

---

# Gödelian sentences, Rosserian sentences and truth

ZIBA ASSADI, *Independent Scholar, Tabriz, Iran.*

SAEED SALEHI, *Department of Mathematical Sciences, University of Tabriz, P.O. Box 51666-16471, Tabriz, Iran.*

## Abstract

There is a long-standing debate in the logico-philosophical community as to if/why the Gödelian sentences of a consistent and sufficiently strong theory are true ( $\gamma$  is a Gödelian sentence of  $T$  when  $\gamma$  is equivalent to the  $T$ -unprovability of  $\gamma$  inside  $T$ ). The prevalent argument seems to be something like the following: *since every one of the Gödelian sentences of such a theory is equivalent to the consistency statement of the theory, even provably so inside the theory, the truth of those sentences follows from the consistency of the theory in question; so, Gödelian sentences of consistent and sufficiently strong theories are true.* In this paper, we critically examine this argument and present necessary and sufficient conditions for the truth of Gödelian sentences (and Rosserian sentences) of consistent and sufficiently strong arithmetical theories.

*Keywords:* Incompleteness theorem, Gödelian sentences, Rosser's trick, Rosserian sentences, soundness, consistency,  $\Sigma_n$ -soundness

## 1 Introduction

By the first incompleteness theorem of Gödel [9], for every consistent and sufficiently strong arithmetical theory there are sentences that are undecidable in the theory. Examples of such undecidable sentences are actually constructed in Gödel's original proof in a way that each of those sentences is equivalent to its own unprovability in the theory; see Definition 3.1 below. A natural question here is that while the theory in question cannot decide the truth of its Gödelian sentences, what about *us* (human beings)? Can we 'see' (or demonstrate) their truth? This question has attracted the attention of many philosophers, physicists, computer scientists, as well as mathematical logicians. As there are numerous papers and books on this subject, it is not possible to cite them all here; see Conclusions for a few.

If our theory (which is an RE set of sentences) is sufficiently strong and *sound*, then it proves the equivalence of each of its Gödelian sentences with the consistency statement of the theory (see Remark 2.3.III below). Since the consistency statement is true, and the theory is sound, then it follows that all the Gödelian sentences of the theory are true. Now, let us see an example of a false Gödelian sentence (of a consistent and sufficiently strong theory). Put  $T$  be a consistent and sufficiently strong RE theory; by Gödel's second incompleteness theorem,  $T$  cannot prove its consistency statement. So, the theory  $T$  plus its *inconsistency* statement is consistent; call it  $U$ . Now,  $U$  proves the inconsistency of  $T$  and so it proves the inconsistency of  $U$ , itself, as well. Therefore,  $U$  proves that the contradiction,  $\perp$ , is  $U$ -provable; so,  $\perp$  is  $U$ -provably equivalent to the  $U$ -unprovability of  $\perp$ . Thus,  $\perp$  is a Gödelian sentence of  $U$  (note that  $U$  is an *unsound* theory; see also [17, 18]).

So, something must be wrong with the argument mentioned in the Abstract (that Gödelian sentences of consistent and sufficiently strong theories are true).

The paper is organized as follows: in Section 2 we present some preliminaries necessary for following the upcoming arguments. In Section 3 we show that Gödelian (I-am-unprovable) sentences constitute all the unprovable sentences in a sense and provide necessary and sufficient conditions for their truth. In Section 4 we study Rosserian sentences; the sentences that express ‘for every proof of me there is a smaller proof of my negation’. We also provide necessary and sufficient conditions for their truth. In Section 5 we conclude the paper with a diagram on the truth (and falsity) of Gödelian and Rosserian sentences by presenting some nice equivalent conditions on the underlying theory.

## 2 Preliminaries

We assume familiarity with the notions of  $\Pi_n$  and  $\Sigma_n$  formulas, Peano’s arithmetic  $\mathbf{PA}$  and its fragments, Robinson’s Arithmetic  $\mathbf{Q}$  and the fact that  $\mathbf{Q}$  is a sound and  $\Sigma_1$ -complete theory (i.e. every  $\mathbf{Q}$ -provable sentence is true and every true  $\Sigma_1$ -sentence is  $\mathbf{Q}$ -provable). By the diagonal lemma of Gödel and Carnap, for every formula  $\Psi(x)$  with the only free variable  $x$ , there exists a sentence  $\theta$  such that  $\theta \leftrightarrow \Psi(\#\theta)$  is true (in the standard model of natural numbers  $\mathbb{N}$ ) and also provable in  $\mathbf{Q}$ ; here  $\#A$  denotes the numeral of the Gödel code of  $A$ , relative to a fixed Gödel numbering (arithmetization) of the syntax. Moreover, if  $\Psi(x)$  is a  $\Pi_n$ -formula, for some  $n \geq 1$ , then  $\theta$  can be taken to be a  $\Pi_n$ -sentence; and if  $\Psi(x)$  is  $\Sigma_n$ , then  $\theta$  can be taken to be  $\Sigma_n$  too. We provide more details in the following:

LEMMA 2.1 (The diagonal lemma).

Let  $n \geq 1$ . For every  $\Pi_n$ -formula  $\Psi(x)$  there exists a  $\Pi_n$ -sentence  $\theta$  such that  $\mathbf{Q} \vdash \theta \leftrightarrow \Psi(\#\theta)$ . And for every  $\Sigma_n$ -formula  $\Psi(x)$  there exists a  $\Sigma_n$ -sentence  $\theta$  with the same property.

PROOF. There is a primitive recursive partial function  $d$  that assigns to a given  $m$ , when  $m$  codes a formula with the only free variable  $x$ , the Gödel code of the sentence that results from substituting  $\bar{m}$  for  $x$  (where  $\bar{m}$  is the numeral of  $m$ , a term in the language of arithmetic representing  $m$ ). There is a  $\Sigma_1$ -formula  $\delta(x, y)$ , in the language of arithmetic, that strongly represents  $d$  in  $\mathbf{Q}$ . This means that  $\mathbf{Q} \vdash \forall y[\delta(\bar{m}, y) \leftrightarrow y = \bar{d}(m)]$  for every  $m \in \mathbb{N}$ .

If  $\Psi(x)$  is a  $\Pi_n$ -formula, then put  $\alpha(x) = \forall y[\delta(x, y) \rightarrow \Psi(y)]$  and let  $\mathbf{a}$  be its Gödel code. Now, let  $\theta = \alpha(\bar{\mathbf{a}})$ ; then  $\theta$  is a  $\Pi_n$ -sentence and we have provably in  $\mathbf{Q}$  that

$$\begin{aligned} \theta &\leftrightarrow \forall y[\delta(\bar{\mathbf{a}}, y) \rightarrow \Psi(y)] \\ &\leftrightarrow \forall y[y = \bar{\mathbf{d}}(\mathbf{a}) \rightarrow \Psi(y)] \\ &\leftrightarrow \forall y[y = \#\theta \rightarrow \Psi(y)] \\ &\leftrightarrow \Psi(\#\theta). \end{aligned}$$

If  $\Psi(x)$  is a  $\Sigma_n$ -formula, then put  $\eta(x) = \exists y[\delta(x, y) \wedge \Psi(y)]$  and let  $\mathbf{e}$  be its Gödel code. Now, let  $\theta = \eta(\bar{\mathbf{e}})$ ; then  $\theta$  is a  $\Sigma_n$ -sentence and we have provably in  $\mathbf{Q}$  that

$$\begin{aligned} \theta &\leftrightarrow \exists y[\delta(\bar{\mathbf{e}}, y) \wedge \Psi(y)] \\ &\leftrightarrow \exists y[y = \bar{\mathbf{d}}(\mathbf{e}) \wedge \Psi(y)] \\ &\leftrightarrow \exists y[y = \#\theta \wedge \Psi(y)] \\ &\leftrightarrow \Psi(\#\theta). \end{aligned} \quad \square$$

The incompleteness theorem is usually stated for recursively enumerable (RE) theories that extend  $\mathbf{Q}$ ; though it also holds for more general theories, see, e.g. [32] or [15]. For us a theory is a set of sentences. If  $T$  is an RE theory, then by [14, Corollary 3.4] the theory  $T$  is  $\Sigma_1$ -definable, in the sense that we have  $T = \{\theta \mid \mathbb{N} \models \sigma(\#\theta)\}$  for a  $\Sigma_1$ -formula  $\sigma(x)$ , where  $\theta$  ranges over the sentences.

