

# Chapter 7

## On the Theory of Self-Reference

By self-reference we basically mean the possibility of talking inside a theory  $T$  about  $T$  itself or related theories. Here we can give merely a glimpse into this recently much advanced area of research; see e.g. [Bu]. We will prove Gödel's second incompleteness theorem, Löb's theorem, and many other results related to self-reference, while further results are discussed only briefly and elucidated by means of applications. All this is of great interest both for epistemology and the foundations of mathematics.

The mountain we first have to climb is the proof of the derivability conditions for PA and related theories in 7.1, and the derivable  $\Sigma_1$ -completeness in 7.2. But anyone contented with leafing through these sections can begin straight away in 7.3; from then on we will just be reaping the fruits of our labor. However, one would forgo a real adventure in doing so, namely the fusion of logic and number theory in the analysis of PA. For a comprehensive understanding of self-reference, the material of 7.1 and 7.2 (partly prepared in Chapter 6) should be studied anyway.

Gödel himself tried to interpret the notion "provable" using a modal operator in the framework of the modal system S4. This attempt reflects some of his own results, though not adequately. Only after 1970, when modal logic was sufficiently advanced, could such a program be successfully carried out. A suitable instrument turned out to be the modal logic denoted by G (or GL). The Kripke semantics for G introduced in 7.4 is an excellent tool for confirming or refuting self-referential statements. Solovay's completeness theorem and the completeness theorem of Kripke semantics for G in 7.5 are fortunately of the kind that allows application without knowing the completeness proof itself, which in both cases are not quite easy and use several technical tricks.

There are several extensions of  $\mathbf{G}$ , for example, the bimodal logic  $\mathbf{GD}$  in 7.6. This logic is related to Hilbert's famous  $\omega$ -rule. A weakening of it can be expressed by the modal operator  $\Box$  of  $\mathbf{GD}$ . A comprehensive survey can be found in [Bu, Chapter VII]; see also [Vi2]. In 7.7 we discuss some questions regarding self-reference in axiomatic set theory.

## 7.1 The Derivability Conditions

Put somewhat simply, Gödel's second incompleteness theorem states that  $\vdash_T \mathbf{Con}_T$  cannot hold for a sufficiently strong and consistent axiomatizable theory  $T$ . Here  $\mathbf{Con}_T$  is a sentence reflecting the metatheoretic statement of consistency of  $T$  *inside*  $T$ , more precisely, inside the (first-order) language  $\mathcal{L}$  of  $T$ . In a popular formulation: *If  $T$  is consistent, then this consistency is unprovable in  $T$ .* As was outlined by Gödel and will be verified in this chapter, the italicized sentence is not only true but also formalizable in  $\mathcal{L}$  and even provable in the framework of  $T$ .

The easiest way to obtain Gödel's theorem is first to prove the *derivability conditions* stated below. Their formulation supposes the arithmetizability of  $T$ , which includes the distinguishing of a sequence  $0, \underline{1}, \dots$  of ground terms; see page 250. Let  $\mathbf{bew}_T(y, x)$  be a formula that represents the recursive predicate  $\mathit{bew}_T$  in  $T$  as in 6.4. For  $\mathbf{bwb}_T(x) = \exists y \mathbf{bew}_T(y, x)$  we write  $\Box(x)$ , and  $\Box\alpha$  is to mean  $\mathbf{bwb}_T \ulcorner \alpha \urcorner$ . We may read  $\Box\alpha$  as “box  $\alpha$ ” or more suggestively “ $\alpha$  is provable in  $T$ ,” because  $\Box\alpha$  reflects the metatheoretic property  $\vdash_T \alpha$  in  $T$ . If  $\Box$  refers to some theory  $T' \neq T$  then  $\Box$  has to be indexed correspondingly. For instance,  $\Box_{\mathbf{ZFC}}\alpha$  for  $\alpha \in \mathcal{L}_\epsilon$  can easily be expressed also in  $\mathcal{L}_{ar}$ . Note that  $\Box\alpha$  is always a sentence, even if  $\alpha$  contains free variables.

Further, set  $\Diamond\alpha := \neg\Box\neg\alpha$  for  $\alpha \in \mathcal{L}$ . If  $\alpha$  is a sentence,  $\Diamond\alpha$  may be read as  *$\alpha$  is compatible with  $T$* , because it formalizes ‘ $\not\vdash_T \neg\alpha$ ’, which is, as we know, equivalent to the consistency of  $T + \alpha$ . First of all, we define  $\mathbf{Con}_T$  in a natural way by

$$\mathbf{Con}_T := \neg\Box\perp \quad (= \neg\mathbf{bwb}_T(\ulcorner \perp \urcorner)),$$

where  $\perp$  is a contradiction,  $0 \neq 0$ , for instance. We shall see in a moment that  $\mathbf{Con}_T$  is independent modulo  $T$  of the choice of  $\perp$ . The mentioned derivability conditions then read as follows:

- $D1: \vdash_T \alpha \Rightarrow \vdash_T \Box\alpha,$   
 $D2: \Box\alpha \wedge \Box(\alpha \rightarrow \beta) \vdash_T \Box\beta,$   
 $D3: \vdash_T \Box\alpha \rightarrow \Box\Box\alpha.$

Here  $\alpha, \beta$  run through all sentences of  $\mathcal{L}$ . These conditions are due to Löb, but they were considered in a slightly different setting already in [HB]. Sometimes  $D2$  is written in the equivalent form  $\Box(\alpha \rightarrow \beta) \vdash_T \Box\alpha \rightarrow \Box\beta$ , and  $D3$  as  $\Box\alpha \vdash_T \Box\Box\alpha$ .

