between these two sentences is that while independence of Rosser sentence is formalizable in PA, viz  $\text{PA} \vdash \text{Con}(\text{PA}) \rightarrow \neg \text{Pr}_{\text{PA}}(\overline{\rho}) \ \& \ \neg \text{Pr}_{\text{PA}}(\overline{\neg \rho})$ , only one half of independence of Gödel sentence can be formalized:  $\text{PA} \vdash \text{Con}(\text{PA}) \rightarrow \neg \text{Pr}_{\text{PA}}(\overline{\nu})$ , but  $\text{PA} \not\vdash \text{Con}(\text{PA}) \rightarrow \neg \text{Pr}_{\text{PA}}(\overline{\neg \nu})$ . Since my point is that properties of self-referential sentences can be obtained using modal logic, no proofs are given here.

*L. Henkin* asked a question whether a sentence saying **I am provable** must be provable. In the equational setting, Henkin's question concerns the equation  $\text{PA} \vdash \kappa \equiv \text{Pr}_{\text{PA}}(\overline{\kappa})$  for an unknown sentence  $\kappa$ . Since this equation has a trivial (provable) solution  $0 = 0$ , Henkin's question should be understood as a question whether *any* Henkin sentence is provable, i.e. whether the Henkin equation has, up to provable equivalence, only one solution. The question brought the problem of uniqueness of self-referential equations to the attention of logicians and was answered positively in the influential paper by [Löb \(1955\)](#).

Let  $\text{PA} \upharpoonright y + x$  denote the theory whose axioms are the sentence  $x$  plus all axioms of PA which are less than  $y$ . Our fourth example of a self-referential

equation is

$$\text{PA} \vdash \eta \equiv \exists y(-\text{Con}(\text{PA} \uparrow y + \bar{\eta}) \& \text{Con}(\text{PA} \uparrow y + \overline{\neg\eta})). \quad (2)$$

Any sentence  $\eta$  satisfying this equation has the property that both theories  $\text{PA} + \eta$  and  $\text{PA} + \neg\eta$  are *interpretable* in PA, in the syntactical sense of [Tarski, Mostowski, and Robinson \(1953\)](#). This construction of a symmetrically interpretable sentence is mentioned e.g. in [Švejdar \(1983\)](#), but is in fact distilled out from [Hájek and Hájková \(1972\)](#); so the sentence  $\eta$  could be called *Hájek sentence*. However it should be noted that other self-referential constructions named by Petr Hájek exist.

## 4 Self-reference and modal logic

A crucial idea, connecting meta-mathematics and modal logic, is interpreting the usual modal symbol  $\Box$ , the necessity operator, by formalized provability expressed by the formula  $\text{Pr}_{\text{PA}}(x)$ . We directly proceed to modal systems with two additional “modalities” that are not met in traditional modal logical studies and that are closely related one to another, the Rosser symbols  $\preceq, \prec$ .

So we consider propositional modal language with propositional atoms, the symbol  $\perp$  for falsity, connectives  $\rightarrow, \&, \vee$ , the unary modality  $\Box$ , and two binary *Rosser* (or *witness comparison*) modalities  $\preceq$  and  $\prec$ . *Modal formulas* are built up from atoms and  $\perp$  using connectives and modalities, with the restriction that  $\preceq, \prec$  are applicable only to formulas starting with  $\Box$ . So  $\Box(\Box p \rightarrow p)$  or  $q \& (\Box p \preceq \Box \perp)$  are examples of modal formulas. We use  $\neg A, \top, A \equiv B$ , and  $\diamond A$  as shorthands for  $A \rightarrow \perp, \perp \rightarrow \perp, (A \rightarrow B) \& (B \rightarrow A)$ , and  $\neg \Box \neg A$  respectively. For omitting parentheses, we assign the symbols  $\preceq, \prec$  higher priority than binary connectives; among connectives, implication  $\rightarrow$  has higher priority than equivalence  $\equiv$ , but lower than conjunction  $\&$  and disjunction  $\vee$ . So  $p \equiv q \vee \Box p \prec \Box q$  is the same formula as  $p \equiv (q \vee (\Box p \prec \Box q))$ .