By Craig's trick [4], every such RE theory can be axiomatized by a  $\Delta_0$ -definable set of sentences (see [32, Lemma 2.4]): write  $\sigma(x) \equiv \exists y \chi(y, x)$  where  $\chi$  is a  $\Delta_0$ -formula and consider the theory  $T^* = \{\theta \wedge (\bar{n} = \bar{n}) \mid \mathbb{N} \models \chi(\bar{n}, \#\theta)\}$ ; now  $T^*$  is equivalent to  $T$ , and is definable by the  $\Delta_0$ -formula  $\tau(x) = \exists y, z \leq x [x = [y \& (\bar{z} = \bar{z})] \wedge \chi(z, y)]$ .

Let  $\tau(x)$  be an arbitrary  $\Delta_0$ -formula. Put  $\text{Th}_\tau = \{\theta \mid \mathbb{N} \models \tau(\#\theta)\}$  to be the theory defined by  $\tau$ , where  $\theta$  ranges over the sentences. By the arithmetization of syntax with respect to a fixed Gödel coding, we can write a  $\Sigma_1$ -formula  $\text{prf}_\tau(y, x)$  stating that 'y is (the code of) a proof of the sentence (with code) x in the theory  $\text{Th}_\tau$ '. Let us note that it suffices for y to be the code of a sequence of formulas, each of which is either a logical axiom or satisfies  $\tau$  (is an axiom of  $\text{Th}_\tau$ ) or is derived from one or two earlier formulas by a logical rule (which could be Modus Ponens or Generalization). Since  $\text{prf}_\tau(y, x)$  is a decidable relation when  $\tau$  is  $\Delta_0$  (which implies that  $\text{Th}_\tau$  is a decidable set of sentences), by [14, Corollary 3.5] it is equivalent to a  $\Pi_1$ -formula  $\pi(y, x)$ , and this equivalence is provable in the theory  $\text{IS}_1$ , a fragment of  $\text{PA}$  in which the axiom scheme of induction is restricted to  $\Sigma_1$ -formulas (with parameters), by [11, Definition I.4.3], this is to say that we have  $\text{IS}_1 \vdash \forall x, y [\text{prf}_\tau(y, x) \leftrightarrow \pi(y, x)]$ . Let us note that the theory  $\text{IS}_1$  is finitely axiomatizable, and 'arithmetization of metamathematics' can be developed in it; see [11].

*Throughout the paper, we consider  $\Delta_0$ -definable theories that extend  $\text{IS}_1$ .*<sup>1</sup>

Thus, for a  $\Delta_0$ -formula  $\tau$  the *proof predicate*  $\text{prf}_\tau$  of  $\tau$  is a  $\text{IS}_1$ -provably  $\Delta_1$ -formula; though it could be a  $\Delta_0$ -formula by the techniques of [11, §V.3]. Let  $\text{Pr}_\tau(x) = \exists y \text{prf}_\tau(y, x)$  be the *provability predicate* of  $\tau$  and  $\text{Con}_\tau = \neg \text{Pr}_\tau(\#[0 \neq 0])$  be its *consistency statement*. Note that  $\text{Pr}_\tau(x)$  is a  $\Sigma_1$ -formula, and  $\text{Con}_\tau$  is a  $\Pi_1$ -sentence. For our coding and arithmetization we expect the following to hold.

CONVENTION 2.2 For every  $\Delta_0$ -formula  $\tau(x)$  and every sentence  $\varphi$  the following hold:

- (C1)  $\text{Th}_\tau \vdash \varphi \iff \mathbf{Q} \vdash \text{prf}_\tau(\bar{m}, \#\varphi)$  for some  $m \in \mathbb{N}$ .
- (C2)  $\text{Th}_\tau \not\vdash \varphi \implies \text{IS}_1 \vdash \neg \text{prf}_\tau(\bar{n}, \#\varphi)$  for every  $n \in \mathbb{N}$ .

Also, the following *derivability conditions* hold for  $\text{Pr}_\tau(x)$ , when  $\tau$  is a  $\Delta_0$ -formula such that  $\text{Th}_\tau \supseteq \text{IS}_1$  and  $\varphi, \psi$  are sentences:

- (D1)  $\text{Th}_\tau \vdash \varphi \iff \mathbf{Q} \vdash \text{Pr}_\tau(\#\varphi)$ .
- (D2)  $\text{IS}_1 \vdash \text{Pr}_\tau(\#[\varphi \rightarrow \psi]) \rightarrow [\text{Pr}_\tau(\#\varphi) \rightarrow \text{Pr}_\tau(\#\psi)]$ .
- (D3)  $\text{IS}_1 \vdash \text{Pr}_\tau(\#\varphi) \rightarrow \text{Pr}_\tau(\#[\text{Pr}_\tau(\#\varphi)])$ . ◇

REMARK 2.3

We will need the following consequences of the derivability conditions in Convention 2.2, where  $\tau$  is a  $\Delta_0$ -formula such that  $\text{Th}_\tau \supseteq \text{IS}_1$  and  $\varphi$  is an arbitrary sentence:

- (I)  $\text{Th}_\tau \vdash \neg \text{Con}_\tau \rightarrow \text{Pr}_\tau(\#\varphi)$ .
- (II) If  $\text{Th}_\tau \vdash \text{Pr}_\tau(\#\varphi) \rightarrow \varphi$ , then  $\text{Th}_\tau \vdash \varphi$ .
- (III) If  $\text{Th}_\tau \vdash \varphi \leftrightarrow \neg \text{Pr}_\tau(\#\varphi)$ , then  $\text{Th}_\tau \vdash \varphi \leftrightarrow \text{Con}_\tau$ . ◇

One can find proofs for Remark 2.3 in [38] and [2]; let us note that Remark 2.3.II is the so-called Löb's Rule.

---

<sup>1</sup>By a result of [40], instead of  $\text{IS}_1$  one can take the slightly weaker theory  $\text{EA} + \text{BS}_1$ , which is equivalent to  $\text{ID}_1$  by a result of [35].

### 3 Gödelian sentences and their truth

Gödel's proof of his incompleteness theorem uses the diagonal lemma (2.1) for the negation of the provability predicate of the  $\Delta_0$ -formula  $\tau$ .

DEFINITION 3.1 (Gödelian Sentences).

A sentence  $\gamma$  is called a Gödelian sentence of the  $\Delta_0$ -formula  $\tau(x)$  when  $\gamma$  is equivalent to its unprovability in the theory defined by  $\tau$ , i.e., we have  $\text{Th}_\tau \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\#\gamma)$ .  $\diamond$

By Remark 2.3.III any two Gödelian sentences of  $\tau$ , when  $\text{Th}_\tau \supseteq \text{IS}_1$ , are  $\text{Th}_\tau$ -provably equivalent; so, many authors talk of *the* Gödel sentence of  $\text{Th}_\tau$ . In some sources, the assumptions and definitions underlying (variants of) the following argument are not always made fully explicit. However, Gödel [9] and several authors after him argue that the Gödelian sentences of a consistent theory are true, since

- (1) they are provably equivalent to their unprovability in the theory, and
- (2) they are indeed unprovable in the theory; and so
- (3) they must be true.