A consequence of  $D1$  and  $D2$  is  $D0: \alpha \vdash_T \beta \Rightarrow \Box\alpha \vdash_T \Box\beta$ . This results from the following chain of implications:

$$\alpha \vdash_T \beta \Rightarrow \vdash_T \alpha \rightarrow \beta \Rightarrow \vdash_T \Box(\alpha \rightarrow \beta) \Rightarrow \vdash_T \Box\alpha \rightarrow \Box\beta \Rightarrow \Box\alpha \vdash_T \Box\beta.$$

From  $D0$  it clearly follows that  $\alpha \equiv_T \beta \Rightarrow \Box\alpha \equiv_T \Box\beta$ . In particular, the choice of  $\perp$  in  $\text{Con}_T$  is arbitrary as long as  $\perp \equiv_T 0 \neq 0$ .

**Remark 1.** Any operator  $\partial: \mathcal{L} \rightarrow \mathcal{L}$  satisfying the conditions  $d1: \vdash_T \alpha \Rightarrow \vdash_T \partial\alpha$  and  $d2: \partial(\alpha \rightarrow \beta) \vdash_T \partial\alpha \rightarrow \partial\beta$  thus satisfies also  $d0: \alpha \vdash_T \beta \Rightarrow \partial\alpha \vdash_T \partial\beta$ , and hence  $d00: \alpha \equiv_T \beta \Rightarrow \partial\alpha \equiv_T \partial\beta$ , for all  $\alpha, \beta \in \mathcal{L}$ . It likewise satisfies  $d\wedge: \partial(\alpha \wedge \beta) \equiv_T \partial\alpha \wedge \partial\beta$ , for  $\alpha \wedge \beta \vdash_T \alpha, \beta$ , hence  $\partial(\alpha \wedge \beta) \vdash_T \partial\alpha, \partial\beta \vdash_T \partial\alpha \wedge \partial\beta$  in view of  $d0$ . The converse direction  $\partial\alpha \wedge \partial\beta \vdash_T \partial(\alpha \wedge \beta)$  readily follows from  $\alpha \vdash_T \beta \rightarrow \alpha \wedge \beta$  by first applying  $d0$  and then  $d2$ .

Whereas  $D2$  and  $D3$  represent sentence schemata in  $T$ , condition  $D1$  is of metatheoretic nature and follows obviously from the representability of  $\text{bew}_T$  in  $T$ . Thus,  $D1$  holds even for weak theories such as  $T = \text{Q}$ . On the other hand, the converse of  $D1$ ,

$$D1^*: \vdash_T \Box\alpha \Rightarrow \vdash_T \alpha, \text{ for all } \alpha \in \mathcal{L}^0,$$

may fail. Fortunately, it holds for all  $\omega$ -consistent axiomatic extensions  $T \supseteq \text{Q}$  such as  $T = \text{PA}$ . Indeed,  $\not\vdash_T \alpha$  implies  $\vdash_T \neg \text{bew}_T(\underline{n}, \ulcorner \alpha \urcorner)$  for all  $n$  (Corollary 6.4.3). Hence,  $\not\vdash_T \exists y \text{bew}_T(y, \ulcorner \alpha \urcorner)$  in view of the  $\omega$ -consistency of  $T$ , that is,  $\not\vdash_T \Box\alpha$ .

Unlike  $D1$ , the properties  $D2$  and  $D3$  are not so easily obtained. The theory  $T$  must be able not only to *speak* about provability in  $T$  (perhaps via arithmetization), but also to *prove* basic properties about provability.  $D3$  is nothing else than *condition D1 formalized within T*, while  $D2$  formalizes (7) from page 230, the closure under MP in arithmetical terms. Let us first realize that  $D2$  holds, provided it has been shown that

$$D2^*: \text{bew}_T(u, x) \wedge \text{bew}_T(v, x \dot{\sim} y) \vdash_T \text{bew}_T(u * v * \langle y \rangle, y),$$

where the p.r. functions  $\tilde{\rightarrow}$ ,  $*$ , and  $y \mapsto \langle y \rangle$  appearing in  $D2^*$  must either be present or definable in  $T$ . Generally speaking,  $f \in \mathbf{F}_n$  is called *definable* in an arithmetizable theory  $T \subseteq \mathcal{L}$  (with respect to the sequence of terms  $(\underline{n})_{n \in \mathbb{N}}$  in  $T$ ) if there is a formula  $\delta(\vec{x}, y) \in \mathcal{L}$  such that

$$(1) \quad (a) \vdash_T \delta(\underline{\vec{x}}, f\underline{\vec{a}}) \text{ for all } \underline{\vec{a}}, \quad (b) \vdash_T \forall \vec{x} \exists! y \delta(\vec{x}, y).$$

