

## GÖDEL'S SECOND INCOMPLETENESS THEOREM: HOW IT IS DERIVED AND WHAT IT DELIVERS

SAEED SALEHI

**Abstract.** The proofs of Gödel (1931), Rosser (1936), Kleene (first 1936 and second 1950), Chaitin (1970), and Boolos (1989) for the first incompleteness theorem are compared with each other, especially from the viewpoint of the second incompleteness theorem. It is shown that Gödel's (first incompleteness theorem) and Kleene's first theorems are equivalent with the second incompleteness theorem, Rosser's and Kleene's second theorems do deliver the second incompleteness theorem, and Boolos' theorem is derived from the second incompleteness theorem in the standard way. It is also shown that none of Rosser's, Kleene's second, or Boolos' theorems is equivalent with the second incompleteness theorem, and Chaitin's incompleteness theorem neither delivers nor is derived from the second incompleteness theorem. We compare (the strength of) these six proofs with one another.

**§1. Introduction and preliminaries.** The first incompleteness theorem states the existence of a  $\Pi_1$ -sentence  $\psi$  for a given sufficiently strong and recursively enumerable (RE) arithmetical theory  $T$  such that if  $T$  is consistent, then  $\mathbb{N} \models \psi$  and  $T \not\vdash \psi$ . If  $T$  is, moreover,  $\Sigma_1$ -sound (i.e., every  $T$ -provable  $\Sigma_1$ -sentence is true in the standard model of natural numbers  $\mathbb{N}$ ), then we also have  $T \not\vdash \neg\psi$  (since if we had  $T \vdash \neg\psi$ , then we would have  $\mathbb{N} \models \neg\psi$  by the  $\Sigma_1$ -soundness of  $T$  and the fact that  $\neg\psi$  is a  $\Sigma_1$ -sentence). The  $\Pi_1$ -sentence  $\psi$  depends on the theory  $T$ , or more precisely, on how  $T$  is presented. An RE theory  $T$  may be presented (given) by, for example, an input-free Turing machine (or a program) that outputs a set of axioms for the theory  $T$  (after running). It is known that a theory  $T$  is RE if and only if it can be defined by a  $\Sigma_1$ -formula; that is, for some  $\Sigma_1$ -formula  $\xi(x)$ , the set  $\{\theta \mid \mathbb{N} \models \xi(\ulcorner\theta\urcorner)\}$  axiomatizes  $T$ , where  $\theta$  ranges over the sentences and  $\ulcorner\theta\urcorner$  denotes the Gödel code of  $\theta$  (see, e.g., [13, Theorem 3.3]). By Craig's trick [8], every such theory can be axiomatized by a  $\Delta_0$ -definable set of axioms (see, e.g., [24, Lemma 2.4] or [29, Section 4.3.4]): if  $\xi(x) = \exists y \zeta(y, x)$  for some  $\Delta_0$ -formula  $\zeta$ , then the  $\Delta_0$ -formula  $\tau(x) = \exists u, v \leq x [x = (u \wedge \bar{v} = \bar{v}) \wedge \zeta(v, u)]$  defines another axiomatization for the theory  $T$ .

---

Received January 21, 2019.

2020 *Mathematics Subject Classification.* 03F40, 03F30, 68Q30, 03D32.

*Key words and phrases.* first incompleteness theorem, second incompleteness theorem, Gödel's proof, Rosser's proof, Kleene's proof, Chaitin's proof, Boolos' proof.

Any given  $\Delta_0$ -formula  $\tau(x)$ , with the only one free variable  $x$ , defines the theory  $\mathcal{Th}_\tau = \{\theta \mid \mathbb{N} \models \tau(\ulcorner \theta \urcorner)\}$ , where  $\theta$  ranges over sentences. By some fixed Gödel coding, one can construct a  $\Delta_0$ -formula  $\text{prf}_\tau(y, x)$ , called *the proof predicate* of  $\tau$ , in the language of arithmetic which states that “ $y$  is (the Gödel code of) a proof in  $\mathcal{Th}_\tau$  of the sentence (with the Gödel code)  $x$ ” (see [1, p. 215] or [3, p. 204]). Then *the provability predicate* of a system  $\tau$  is the  $\Sigma_1$ -formula  $\text{Pr}_\tau(x) = \exists y \text{prf}_\tau(y, x)$ , and  $\text{Con}_\tau = \neg \text{Pr}_\tau(\ulcorner 0 \neq 0 \urcorner)$  is *the consistency statement* of  $\tau$ . Let us note that a theory may have different axiomatizations, and even one single axiomatization of it may have different defining formulas, and so different proof (and provability) predicates, and different consistency statements.

Let us fix a *Base Theory*  $\mathfrak{B}$ , which is an RE theory such that:

- The theory  $\mathfrak{B}$  is a *sound* extension of Robinson’s arithmetic (therefore, the theory  $\mathfrak{B}$  is  $\Sigma_1$ -complete, i.e., can prove all the true  $\Sigma_1$ -sentences, and can *strongly represent* all the recursive functions; see [28]).
- The theory  $\mathfrak{B}$  can prove the *Derivability* (or *Provability*) *Conditions* of Gödel, Hilbert, Bernays, and Löb (see p. 3 below, and cf. [6, 27]).

Of course, Peano’s Arithmetic could be taken for  $\mathfrak{B}$ , though it is too strong for that. However, Robinson’s arithmetic ( $\mathbb{Q}$  or  $\mathbb{R}$ ) seems too weak to be such a base theory (though, we have no concrete proof for, e.g.,  $\mathbb{Q}$ ’s weakness at hand). By [1, Proposition 16] the Elementary Arithmetic EA may suffice for us (cf. [29, Remark 6.7] where it is argued that one needs  $\text{EA} + \text{B}\Sigma_1$ , or equivalently  $\text{I}\Delta_1$  by [26], for handling the  $\Sigma_1$ -formulas  $\text{Pr}_\tau$ ). To stay on the safe side one can take for  $\mathfrak{B}$  the finitely axiomatizable theory  $\text{I}\Sigma_1$  (the fragment of Peano’s arithmetic where the induction axiom scheme is restricted to  $\Sigma_1$ -formulas). One good reason (other than the ability of  $\text{I}\Sigma_1$  to arithmetize the syntax and prove the basic propositions of it, see [11]) is that we will need a variant of the proof predicate, denoted  $\overline{\text{prf}}_\tau(y, x)$ , whose all Rosserian sentences are equivalent with each other (in the base theory); and for that Primitive Recursive Arithmetic PRA suffices (see [27, Chapter 6, Theorem 3.6]).

By *a system*, we mean a  $\Delta_0$ -formula  $\tau(x)$ , with the only one free variable  $x$ , such that  $\mathcal{Th}_\tau \vdash \mathfrak{B}$ . A system is consistent (or  $\Sigma_1$ -sound) when  $\mathcal{Th}_\tau$  is a consistent (or  $\Sigma_1$ -sound) theory.

A mapping  $\mathcal{F}: \tau \mapsto \mathcal{F}_\tau$  which assigns a  $\Pi_1$ -sentence  $\mathcal{F}_\tau$  to any given system  $\tau$  is called a  $\Pi_1$ -*incompleteness witness* when for every consistent system  $\tau$  we have  $\mathbb{N} \models \mathcal{F}_\tau$  and  $\mathcal{Th}_\tau \not\vdash \mathcal{F}_\tau$ . In this paper, we investigate the  $\Pi_1$ -incompleteness witnesses of Gödel [9], Rosser [22], Kleene (first [15] and second [16]), Chaitin [7] and Boolos [5]. Our purpose is comparing those  $\Pi_1$ -incompleteness witnesses with each other, and with *Gödel’s second incompleteness theorem*, which is the following statement:

If  $\tau$  is a consistent system, then  $\mathcal{Th}_\tau \not\vdash \text{Con}_\tau$ . ( $\mathbb{G}_2$ )

Let  $\mathcal{F}$  and  $\mathcal{H}$  be two  $\Pi_1$ -incompleteness witnesses. We say that  $\mathcal{F}$  is *derived from*  $\mathcal{H}$ , or  $\mathcal{H}$  *delivers*  $\mathcal{F}$  (in the standard way) denoted  $\mathcal{F} \preceq \mathcal{H}$  when for

every system  $\tau$  we have  $\mathcal{T}h_\tau \vdash \mathcal{F}_\tau \rightarrow \mathcal{H}_\tau$ . We note that then the unprovability of  $\mathcal{F}_\tau$  (in  $\mathcal{T}h_\tau$ ) follows from the unprovability of  $\mathcal{H}_\tau$ .

Some results of this paper are summarized in the following table:

$\Pi_1$ -incompleteness	Does deliver $\mathbb{G}_2$	Is derived from $\mathbb{G}_2$
Gödel <sub>1</sub> (1931) [9]	✓	✓
Kleene <sub>1</sub> (1936) [15]	✓	✓
Rosser (1936) [22]	✓	✗
Kleene <sub>2</sub> (1950) [16]	✓	✗
Chaitin (1970) [7]	✗	✗
Boolos (1989) [5]	✗	✓

Let us denote Gödel's (respectively, Kleene's) first  $\Pi_1$ -incompleteness witness by  $\mathbb{G}$  (respectively,  $\mathbb{K}$ ); let  $\overline{\mathbb{R}}$  (respectively,  $\overline{\mathbb{K}'}$ ) denote an alternative version of Rosser's (respectively, Keelen's second)  $\Pi_1$ -incompleteness witness (which will be rigorously determined later). Let  $\mathbb{C}$  denote (one of the infinitely many  $\Pi_1$ -sentences that) Chaitin's incompleteness witness (proves to be true and unprovable in consistent systems). Finally, we denote by  $\tilde{\mathbb{B}}$  a  $\Pi_1$ -incompleteness witness that is very similar to that of Boolos, but substantially different from his original formulation. Some results of our comparison are depicted in the following diagram (where  $\mathcal{F} \cong \mathcal{H}$  abbreviates  $\mathcal{F} \preceq \mathcal{H} \preceq \mathcal{F}$  and  $\mathcal{F} \approx \mathcal{H}$  abbreviates  $\mathcal{F} \preceq \mathcal{H} \not\preceq \mathcal{F}$ ):

$$\mathbb{C} \approx \tilde{\mathbb{B}} \approx \mathbb{G}_2 \cong \mathbb{K} \cong \mathbb{G} \approx \overline{\mathbb{R}} \cong \overline{\mathbb{K}'}$$

**§2. Proofs of Gödel, Rosser, and Kleene.** For  $n \in \mathbb{N}$ , let  $\bar{n}$  denote its numeral (the term representing  $n$  in the language of arithmetic). A function  $f : \mathbb{N} \rightarrow \mathbb{N}$  is said to be *strongly representable* in  $\mathfrak{B}$  when for a formula  $\langle\langle f(x) = y \rangle\rangle$  in its language, with the only free variables  $x$  and  $y$ , we have

$$\text{for all } m, n \in \mathbb{N}, \text{ if } f(m) = n, \text{ then } \mathfrak{B} \vdash \forall y [\langle\langle f(\bar{m}) = y \rangle\rangle \leftrightarrow y = \bar{n}].$$

The derivability conditions are the following (for a system  $\tau$ ):

- (C<sub>1</sub>)  $\mathfrak{B} \vdash \text{Pr}_\tau(\ulcorner \theta \urcorner)$  if and only if  $\mathcal{T}h_\tau \vdash \theta$ , for all sentences  $\theta$ .
- (C<sub>2</sub>)  $\mathfrak{B} \vdash \text{Pr}_\tau(\ulcorner \theta \rightarrow \eta \urcorner) \rightarrow [\text{Pr}_\tau(\ulcorner \theta \urcorner) \rightarrow \text{Pr}_\tau(\ulcorner \eta \urcorner)]$ , for all sentences  $\theta$  and  $\eta$ .
- (C<sub>3</sub>)  $\mathfrak{B} \vdash \sigma \rightarrow \text{Pr}_\tau(\ulcorner \sigma \urcorner)$ , for all  $\Sigma_1$ -formulas  $\sigma$ .