It is argued in [17] that this line of reasoning does not demonstrate the truth of Gödelian sentences, and indeed some  $\Sigma_1$ -unsound theories may have false Gödelian sentences. In fact, step (2) in the above argument is redundant:

LEMMA 3.2

Suppose that for the  $\Delta_0$ -formula  $\tau(x)$ , the theory  $\text{Th}_\tau$  is consistent and contains  $\text{IS}_1$ . For every sentence  $\varphi$ , if  $\text{Th}_\tau \vdash \varphi \rightarrow \neg \text{Pr}_\tau(\#\varphi)$ , then  $\text{Th}_\tau \not\vdash \varphi$ .

PROOF. Because  $\text{Th}_\tau \vdash \varphi$  would imply on the one hand  $\text{Th}_\tau \vdash \neg \text{Pr}_\tau(\#\varphi)$  by the assumption, and on the other hand  $\text{Th}_\tau \vdash \text{Pr}_\tau(\#\varphi)$  by Convention 2.2.D1.  $\square$

So, the question of the validity of the above reasoning for the truth of Gödelian sentences boils down to the following question:

*Does  $\text{Th}_\tau \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\#\gamma)$  imply  $\mathbb{N} \models \gamma$ ,  
for a  $\Delta_0$ -formula  $\tau$  with consistent  $\text{Th}_\tau \supseteq \text{IS}_1$ ?*

Put in another way,

*under which conditions are all the Gödelian sentences of  $\tau$  true?*

We answer this question in the present section, and in the next section we answer a similar question for the Rosserian sentences (of arithmetical theories). Let us start with an amusing result (cf. [18, Theorem 1]):

PROPOSITION 3.3 (Characterizing Gödelian sentences of super-theories).

Let  $\tau(x)$  be a  $\Delta_0$ -formula such that  $\text{Th}_\tau \supseteq \text{IS}_1$ . The following are equivalent for a sentence  $\varphi$ :

- (1)  $\varphi$  is unprovable in  $\text{Th}_\tau$ , i.e.  $\text{Th}_\tau \not\vdash \varphi$ ;
- (2)  $\varphi$  is a Gödelian sentence of some consistent extension of  $\text{Th}_\tau$ ;
- (3)  $\text{Th}_\tau + [\varphi \leftrightarrow \neg \text{Pr}_\tau(\#\varphi)]$  is consistent.

PROOF. (1  $\Rightarrow$  2): By Lemma 2.1, for a sentence  $\xi$  we have

$$\mathbf{Q} \vdash \xi \leftrightarrow [\varphi \leftrightarrow \neg \text{Pr}_{\tau'}(\#\varphi)]$$

where  $\tau'(x) = \tau(x) \vee (x = \#\xi)$ . Then  $\text{Th}_{\tau'} \vdash \varphi \leftrightarrow \neg \text{Pr}_{\tau'}(\#\varphi)$  and it remains to show that the theory  $\text{Th}_{\tau'}$  (which is  $\text{Th}_{\tau} + \xi$ ) is consistent. If not, then  $\text{Th}_{\tau} \vdash \neg\xi$ . So, on the one hand we have (i)  $\text{Th}_{\tau} \vdash \neg[\varphi \leftrightarrow \neg \text{Pr}_{\tau'}(\#\varphi)]$ , and on the other hand  $\text{Th}_{\tau'} \vdash \varphi$ , which implies (ii)  $\mathcal{Q} \vdash \text{Pr}_{\tau'}(\#\varphi)$  by Convention 2.2.D1. Now, (i) and (ii) imply that  $\text{Th}_{\tau} \vdash \varphi$ , contradicting the assumption.

(2  $\Rightarrow$  3): Suppose that the theory  $\text{Th}_{\tau} + [\varphi \leftrightarrow \neg \text{Pr}_{\tau}(\#\varphi)]$  is not consistent; then we have  $\text{Th}_{\tau} \vdash \neg[\varphi \leftrightarrow \neg \text{Pr}_{\tau}(\#\varphi)]$ , and so  $\text{Th}_{\tau} \vdash \text{Pr}_{\tau}(\#\varphi) \rightarrow \varphi$ , which implies  $\text{Th}_{\tau} \vdash \varphi$  by Löb's Rule (Remark 2.3.II). So, for every extension  $\text{Th}_{\tau'}$  of  $\text{Th}_{\tau}$  we have  $\text{Th}_{\tau'} \vdash \varphi$ , and so by Convention 2.2.D1, we have  $\mathcal{Q} \vdash \text{Pr}_{\tau'}(\#\varphi)$ . Therefore, for every such  $\text{Th}_{\tau'}$  we have  $\text{Th}_{\tau'} \vdash \neg[\varphi \leftrightarrow \neg \text{Pr}_{\tau'}(\#\varphi)]$ , which contradicts the assumption.

(3  $\Rightarrow$  1): If  $\text{Th}_{\tau} \vdash \varphi$ , then, by Convention 2.2.D1, we have  $\mathcal{Q} \vdash \text{Pr}_{\tau}(\#\varphi)$ , and so we should have also  $\text{Th}_{\tau} \vdash \neg[\varphi \leftrightarrow \neg \text{Pr}_{\tau}(\#\varphi)]$ .  $\square$

We now provide a necessary and sufficient condition for the truth of all the Gödelian  $\Pi_1$ -sentences (cf. [17, Theorem 3.4]):

**THEOREM 3.4** (On the truth and independence of Gödelian  $\Pi_1$ -sentences).

We assume that for the  $\Delta_0$ -formula  $\tau(x)$  we have  $\text{Th}_{\tau} \supseteq \text{IS}_1$ .

If  $\text{Th}_{\tau} \vdash \neg \text{Con}_{\tau}$ , then every false  $\Pi_1$ -sentence is a Gödelian sentence of  $\tau$ , and no Gödelian sentence of  $\tau$  is independent from  $\text{Th}_{\tau}$ .

If  $\text{Th}_{\tau} \not\vdash \neg \text{Con}_{\tau}$ , then all the Gödelian  $\Pi_1$ -sentences of  $\tau$  are true, and all the Gödelian sentences of  $\tau$  are independent from  $\text{Th}_{\tau}$ .

**PROOF.** If  $\text{Th}_{\tau} \vdash \neg \text{Con}_{\tau}$ , then by Remark 2.3.I we have  $\text{Th}_{\tau} \vdash \text{Pr}_{\tau}(\#\varphi)$  for every sentence  $\varphi$ . So, for every Gödelian sentence  $\gamma$  of  $\tau$  we have  $\text{Th}_{\tau} \vdash \neg\gamma$ ; thus no Gödelian sentence of  $\tau$  can be independent from  $\text{Th}_{\tau}$ . Now, let  $\phi$  be an arbitrary false  $\Pi_1$ -sentence; then  $\neg\phi$  is a true  $\Sigma_1$ -sentence, and so provable in  $\mathcal{Q}$ . Thus,  $\text{Th}_{\tau} \vdash \neg\phi$ ; and so from  $\text{Th}_{\tau} \vdash \text{Pr}_{\tau}(\#\phi)$  we have  $\text{Th}_{\tau} \vdash \phi \leftrightarrow \neg \text{Pr}_{\tau}(\#\phi)$ , which means that  $\phi$  is a (false) Gödelian  $\Pi_1$ -sentence of  $\tau$ .