To define the *arithmetical semantics* for modal formulas with Rosser modalities we need two auxiliary notions, standard proof predicate and arithmetical evaluation. A *standard proof predicate* is an arithmetical formula  $\text{Prf}(x, y)$  satisfying  $\text{PA} \vdash \forall x(\exists y \text{Prf}(x, y) \equiv \text{Pr}_{\text{PA}}(x))$  and such that both  $\text{Prf}(x, y)$  and  $\neg \text{Prf}(x, y)$  are PA-equivalent to a  $\Sigma_1$ -formula (i.e.,  $\text{Prf}(x, y)$  is a  $\Delta_1(\text{PA})$ -formula). The **formula**  $\text{Proof}_{\text{PA}}(x, y)$  is an example of a standard proof predicate. It is evident that if  $\text{Prf}(x, y)$  is a standard proof predicate then the formula  $\exists y \text{Prf}(x, y)$ , which may be called a provability predicate associated with that proof predicate, satisfies the **derivability conditions D1–D3**. A function  $e$  from modal formulas to arithmetical sentences is an (*arithmetical*) *evaluation* based on a standard proof predicate  $\text{Prf}(x, y)$  if it commutes with logical connectives ( $e(A \& B) = e(A) \& e(B)$ , etc.) and satisfies

$e(\perp) = (0 = S(0))$  and

$$\begin{aligned} e(\Box A) &= \exists y \text{Prf}(\overline{e(A)}, y), \\ e(\Box A \preceq \Box B) &= \exists y (\text{Prf}(\overline{e(A)}, y) \ \& \ \forall v < y \neg \text{Prf}(\overline{e(B)}, v)), \\ e(\Box A \prec \Box B) &= \exists y (\text{Prf}(\overline{e(A)}, y) \ \& \ \forall v \leq y \neg \text{Prf}(\overline{e(B)}, v)). \end{aligned}$$

Thus the modal formulas  $\Box A \preceq \Box B$  and  $\Box A \prec \Box B$  can be read “ $A$  has a proof which is less than or equal to (or less than, respectively) any possible proof of  $B$ ”. It is natural to read the formula  $\Box A$  as “ $A$  is provable” rather than “ $A$  is necessary”.

A modal formula is a PA-*tautology* if  $\text{PA} \vdash e(A)$  for each evaluation  $e$  based on any standard proof predicate  $\text{Prf}(x, y)$ .

Note that an evaluation is fully determined by its values on propositional atoms and that a value  $e(A)$  of a formula  $A$  depends on only those atoms that occur in  $A$ . So if  $A$  is a formula like  $\neg \Box \perp$ , containing no atoms,  $e(A)$  is the same sentence for each  $e$ . Also note that if  $A$  contains no occurrences of  $\preceq, \prec$  then  $e(A)$  is equivalent to a sentence constructed from values of atoms using only logical connectives and the formula  $\text{Pr}_{\text{PA}}(x)$ ; in this sense  $e(A)$  is independent of the choice of the proof predicate  $\text{Prf}(x, y)$ . This follows from the fact that if  $\Box B$  is any subformula of  $A$  then  $e(\Box B)$  is **PA-equivalent** to  $\text{Pr}_{\text{PA}}(\overline{e(B)})$ .