Let us start with a straightforward observation about the consistency statements:

LEMMA 2.1. *For a system  $\tau$  and a sentence  $\theta$ , we have*

- (1)  $\mathfrak{B} \vdash \neg \text{Con}_\tau \rightarrow \text{Pr}_\tau(\ulcorner \theta \urcorner)$  and
- (2)  $\mathfrak{B} \vdash \text{Pr}_\tau(\ulcorner \neg \theta \urcorner) \wedge \text{Pr}_\tau(\ulcorner \theta \urcorner) \rightarrow \neg \text{Con}_\tau$ .

PROOF. Both parts of the lemma follow from the derivability conditions by using the tautologies  $0 \neq 0 \rightarrow \theta$  for (1), and  $\neg \theta \rightarrow (\theta \rightarrow 0 \neq 0)$  for (2).  $\dashv$

In the following proposition, we show some necessary and sufficient conditions for delivering  $\mathbb{G}_2$ , and being derivable from  $\mathbb{G}_2$ , in the standard way (cf. [19]):

**PROPOSITION 2.2.** *For every system  $\tau$  and every  $\Pi_1$ -sentence  $\psi$ , we have*

- (1)  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \psi$  if and only if  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \psi$ ;
- (2)  $\mathcal{Th}_\tau \vdash \psi \rightarrow \text{Con}_\tau$  if and only if  $\mathcal{Th}_\tau \vdash \psi \rightarrow \neg \text{Pr}_\tau(\ulcorner \psi \urcorner)$ .

**PROOF.** (1 $\Rightarrow$ ): Suppose that  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \psi$  holds. By Lemma 2.1(1) we have  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \text{Con}_\tau$ ; so  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \psi$  follows.

(1 $\Leftarrow$ ): Now, suppose  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \psi$ ; so,  $\mathcal{Th}_\tau \vdash \neg \psi \rightarrow \text{Pr}_\tau(\ulcorner \psi \urcorner)$  holds. Also, since  $\neg \psi$  is  $\Sigma_1$ , by (C<sub>3</sub>), we have  $\mathcal{Th}_\tau \vdash \neg \psi \rightarrow \text{Pr}_\tau(\ulcorner \neg \psi \urcorner)$ . Thus, by Lemma 2.1(2) we get  $\mathcal{Th}_\tau \vdash \neg \psi \rightarrow \neg \text{Con}_\tau$ , and so  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \psi$ .

(2 $\Rightarrow$ ): Suppose that we have  $\mathcal{Th}_\tau \vdash \psi \rightarrow \text{Con}_\tau$ ; then by (C<sub>1</sub>) and (C<sub>2</sub>), it follows that  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \text{Con}_\tau \urcorner) \rightarrow \neg \text{Pr}_\tau(\ulcorner \psi \urcorner)$ . Now, by Löb's theorem (or, formalized  $\mathbb{G}_2$ , see [6]) we have  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \neg \text{Pr}_\tau(\ulcorner \text{Con}_\tau \urcorner)$ . Therefore, we have the desired conclusion  $\mathcal{Th}_\tau \vdash \psi \rightarrow \neg \text{Pr}_\tau(\ulcorner \psi \urcorner)$ .

(2 $\Leftarrow$ ): Now, suppose that  $\mathcal{Th}_\tau \vdash \psi \rightarrow \neg \text{Pr}_\tau(\ulcorner \psi \urcorner)$  holds. Then by Lemma 2.1(1), that  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \text{Con}_\tau$ , we get  $\mathcal{Th}_\tau \vdash \psi \rightarrow \text{Con}_\tau$ .  $\dashv$

Let us note that the assumption “ $\psi$  is a  $\Pi_1$ -sentence” is used only in the proof of (1 $\Leftarrow$ ): if  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \psi$ , then  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \psi$ . The rest of the implications in Proposition 2.2 hold for arbitrary  $\psi$ . That  $\psi$  should be  $\Pi_1$  in (1 $\Leftarrow$ ) can be seen by the following example: Let  $\tau$  be a system such that  $\mathcal{Th}_\tau \vdash \text{Pr}_\tau(\neg \text{Con}_\tau)$  and  $\mathcal{Th}_\tau \not\vdash \neg \text{Con}_\tau$ ; one can take  $\tau(x)$  to be the system  $\zeta(x) \vee [x = \text{Pr}_\zeta(\neg \text{Con}_\zeta)]$  for a system  $\zeta$  with  $\mathcal{Th}_\zeta \not\vdash \neg \text{Con}_\zeta$ , see [6] or [12, Theorem 36]. Then, for the  $\Sigma_1$ -sentence  $\psi = \neg \text{Con}_\tau$  we have  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \psi \urcorner) \rightarrow \psi$ , but  $\mathcal{Th}_\tau \not\vdash \text{Con}_\tau \rightarrow \psi$ .

**2.1. Gödel's proof.** The proof of Gödel constructs a  $\Pi_1$ -sentence  $\gamma$ , for a given system  $\tau$ , such that (\*)  $\mathfrak{B} \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  holds. If  $\mathcal{Th}_\tau$  is consistent, then the unprovability of  $\gamma$  from  $\mathcal{Th}_\tau$  follows from the derivability conditions (if  $\mathcal{Th}_\tau \vdash \gamma$ , then  $\mathcal{Th}_\tau \vdash \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  by (C<sub>1</sub>), and also  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  by (\*), contradicting the consistency of  $\mathcal{Th}_\tau$ ), and the truth of  $\gamma$  follows from the soundness of  $\mathfrak{B}$  (since we have  $\mathbb{N} \models \gamma \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  by (\*), and  $\mathbb{N} \models \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  by  $\mathcal{Th}_\tau \not\vdash \gamma$ , then we have  $\mathbb{N} \models \gamma$ ). Proposition 2.2 immediately implies that  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \leftrightarrow \gamma$  holds for every  $\Pi_1$ -sentence  $\gamma$  that satisfies (\*). Essentially, the same proof can show that ( $\ddagger$ )  $\mathfrak{B} \vdash \text{Con}_\tau \leftrightarrow \gamma$  holds for all such sentences, and conversely, that if ( $\ddagger$ ) holds, then (\*) holds too.

**THEOREM 2.3.** *For every system  $\tau$  and  $\Pi_1$ -sentence  $\gamma$  we have*

$$\mathfrak{B} \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner) \quad \text{if and only if} \quad \mathfrak{B} \vdash \text{Con}_\tau \leftrightarrow \gamma. \quad \dashv$$

This theorem is a special case of the theorem(s) of D. de Jongh, G. Sambin, and C. Bernardi (see [6, Chapter 8]).

So, we can define *the* Gödel sentence  $\mathbb{G}_\tau$  of a system  $\tau$  to be any of the  $\Pi_1$ -sentences  $\gamma$  that satisfy  $\mathfrak{B} \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$ ; noting that  $\mathbb{G}_\tau$  is unique up to equivalence, even provably so in  $\mathfrak{B}$ .

**2.2. Rosser's proof.** The proof of Rosser constructs, by using the diagonal lemma, a  $\Pi_1$ -sentence  $\rho$ , for a given system  $\tau$ , such that

$$\mathfrak{B} \vdash \rho \leftrightarrow \forall x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \rightarrow \exists y < x \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)]$$

holds. Let us call any such  $\Pi_1$ -sentence  $\rho$ , a *Rosserian* sentence of the system  $\tau$ .

The classical proof of Rosser shows the independence of  $\rho$  from  $\mathcal{T}\mathcal{h}_\tau$ , if  $\mathcal{T}\mathcal{h}_\tau$  is (only) consistent (and not necessarily  $\Sigma_1$ -sound): if  $\mathcal{T}\mathcal{h}_\tau \vdash \rho$ , then  $\mathcal{T}\mathcal{h}_\tau \vdash \forall x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \rightarrow \exists y < x \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)]$ , and also, by  $(C_1)$ ,  $\mathcal{T}\mathcal{h}_\tau \vdash \text{prf}_\tau(\bar{n}, \ulcorner \rho \urcorner)$  for some  $n \in \mathbb{N}$ ; whence,  $\mathcal{T}\mathcal{h}_\tau \vdash \exists y < \bar{n} \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)$ , but the  $\Delta_0$ -sentence  $\exists y < \bar{n} \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)$  would be false by the consistency of  $\mathcal{T}\mathcal{h}_\tau$ , and so would be refutable in  $\mathfrak{B}$ , contradiction. Also, the assumption  $\mathcal{T}\mathcal{h}_\tau \vdash \neg \rho$ , would imply  $\mathcal{T}\mathcal{h}_\tau \vdash \exists x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \wedge \forall y < x \neg \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)]$  on the one hand, and  $\mathfrak{B} \vdash \text{prf}_\tau(\bar{m}, \ulcorner \neg \rho \urcorner)$  for some  $m \in \mathbb{N}$ , by  $(C_1)$ , on the other hand; thus,  $\mathcal{T}\mathcal{h}_\tau \vdash \exists x \leq \bar{m} \text{prf}_\tau(x, \ulcorner \rho \urcorner)$ , but by the consistency of  $\mathcal{T}\mathcal{h}_\tau$ , the  $\Delta_0$ -sentence  $\exists x \leq \bar{m} \text{prf}_\tau(x, \ulcorner \rho \urcorner)$  would be false, and so refutable in  $\mathfrak{B}$ , contradiction.