If  $\text{Th}_{\tau} \not\vdash \neg \text{Con}_{\tau}$ , then Remark 2.3.III implies that for every Gödelian sentence  $\gamma$  of  $\tau$  we have  $\text{Th}_{\tau} \not\vdash \neg\gamma$ ; thus,  $\gamma$  is independent from  $\text{Th}_{\tau}$  (noting that  $\text{Th}_{\tau}$  is consistent and so we also have  $\text{Th}_{\tau} \not\vdash \gamma$  by Lemma 3.2). If a Gödelian  $\Pi_1$ -sentence  $\gamma$  of  $\tau$  is not true, then  $\neg\gamma$  is a true  $\Sigma_1$ -sentence, and so should be  $\mathcal{Q}$ -provable; a contradiction with the  $\text{Th}_{\tau}$ -independence of  $\gamma$ , proved above.  $\square$

If the theory  $\text{Th}_{\tau}$  is  $\Sigma_1$ -sound, then we have  $\text{Th}_{\tau} \not\vdash \neg \text{Con}_{\tau}$ . If  $\text{Th}_{\tau}$  is inconsistent or we have  $\tau(x) = \vartheta(x) \vee (x = \#[\neg \text{Con}_{\vartheta}])$  for a  $\Delta_0$ -formula  $\vartheta$  such that  $\text{Th}_{\vartheta}$  is a consistent extension of  $\text{IS}_1$ , then  $\text{Th}_{\tau} \vdash \neg \text{Con}_{\tau}$  (noting that  $\text{Th}_{\tau} = \text{Th}_{\vartheta} + \neg \text{Con}_{\vartheta}$ ); in the latter case  $\text{Th}_{\tau}$  is consistent by Gödel's second incompleteness theorem. Thus, by Theorem 3.4, a necessary and sufficient condition for the truth of all the Gödelian  $\Pi_1$ -sentences of  $\tau$  is the consistency of  $\text{Th}_{\tau}$  with  $\text{Con}_{\tau}$ , a condition obviously implied by  $\omega$ -consistency, although this condition is stronger than the mere consistency of  $\text{Th}_{\tau}$ ; see [13, Theorem 36].

For investigating on the truth of Gödelian  $\Pi_{n+1}$ -sentences (and  $\Sigma_{n+1}$ -sentences) we make a definition and an observation. Before that let us note that no Gödelian  $\Sigma_1$ -sentence of a consistent  $\Delta_0$ -definable extension of  $\text{IS}_1$  can be true:

**PROPOSITION 3.5** (On the truth of Gödelian  $\Sigma_1$ -sentences).

For every  $\Delta_0$ -formula  $\tau(x)$ , no Gödelian  $\Sigma_1$ -sentence of  $\tau$  can be true if  $\text{Th}_{\tau}$  is consistent and contains  $\text{IS}_1$ .

PROOF. If a Gödelian  $\Sigma_1$ -sentence of  $\tau$  were true, then it would have been provable in  $\mathcal{Q}$ , and this would have contradicted Lemma 3.2 for consistent  $\text{Th}_\tau$ .  $\square$

DEFINITION 3.6 ( $\mathcal{Y}$ -Soundness).

Let  $\mathcal{Y}$  be a class of sentences. A theory  $S$  is called  $\mathcal{Y}$ -sound when every  $S$ -provable  $\mathcal{Y}$ -sentence is true.  $\diamond$

The following lemma has been proved for  $\mathcal{Y} = \Sigma_1, \Sigma_2$  in [13, Theorems 25, 27, 30 and 32]:

LEMMA 3.7 (On extensions of  $\mathcal{Y}$ -sound theories).

Let  $\mathcal{Y}$  be a class of sentences that is closed under disjunction. If  $T$  is a  $\mathcal{Y}$ -sound theory, then for every sentence  $\varphi$ , either  $T + \varphi$  or  $T + \neg\varphi$  is  $\mathcal{Y}$ -sound.

PROOF. If neither  $T + \varphi$  nor  $T + \neg\varphi$  is  $\mathcal{Y}$ -sound, then for some false  $\mathcal{Y}$ -sentences  $\xi$  and  $\xi'$  we have  $T + \varphi \vdash \xi$  and  $T + \neg\varphi \vdash \xi'$ . Thus,  $T \vdash \xi \vee \xi'$ , and  $\xi \vee \xi'$  is a false  $\mathcal{Y}$ -sentence; a contradiction with the  $\mathcal{Y}$ -soundness of  $T$ .  $\square$

One of our main results is the following necessary and sufficient condition for the truth of Gödelian ( $\Pi_{n+1}$ - and  $\Sigma_{n+1}$ -) sentences:

THEOREM 3.8 (On the truth of Gödelian  $\Pi_{n+1}$ - and  $\Sigma_{n+1}$ -sentences).

Let  $n \geq 1$ , and let  $\tau$  be a  $\Delta_0$ -formula such that  $\text{Th}_\tau \supseteq \text{IS}_1$ .

All the Gödelian  $\Pi_{n+1}$ -sentences of  $\tau$  are true if and only if  $\text{Th}_\tau$  is  $\Pi_{n+1}$ -sound.

All the Gödelian  $\Sigma_{n+1}$ -sentences of  $\tau$  are true if and only if  $\text{Th}_\tau$  is  $\Sigma_{n+1}$ -sound.

PROOF. Let  $\mathcal{Y}$  be either of the two classes of sentences (either  $\Pi_{n+1}$  or  $\Sigma_{n+1}$ ).

First, suppose that  $\text{Th}_\tau$  is  $\mathcal{Y}$ -sound, and let  $\gamma$  be a Gödelian  $\mathcal{Y}$ -sentence of  $\tau$ . By Lemma 3.2 and Convention 2.2.D1 we have  $\mathbb{N} \models \neg \text{Pr}_\tau(\#\gamma)$ , and so  $\text{Pr}_\tau(\#\gamma)$  is a false  $\Sigma_1$ -sentence. Now,  $\text{Th}_\tau + \neg\gamma \vdash \text{Pr}_\tau(\#\gamma)$ , and so  $\text{Th}_\tau + \neg\gamma$  is not  $\Sigma_1$ -sound; hence, it is not  $\mathcal{Y}$ -sound either. Thus, by Lemma 3.7, the theory  $\text{Th}_\tau + \gamma$  should be  $\mathcal{Y}$ -sound. Therefore,  $\gamma$  must be true.<sup>2</sup>

Now, suppose that all the Gödelian  $\mathcal{Y}$ -sentences of  $\tau$  are true. We show that the theory  $\text{Th}_\tau$  is  $\mathcal{Y}$ -sound. Assume that  $\text{Th}_\tau \vdash \xi$  for a  $\mathcal{Y}$ -sentence  $\xi$ . We prove that  $\xi$  is true. By Lemma 2.1 there exists a  $\mathcal{Y}$ -sentence  $\gamma$  such that  $\mathcal{Q} \vdash \gamma \leftrightarrow [\xi \wedge \neg \text{Pr}_\tau(\#\gamma)]$ . Thus, from  $\text{Th}_\tau \vdash \xi$  we have  $\text{Th}_\tau \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\#\gamma)$ , and so  $\gamma$  is a Gödelian  $\mathcal{Y}$ -sentence of  $\tau$ . Hence,  $\gamma$  is true, and so, by the soundness of  $\mathcal{Q}$ , we have  $\mathbb{N} \models \xi$ .  $\square$

Hence, all the Gödelian sentences of a theory are true if and only if the theory is sound; cf. [36, Theorem 24.7].

REMARK 3.9 (On the hierarchy of  $\Pi_n$ - and  $\Sigma_n$ -soundness).

Let us note that an extension of  $\mathcal{Q}$  is consistent if and only if it is  $\Pi_1$ -sound: indeed, no consistent extension of  $\mathcal{Q}$  can prove a false  $\Pi_1$ -sentence, since the negation of such a sentence would be a true  $\Sigma_1$ -sentence and so would be provable in  $\mathcal{Q}$ .