The formula  $\Box(\perp \rightarrow p)$  is an example of a PA-tautology because PA knows (can prove) that *ex falso quodlibet*. The formula  $\neg \Box \perp$  is *not* a PA-tautology because its value  $e(A)$  is (for each  $e$ ) the sentence  $\text{Con}(\text{PA})$  which, by Gödel Second Incompleteness Theorem, is not provable in PA. Now recall that  $e(\perp)$  is  $0 = S(0)$  for each  $e$ , and that  $e(\top)$  is provable and  $e(\perp)$  is refutable in PA. So if  $e$  is based on a standard proof predicate  $\text{Prf}(x, y)$  then  $\mathbf{N} \models \exists y \text{Prf}(\overline{e(\top)}, y)$  and  $\mathbf{N} \not\models \exists y \text{Prf}(\overline{e(\perp)}, y)$ . Hence  $\mathbf{N} \models \text{Prf}(\overline{e(\top)}, \overline{m})$  for some  $\overline{m}$  and  $\mathbf{N} \models \neg \text{Prf}(\overline{e(\perp)}, \overline{k})$  for each  $k$ . Since  $\Sigma$ -completeness is applicable to both formulas  $\text{Prf}$  and  $\neg \text{Prf}$ , we have  $\text{PA} \vdash \text{Prf}(\overline{e(\top)}, \overline{m})$  and  $\text{PA} \vdash \forall v \leq \overline{m} \neg \text{Prf}(\overline{e(\perp)}, v)$ . Thus

$$\text{PA} \vdash \exists y (\text{Prf}(\overline{e(\top)}, y) \ \& \ \forall v \leq y \neg \text{Prf}(\overline{e(\perp)}, v)).$$

This argument shows that  $\Box \top \prec \Box \perp$  is another example of a PA-tautology; PA knows that the trivially true sentence  $\neg(0 = S(0))$  has a proof less than any proof of contradiction without discussing the *existence* of a proof of contradiction. Now consider, as our final example, the two trivially provable sentences  $\neg(0 = S(0))$  and  $\text{Pr}_{\text{PA}}(\overline{\neg(0 = S(0))})$  which are the values  $e(\top)$  and  $e(\Box \top)$  of the modal formulas  $\top$  and  $\Box \top$  under any evaluation  $e$ . It is possible to choose a standard proof predicate  $\text{Prf}(x, y)$  and an evaluation  $e$  based on it such that, in the sense of  $\text{Prf}(x, y)$ , the first proof of  $e(\Box \top)$  is less than the first proof of  $e(\top)$ . Then  $\text{PA} \not\vdash e(\Box \top \preceq \Box \Box \top)$  and so  $\Box \top \preceq \Box \Box \top$  is not a

PA-tautology. The same argument together with a choice of another proof predicate shows that neither  $\Box\Box\top \prec \Box\top$  is a PA-tautology. This example demonstrates why the notion of standard proof predicate is useful: letting the proof predicate vary, i.e. opening the possibility of artificial reordering proofs, is a means how to get rid of irrelevant facts (possibly depending on details of coding of syntax) about ordering of proofs.

The *modal logic R* of [Guaspari and Solovay \(1979\)](#), the classical theory of witness comparison, has the following axiom schemas A1–A4, B1–B4, and P and deduction rules MP, Nec, and Un:

- A1: all propositional tautologies,  
A2:  $\Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$ ,  
A3:  $\Box A \rightarrow \Box\Box A$ ,  
A4:  $\Box(\Box A \rightarrow A) \rightarrow \Box A$ ,  
MP:  $A, A \rightarrow B / B$ ,                      Nec:  $A / \Box A$ ,  
Un:  $\Box A / A$ ,  
B1:  $\Box A \preceq \Box B \rightarrow \Box A$ ,  
B2:  $\Box A \preceq \Box B \ \& \ \Box B \preceq \Box C \rightarrow \Box A \preceq \Box C$ ,  
B3:  $\Box A \vee \Box B \rightarrow \Box A \preceq \Box B \vee \Box B \prec \Box A$ ,  
B4:  $\Box A \prec \Box B \equiv \Box A \preceq \Box B \ \& \ \neg(\Box B \preceq \Box A)$ ,  
P:  $\Box A \preceq \Box B \rightarrow \Box(\Box A \preceq \Box B)$ ,               $\Box A \prec \Box B \rightarrow \Box(\Box A \prec \Box B)$ .