Rosser's theorem is not derivable from  $\mathbb{G}_2$  in the standard way (though, it does deliver  $\mathbb{G}_2$ ): For a consistent system  $\tau$ , let  $\mathfrak{g}(x) = \tau(x) \vee [x = \neg \text{Con}_\tau]$ ; then  $\mathcal{T}\mathcal{h}_\mathfrak{g}$  is consistent by  $\mathbb{G}_2$ , and  $\mathcal{T}\mathcal{h}_\mathfrak{g} \vdash \neg \text{Con}_\mathfrak{g}$ . So, for every Rosserian sentence  $\rho$  of  $\mathfrak{g}$ , we have  $\mathcal{T}\mathcal{h}_\mathfrak{g} \not\vdash \rho \rightarrow \text{Con}_\mathfrak{g}$ , since otherwise  $\mathcal{T}\mathcal{h}_\mathfrak{g} \vdash \neg \rho$  would hold, contradicting Rosser's theorem. Below we show a stronger result.<sup>1</sup>

**THEOREM 2.4.** *Let  $\tau$  be a system, and  $\rho$  be a Rosserian sentence of  $\tau$ . Then  $\mathfrak{B} \vdash \text{Con}_\tau \rightarrow \rho$  holds; and if  $\mathcal{T}\mathcal{h}_\tau$  is consistent, then  $\mathcal{T}\mathcal{h}_\tau \not\vdash \rho \rightarrow \text{Con}_\tau$ .*

**PROOF.** By  $\mathfrak{B} \vdash \neg \rho \leftrightarrow \exists x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \wedge \forall y < x \neg \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)]$ , we have that  $\mathfrak{B} \vdash \neg \rho \rightarrow \text{Pr}_\tau(\ulcorner \rho \urcorner)$ . Thus, Proposition 2.2(1) implies  $\mathfrak{B} \vdash \text{Con}_\tau \rightarrow \rho$ .

For showing the second part, we first show a formalized version of Rosser's theorem,  $\mathfrak{B} \vdash \text{Con}_\tau \rightarrow \neg \text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$ . Reason inside  $\mathfrak{B} + \text{Con}_\tau + \text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$ :

We have  $\neg \rho \leftrightarrow \exists x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \wedge \forall y < x \neg \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)]$ , and so  $\text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$  implies  $(\ddagger) \text{Pr}_\tau(\ulcorner \exists b[\text{prf}_\tau(b, \ulcorner \rho \urcorner) \wedge \forall y < b \neg \text{prf}_\tau(y, \ulcorner \neg \rho \urcorner)] \urcorner)$  by  $(C_1, C_2)$ . Also,  $\text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$  implies the existence of some  $a$  such that  $(\ddot{\ddagger}) \text{prf}_\tau(a, \ulcorner \neg \rho \urcorner)$  holds. Now, by  $(\ddagger)$ ,  $(\ddot{\ddagger})$ , and  $\text{Con}_\tau$ , we have  $(*) \text{Pr}_\tau(\ulcorner \exists b \leq a \text{prf}_\tau(b, \ulcorner \rho \urcorner) \urcorner)$ . On the other hand, by the assumption  $\text{Con}_\tau + \text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$ , we have  $\neg \text{Pr}_\tau(\ulcorner \rho \urcorner)$ , and so the  $\Delta_0$ -formula  $\forall x \leq a \neg \text{prf}_\tau(x, \ulcorner \rho \urcorner)$  is true; whence, by  $(C_3)$ , we have  $(**) \text{Pr}_\tau(\ulcorner \forall x \leq a \neg \text{prf}_\tau(x, \ulcorner \rho \urcorner) \urcorner)$ . Now,  $(*)$  and  $(**)$  contradict  $\text{Con}_\tau$ .

Assume now that  $\mathcal{T}\mathcal{h}_\tau$  is consistent; and assume (for the sake of a contradiction) that  $\mathcal{T}\mathcal{h}_\tau \vdash \rho \rightarrow \text{Con}_\tau$  holds. Then by  $\mathcal{T}\mathcal{h}_\tau \vdash \text{Con}_\tau \rightarrow \neg \text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$ , proved above, we have  $\mathcal{T}\mathcal{h}_\tau + \rho \vdash \neg \text{Pr}_\tau(\ulcorner \neg \rho \urcorner)$ . This contradicts

<sup>1</sup>Theorem 2.4 was first proved in [18, p. 16] where it is stated that this result was “implicit in the papers of Gödel [9] and Rosser [22]” (see also the end of [17]). It is also proved in [2] (see the Lemma on page 405) in which it is stated that although this result “was certainly known before,” the author “was unable to find a proof of it in the literature.”

$\mathbb{G}_2$ , unless  $\mathcal{Th}_\tau + \rho$  is inconsistent, or  $\mathcal{Th}_\tau \vdash \neg\rho$ ; and this contradicts Rosser's theorem, unless  $\mathcal{Th}_\tau$  is inconsistent. So, we showed that if  $\mathcal{Th}_\tau$  is consistent, then  $\mathcal{Th}_\tau \not\vdash \rho \rightarrow \text{Con}_\tau$ .  $\dashv$

It follows from Theorem 2.3 that all the  $\Pi_1$ -sentences  $\gamma$  for which we have that  $\mathfrak{B} \vdash \gamma \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \gamma \urcorner)$  are  $\mathfrak{B}$ -provably equivalent with (each other and with)  $\text{Con}_\tau$ . That all the Rosserian sentences of  $\tau$ , that is, all the  $\Pi_1$ -sentences  $\rho$  for which we have  $\mathfrak{B} \vdash \rho \leftrightarrow \forall x[\text{prf}_\tau(x, \ulcorner \rho \urcorner) \rightarrow \exists y < x \text{prf}_\tau(y, \ulcorner \neg\rho \urcorner)]$ , are equivalent with each other (in  $\mathcal{Th}_\tau$ ) was posed as an open question in [18, p. 16]. This was answered in [10] as follows: there are standard proof predicates<sup>2</sup> for which all the Rosserian sentences are equivalent (see [10, Theorem 6.2], and [30] for a correct proof); and there are standard proof predicates for which there are nonequivalent Rosserian sentences (see [10, Theorem 6.1]). For a given system  $\tau$ , let  $\overline{\text{prf}}_\tau(y, x)$  be a proof predicate of  $\tau$  all of whose Rosserian sentences are equivalent (see [4]). Then, (a variant of) the Rosser sentence  $\overline{\text{R}}_\tau$  of  $\tau$  can be defined to be any of the  $\Pi_1$ -sentences  $\rho$  that satisfy  $\mathfrak{B} \vdash \rho \leftrightarrow \forall x[\overline{\text{prf}}_\tau(x, \ulcorner \rho \urcorner) \rightarrow \exists y < x \overline{\text{prf}}_\tau(y, \ulcorner \neg\rho \urcorner)]$ .

So far, we have noticed that for every system  $\tau$ ,

$\mathfrak{B} \vdash \mathbb{G}_\tau \leftrightarrow \text{Con}_\tau$  and  $\mathfrak{B} \vdash \text{Con}_\tau \rightarrow \overline{\text{R}}_\tau$ ; but  $\mathcal{Th}_\tau \not\vdash \overline{\text{R}}_\tau \rightarrow \text{Con}_\tau$  for consistent  $\mathcal{Th}_\tau$ .

Let us note that  $\overline{\text{R}}_\tau$  is sensitive to implementation details modulo provable equivalence in the ambient theory, while  $\mathbb{G}_\tau$  is not so.

**2.3. Kleene's proof(s).** The proofs of Gödel and Rosser use the diagonal (aka self-referential) lemma for constructing the unprovable sentences. Kleene's (both first and second) proof can be considered diagonal-free in a sense, since it does not use this lemma directly, and avoids self-referentiality. However, as will be seen below, one can still argue that there could be some (at least, hidden) circularity in the proof. For stating Kleene's proofs, let us fix the notation.

**DEFINITION 2.5.** Let all the *unary* partial recursive functions be effectively (recursively) listed as  $\varphi_0, \varphi_1, \varphi_2, \dots$ . For  $m, n \in \mathbb{N}$ , if  $\varphi_m$  is defined at  $n$  (i.e.,  $\varphi_m(n)$  exists), then we write  $\varphi_m(n) \downarrow$ , and say that  $\varphi_m$  halts on  $n$ ; likewise,  $\varphi_m(n) \uparrow$  means that the function  $\varphi_m$  is not defined at  $n$ . Let us take  $\langle\langle \varphi_{\overline{m}}(\overline{n}) \uparrow \rangle\rangle$  to be the formula, in the language of arithmetic, that expresses  $\varphi_m(n) \uparrow$ .  $\dashv$

Let us note that  $\langle\langle \varphi_{\overline{m}}(\overline{n}) \uparrow \rangle\rangle$  is a  $\Pi_1$ -sentence. Kleene's proof [15] shows the existence of some  $k \in \mathbb{N}$ , for a given system  $\tau$ , such that the sentence  $\langle\langle \varphi_{\overline{k}}(\overline{k}) \uparrow \rangle\rangle$  is (true but) unprovable in  $\mathcal{Th}_\tau$ , if  $\mathcal{Th}_\tau$  is consistent. The nonconstructive version of the proof goes as follows: since the set  $\{n \in \mathbb{N} \mid \mathcal{Th}_\tau \vdash \langle\langle \varphi_{\overline{n}}(\overline{n}) \uparrow \rangle\rangle\}$  is RE, but the set  $\{n \in \mathbb{N} \mid \mathbb{N} \models \langle\langle \varphi_{\overline{n}}(\overline{n}) \uparrow \rangle\rangle\}$  is not, and the former is contained in the latter for consistent system  $\mathcal{Th}_\tau$ , so there must exist some  $k$  such that we have  $\mathbb{N} \models \langle\langle \varphi_{\overline{k}}(\overline{k}) \uparrow \rangle\rangle$  but  $\mathcal{Th}_\tau \not\vdash \langle\langle \varphi_{\overline{k}}(\overline{k}) \uparrow \rangle\rangle$ . Indeed, any RE and (effectively) undecidable set could be used for the proof; so, the existence of an RE and undecidable set implies Gödel's first incompleteness theorem,

<sup>2</sup>A proof predicate  $\lambda(y, x)$  is called *standard*, when its provability predicate, defined as  $\Lambda(x) = \exists y \lambda(y, x)$ , satisfies the derivability conditions.

by this argument of Kleene. It suffices for  $\varphi_k$  to have the property  $\varphi_k(n) \downarrow \iff \mathcal{T}\mathcal{H}_\tau \vdash \langle\langle \varphi_{\bar{n}}(\bar{n}) \uparrow \rangle\rangle$ , for every  $n \in \mathbb{N}$  (see [25, Theorem 2.2]); it is worth noting that there are indeed infinitely many such  $k$ 's. Now, the (true and) unprovable sentence of Kleene's first proof, for a given system  $\tau$ , can be constructed in a diagonal-free way as follows.