One can also show that a theory is  $\Sigma_n$ -sound if and only if it is  $\Pi_{n+1}$ -sound: if the theory  $S$  is  $\Sigma_n$ -sound and  $S \vdash \pi$ , where  $\pi$  is a  $\Pi_{n+1}$ -sentence, then write  $\pi = \forall x \sigma(x)$  for a  $\Sigma_n$ -formula  $\sigma$ ; since

<sup>2</sup>Another proof (without appeal to Lemma 3.7): If  $\gamma$  is a Gödelian  $\mathcal{Y}$ -sentence of  $\tau$ , then  $\text{Th}_\tau \vdash \gamma \vee \text{Pr}_\tau(\#\gamma)$  and so  $\mathbb{N} \models \gamma \vee \text{Pr}_\tau(\#\gamma)$  since  $\gamma \vee \text{Pr}_\tau(\#\gamma)$  is a  $\mathcal{Y}$ -sentence and  $\text{Th}_\tau$  is  $\mathcal{Y}$ -sound; as  $\mathbb{N} \not\models \text{Pr}_\tau(\#\gamma)$  by Lemma 3.2 and Convention 2.2.D1, we should have  $\mathbb{N} \models \gamma$ . *QED*

for every  $k \in \mathbb{N}$  we have  $S \vdash \sigma(\bar{k})$ , and  $\sigma(\bar{k})$  is a  $\Sigma_n$ -sentence, then  $\mathbb{N} \models \sigma(\bar{k})$  for every  $k \in \mathbb{N}$ , so  $\mathbb{N} \models \forall x \sigma(x) = \pi$ .

The hierarchy of  $\Sigma_n$ -sound theories is strict, since there exist some  $\Sigma_n$ -sound theories, which are not  $\Sigma_{n+1}$ -sound; see, e.g. [32, Theorem 2.5] or [15, Theorem 4.8]. Therefore, the truth of (even all) the Gödelian  $\Pi_{n+1}$ -sentences (respectively,  $\Sigma_{n+1}$ -sentences) of a theory does not necessarily imply the truth of its Gödelian  $\Pi_{n+2}$ -sentences (respectively,  $\Sigma_{n+2}$ -sentences).  $\diamond$

#### 4 Rosserian sentences and their truth

In Theorem 3.4 we saw that (all of the) Gödelian sentences of some theories could be refutable in those theories (though, they are always unprovable in consistent theories, see Lemma 3.2). Rosser's trick [31] constructs an independent sentence for a given theory when it is consistent (recall that our theories are RE extensions of  $\mathcal{Q}$ ). Before going into Rosser's construction, let us note that no construction similar to Gödel's can result in an independent sentence.

DEFINITION 4.1 (Pseudo-Gödelian sentences).

Let  $\tau$  be a  $\Delta_0$ -formula. Let us call the sentence  $\psi$  a *pseudo-Gödelian sentence* of  $\tau$  when there exist some propositional formulas  $C_1(p), \dots, C_n(p)$ , over the propositional variable  $p$ , and there exists a propositional formula  $B(p_1, \dots, p_n)$ , over the propositional variables  $p_1, \dots, p_n$ , such that we have  $\text{Th}_\tau \vdash \psi \leftrightarrow B(\text{Pr}_\tau[\#C_1(\psi)], \dots, \text{Pr}_\tau[\#C_n(\psi)])$ .  $\diamond$

For example, the sentences  $\mathcal{P}$  and  $\mathcal{R}$  for which we have

$$\text{Th}_\tau \vdash \mathcal{P} \leftrightarrow [\neg \text{Pr}_\tau(\#\mathcal{P}) \wedge \neg \text{Pr}_\tau(\#[\neg\mathcal{P}])]$$

and

$$\text{Th}_\tau \vdash \mathcal{R} \leftrightarrow [\text{Pr}_\tau(\#\mathcal{R}) \rightarrow \text{Pr}_\tau(\#[\neg\mathcal{R}])]$$

are both some pseudo-Gödelian sentences of  $\tau$ .

For a  $\Delta_0$ -formula  $\tau$  such that  $\text{Th}_\tau \supseteq \text{IS}_1$  is consistent, let  $\nu(x) = \tau(x) \vee (x = \#[\neg\text{Con}_\tau])$ . The theory  $\text{Th}_\nu$ , which is  $\text{Th}_\tau + \neg\text{Con}_\tau$ , is consistent by Gödel's second incompleteness theorem. Now, from  $\text{Th}_\nu \vdash \neg\text{Con}_\nu$ , and Remark 2.3.I, we have  $\text{Th}_\nu \vdash \text{Pr}_\nu(\#\theta)$  for every sentence  $\theta$ . Hence,  $\text{Th}_\nu$  decides every pseudo-Gödelian sentence, and so we have the following Proposition (4.2); cf. [38, Exercise 1, p.149].

PROPOSITION 4.2 (On the decidability of pseudo-Gödelian sentences).

Let  $\tau$  be a  $\Delta_0$ -formula, and suppose that the theory  $\text{Th}_\tau$  is consistent and contains  $\text{IS}_1$ . Let  $\nu(x) = \tau(x) \vee (x = \#[\neg\text{Con}_\tau])$ . Then, no pseudo-Gödelian sentence of  $\nu$  can be independent from the theory  $\text{Th}_\nu$ .  $\square$

In the above examples, it can be seen that  $\text{Th}_\nu \vdash \neg\mathcal{P}$  and  $\text{Th}_\nu \vdash \mathcal{R}$  (noting that we have  $\text{Th}_\nu = \text{Th}_\tau + \neg\text{Con}_\tau$ ). Thus, for getting independent sentences (of consistent theories) one should go beyond the (pseudo-)Gödelian sentences.

DEFINITION 4.3 (Rosserian provability and Rosserian sentences).

The Rosserian provability predicate of a  $\Delta_0$ -formula  $\tau$  is

$$\text{R.Pr}_\tau(x) = \exists y [\text{prf}_\tau(y, x) \wedge \forall z < y \neg \text{prf}_\tau(z, \neg x)].$$

If  $\text{Th}_\tau \vdash \rho \leftrightarrow \neg \text{R.Pr}_\tau(\#\rho)$ , then  $\rho$  is called a Rosserian sentence of  $\tau$ .  $\diamond$

Let us note that  $R.Pr_\tau(x)$  is an  $I\Sigma_1$ -provably  $\Sigma_1$ -formula, when  $\tau(x)$  is  $\Delta_0$ ; so,  $\rho$  is an  $I\Sigma_1$ -provably  $\Pi_1$ -sentence, cf. [11, Remark III.4.19]. The independence of the Rosserian sentences (from the theory in question) follows from the following basic properties of the Rosserian provability:

LEMMA 4.4

If  $Th_\tau$  is consistent and contains  $I\Sigma_1$  for a  $\Delta_0$ -formula  $\tau$ , then for every sentence  $\varphi$  we have

- (1)  $Th_\tau \vdash \varphi \iff I\Sigma_1 \vdash R.Pr_\tau(\#\varphi)$ .
- (2)  $Th_\tau \vdash \neg\varphi \implies I\Sigma_1 \vdash \neg R.Pr_\tau(\#\varphi)$ .

PROOF. For (1) it suffices to note that for consistent  $Th_\tau$  we have  $Th_\tau \vdash \varphi$  if and only if the  $I\Sigma_1$ -provably  $\Sigma_1$ -sentence  $R.Pr_\tau(\#\varphi)$  is true. For (2) suppose that  $Th_\tau \vdash \neg\varphi$ ; then by Convention 2.2.C1 we have  $\mathcal{Q} \vdash \text{prf}_\tau(\bar{m}, \#[\neg\varphi])$  for some  $m$ . Now, reason inside  $I\Sigma_1$ :

For any  $y$  with  $\text{prf}_\tau(y, \#\varphi)$  we have  $y > \bar{m}$ , since no  $i \leq \bar{m}$  (which are  $i = 0, \dots, \bar{m}$ ) could satisfy  $\text{prf}_\tau(i, \#\varphi)$  by Convention 2.2.C2, and so for some  $z < y$ , which is  $z = \bar{m}$ , we have  $\text{prf}_\tau(z, \#[\neg\varphi])$ . Thus,  $\forall y[\text{prf}_\tau(y, \#\varphi) \rightarrow \exists z < y \text{prf}_\tau(z, \#[\neg\varphi])]$  holds, and so  $\neg R.Pr_\tau(\#\varphi)$ .  $\square$

Now, we can characterize the Rosserian sentences of super-theories:

PROPOSITION 4.5 (Characterizing Rosserian sentences of super-theories).