The restriction of the logic R to the usual modal language, with  $\Box$  as the only modality, axioms A1–A4, and rules modus ponens MP and necessitation Nec, is *provability logic* GL. Nowadays, there is a rich literature about provability logic. The logic R as well as the other systems mentioned below are conservative extensions of provability logic. There is no need to add the unnecessitation rule Un to provability logic because Un is a derived rule in that logic. The schema P could be called a *persistence schema*. We keep it a little bit apart from the *basic axioms B1–B4* about witness comparison because we will discuss some alternatives to it.

The logic R is complete with respect to the arithmetical semantics. It also has a satisfactory Kripke semantics and is decidable; which is quite surprising in the light of the highly inefficient nature of the arithmetical semantics. Out of these facts, it is rather straightforward to verify soundness of all axioms and rules, perhaps with an exception of A4, the Löb axiom schema. Note that soundness of Nec, A2, and A3 follows immediately from the *derivability conditions*, while soundness of the schema P follows from the formalized  $\Sigma$ -*completeness* theorem.

Consider the modal formula

$$\Box(p \equiv \Box\neg p \prec \Box p) \rightarrow (\neg\Box\perp \rightarrow \neg\Box p \ \& \ \neg\Box\neg p). \quad (3)$$

It says that “if  $p$  is provably equivalent to the statement saying that there is a proof of  $\neg p$  which is less than any proof of  $p$ , then, under the assumption of consistency, neither  $p$  nor  $\neg p$  is provable”. Formula (3) is an example of a formula provable in R; this fact is left as an exercise because we give a similar proof in our logic SR below. Then, by arithmetical soundness of the logic R, the arithmetical value of formula (3) is provable in PA under any evaluation. Let  $e$  be an evaluation sending  $p$  to the Rosser sentence  $\rho$ . For this evaluation we also have  $\text{PA} \vdash e(\Box(p \equiv \Box\neg p \prec \Box p))$  and so  $\text{PA} \vdash e(\neg\Box\perp \rightarrow \neg\Box p \ \& \ \neg\Box\neg p)$ . Since  $e(\neg\Box\perp \rightarrow \neg\Box p \ \& \ \neg\Box\neg p)$  is the sentence  $\text{Con}(\text{PA}) \rightarrow \neg\text{Pr}_{\text{PA}}(\bar{\rho}) \ \& \ \neg\text{Pr}_{\text{PA}}(\overline{\neg\rho})$ , we have proved the fact mentioned above when speaking about Gödel and Rosser sentences: independence of the Rosser sentence  $\rho$  is formalizable in PA.

The reasoning in the previous paragraph is an example of a “direct” application of modal logic: properties of a particular self-referential sentence can be obtained by proving appropriate modal formulas in an appropriate modal system. Here is another example of a direct application of modal logic in meta-mathematics: from provability of the modal formula  $\Box(p \equiv \neg\Box p) \rightarrow (\neg\Box\perp \rightarrow \neg\Box p)$  in R (and in GL) we get that PA proves the sentence  $\text{Con}(\text{PA}) \rightarrow \neg\text{Pr}_{\text{PA}}(\bar{v})$  expressing formalized unprovability of the Gödel sentence  $v$ . An example of a very nice and less direct application of modal logic is this: a sentence like  $\rho$ , whose full independence (unprovability and unrefutability) is formalizable in PA, cannot be obtained by *Gödelian self-reference*, i.e. by writing down an equation  $\text{PA} \vdash \varphi \equiv \psi(\bar{\varphi})$  with  $\psi$  speaking about provability only, not employing things like the Rosser trick. Out of our examples, Gödel and Henkin sentences are obtained by Gödelian self-reference, Rosser and Hájek sentences are not.