**DEFINITION 2.6.** For a system  $\tau$ , let  $t \in \mathbb{N}$  be an index (out of the infinitely many indexes) for the recursive function  $n \mapsto \mu z: \text{prf}_\tau(z, \ulcorner \langle\langle \varphi_{\bar{n}}(\bar{n}) \uparrow \rangle\rangle \urcorner)$ . Let  $\mathbb{K}_\tau = \langle\langle \varphi_{\bar{t}}(\bar{t}) \uparrow \rangle\rangle$  be *Kleene's (first) sentence* for system  $\tau$ .  $\dashv$

Let us note that by the definition of  $t$ ,  $\varphi_t(t) = \mu z: \text{prf}_\tau(z, \ulcorner \langle\langle \varphi_{\bar{t}}(\bar{t}) \uparrow \rangle\rangle \urcorner)$ , and so  $\varphi_t(t) \uparrow \iff \mathcal{T}\mathcal{H}_\tau \not\vdash \langle\langle \varphi_{\bar{t}}(\bar{t}) \uparrow \rangle\rangle$ , or equivalently  $\mathbb{K}_\tau \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \mathbb{K}_\tau \urcorner)$  (even provably in  $\mathfrak{B}$ ), which resembles Gödel's equivalence  $\mathbb{G}_\tau \leftrightarrow \neg \text{Pr}_\tau(\ulcorner \mathbb{G}_\tau \urcorner)$ . So, the equivalence of Kleene's first incompleteness witness with  $\mathbb{G}_2$  is a consequence of Theorem 2.3 (which also shows that  $\mathbb{K}_\tau$  is not sensitive to various implementation details):

**COROLLARY 2.7.** For every system  $\tau$  we have  $\mathfrak{B} \vdash \text{Con}_\tau \leftrightarrow \mathbb{K}_\tau$ .  $\dashv$

**2.4. Kleene's symmetric proof.** Kleene's first sentence is not independent from the system if the system is consistent (see [25, Theorem 2.3]); and so Kleene [16] gave another proof for the Gödel–Rosser theorem, which was called by him “a symmetric form” of Gödel's (incompleteness) theorem; see also [23, Theorem 14].

**DEFINITION 2.8.** Let us effectively list all the *binary* partial recursive functions as  $\phi_0, \phi_1, \phi_2, \dots$ . Let  $\phi_n(k, l) \downarrow_m$  mean that the binary partial recursive function  $\phi_n$  is defined at  $(k, l)$  and its value can be computed (in a fixed programming language) in  $m$  steps (or less). As before, its formalization in the language of arithmetic is denoted by  $\langle\langle \phi_{\bar{n}}(\bar{k}, \bar{l}) \downarrow_{\bar{m}} \rangle\rangle$ .  $\dashv$

Let us note that  $\langle\langle \phi_x(u, v) \downarrow_y \rangle\rangle$  can be written by a  $\Delta_0$ -formula of the free variables  $u, v, x$ , and  $y$ . For a system  $\tau$ , let  $f_\tau: \mathbb{N} \rightarrow \mathbb{N}$  be the  $\tau$ -proof search function,  $f_\tau(u) = \mu z: \overline{\text{prf}}_\tau(z, u)$ , whose algorithm is as follows:

input  $u$ , put  $i := 0$ , and run the sub-program  $\sharp_i$ ;  
 $\sharp_i$ : check if  $i$  is a  $\tau$ -proof (if for some  $j \leq i$ ,  $\overline{\text{prf}}_\tau(i, j)$  holds); if not, then put  $i := i + 1$  and repeat  $\sharp_i$ ; if yes, then compute what  $i$  proves (the above  $j$ ); if  $j \neq u$ , then put  $i := i + 1$  and repeat  $\sharp_i$ ; if  $j = u$ , then output  $i$  and halt.

**DEFINITION 2.9.** For numbers  $m, n \in \mathbb{N}$ , let  $\mathbb{J}_n^m$  be the following sentence:  $\forall x[\langle\langle \phi_{\bar{m}}(\bar{m}, \bar{n}) \downarrow_x \rangle\rangle \rightarrow \exists y < x \langle\langle \phi_{\bar{n}}(\bar{m}, \bar{n}) \downarrow_y \rangle\rangle]$ . For a given system  $\tau$ , let  $r$  and  $s$  be some indexes for the (binary recursive) functions  $(m, n) \mapsto f_\tau(\ulcorner \mathbb{J}_n^m \urcorner)$  and  $(m, n) \mapsto f_\tau(\ulcorner \neg \mathbb{J}_n^m \urcorner)$ , respectively. For a given system  $\tau$ , let  $\mathbb{K}'_\tau = \mathbb{J}_s^r$  be *Kleene's second sentence* for  $\tau$ .  $\dashv$

Thus, for a system  $\tau$ ,  $\phi_{\bar{r}}(m, n) = f_\tau(\ulcorner \mathbb{J}_n^m \urcorner)$  and  $\phi_{\bar{s}}(m, n) = f_\tau(\ulcorner \neg \mathbb{J}_n^m \urcorner)$ ; also  $\mathbb{K}'_\tau = \forall x[\langle\langle \phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_x \rangle\rangle \rightarrow \exists y < x \langle\langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_y \rangle\rangle]$ , which is a  $\Pi_1$ -sentence.

For a consistent system  $\tau$ , the independence of  $\mathbb{K}'_\tau$  from  $\mathcal{T}\mathcal{H}_\tau$  can be shown along the lines of Rosser's proof: If  $\mathcal{T}\mathcal{H}_\tau \vdash \mathbb{K}'_\tau$ , then  $\phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_m$  for some  $m \in \mathbb{N}$ ,



and so  $\mathcal{Th}_\tau \vdash \exists y < \bar{m} \langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_y \rangle \rangle$ . But by the consistency of  $\mathcal{Th}_\tau$  we have  $\mathcal{Th}_\tau \not\vdash \neg \overline{\mathbb{K}'_\tau}$  and so  $\phi_{\bar{s}}(\bar{r}, \bar{s}) \uparrow$ ; whence, the  $\Delta_0$ -sentence  $\exists y < \bar{m} \langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_y \rangle \rangle$  is false and so should be refutable in  $\mathcal{Th}_\tau$ , contradiction. Also, if  $\mathcal{Th}_\tau \vdash \neg \overline{\mathbb{K}'_\tau}$ , then  $\phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_n$  for some  $n \in \mathbb{N}$ , and so  $\mathcal{Th}_\tau \vdash \exists x \leq \bar{n} \langle \langle \phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_x \rangle \rangle$ . But the  $\Delta_0$ -sentence  $\exists x \leq \bar{n} \langle \langle \phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_x \rangle \rangle$  is false by the consistency of  $\mathcal{Th}_\tau$  (which implies  $\mathcal{Th}_\tau \not\vdash \overline{\mathbb{K}'_\tau}$ , thus  $\phi_{\bar{r}}(\bar{r}, \bar{s}) \uparrow$ ) and so it should be refutable in  $\mathcal{Th}_\tau$ ; contradiction again.

Also, very similarly to the proof of Theorem 2.4, it can be shown that for every system  $\tau$  we have  $\mathcal{Th}_\tau \vdash \text{Con}_\tau \rightarrow \overline{\mathbb{K}'_\tau}$ ; and  $\mathcal{Th}_\tau \not\vdash \overline{\mathbb{K}'_\tau} \rightarrow \text{Con}_\tau$  if  $\mathcal{Th}_\tau$  is consistent. We show, more generally, that  $\overline{\mathbb{K}'_\tau}$  is a Rosserian sentence of  $\tau$ , and so this follows directly from Theorem 2.4.

If  $x = f_\tau(u)$  for a system  $\tau$ , and some  $u$ , then let  $\hat{x}$  be the least  $s$  such that  $f_\tau(u) \downarrow_s$  holds. We note that, for every  $u, v, x, y$ , if  $f_\tau(u) = x < y = f_\tau(v)$ , then  $\hat{x} < \hat{y}$ : this is because the algorithm of  $f_\tau$ , on the input  $v$ , has already checked  $x$ , before halting at the step  $\hat{y}$ , to see if it is a  $\tau$ -proof, and if (yes, then) it is a  $\tau$ -proof of  $v$ . The algorithm has noticed that  $x$  is a  $\tau$ -proof of some  $u \neq v$ , and then has gone to the next step (to check  $x+1$  and so on). So, the number of steps that  $f_\tau$  needs to calculate  $y$  (on the input  $v$ ) is greater than the number of steps that it needs to calculate  $x$  (on the input  $u$ ); thus  $\hat{x} < \hat{y}$ .

**THEOREM 2.10.** *For every system  $\tau$ , the sentence  $\overline{\mathbb{K}'_\tau}$  of  $\tau$  is a Rosserian sentence:  $\mathfrak{B} \vdash \overline{\mathbb{K}'_\tau} \leftrightarrow \forall x [\overline{\text{prf}}_\tau(x, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner) \rightarrow \exists y < x \overline{\text{prf}}_\tau(y, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner)]$ .*

**PROOF.** Reason inside  $\mathfrak{B}$ :

Suppose that  $\overline{\mathbb{K}'_\tau}$  holds, and for some  $a$ , we have  $\overline{\text{prf}}_\tau(a, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner)$ . Let  $\alpha$  be the minimum  $x$  with  $\overline{\text{prf}}_\tau(x, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner)$ ; then  $a \geq \alpha$ , and  $\langle \langle \phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_{\hat{\alpha}} \rangle \rangle$  holds. Thus, by  $\overline{\mathbb{K}'_\tau}$ , there exists some  $\beta < \hat{\alpha}$  such that  $\langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_\beta \rangle \rangle$  holds. So, for  $b = \phi_{\bar{s}}(\bar{r}, \bar{s})$  we have  $\overline{\text{prf}}_\tau(b, \ulcorner \neg \overline{\mathbb{K}'_\tau} \urcorner)$ . We show that  $b < a$ ; if on the contrary  $b \geq a$ , then we would have  $b \geq \alpha$ , and so  $\beta \geq \hat{b} \geq \hat{\alpha}$ ; contradiction with  $\beta < \hat{\alpha}$ . Thus, if  $\overline{\mathbb{K}'_\tau}$  holds, then for all  $a$  with  $\overline{\text{prf}}_\tau(a, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner)$  there exists some  $b < a$  with  $\overline{\text{prf}}_\tau(b, \ulcorner \neg \overline{\mathbb{K}'_\tau} \urcorner)$ .