Let  $\tau$  be a  $\Delta_0$ -formula such that  $Th_\tau \supseteq I\Sigma_1$ . The following are equivalent for a sentence  $\varphi$ :

- (1)  $\varphi$  is independent from  $Th_\tau$ , i.e.  $Th_\tau \not\vdash \varphi$  and  $Th_\tau \not\vdash \neg\varphi$ ;
- (2)  $\varphi$  is a Rosserian sentence of some consistent extension of  $Th_\tau$ ;

and are implied by the following:

- (3)  $Th_\tau + [\varphi \leftrightarrow \neg R.Pr_\tau(\#\varphi)]$  is consistent.

PROOF. First we show the equivalence of (1) and (2).

(1  $\Rightarrow$  2): By Lemma 2.1 we have  $\mathcal{Q} \vdash \xi \leftrightarrow [\varphi \leftrightarrow \neg R.Pr_{\tau'}(\#\varphi)]$  for some sentence  $\xi$  where  $\tau'(x) = \tau(x) \vee (x = \#\xi)$ . Then  $Th_{\tau'} \vdash \varphi \leftrightarrow \neg R.Pr_{\tau'}(\#\varphi)$ , which shows that  $\varphi$  is a Rosserian sentence of  $\tau'$ . We show that the theory  $Th_{\tau'}$  is consistent. If not, then  $Th_\tau \vdash \neg\xi$ . Thus, we have (\*)  $Th_\tau \vdash \neg[\varphi \leftrightarrow \neg R.Pr_{\tau'}(\#\varphi)]$ . Also,  $Th_{\tau'} \vdash \varphi$  and  $Th_{\tau'} \vdash \neg\varphi$ , and so by Convention 2.2.C1 there are  $m, n \in \mathbb{N}$  such that  $\mathcal{Q} \vdash \text{prf}_{\tau'}(\bar{m}, \#\varphi)$  and  $\mathcal{Q} \vdash \text{prf}_{\tau'}(\bar{n}, \#[\neg\varphi])$ ; we can assume that  $m$  and  $n$  are the least such numbers.

(i) If  $m \leq n$ , then  $I\Sigma_1 \vdash \text{prf}_{\tau'}(\bar{m}, \#\varphi) \wedge \forall z < \bar{m} \neg \text{prf}_{\tau'}(z, \#[\neg\varphi])$  and so  $I\Sigma_1 \vdash R.Pr_{\tau'}(\#\varphi)$ , which implies by (\*) that  $Th_\tau \vdash \varphi$ , contradicting (1).

(ii) If  $n < m$ , then we have  $I\Sigma_1$ -provably that for every  $y$ :

$$\begin{aligned} & \text{prf}_{\tau'}(y, \#\varphi) \\ \rightarrow & y \geq \bar{m} && \text{since } m \text{ is the least with } \text{prf}_{\tau'}(m, \#\varphi) \\ \rightarrow & \exists z < y \text{prf}_{\tau'}(z, \#[\neg\varphi]) && \text{since one can take } z = \bar{n} (< \bar{m} \leq y). \end{aligned}$$

So,  $I\Sigma_1 \vdash \neg R.Pr_{\tau'}(\#\varphi)$ , which implies by (\*) that  $Th_\tau \vdash \neg\varphi$ , and this contradicts (1).

Thus,  $Th_{\tau'}$  must be consistent.

(2  $\Rightarrow$  1): Suppose that  $Th_{\tau'}$  is a consistent extension of  $Th_\tau$  such that  $\varphi$  is a Rosserian sentence of it. It suffices to show that  $\varphi$  is independent from  $Th_\tau$ . If  $Th_{\tau'} \vdash \varphi$ , then we should have on the one hand  $Th_{\tau'} \vdash R.Pr_{\tau'}(\#\varphi)$  by Lemma 4.4.1 and on the other hand  $Th_{\tau'} \vdash \neg R.Pr_{\tau'}(\#\varphi)$  by Definition 4.3; contradicting the consistency of  $Th_{\tau'}$ . Also,  $Th_{\tau'} \vdash \neg\varphi$  would imply on the one hand  $Th_{\tau'} \vdash \neg R.Pr_{\tau'}(\#\varphi)$  by Lemma 4.4.2 and on the other hand  $Th_{\tau'} \vdash R.Pr_{\tau'}(\#\varphi)$  by Definition 4.3; contradicting  $Th_{\tau'}$ 's consistency again.

Now, we show that (3) implies (1), and so (2) too.

(3  $\Rightarrow$  1): Note that  $\text{Th}_\tau$  is consistent by the assumption. If  $\text{Th}_\tau \vdash \varphi$ , then  $\text{Th}_\tau \vdash \text{R}.\text{Pr}_\tau(\#\varphi)$  by Lemma 4.4.1, and so  $\text{Th}_\tau \vdash \neg[\varphi \leftrightarrow \neg\text{R}.\text{Pr}_\tau(\#\varphi)]$ . If  $\text{Th}_\tau \vdash \neg\varphi$ , then Lemma 4.4.2 would imply that  $\text{Th}_\tau \vdash \neg\text{R}.\text{Pr}_\tau(\#\varphi)$ , and so  $\text{Th}_\tau \vdash \neg[\varphi \leftrightarrow \neg\text{R}.\text{Pr}_\tau(\#\varphi)]$  would hold again.  $\square$

REMARK 4.6 (Löb's rule for Rosserian provability).

Let us note that the contraposition of the implication (1  $\Rightarrow$  3) in Proposition 4.5 says that if  $\text{Th}_\tau \vdash \varphi \leftrightarrow \text{R}.\text{Pr}_\tau(\#\varphi)$ , i.e. if  $\varphi$  is a *Rosser-type Henkin sentence*,<sup>3</sup> so-called in [16], then  $\varphi$  is not independent from  $T$ . Actually, it is shown in [16] that there are *standard* proof predicates (i.e. those satisfying Convention 2.2), which have independent Rosser-type Henkin sentences, and there are standard proof predicates none of whose Rosser-type Henkin sentences are independent. The latter proof predicates satisfy (1  $\Rightarrow$  3) in Proposition 4.5 and satisfy a Löb-like rule for Rosserian provability, while the former ones do not satisfy (1  $\Rightarrow$  3) in Proposition 4.5 and do not satisfy any Löb-like rule for Rosserian provability. So, the implication (1  $\Rightarrow$  3) in Proposition 4.5 depends on  $\text{prf}_\tau(y, x)$ , and is not robust.  $\diamond$

Unlike Gödelian  $\Pi_1$ -sentences, all the Rosserian  $\Pi_1$ -sentences of consistent theories are true, and like Gödelian  $\Sigma_1$ -sentence, all of their Rosserian  $\Sigma_1$ -sentences are false:

THEOREM 4.7 (On the truth of Rosserian  $\Pi_1$ - and  $\Sigma_1$ -sentences).

For an arbitrary  $\Delta_0$ -formula  $\tau(x)$ , every Rosserian  $\Pi_1$ -sentence of  $\tau$  is true and every Rosserian  $\Sigma_1$ -sentence of  $\tau$  is false, if  $\text{Th}_\tau$  is consistent and contains  $\text{IS}_1$ .