## 5 Alternative systems with Rosser modalities

The Hájek sentence (2) says that a number  $y$  such that  $\neg\text{Con}(\text{PA} \upharpoonright y + \bar{\eta})$  appears before any possible  $v$  such that  $\neg\text{Con}(\text{PA} \upharpoonright v + \overline{\neg\eta})$ ; so it in fact also uses the Rosser trick. Moreover, the formula  $\exists y \neg\text{Con}(\text{PA} \upharpoonright y + \neg x)$  is equivalent to  $\text{Pr}_{\text{PA}}(x)$ ; so it can be viewed as a *generalized proof predicate*, and a  $y$  such that  $\neg\text{Con}(\text{PA} \upharpoonright y + \neg x)$  can be viewed as a proof of  $x$  in a generalized sense. Put together, the Hájek sentence is a sort of a Rosser sentence. However, the modal system R does not apply to it because the sentence is not a  $\Sigma$ -sentence. Speaking otherwise, if evaluations  $e$  based on the generalized proof predicate are allowed, then the formalized  $\Sigma$ -completeness theorem is not applicable to sentences of the form  $e(\Box A \preceq \Box B)$  and  $e(\Box A \prec \Box B)$ , and thus the schema P is not arithmetically sound.

A modal system applicable to the Hájek sentence was proposed in Švejdar (1983). The system is denoted Z there; it results from replacing the schema P



by another schema Z:

$$\text{Z: } \Box A \rightarrow \Box(\neg B \rightarrow \Box A \prec \Box B).$$

This schema says that “if  $A$  is provable then it can be proved that it has a proof smaller than any proof of any false statement  $B$ ”, or more briefly, “false statements do not have small proofs”. It is not difficult to verify that [formula \(3\)](#), the modalized Rosser theorem, is provable also in the logic Z and that the schema Z is provable in the logic R. Soundness of the schema Z w.r.t semantics where evaluations can be based also on the generalized proof predicate  $\neg\text{Con}(\text{PA} \upharpoonright y + \neg x)$  follows from [Theorem 1](#) below. So (formalized or not) independence of the Hájek sentence is in fact a consequence of modal considerations. It is also known ([Švejdar, 1983](#)) that the logic Z has a Kripke completeness theorem and is decidable.

Independence is not the only important property of the Hájek sentence by far. But since the system Z cannot prove any formula which is not provable in R, it offers no modal explanation why the Hájek sentence [\(2\)](#) should be important. The goal of this paper is to suggest a system stronger than Z that could express what Hájek sentence and similar constructions have in comparison to the usual Rosser sentence.

Let a (*modal*)  $\Sigma$ -*formula* be any formula obtained from formulas starting with  $\Box$  and from  $\perp$  and  $\top$  using only conjunctions and disjunctions. Let S be the schema

$$\text{S: } \Box(E \rightarrow A) \rightarrow \Box(E \& \neg B \rightarrow \Box A \prec \Box B),$$

where  $E$  is a  $\Sigma$ -formula and  $A$  and  $B$  are arbitrary modal formulas. Let SR be the system like Z, but with the schema Z replaced by the schema S. The notation is somewhat tentative: The letter “S” refers to Greek  $\Sigma$  because the logic SR is a Rosserian system with a special attention to arithmetical  $\Sigma$ -sentences, or “SR” can mean “semi-reflexive”.

**Theorem 1** *The logic SR is sound w.r.t. the arithmetical semantics with evaluations based on the generalized proof predicate  $\neg\text{Con}(\text{PA} \upharpoonright y + \neg x)$ .*

**Proof** We show arithmetical soundness of the schema S. A key step is to use the *essential reflexivity* of PA: the implication  $\chi \rightarrow \text{Con}(\text{PA} \upharpoonright \bar{k} + \bar{\chi})$  is provable in PA for each  $m$  and each sentence  $\chi$ . So let modal formulas  $A$  and  $B$ , modal  $\Sigma$ -formula  $E$  and an evaluation  $e$  based on the generalized proof predicate be given. Let  $\varphi$ ,  $\psi$ , and  $\sigma$  be the arithmetical sentences that are the values of formulas  $A$ ,  $B$ , and  $E$  (respectively) under  $v$ . We have to prove in PA that if  $\sigma \rightarrow \varphi$  is provable then one can prove, using the assumptions  $\sigma$  and  $\neg\psi$ , that  $\bar{\varphi}$  has a generalized proof smaller than any generalized proof of  $\bar{\psi}$ . We shift the speech one level up, giving a meta-mathematical argument. The

real proof of the theorem is then obtained by formalizing the considerations in PA, i.e. essentially by rewriting the argument using the **sans-serif font**. So let  $\text{PA} \vdash \sigma \rightarrow \varphi$ . We have to show

$$\text{PA} \vdash \sigma \ \& \ \neg\psi \rightarrow \exists y(\neg\text{Con}(\text{PA} \upharpoonright y + \overline{\neg\varphi}) \ \& \ \forall v \leq y \text{Con}(\text{PA} \upharpoonright v + \overline{\neg\psi})). \quad (4)$$