Now, suppose that  $\neg \overline{\mathbb{K}'_\tau}$  holds. So, there exists some  $p$  such that  $\langle \langle \phi_{\bar{r}}(\bar{r}, \bar{s}) \downarrow_p \rangle \rangle$  holds, and for no  $q < p$  can  $\langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_q \rangle \rangle$  hold. Then, let  $a = \phi_{\bar{r}}(\bar{r}, \bar{s})$ ; so,  $p \geq \hat{a}$ , and we have  $\overline{\text{prf}}_\tau(a, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner)$ . We show that for no  $b < a$  can  $\overline{\text{prf}}_\tau(b, \ulcorner \neg \overline{\mathbb{K}'_\tau} \urcorner)$  hold. Assume, on the contrary, that for some  $b < a$  we have  $\overline{\text{prf}}_\tau(b, \ulcorner \neg \overline{\mathbb{K}'_\tau} \urcorner)$ . Then,  $\langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow \rangle \rangle$ ; let  $d = \phi_{\bar{s}}(\bar{r}, \bar{s})$ . So, we have  $\langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_{\hat{d}} \rangle \rangle$ , and we also have  $d \leq b < a$ . Thus,  $d < a$  holds, and so we have  $\hat{d} < \hat{a} \leq p$ ; this is a contradiction with the property of  $p$  (that  $\forall q < p \neg \langle \langle \phi_{\bar{s}}(\bar{r}, \bar{s}) \downarrow_q \rangle \rangle$ ). Thus, if  $\neg \overline{\mathbb{K}'_\tau}$  holds, then there exists some  $a$  such that  $\overline{\text{prf}}_\tau(a, \ulcorner \overline{\mathbb{K}'_\tau} \urcorner) \wedge \forall b < a \neg \overline{\text{prf}}_\tau(b, \ulcorner \neg \overline{\mathbb{K}'_\tau} \urcorner)$ .  $\dashv$

Therefore, under some certain considerations (and the choice of  $\overline{\text{prf}}_\tau(y, x)$  for proof predicate), Kleene's second  $\Pi_1$ -incompleteness witness is equivalent to Rosser's  $\Pi_1$ -incompleteness witness.



**COROLLARY 2.11.** *For every system  $\tau$ , we have  $\mathfrak{B} \vdash \overline{\mathbb{K}'_\tau} \leftrightarrow \overline{\mathbb{R}_\tau}$ . So, we have  $\mathfrak{B} \vdash \text{Con}_\tau \rightarrow \overline{\mathbb{K}'_\tau}$ ; and  $\mathcal{T}\mathfrak{h}_\tau \not\vdash \overline{\mathbb{K}'_\tau} \rightarrow \text{Con}_\tau$  if  $\mathcal{T}\mathfrak{h}_\tau$  is consistent*  $\dashv$

Thus far, we have shown that Gödel's first theorem is equivalent with  $\mathfrak{G}_2$  and with Kleene's first theorem; (a variant of) Rosser's theorem is equivalent with (a variant of) Kleene's second theorem, and it does deliver  $\mathfrak{G}_2$ , but is not derivable from  $\mathfrak{G}_2$ . In picture:

$$(\text{Gödel}_2 \cong \text{Kleene} \cong \text{Gödel}_1) \overset{\sim}{\approx} (\overline{\text{Rosser}} \cong \overline{\text{Kleene}}).$$

**§3. Proofs of Chaitin and Boolos.** We say that the  $\Pi_1$ -incompleteness witness  $\mathcal{F}$  is *constructive* (or *effective*), when  $\mathcal{F}_\tau$  can be effectively (computably) constructed from  $\tau$ . We say that the  $\Pi_1$ -incompleteness witness  $\mathcal{F}$  has *the Rosser property*, when  $\mathcal{F}_\tau$  is independent from  $\mathcal{T}\mathfrak{h}_\tau$ , for every consistent system  $\tau$ . Constructivity and the Rosser property of the  $\Pi_1$ -incompleteness theorems of Gödel, Rosser, Kleene, Chaitin, and Boolos were studied in [25] (see the table on its p. 579). There, it was shown that the  $\Pi_1$ -incompleteness theorems of Gödel, Rosser, and (both theorems of) Kleene are constructive, and the theorems of Chaitin and Boolos are not; also none of the theorems of Gödel, Kleene's first, or Boolos have the Rosser property, while the theorems of Rosser, Kleene's second, and (a variant of) Chaitin do have the Rosser property.

**3.1. Chaitin's proof.** There are several variants of Chaitin's theorem. Here, we consider one of the simple ones.

**DEFINITION 3.1.** The *Kolmogorov–Chaitin Complexity* function is defined to be the mapping  $\mathcal{K}(w) = \mu e : [\varphi_e(0) = w]$ , for  $w \in \mathbb{N}$ . Let  $\langle\langle \mathcal{K}(x) > y \rangle\rangle$  denote the  $\Pi_1$ -formula  $\forall v \leq y [\langle\langle \varphi_v(0) \downarrow \rangle\rangle \rightarrow \langle\langle \varphi_v(0) \neq x \rangle\rangle]$ , in the language of arithmetic, with the free variables  $x$  and  $y$ , that expresses  $\mathcal{K}(x) > y$ .  $\dashv$

Chaitin's theorem [7] shows the existence of some  $c_\tau \in \mathbb{N}$ , for a given system  $\tau$ , such that for every  $w, e \in \mathbb{N}$  with  $e \geq c_\tau$  we have  $\mathcal{T}\mathfrak{h}_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) > \overline{e} \rangle\rangle$  if  $\mathcal{T}\mathfrak{h}_\tau$  is consistent. It suffices to put for a system  $\tau$ ,

$$\varphi_{c_\tau}(x) = \pi_1[\mu z : \text{prf}_\tau(\pi_2[z], \ulcorner \langle\langle \mathcal{K}(\pi_1[z]) > x + c_\tau \rangle\rangle \urcorner)],$$

which is possible by Kleene's recursion (or fixed point) theorem; here,  $\pi_1, \pi_2$  are the projection functions of a fixed enumeration of ordered pairs (i.e., if we have a bijective mapping  $(a, b) \mapsto \langle a, b \rangle$  between ordered pairs of numbers and numbers, then we have  $\pi_1(\langle a, b \rangle) = a$  and  $\pi_2(\langle a, b \rangle) = b$ ).

Chaitin's proof goes as follows: If  $\mathcal{T}\mathfrak{h}_\tau$  is consistent, and  $\mathcal{T}\mathfrak{h}_\tau \vdash \langle\langle \mathcal{K}(\overline{w}) > \overline{e} \rangle\rangle$  for some  $w, e \in \mathbb{N}$  with  $e \geq c_\tau$ , then let  $z = \langle u, p \rangle$  be the minimum ordered pair such that  $\text{prf}_\tau(p, \ulcorner \langle\langle \mathcal{K}(\overline{u}) > \overline{c}_\tau \rangle\rangle \urcorner)$  holds. Then,  $\mathcal{T}\mathfrak{h}_\tau \vdash \langle\langle \mathcal{K}(\overline{u}) > \overline{c}_\tau \rangle\rangle$  and also  $\varphi_{c_\tau}(0) = u$  holds; thus,  $\mathcal{K}(u) \leq c_\tau$ . Whence, the  $\Sigma_1$ -sentence  $\langle\langle \mathcal{K}(\overline{u}) \leq \overline{c}_\tau \rangle\rangle$  is true, and so should be provable in  $\mathcal{T}\mathfrak{h}_\tau$ ; a contradiction with the consistency of  $\mathcal{T}\mathfrak{h}_\tau$  (cf. [25, Theorem 3.3]). Let us note that there are cofinitely many  $w$ 's such that  $\mathbb{N} \models \langle\langle \mathcal{K}(\overline{w}) > \overline{c}_\tau \rangle\rangle$ ; so it is tempting to fix one of those  $w$ 's as  $w_\tau$ , and define *the Chaitin sentence of system  $\tau$*  as

$\mathfrak{C}_\tau = \langle\langle \mathcal{K}(\overline{w}_\tau) \rangle \overline{c}_\tau \rangle$ . For technical reasons, we will postpone this till the next subsection (Definition 3.12); for the moment, we would like our results on Chaitin's proof to be as general as possible. We now show a Rosserian form of Chaitin's theorem (see [25, Theorem 3.9]).

**THEOREM 3.2.** *For every consistent system  $\tau$  and for every  $e \geq c_\tau$ , there are cofinitely many  $w$ 's such that  $\langle\langle \mathcal{K}(\overline{w}) \rangle \overline{e} \rangle$  is independent from  $\mathcal{F}h_\tau$ .*

**PROOF.** Fix an  $e \geq c_\tau$ ; we showed that for no  $w$  can  $\mathcal{F}h_\tau \vdash \langle\langle \mathcal{K}(\overline{w}) \rangle \overline{e} \rangle$  hold. We now show that  $\mathcal{F}h_\tau \vdash \langle\langle \mathcal{K}(\overline{w}) \leq \overline{e} \rangle \rangle$  can hold for at most  $(e+1)$ -many  $w$ 's. First let us note that the  $\Delta_0$ -sentence  $\forall \{x_i < \overline{n}\}_{i \leq n} (\bigvee_{i < j \leq n} x_i = x_j)$ , which is a version of the Pigeonhole Principle, is true for each  $n \in \mathbb{N}$ , and thus is provable in  $\mathfrak{B}$ . Reason inside  $\mathcal{F}h_\tau$ :

If for some distinct  $\overline{w}_0, \overline{w}_1, \dots, \overline{w}_{e+1}$  we have  $\langle\langle \mathcal{K}(\overline{w}_i) \leq \overline{e} \rangle \rangle$ , then for each  $i \leq e+1$  we have  $\langle\langle \varphi_{\overline{z}_i}(0) = \overline{w}_i \rangle \rangle$  for some  $z_i \leq e$ . By the Pigeonhole Principle there should exist some  $j < k \leq e+1$  such that  $\overline{z}_j = \overline{z}_k$ . Thus, we should have  $\overline{w}_j = \overline{w}_k$ , contradicting the distinctness of  $\overline{w}_i$ 's.

So, there are cofinitely many  $w$ 's for which we have both  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) \rangle \overline{e} \rangle$  and  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) \leq \overline{e} \rangle \rangle$ .  $\dashv$

We note that if  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) \leq \overline{e} \rangle \rangle$ , then  $\mathbb{N} \models \langle\langle \mathcal{K}(\overline{w}) \rangle \overline{e} \rangle$ . So, it would be more tempting to fix as  $w_\tau$  one of the cofinitely many  $w$ 's with  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) \leq \overline{e} \rangle \rangle$ ; since then the Chaitin sentence  $\mathfrak{C}_\tau$  of  $\tau$  will be independent from  $\mathcal{F}h_\tau$ . In the following theorem, we show that Chaitin's theorem can deliver no constructive  $\Pi_1$ -incompleteness witness. Let us note that adding a true  $\Pi_1$ -sentence to a  $\Sigma_1$ -sound theory results in a  $\Sigma_1$ -sound theory; and the union of a chain of  $\Sigma_1$ -sound theories is also a  $\Sigma_1$ -sound theory.