PROOF. If a Rosserian  $\Pi_1$ -sentence of  $\tau$  were false, then its negation would be a true  $\Sigma_1$ -sentence, and so would be provable in  $\mathcal{Q}$ ; contradicting Rosser's theorem on the independence of Rosserian sentences (see Proposition 4.5). If a Rosserian  $\Sigma_1$ -sentence were true, then it would be provable in  $\mathcal{Q}$ ; contradicting the unprovability of Rosserian sentences.  $\square$

However, for  $n \geq 1$ , the truth of all the Gödelian  $\Pi_{n+1}$ -sentences is equivalent to the truth of all the Rosserian  $\Pi_{n+1}$ -sentences and the truth of all the Gödelian  $\Sigma_{n+1}$ -sentences is equivalent to the truth of all the Rosserian  $\Sigma_{n+1}$ -sentences:

THEOREM 4.8 (On the truth of Rosserian  $\Pi_{n+1}$ - and  $\Sigma_{n+1}$ -sentences).

Let  $n \geq 1$ , and let  $\tau$  be a  $\Delta_0$ -formula such that  $\text{Th}_\tau \supseteq \text{IS}_1$ .

All the Rosserian  $\Pi_{n+1}$ -sentences of  $\tau$  are true iff  $\text{Th}_\tau$  is  $\Pi_{n+1}$ -sound.

All the Rosserian  $\Sigma_{n+1}$ -sentences of  $\tau$  are true iff  $\text{Th}_\tau$  is  $\Sigma_{n+1}$ -sound.

PROOF. Let  $\mathcal{Y}$  be either of the two classes of sentences (either  $\Pi_{n+1}$  or  $\Sigma_{n+1}$ ). If  $\text{Th}_\tau$  is  $\mathcal{Y}$ -sound and  $\rho$  is a Rosserian  $\mathcal{Y}$ -sentence of  $\tau$ , then  $\text{R}.\text{Pr}_\tau(\#\rho)$  is a false  $\text{IS}_1$ -provably  $\Sigma_1$ -sentence by Proposition 4.5 and Lemma 4.4.1. Since  $\text{Th}_\tau \vdash \rho \vee \text{R}.\text{Pr}_\tau(\#\rho)$  and  $\rho \vee \text{R}.\text{Pr}_\tau(\#\rho)$  is a  $\mathcal{Y}$ -sentence, then  $\mathbb{N} \models \rho \vee \text{R}.\text{Pr}_\tau(\#\rho)$ , and so  $\rho$  is true. Now, suppose that all the Rosserian  $\mathcal{Y}$ -sentence of  $\tau$  are true and  $\text{Th}_\tau \vdash \xi$ , where  $\xi$  is a  $\mathcal{Y}$ -sentence. By Lemma 2.1 there is a  $\mathcal{Y}$ -sentence  $\rho$  such that  $\text{IS}_1 \vdash \rho \leftrightarrow [\xi \wedge \neg\text{R}.\text{Pr}_\tau(\#\rho)]$ . So,  $\rho$  is a Rosserian  $\mathcal{Y}$ -sentence of  $\tau$ ; thus, it is true by the assumption. Therefore, by the soundness of  $\text{IS}_1$  the sentence  $\xi$  is true too. Hence,  $\text{Th}_\tau$  is  $\mathcal{Y}$ -sound.  $\square$

Therefore, all the Rosserian sentences of  $\tau$  are true if and only if  $\text{Th}_\tau$  is sound; cf. also [36, Theorem 24.7].

---

<sup>3</sup>The sentence  $\psi$  is a *Henkin sentence* (of  $\tau$ ) when it is equivalent to its own provability in the theory, i.e. when we have  $\text{Th}_\tau \vdash \psi \leftrightarrow \text{Pr}_\tau(\#\psi)$ . A *Rosser-type Henkin sentence*  $\varphi$  is equivalent to its own Rosserian provability in the theory, i.e. we have that  $\text{Th}_\tau \vdash \varphi \leftrightarrow \text{R}.\text{Pr}_\tau(\#\varphi)$ .

### 5 Conclusions

The first one who talked about the truth of Gödelian sentences was Gödel himself [9]. This turned into a serious debate with [10] in which (what we call now) the *Gödel Disjunction* was announced; see [8] and [12], and the references therein. The so called *Anti-Mechanism Thesis*, or the *Lucas-Penrose Argument*, started with [19] and was popularized by [25]; see also [24] and [29]. After that, there has been a large discussion on the truth of Gödelian sentences; see, e.g. [7], [37, 38], [1], [34], [20, 21], [39], [30], [26], [22], [33], [13], [3], [5, 6], [27, 28] and [23].

As argued above, the consistency of a theory does not imply the truth of (all of) its Gödelian  $\Pi_1$ -sentences, but does imply the truth of its all Rosserian  $\Pi_1$ -sentences. One may wonder why the proponents of the anti-mechanism thesis have not used the Rosserian  $\Pi_1$ -sentences in their arguments, given that the truth of those sentences is straightforward (and immediately follows from the consistency of the theory). Though, the opponents have argued that actually for ‘seeing’ the truth of Gödelian  $\Pi_1$ -sentences one should ‘see’ (at least) the consistency of the theory (and indeed, more than that).

The following diagram summarizes our old and new results. Note that the conditions get (strictly) stronger from bottom to top. As the diagram shows, if one faces the question as to whether a given Gödelian sentence  $\gamma$  of a consistent and sufficiently strong RE theory  $T$  is true or not, then one should consider the complexity of the sentence: if  $\gamma$  is  $\Sigma_1$ , then it is false; if  $\gamma$  is  $\Pi_1$ , then it is true when  $T$  is consistent with the consistency statement of  $T$ ; if  $\gamma$  is  $\Pi_2$ , then it is true when  $T$  is  $\Sigma_1$ -sound; and, finally, if  $\gamma$  is  $\Sigma_{n+1}$  or  $\Pi_{n+2}$  for some  $n \geq 1$ , then it is true when  $T$  is  $\Sigma_{n+1}$ -sound. Let  $\rho$  be a Rosserian sentence of such a theory  $T$ ; if  $\rho$  is  $\Sigma_1$ , then it is false; if  $\rho$  is  $\Pi_1$ , then it is true; if  $\rho$  is  $\Pi_2$ , then it is true when  $T$  is  $\Sigma_1$ -sound; and if  $\rho$  is  $\Sigma_{n+1}$  or  $\Pi_{n+2}$  for some  $n \geq 1$ , then it is true when  $T$  is  $\Sigma_{n+1}$ -sound.

Soundness	≡	Truth of all the Gödelian sentences
	≡	Truth of all the Rosserian sentences
.....	.	.....
$(n \geq 1)$ $\Sigma_{n+1}$ -soundness	≡	$\Pi_{n+2}$ -soundness
	≡	Truth of all the Gödelian $\Sigma_{n+1}$ -sentences
	≡	Truth of all the Gödelian $\Pi_{n+2}$ -sentences
	≡	Truth of all the Rosserian $\Sigma_{n+1}$ -sentences
	≡	Truth of all the Rosserian $\Pi_{n+2}$ -sentences
.....	.	.....
.....	.	.....
$\Sigma_1$ -soundness	≡	$\Pi_2$ -soundness
	≡	Truth of all the Gödelian $\Pi_2$ -sentences
	≡	Truth of all the Rosserian $\Pi_2$ -sentences
.....	.	.....
Consistency with $\text{Con}_T$	≡	Truth of all the Gödelian $\Pi_1$ -sentences
.....	.	.....
Consistency	≡	$\Pi_1$ -soundness
	≡	Truth of all the Rosserian $\Pi_1$ -sentences
	≡	Falsity of all the Gödelian $\Sigma_1$ -sentences
	≡	Falsity of all the Rosserian $\Sigma_1$ -sentences

## Acknowledgements

This research is supported by the grant 98013437 of the Iran National Science Foundation (INSF). The authors warmly thank Kaave Lajevardi for the most helpful discussions and comments. This is a continuation of a project that he started a while ago, which have resulted in e.g. [17, 18].