Note that  $\forall v \leq y \text{Con}(\text{PA} \upharpoonright v + \overline{\neg\psi})$  is the same as  $\text{Con}(\text{PA} \upharpoonright y + \overline{\neg\psi})$ . A fixed proof of  $\sigma \rightarrow \varphi$  uses only a finite number of axioms of PA. So if  $m_1$  is sufficiently big then the axioms of PA less than  $m_1$ , together with the sentence  $\neg(\sigma \rightarrow \varphi)$ , constitute a contradictory theory. Thus, in fact by  $\Sigma$ -completeness, we have  $\text{PA} \vdash \neg\text{Con}(\text{PA} \upharpoonright \overline{m_1} + \overline{\neg(\sigma \rightarrow \varphi)})$ . An inspection of the proof of formalized  $\Sigma$ -completeness yields an  $m_0$  such that  $\text{PA} \vdash \sigma \rightarrow \neg\text{Con}(\text{PA} \upharpoonright \overline{m_0} + \overline{\neg\sigma})$  for (each and thus) our sentence  $\sigma$ . So if  $m = \max\{m_0, m_1\}$  we have

$$\text{PA} \vdash \sigma \rightarrow \neg\text{Con}(\text{PA} \upharpoonright \overline{m} + \overline{\neg\varphi}), \quad (5)$$

$m$  is a generalized proof of  $\varphi$ . Essential reflexivity applied to the sentence  $\psi$  yields  $\text{PA} \vdash \neg\psi \rightarrow \text{Con}(\text{PA} \upharpoonright \overline{m} + \overline{\neg\psi})$ . This and (5) implies (4), q.e.d. ■

Note the difference between the proof predicate  $\neg\text{Con}(\text{PA} \upharpoonright y + \neg x)$  and the standard proof predicate  $\text{Proof}_{\text{PA}}(x, y)$ : there is no common  $m_0$  such that  $\sigma \rightarrow \exists v \leq \overline{m_0} \text{Proof}_{\text{PA}}(\overline{\sigma}, v)$  is provable for each  $\Sigma$ -sentence  $\sigma$ .

A substitution of  $\Box\perp$ ,  $\perp$ , and  $\Box\perp$  for  $A$ ,  $B$ , and  $E$  in the schema S yields  $\Box(\Box\perp \rightarrow \Box\Box\perp \prec \Box\perp)$ . Then, using the rule Un, we get  $\Box\perp \rightarrow \Box\Box\perp \prec \Box\perp$ . This formula, as well as the formula  $\Box\Box\perp \rightarrow \Box\Box\perp \prec \Box\perp$ , are examples of formulas provable in the logic SR but unprovable in R. As an exercise, the reader may try to derive the latter from the former by distinguishing cases  $\Box\perp$  and  $\neg\Box\perp$ . A further example of a formula provable in SR is

$$\Box(p \equiv \Box\neg p \prec \Box p) \rightarrow (\Box(E \rightarrow p) \vee \Box(E \rightarrow \neg p) \rightarrow \Box\neg E), \quad (6)$$

for any  $\Sigma$ -formula  $E$ . For, let  $D$  be  $\Box(p \equiv \Box\neg p \prec \Box p)$ . Then the proof of (6) in SR may proceed like this:

$$\begin{aligned} & \Box(E \rightarrow p) \rightarrow \Box(E \ \& \ p \rightarrow \Box p \prec \Box\neg p) && ; \text{S} \\ & \Box(E \rightarrow p) \rightarrow \Box(E \rightarrow \Box p \prec \Box\neg p) \\ & \Box(E \rightarrow p) \rightarrow \Box(E \rightarrow \Box p \preceq \Box\neg p) && ; \text{B4} \\ & \Box(E \rightarrow p) \rightarrow \Box(E \rightarrow \neg(\Box\neg p \prec \Box p)) && ; \text{B4} \\ & D \rightarrow \Box(p \rightarrow \Box\neg p \prec \Box p) \\ & D \ \& \ \Box(E \rightarrow p) \rightarrow \Box(E \rightarrow \Box\neg p \prec \Box p) \\ & D \ \& \ \Box(E \rightarrow p) \rightarrow \Box\neg E, \end{aligned}$$

where the seventh line follows from the sixth and fourth. The proof of the second half,  $D \ \& \ \Box(E \rightarrow \neg p) \rightarrow \Box\neg E$ , is similar.

A substitution of  $\top$  for  $E$  in (6) yields formula (3); so (6) is a stronger version of modalized Rosser theorem. Formula (6), as well as formula (3), speaks about properties of the Rosser sentence, i.e. sentence  $p$  provably equivalent to the statement there is a proof of  $\neg p$  which is less than any proof of  $p$ . While formula (3) says that the Rosser sentence is independent provided the base theory is consistent, formula (6) says that neither  $p$  nor  $\neg p$  can be proved using any assumption expressed by a consistent  $\Sigma$ -sentence. It is known that, in PA, the unprovability of a sentence  $\varphi$  from any consistent  $\Sigma$ -sentence is the same as interpretability of  $\text{PA} + \neg\varphi$  in PA. This fact is due to Petr Hájek. So provability of formula (6) in SR together with Theorem 1 yields the symmetric interpretability of the Hájek sentence.

Speaking precisely, there is a little gap in the conclusion that the symmetric interpretability of Hájek sentence can be obtained modally. Our considerations only show that neither Hájek sentence  $\eta$  nor its negation  $\neg\eta$  can be proved from a consistent assumption *which is a value of a modal  $\Sigma$ -formula*. Note that values of modal  $\Sigma$ -formulas constitute a proper subset of all  $\Sigma$ -sentences. To get the full result, one would modify the logic SR as follows: (i) consider two sorts of propositional atoms, (normal) atoms  $p, q, \dots$  and  $\Sigma$ -atoms  $\sigma, \tau, \dots$ , and (ii) prove a formula like (6), but with  $E$  replaced by a  $\Sigma$ -atom  $\sigma$ . We did not do this modification to keep things simpler.

To finish, it must be admitted that we are still rather far from a well-rounded modal system for Rosser modalities. Nevertheless, I believe that the schemas Z and S do formalize arguments that are quite frequent when dealing with self-referential sentences that use some version of the Rosser trick.

Vítězslav Švejdar

Department of Logic, Charles University

Palachovo nám. 2, 116 38 Praha 1, Czech Republic

vitezslavdotsvejdaratcunidotcz, <http://www1.cuni.cz/~svejdar/>

## References

- Gödel, K. (1931). Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik*, 37, 349–360.
- Guaspari, D., & Solovay, R. M. (1979). Rosser sentences. *Annals of Math. Logic*, 16, 81–99.
- Hájek, P., & Hájková, M. (1972). On interpretability in theories containing arithmetic. *Fundamenta Mathematicae*, 76, 131–137.
- Löb, M. H. (1955). Solution of a problem of Leon Henkin. *J. Symbolic Logic*, 20, 115–118.
- Smoryński, C. (1985). *Self-reference and modal logic*. New York: Springer.

- Solovay, R. M. (1976). Provability interpretations of modal logic. *Israel J. Math.*, 25, 287–304.
- Švejdar, V. (1983). Modal analysis of generalized Rosser sentences. *J. Symbolic Logic*, 48(4), 986–999.
- Tarski, A., Mostowski, A., & Robinson, R. M. (1953). *Undecidable theories*. Amsterdam: North-Holland.