**THEOREM 3.3.** *For every constructive  $\Pi_1$ -incompleteness witness  $\mathcal{F}$ , and for every  $\Sigma_1$ -sound system  $\tau$  there exists some  $\Sigma_1$ -sound super-system  $\varrho$  of  $\tau$  such that we have  $\mathcal{F}h_\varrho \not\vdash \mathcal{F}_\varrho \rightarrow \langle\langle \mathcal{K}(\overline{w}) \rangle \overline{c}_\varrho \rangle \rangle$  for all  $w$ .*

**PROOF.** Let  $\tau^0 = \tau$ , and inductively,  $\tau^{n+1}(x) = \tau^n(x) \vee [x = \mathcal{F}_{\tau^n}]$ . We note that all the  $\{\mathcal{F}h_{\tau^n}\}_n$ 's and also  $\mathcal{F}h_{\tau^\infty} = \bigcup_n \mathcal{F}h_{\tau^n}$  are  $\Sigma_1$ -sound. We show that there is some  $n$  such that  $\mathcal{F}h_{\tau^n} \not\vdash \mathcal{F}_{\tau^n} \rightarrow \langle\langle \mathcal{K}(\overline{w}) \rangle \overline{c}_{\tau^n} \rangle \rangle$  for all  $w$ ; thus proving the theorem. Assume, for the sake of a contradiction, that for every  $n$  there exists some  $w_n$  such that  $\mathcal{F}h_{\tau^n} \vdash \mathcal{F}_{\tau^n} \rightarrow \langle\langle \mathcal{K}(\overline{w}_n) \rangle \overline{c}_{\tau^n} \rangle \rangle$ . There are two cases:

- (1) For some  $n$ ,  $c_{\tau^n} \geq c_{\tau^{n+1}}$ . Then, from  $\mathcal{F}h_{\tau^{n+1}} \vdash \langle\langle \mathcal{K}(\overline{w}_n) \rangle \overline{c}_{\tau^n} \rangle \rangle$  we have  $\mathcal{F}h_{\tau^{n+1}} \vdash \langle\langle \mathcal{K}(\overline{w}_n) \rangle \overline{c}_{\tau^{n+1}} \rangle \rangle$ , contradicting Chaitin's theorem for  $\tau^{n+1}$ .
- (2) For every  $n$ ,  $c_{\tau^n} < c_{\tau^{n+1}}$ . By the constructivity of  $\mathcal{F}$  the system  $\mathcal{F}h_{\tau^\infty}$  is RE. So, there should exist some  $m$  such that  $c_{\tau^\infty} \leq c_{\tau^m}$ . Thus, from  $\mathcal{F}h_{\tau^{m+1}} \vdash \langle\langle \mathcal{K}(\overline{w}_m) \rangle \overline{c}_{\tau^m} \rangle \rangle$ , we get  $\mathcal{F}h_{\tau^\infty} \vdash \langle\langle \mathcal{K}(\overline{w}_m) \rangle \overline{c}_{\tau^\infty} \rangle \rangle$ , which contradicts Chaitin's theorem for the consistent  $\tau^\infty$ .  $\dashv$

A mapping  $\tau \mapsto w_\tau$ , which assigns  $w_\tau \in \mathbb{N}$  to a given system  $\tau$ , is called a *Chaitin mapping* when  $\mathbb{N} \models \langle\langle \mathcal{K}(\overline{w}_\tau) \rangle \overline{c}_\tau \rangle \rangle$  holds for every consistent system  $\tau$ . Let us call it a *Rosser–Chaitin mapping* when  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}_\tau) \leq \overline{c}_\tau \rangle \rangle$  holds for every consistent system  $\tau$ .

COROLLARY 3.4. *No Chaitin mapping can be constructive ([25, Theorem 3.5]).*

*Chaitin's  $\Pi_1$ -incompleteness theorem can deliver none of the  $\Pi_1$ -incompleteness theorems of Gödel (first and second), Kleene (first and second), or Rosser.*  $\dashv$

We will see below that Chaitin's  $\Pi_1$ -incompleteness witness does deliver a variant of the  $\Pi_1$ -incompleteness witness of Boolos; thus, Boolos' theorem is not constructive (see [25, Theorem 4.5]). As the last result on Chaitin's theorem, we show that essentially no  $\Pi_1$ -incompleteness witness may deliver Chaitin's  $\Pi_1$ -incompleteness witness in the standard way:

COROLLARY 3.5. *For every consistent system  $\tau$  and every sentence  $\psi$  with  $\mathcal{T}\mathcal{H}_\tau \not\vdash \psi$ , there are cofinitely many  $w$ 's with  $\mathcal{T}\mathcal{H}_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}) > \overline{c}_\tau \rangle\rangle \rightarrow \psi$ .*

PROOF. Since  $\mathcal{T}\mathcal{H}_\tau + \neg\psi$  is consistent, then by (the proof of) Theorem 3.2 for cofinitely many  $w$ 's, we have  $\mathcal{T}\mathcal{H}_\tau + \neg\psi \not\vdash \langle\langle \mathcal{K}(\overline{w}) \leq \overline{c}_\tau \rangle\rangle$ .  $\dashv$

As a consequence, for every given Rosser–Chaitin mapping  $\tau \mapsto v_\tau$ , there exists another Rosser–Chaitin mapping  $\tau \mapsto w_\tau$  such that for every consistent system  $\tau$  we have  $\mathcal{T}\mathcal{H}_\tau \not\vdash \langle\langle \mathcal{K}(\overline{w}_\tau) > \overline{c}_\tau \rangle\rangle \rightarrow \langle\langle \mathcal{K}(\overline{v}_\tau) > \overline{c}_\tau \rangle\rangle$ .

**3.2. Boolos' proof.** Finally, we consider the theorem of Boolos, for which we make the following convention.

CONVENTION. *All the variables are  $\vartheta, \vartheta', \vartheta'', \vartheta''', \dots$  whose lengths are 1, 2, 3, 4, ..., respectively.*  $\dashv$

The length of an expression is the number of symbols in it. By the above convention (3.6), for any natural number  $n \in \mathbb{N}$ , there are at most finitely many formulas with length  $n$ ; without this convention, for variables  $x, y, z, \dots$ , all of the formulas  $x=0, y=0, z=0, \dots$  would be length of 3.

DEFINITION 3.7. A number  $n \in \mathbb{N}$  is *definable* in the theory  $T$ , by the formula  $\theta(\vartheta)$  with the only free variable  $\vartheta$ , when  $T \vdash \forall \vartheta [\theta(\vartheta) \leftrightarrow \vartheta = \overline{n}]$  holds. Let  $\delta(\ulcorner \theta \urcorner, n)$  denote the Gödel code  $\ulcorner \forall \vartheta [\theta(\vartheta) \leftrightarrow \vartheta = \overline{n}] \urcorner$ .  $\dashv$

Suppose that the formula  $F_1(x)$  in the language of arithmetic states that “ $x$  is (the Gödel code of) a formula which has  $\vartheta$  as its only free variable,” and the formula  $L^{<y}(x)$  states that “the formula (with Gödel code)  $x$  has length less than  $y$ .” Indeed, there are such  $\Sigma_1$ -formulas in the language of arithmetic, whose existence can be shown by the techniques of Gödel's arithmetization.

DEFINITION 3.8. For a system  $\tau$ , let  $\mathcal{D}_\tau^{<y}(x) = \exists \xi [F_1(\xi) \wedge L^{<y}(\xi) \wedge \text{Pr}_\tau(\delta(\xi, x))]$  be a formula, in the language of arithmetic, stating that “ $x$  is definable in theory  $\mathcal{T}\mathcal{H}_\tau$  by a formula with length  $< y$ .” Let  $\mathcal{B}_\tau^{<y}(x) = \neg \mathcal{D}_\tau^{<y}(x) \wedge \forall z < x \mathcal{D}_\tau^{<y}(z)$  be the formula which states that “ $x$  is the least number not definable (in  $\mathcal{T}\mathcal{H}_\tau$ ) by any formula with length less than  $y$ .” Let  $\ell_\tau$  be the length of  $\mathcal{B}_\tau^{<\vartheta''}(\vartheta)$ .  $\dashv$

Let  $k$  be any natural number nonsmaller than 10 (Boolos [5] originally takes it to be 10). For a system  $\tau$ , let  $k_\tau = \bar{k} \cdot \bar{\ell}_\tau$  be a term representing  $k\ell_\tau$ . Let  $\beta_\tau(\vartheta) = \exists \vartheta' [\vartheta' = k_\tau \wedge \mathcal{B}_\tau^{<\vartheta'}(\vartheta)]$  be the formula stating that “ $\vartheta$  is the least number that is not definable by any formula with length less than  $k\ell_\tau$ .” It can be shown that the length of  $\beta_\tau$  is less than  $k\ell_\tau$  (cf. the proof of Theorem 4.3 in [25]). Let  $b_\tau$  be the least number (if any) that is not definable (in  $\mathcal{Th}_\tau$ ) by any formula with length less than  $k\ell_\tau$ . Boolos’ original theorem is the following:

**THEOREM 3.9.** *If  $\tau$  is a consistent system, then  $\beta_\tau(\bar{b}_\tau)$  is a true sentence that is not provable in  $\mathcal{Th}_\tau$ .*

**PROOF.** The truth of  $\beta_\tau(\bar{b}_\tau)$  follows from the definition of  $b_\tau$ . Assume, for the sake of a contradiction, that  $\mathcal{Th}_\tau \vdash \beta_\tau(\bar{b}_\tau)$ ; thus,  $\mathcal{Th}_\tau \vdash \mathcal{B}_\tau^{<k_\tau}(\bar{b}_\tau)$ . Then,  $\mathfrak{B} \vdash \forall x [\mathcal{B}_\tau^{<k_\tau}(x) \leftrightarrow x = \bar{b}_\tau]$ ,<sup>3</sup> and so  $b_\tau$  is definable (in  $\mathcal{Th}_\tau$ ) by the formula  $\mathcal{B}_\tau^{<k_\tau}(\vartheta)$  whose length is less than  $k\ell_\tau$ ; thus,  $\mathcal{Th}_\tau$  is inconsistent.  $\dashv$