## References

- [1] G. Boolos. On “seeing” the truth of the Gödel sentence. *Behavioral and Brain Sciences*, **13**, 655–656, 1990. Reprinted in: R. Jeffrey, ed., G. Boolos, *Logic, Logic, and Logic*, pp. 389–391. Harvard University Press, 1999.
- [2] G. Boolos. *The Logic of Provability*. Cambridge University Press, 1994.
- [3] J. Boyer and G. Sandu. Between proof and truth. *Synthese*, **187**, 821–832. Erratum: *ibid*, 973–974, 2012.
- [4] W. Craig. On axiomatizability within a system. *The Journal of Symbolic Logic*, **18**, 30–32, 1953.
- [5] V. Drăghici. How do we know that G is true? *Logos Architekton, Journal of Logic and Philosophy of Science*, **6**, 39–65, 2012.
- [6] V. Drăghici. Is G true by Gödel’s theorem? *Logos Architekton, Journal of Logic and Philosophy of Science*, **7**, 53–59, 2013.
- [7] M. Dummett. The philosophical significance of Gödel’s theorem. *Ratio*, **5**, 140–155, 1963. Reprinted in *Truth and Other Enigmas*, M. Dummett, pp. 186–201. Harvard University Press (1978, 6th print), 1996.
- [8] S. Feferman. Are there absolutely unsolvable problems? Gödel’s dichotomy. *Philosophia Mathematica*, **14**, 134–152, 2006.
- [9] K. Gödel. Über Formal Unentscheidbare Sätze der Principia Mathematica und Verwandter Systeme I. *Monatshefte für Mathematik und Physik*, **38**, 173–198, 1931 (in German). English Translation: On Formally Undecidable Propositions of Principia Mathematica and Related Systems I. In Kurt Gödel Collected Works, Volume I: Publications 1929–1936, S. Feferman, et al., eds., pp. 135–152. Oxford University Press, 1986.
- [10] K. Gödel. Some basic theorems on the foundations of mathematics and their implications. In *Kurt Gödel Collected Works, Volume III: Unpublished Essays and Lectures*, S. Feferman, et al., eds., pp. 304–323. Oxford University Press (1995), 1951.
- [11] P. Hájek and P. Pudlák. *Metamathematics of First-Order Arithmetic*. Springer, 1993.
- [12] L. Horsten and P. Welch, eds., *Gödel’s Disjunction: The Scope and Limits of Mathematical Knowledge*. Oxford University Press, 2016.
- [13] D. Isaacson. Necessary and sufficient conditions for undecidability of the Gödel sentence and its truth. In *Logic, Mathematics, Philosophy, Vintage Enthusiasms*, D. DeVidi, M. Hallett and P. Clarke, eds., pp. 135–152. Springer, 2011.
- [14] R. Kaye. *Models of Peano Arithmetic*. Oxford University Press, 1991.
- [15] M. Kikuchi and T. Kurahashi. Generalizations of Gödel’s incompleteness theorems for  $\Sigma_n$ -definable theories of arithmetic. *The Review of Symbolic Logic*, **10**, 603–616, 2017.
- [16] T. Kurahashi. Henkin sentences and local reflection principles for Rosser provability. *Annals of Pure and Applied Logic*, **167**, 73–94, 2016.
- [17] K. Lajevardi and S. Salehi. On the arithmetical truth of self-referential sentences. *Theoria: A Swedish Journal of Philosophy*, **85**, 8–17, 2019.
- [18] K. Lajevardi and S. Salehi. There may be many arithmetical Gödel Sentences. *Philosophia Mathematica*, **29**, 278–287, 2021.

- [19] J. R. Lucas. Minds, machines and Gödel. *Philosophy*, **36**, 112–127, 1961. Reprinted in: *The Modeling of Mind*, K.M. Sayre and F.J. Crosson, eds., pp. 255–271. University of Notre Dame Press, 1963; and in *Minds and Machines*, A.R. Anderson, ed., pp. 43–59. Prentice Hall, 1964.
- [20] S. McCall. Can a Turing machine know that the Gödel sentence is true? *The Journal of Philosophy*, **96**, 525–532, 1999. Reprinted in: *The Consistency of Arithmetic, and Other Essays*, S. McCall, pp. 28–35. Oxford University Press, 2014.
- [21] S. McCall. On “seeing” the truth of the Gödel sentence. *Facta Philosophica*, **3**, 25–29, 2001. Reprinted in: *The Consistency of Arithmetic, and Other Essays*, S. McCall, pp. 36–40. Oxford University Press, 2014.
- [22] P. Milne. On Gödel sentences and what they say. *Philosophia Mathematica*, **15**, 193–226, 2007.
- [23] E. Moriconi. Some remarks on true undecidable sentences. In *Truth, Existence and Explanation*, M. Piazza and G. Pulcini, eds., pp. 3–15. Springer, 2018.
- [24] E. Nagel and J. R. Newman. *Gödel’s Proof*. New York University Press, 1958 (revised: 3rd edn. Routledge, 2005).
- [25] R. Penrose. *The Emperor’s New Mind: Concerning Computers, Minds, and The Laws of Physics*. Oxford University Press, 1989 (republished with a new preface, 1999).
- [26] J. Peregrin. Kurt Gödel, completeness, incompleteness. *Journal of Physics: Conference Series*, **82**, pp. 1–4, 2007.
- [27] M. Piazza and G. Pulcini. A deflationary account of the truth of the Gödel sentence  $\mathcal{G}$ . In *From Logic to Practice*, G. Lolli, M. Panza and G. Venturi, eds., pp. 71–90. Springer, 2015.
- [28] M. Piazza and G. Pulcini. What’s so special about the Gödel sentence  $\mathcal{G}$ ? In *Objectivity, Realism, and Proof*, F. Boccuni and A. Sereni, eds., pp. 245–263. Springer, 2016.
- [29] H. Putnam. Minds and machines. In *Dimensions of Mind: A Symposium*, S. Hood, ed., pp. 138–164. New York University Press, 1960. Reprinted in: *Philosophical Papers, Volume 2: Mind, Language and Reality*, H. Putnam, pp. 362–385. Harvard University Press, 1975.
- [30] P. Raatikainen. On the philosophical relevance of Gödel’s incompleteness theorems. *Revue Internationale de Philosophie*, **59**, 513–534, 2005.
- [31] B. Rosser. Extensions of some theorems of Gödel and Church. *The Journal of Symbolic Logic*, **1**, 87–91, 1936.
- [32] S. Salehi and P. Seraji. Gödel–Rosser’s incompleteness theorem, generalized and optimized for definable theories. *Journal of Logic and Computation*, **27**, 1391–1397, 2017.
- [33] G. Serény. How do we know that the Gödel sentence of a consistent theory is true? *Philosophia Mathematica*, **19**, 47–73, 2011.
- [34] S. Shapiro. Induction and indefinite extensibility: the Gödel sentence is true, but did someone change the subject? *Mind*, **107**, 597–624, 1998.
- [35] T. A. Slaman.  $\Sigma_n$ -bounding and  $\Delta_n$ -induction. *Proceedings of The American Mathematical Society*, **132**, 2449–2456, 2004.
- [36] P. Smith. *An Introduction to Gödel’s Theorems*, 2nd edn. Cambridge University Press, 2013.
- [37] C. Smoryński. The incompleteness theorems. In *Handbook of Mathematical Logic*, J. Barwise, ed., pp. 821–865. North Holland, 1977.
- [38] C. Smoryński. *Self-Reference and Modal Logic*. Springer, 1985.
- [39] N. Tennant. On Turing machines knowing their own Gödel sentences. *Philosophia Mathematica*, **9**, 72–79, 2001.
- [40] A. Visser. Another look at the second incompleteness theorem. *The Review of Symbolic Logic*, **13**, 269–295, 2020.