Let us note that  $\beta_\tau(\bar{b}_\tau)$  is not a  $\Pi_1$ -sentence; though we have

$$\mathfrak{B} \vdash \beta_\tau(\bar{b}_\tau) \equiv \mathcal{B}_\tau^{<k_\tau}(\bar{b}_\tau) \equiv \neg \mathcal{D}_\tau^{<k_\tau}(\bar{b}_\tau) \wedge \forall z < \bar{b}_\tau \mathcal{D}_\tau^{<k_\tau}(z) \equiv \neg \mathcal{D}_\tau^{<k_\tau}(\bar{b}_\tau),$$

because  $\forall z < \bar{b}_\tau \mathcal{D}_\tau^{<k_\tau}(z)$  is a true  $\Sigma_1$ -sentence, for a consistent system  $\tau$ , and so it is  $\mathfrak{B}$ -provable. Whence, the essence of Boolos’ theorem is the truth and  $\mathcal{Th}_\tau$ -unprovability of the  $\Pi_1$ -sentence  $\neg \mathcal{D}_\tau^{<k_\tau}(\bar{b}_\tau)$ . Indeed,  $\mathbb{G}_2$  can deliver this, and much more:

**THEOREM 3.10.** *For every consistent system  $\tau$  and every numbers  $m, n \in \mathbb{N}$  with  $m > 3$  we have  $\mathcal{Th}_\tau \not\vdash \neg \mathcal{D}_\tau^{<m}(\bar{n})$ .*

**PROOF.** If  $\mathcal{Th}_\tau \vdash \neg \mathcal{D}_\tau^{<m}(\bar{n})$ , then  $\mathcal{Th}_\tau \vdash \forall \xi [F_1(\xi) \wedge L^{<m}(\xi) \rightarrow \neg \text{Pr}_\tau(\delta(\xi, \bar{n}))]$ . There is some formula  $\zeta(\vartheta)$  with the only free variable  $\vartheta$ , such as  $\vartheta = 0$ , whose length is less than  $m$ ; thus,  $\mathcal{Th}_\tau \vdash \neg \text{Pr}_\tau(\delta(\ulcorner \zeta \urcorner, \bar{n}))$ . So, by Lemma 2.1(1), we should have  $\mathcal{Th}_\tau \vdash \text{Con}_\tau$ , which contradicts  $\mathbb{G}_2$ .  $\dashv$

Chaitin’s theorem can deliver Boolos’ theorem too, even in a more general form:

**THEOREM 3.11.** *For every consistent system  $\tau$  there is a number  $\xi_\tau$  such that for every system  $\varrho$  and for every  $m, n$  with  $m \geq \xi_\tau$  we have  $\mathcal{Th}_\tau \not\vdash \neg \mathcal{D}_\varrho^{<m}(\bar{n})$ .*

**PROOF.** Let  $\xi_\tau$  be a number that is greater than the lengths of all the formulas  $\langle\langle \varphi_{\bar{e}}(0) = \vartheta \rangle\rangle$  with  $e \leq c_\tau$ . Fix  $m, n \in \mathbb{N}$  with  $m \geq \xi_\tau$ . We first show (even inside  $\mathfrak{B}$ ) that (for an arbitrary system  $\varrho$ ) if  $\neg \mathcal{D}_\varrho^{<m}(n)$ , then  $\mathcal{K}(n) > c_\tau$ : If, on the contrary, we had  $\mathcal{K}(n) \leq c_\tau$ , then for some  $e \leq c_\tau$  we would have  $\varphi_e(0) = n$ . Since  $\mathfrak{B}$  can strongly represent all the recursive functions, then  $n$  would be definable in  $\mathfrak{B}$  (and so in  $\mathcal{Th}_\varrho$ ) by the formula  $\langle\langle \varphi_{\bar{e}}(0) = \vartheta \rangle\rangle$  whose length is less than  $\xi_\tau \leq m$ . Thus,  $\mathcal{D}_\varrho^{<\xi_\tau}(n)$  and so  $\mathcal{D}_\varrho^{<m}(n)$  would hold;

<sup>3</sup>Here, we use the fact that  $\mathfrak{B} \vdash \forall x (x < \bar{n} \vee x = \bar{n} \vee \bar{n} < x)$  for each  $n \in \mathbb{N}$  (see the proof of Theorem 4.3 in [25]).

a contradiction. Therefore,  $\mathcal{F}h_\tau \not\vdash \neg D_\varrho^{<\bar{m}}(\bar{n})$  follows from Chaitin's theorem that  $\mathcal{F}h_\tau \not\vdash \langle\langle \mathcal{K}(\bar{n}) > \bar{c}_\tau \rangle\rangle$ .  $\dashv$

Now we can define the Chaitin and the Boolos sentence(s) of a system.

**DEFINITION 3.12.** For a system  $\tau$ , if  $\mathcal{F}h_\tau$  is not consistent, then let  $w_\tau = 0$ . If  $\mathcal{F}h_\tau$  is consistent, then let  $w_\tau$  be one of the cofinitely many  $w$ 's such that

- (i)  $\mathcal{F}h_\tau + \neg \bar{\mathbb{R}}_\tau \not\vdash \langle\langle \mathcal{K}(\bar{w}) \leq \bar{c}_\tau \rangle\rangle$ , and
- (ii)  $\mathbb{N} \models \neg D_\tau^{<\kappa_\tau}(\bar{w})$ ;

where,  $\kappa_\tau$  is the least number  $k$  such that  $k \geq 10 \cdot \ell_\tau$ , and also  $k$  is greater than the lengths of all the formulas  $\langle\langle \varphi_e(0) = \vartheta \rangle\rangle$  with  $e \leq c_\tau$ .

For a consistent system  $\tau$ , let  $\mathbb{C}_\tau = \langle\langle \mathcal{K}(\bar{w}_\tau) > \bar{c}_\tau \rangle\rangle$  be the Chaitin sentence of  $\tau$ , and let  $\bar{\mathbb{B}}_\tau = \neg D_\tau^{<\kappa_\tau}(\bar{w}_\tau)$  be (a variant of) the Boolos sentence of  $\tau$ .  $\dashv$

Let us note that for every consistent system  $\tau$ , the  $\Pi_1$ -sentences  $\mathbb{C}_\tau$  and  $\bar{\mathbb{B}}_\tau$  are both true and unprovable in  $\mathcal{F}h_\tau$ .

- COROLLARY 3.13.** (1) For every consistent system  $\tau$ ,  $\mathcal{F}h_\tau \vdash \bar{\mathbb{B}}_\tau \rightarrow \mathbb{C}_\tau$ ;  
 (2) For every consistent system  $\tau$ ,  $\mathcal{F}h_\tau \vdash \bar{\mathbb{B}}_\tau \rightarrow \text{Con}_\tau$ ;  
 (3) Boolos'  $\Pi_1$ -incompleteness theorem is not constructive.

**PROOF.** (1) The deduction  $\mathfrak{B} \vdash \neg D_\tau^{<\kappa_\tau}(\bar{w}_\tau) \rightarrow \langle\langle \mathcal{K}(\bar{w}_\tau) > \bar{c}_\tau \rangle\rangle$  was shown in the proof of Theorem 3.11. (2) Follows from the proof of Theorem 3.10. (3) Follows from Theorem 3.3 and the item (1) above.  $\dashv$

**COROLLARY 3.14.** For every consistent system  $\tau$ ,

- (1)  $\mathcal{F}h_\tau \not\vdash \mathbb{C}_\tau \rightarrow \bar{\mathbb{R}}_\tau$ ;
- (2)  $\mathcal{F}h_\tau \not\vdash \mathbb{C}_\tau \rightarrow \text{Con}_\tau$ ;
- (3)  $\mathcal{F}h_\tau \not\vdash \mathbb{C}_\tau \rightarrow \bar{\mathbb{B}}_\tau$ ;
- (4)  $\mathcal{F}h_\tau \not\vdash \bar{\mathbb{R}}_\tau \rightarrow \bar{\mathbb{B}}_\tau$ .

**PROOF.** (1) By Definition 3.12. (2) By Theorem 2.4 and (1) above. (3) By Corollary 3.13(2) and (2) above. (4) By Theorem 2.4 and Corollary 3.13(2).  $\dashv$

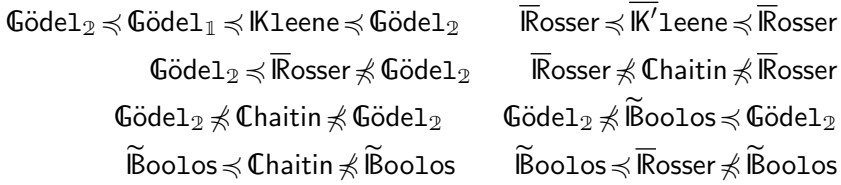
**COROLLARY 3.15.** For every  $\Sigma_1$ -sound system  $\tau$  there exists some  $\Sigma_1$ -sound super-system  $\varrho$  of  $\tau$  such that

- (1)  $\mathcal{F}h_\varrho \not\vdash \text{Con}_\varrho \rightarrow \mathbb{C}_\varrho$ ;
- (2)  $\mathcal{F}h_\varrho \not\vdash \bar{\mathbb{R}}_\varrho \rightarrow \mathbb{C}_\varrho$ ;
- (3)  $\mathcal{F}h_\varrho \not\vdash \text{Con}_\varrho \rightarrow \bar{\mathbb{B}}_\varrho$ .

**PROOF.** (1) By Theorem 3.3. (2) By Theorem 2.4 and (1) above. (3) By Corollary 3.13(1) and (1) above.  $\dashv$

**§4. Conclusions.** We examined the incompleteness theorems of five great minds of symbolic logic, namely Gödel, Rosser, Kleene, Chaitin, and Boolos. We compared their proofs with each other, putting Gödel's second incompleteness theorem at the center of our attention, which resulted in the

following diagram (where  $\mathcal{F} \preceq \mathcal{H}$  means that  $\mathcal{F}$  is derived from  $\mathcal{H}$  in the standard way and thus  $\mathcal{F} \not\preceq \mathcal{H}$  means that  $\mathcal{F}$  is not derivable from  $\mathcal{H}$  in the standard way):



The lines over  $\overline{\text{R}}$  and  $\overline{\text{K}}$  indicate that some alternative versions of the sentences of Rosser and (the second) Kleene have been considered, and the tilde over  $\widetilde{\text{B}}$  indicates that a substantial variant of the sentence of Boolos is considered.

As the diagram shows, Boolos’ theorem is indeed the weakest among the other theorems, since it is derivable from all of them. Rosser’s theorem is the strongest in a sense, since it delivers all of the other theorems except Chaitin’s incompleteness theorem. Chaitin’s is the most neutral one, since it is not derived from any other theorem, and it delivers no other theorem, except Boolos’. Here, we did not study the incompleteness proofs whose unprovable sentences are not  $\Pi_1$ ; one prominent example is Kripke’s proof [21] for the incompleteness theorem, which shows the  $\Pi_2$ -incompleteness of  $\Sigma_2$ -sound and (RE) extensions of  $\mathfrak{B}$ .

Let us examine Boolos’ original proof more closely from [5] to see the little and amusing point that Boolos’ theorem is derivable from  $\mathbb{G}_2$ : His formula  $B(x, y)$  is our  $\mathcal{D}_\tau^{<y}(x)$ , stating that “ $x$  is definable (namable) by a formula with length  $< y$  in theory  $\mathcal{T}\mathfrak{h}_\tau$ .” His  $A(x, y)$  is our  $\mathcal{B}_\tau^{<y}(x)$ , stating that “ $x$  is the least number not definable (not named) by any formula with length  $< y$  in  $\mathcal{T}\mathfrak{h}_\tau$ .” Boolos’  $k$  is our  $\ell_\tau$ , the length of  $B(x, y)$ , and his  $F(x) = \exists y (y = [10] \times [k] \wedge A(x, y))$  is our  $\beta_\tau(x)$ . Boolos notes that the length of  $F(x)$  is less than  $10k$ , and if  $n$  (our  $b_\tau$ ) is the least number not definable by a formula with length less than  $10k$ , then  $\forall x (F(x) \leftrightarrow x = [n])$  is true but not in the output of  $M$  (unprovable in our  $\mathcal{T}\mathfrak{h}_\tau$ ). This sentence is not  $\Pi_1$ , but it is equivalent with  $F([n])$ , and this is equivalent with  $A([n], [10] \times [k])$ . This is not  $\Pi_1$  either, but it is equivalent in  $M$  (or our  $\mathcal{T}\mathfrak{h}_\tau$ ) with  $\neg B([n], [10] \times [k])$ , which is a  $\Pi_1$ -sentence. This sentence says that  $n$  is not definable by any formula with length less than  $10k$ , and in particular it is not definable by the formula  $F(x)$ . Thus,  $\neg B([n], [10] \times [k])$  implies the unprovability of  $F([n])$  in  $\mathcal{T}\mathfrak{h}_\tau$  (that  $F([n])$  is not in the output list of  $M$ ), so it implies  $\text{Con}_\tau$ , the consistency of  $\mathcal{T}\mathfrak{h}_\tau$  (that  $M$  does not output contradictory statements). Whence,  $\forall x (F(x) \leftrightarrow x = [n])$  is not provable in  $\mathcal{T}\mathfrak{h}_\tau$ , because  $\text{Con}_\tau$  is not  $\mathcal{T}\mathfrak{h}_\tau$ -provable by  $\mathbb{G}_2$ . So, the unprovability of the Boolos sentence follows from  $\mathbb{G}_2$  (Boolos continues his argument in [5] and shows the unprovability of  $F([n])$  by an argument similar to Berry’s paradox; see the proof of Theorem 3.9).

One could read in the literature that  $\mathbb{G}_2$  follows from the first incompleteness theorem; this is said (and is true) for Gödel’s proof, and we showed that it is true also for Rosser’s proof and Kleene’s proof(s). As the history goes,



the ground breaking paper [9] of Gödel was the first part, as its title shows. The second part never appeared, as Gödel felt that people could derive  $\mathbb{G}_2$  (which was promised to be proven in a sequel paper) by themselves from the first theorem; so he did not even attempt to write it. On some other proofs for the first incompleteness theorem, one may read the opposite; for example, the authors of [14] write that Maehara [20] “insists that Boolos’ theorem is different from Gödel’s one,” one reason being that “we cannot obtain the second theorem from Boolos’ theorem in the standard way.” We gave a rigorous proof for this insight in Corollary 3.15(3), and showed, moreover, that one cannot obtain  $\mathbb{G}_2$  from Chatin’s theorem, in the standard way, either.

**Acknowledgments.** This research was partially supported by the grant No. 96030030 of the Institute for Research in Fundamental Sciences (IPM), Tehran, Iran. The most helpful comments and suggestions of the referee remarkably improved the presentation and the strength of the results; this is highly appreciated.

#### REFERENCES

- [1] S. ARTEMOV and L. BEKLEMISHEV, *Provability logic*, **Handbook of Philosophical Logic** (D. Gabbay and F. Guenther, editors), second ed., Springer, New York, 2005, pp. 189–360.
- [2] A. AVRON, *A note of provability, truth and existence*, **Journal of Philosophical Logic**, vol. 20 (1991), no. 4, pp. 403–409.
- [3] L. BEKLEMISHEV, *Reflection principles and provability algebras in formal arithmetic*, **Russian Mathematical Surveys**, vol. 60 (2005), no. 2, pp. 197–268.
- [4] R. BLANCK, *On Rosser sentences and proof predicates*, M.A. thesis, Department of Philosophy, University of Göteborg, Sweden, 2006. Available at <http://bit.do/fxUqf>.
- [5] G. BOLOS, *A new proof of the Gödel incompleteness theorem*, **Notices of the American Mathematical Society**, vol. 36 (1989), no. 4, pp. 388–390. *A Letter from George Boolos*, *ibid* 36 (1989), no. 6, p. 676.
- [6] ———, **The Logic of Provability**, Cambridge University Press, Cambridge, 1994.
- [7] G. CHAITIN, *Computational complexity and Gödel’s incompleteness theorem*, **SIGACT News**, vol. 9 (1971), pp. 11–12. Abstract in **Notices of the American Mathematical Society**, vol. 17 (1970), no. 6, p. 672.
- [8] W. CRAIG, *On axiomatizability within a system*, **The Journal of Symbolic Logic**, vol. 18 (1953), no. 1, pp. 30–32.
- [9] K. GÖDEL, *Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme. I.*, **Monatshefte für Mathematik und Physik**, vol. 38 (1931), no. 1, pp. 173–198 (in German). Translated into English as: *On formally undecidable propositions of Principia Mathematica and Related Systems, I*, **Kurt Gödel, Collected Works, Volume I: Publications 1929–1936** (S. Feferman, J. W. Dawson Jr., S. C. Kleene, G. H. Moore, R. M. Solovay, and J. van Heijenoort, editors), Oxford University Press, Oxford, 1986, pp. 135–152.
- [10] D. GUASPARI and R. SOLOVAY, *Rosser sentences*, **Annals of Mathematical Logic**, vol. 16 (1979), no. 1, pp. 81–99.
- [11] P. HÁJEK and P. PUDLÁK, **Metamathematics of First-Order Arithmetic**, second ed., Springer-Verlag, New York, 1998.
- [12] D. ISAACSON, *Necessary and sufficient conditions for undecidability of the Gödel sentence and its truth*, **Logic, Mathematics, Philosophy: Vintage Enthusiasms—Essays in Honour of John L. Bell** (D. DeVidi, M. Hallett, and P. Clarke, editors), Springer, New York, 2011, pp. 135–152.



- [13] R. KAYE, *Models of Peano Arithmetic*, Oxford University Press, Oxford, 1991.
- [14] M. KIKUCHI, T. KURAHASHI, and H. SAKAI, *On proofs of the incompleteness theorems based on Berry's paradox by Vopěnka, Chaitin, and Boolos*. *Mathematical Logic Quarterly*, vol. 58 (2012), no. 4–5, pp. 307–316.
- [15] S. KLEENE, *General recursive functions of natural numbers*. *Mathematische Annalen*, vol. 112 (1936), no. 1, pp. 727–742.
- [16] ———, *A symmetric form of Gödel's theorem*. *Indagationes Mathematicae*, vol. 12 (1950), pp. 244–246.
- [17] G. KREISEL, *On weak completeness of intuitionistic predicate logic*. *The Journal of Symbolic Logic*, vol. 27 (1962), no. 2, pp. 139–158.
- [18] G. KREISEL and G. TAKEUTI, *Formally self-referential propositions for cut free classical analysis and related systems*. *Dissertationes Mathematicæ*, vol. 118 (1974), pp. 1–50.
- [19] A. MACINTYRE and H. SIMMONS, *Gödel's diagonalization technique and related properties of theories*. *Colloquium Mathematicum*, vol. 28 (1973), pp. 165–180.
- [20] S. MAEHARA, *Boolos Shi No Genkou Wo Mite*. *Gendai Shisou* (December 1989), pp. 80–92.
- [21] H. PUTNAM, *Nonstandard models and Kripke's proof of the Gödel theorem*. *Notre Dame Journal of Formal Logic*, vol. 41 (2000), no. 1, pp. 53–58.
- [22] B. ROSSER, *Extensions of some theorems of Gödel and Church*. *The Journal of Symbolic Logic*, vol. 1 (1936), no. 3, pp. 87–91.
- [23] S. SALEHI, *Gödel's incompleteness phenomenon—Computationally*. *Philosophia Scientiæ*, vol. 18 (2014), no. 3, pp. 22–37.
- [24] S. SALEHI and P. SERAJI, *Gödel–Rosser's incompleteness theorem, generalized and optimized for definable theories*. *Journal of Logic and Computation*, vol. 27 (2017), no. 5, pp. 1391–1397.
- [25] ———, *On constructivity and the Rosser property: A closer look at some Gödelean proofs*. *Annals of Pure and Applied Logic*, vol. 169 (2018), no. 10, pp. 971–980.
- [26] T. A. SLAMAN,  $\Sigma_n$ -bounding and  $\Delta_n$ -induction. *Proceedings of The American Mathematical Society*, vol. 132 (2004), no. 8, pp. 2449–2456.
- [27] C. SMORYŃSKI, *Self-Reference and Modal Logic*, Springer, New York, 1985.
- [28] A. TARSKI, A. MOSTOWSKI, and R. M. ROBINSON, *Undecidable Theories*, North-Holland, Netherlands, 1953. Reprinted by Dover Publications, 2010.
- [29] A. VISSER, *Another look at the second incompleteness theorem*. *The Review of Symbolic Logic*, vol. 13 (2020) no. 2, pp. 269–295.
- [30] C. VON BÜLOW, *A remark on equivalent Rosser sentences*. *Annals of Pure and Applied Logic*, vol. 151 (2008), no. 1, pp. 62–67.

RESEARCH INSTITUTE FOR FUNDAMENTAL SCIENCES

UNIVERSITY OF TABRIZ

29 BAHMAN BOULEVARD, P.O. BOX 51666-17766, TABRIZ, IRAN

and

SCHOOL OF MATHEMATICS

INSTITUTE FOR RESEARCH IN FUNDAMENTAL SCIENCES

P.O. BOX 19395-5746, TEHRAN, IRAN

E-mail: [root@saeedsalehi.ir](mailto:root@saeedsalehi.ir)

URL: <http://saeedsalehi.ir/>