Clearly,  $f$  is then also represented by  $\delta(\vec{x}, y)$ . For  $T = \text{PA}$  and related theories, (1) means that  $f$  is explicitly definable in  $T$  in the sense of **2.6** and may be introduced in  $T$  (using a corresponding symbol). From now on we will no longer distinguish between  $T$  and its definitorial extensions and apply  $\vdash_T y = f\vec{x} \leftrightarrow \delta(\vec{x}, y)$  without comment. This and (1) easily imply  $\vdash_T f\underline{\vec{a}} = f\underline{\vec{b}} = \underline{\vec{a}} \tilde{\rightarrow} \underline{\vec{b}}$ , e.g.  $\vdash_T \underline{a} \tilde{\rightarrow} \underline{b} = \underline{a} \tilde{\rightarrow} \underline{b}$ . With  $\ulcorner \alpha \urcorner, \ulcorner \beta \urcorner$  for  $x, y$ , we thus obtain from  $D2^*$  in view of  $\ulcorner \alpha \rightarrow \beta \urcorner = \underline{\dot{\alpha}} \tilde{\rightarrow} \underline{\dot{\beta}} = \underline{\dot{\alpha}} \tilde{\rightarrow} \underline{\dot{\beta}} = \ulcorner \alpha \urcorner \tilde{\rightarrow} \ulcorner \beta \urcorner$ ,

$$\text{bew}_T(u, \ulcorner \alpha \urcorner) \wedge \text{bew}_T(v, \ulcorner \alpha \rightarrow \beta \urcorner) \vdash_T \text{bew}_T(u * v * \langle \ulcorner \beta \urcorner \rangle, \ulcorner \beta \urcorner).$$

Particularization yields  $D2$ . But the real work, the definability of the functions appearing in  $D2^*$  in theories like  $T = \text{PA}$ , still lies ahead.

In order to better keep track of things, we restrict our considerations to the theories ZFC and PA, which are of central interest in nearly all foundational questions. ZFC is only briefly discussed. Here the proofs of  $D2$  and  $D3$  (with  $\square = \square_{\text{ZFC}}$ ) are much easier than in PA and need only a few lines as follows:  $D2^*$  and hence  $D2$  are clear, because the naive proof of  $D2^*$  above with  $\text{bew}_T = \text{bew}_{\text{ZFC}}$  can easily be formalized *inside* ZFC. This includes the definability of all functions occurring in  $D2^*$ , for we *did* define them; for instance, the operation  $*$  on page 224 may be defined by setting  $a * b = \emptyset$  if  $a \notin \omega$  or  $b \notin \omega$ . We arithmetize  $\mathcal{L}_\epsilon$  according to the pattern in **6.2**, encoding formulas with Gödel numbers,<sup>1</sup> so that  $\mathcal{L}_\epsilon$ -formulas are encoded within ZFC by certain  $\omega$ -terms, defined in **3.4**. Formulas from  $\mathcal{L}_{ar}$  are identified with their  $\omega$ -relativized in  $\mathcal{L}_\epsilon$ , called the *arithmetical formulas* of  $\mathcal{L}_\epsilon$ . Moreover, the arithmetical predicate  $\text{bew}_{\text{ZFC}}$  is certainly representable in ZFC by Theorem 6.4.2, since this theorem can be viewed, just like every theorem in this book, as a theorem *within* ZFC. Thus, the naive proof of  $D1$  based on this theorem (up to Corollary 6.4.3) can as a whole be carried out in ZFC, and so  $D3$  is proved.

<sup>1</sup> This is not actually necessary, since in ZFC one can talk directly about finite sequences and hence about  $\mathcal{L}_\epsilon$ -formulas (Remark 2 in **6.6**), but we do so in order to maintain coherence with the exposition in **6.2**.

Roughly speaking,  $D2$  and  $D3$  hold for ZFC because ordinary mathematics, in particular the material in Chapter 6, is formalizable in ZFC. In all of the above, no special set-theoretic constructions such as transfinite recursion are needed. Only relatively simple combinatorial facts are required. Hence there is some hope that the proofs of  $D2$  and  $D3$  can also be carried out in sufficiently strong arithmetical theories like PA. This is indeed so. The proof of  $D3$  for PA will need the most effort and will be completed only in 7.2. Our first goal will be to show that the p.r. functions occurring in  $D2^*$ , and in fact *all* p.r. functions, are explicitly definable in PA.<sup>2</sup> They turn out to be definable even in a sense stronger than required by (1) from the previous page.

**Definition.** An  $n$ -ary recursive function  $f$  is called *provably recursive* or  $\Sigma_1$ -*definable* in PA if there is a  $\Sigma_1$ -formula  $\delta_f(\vec{x}, y)$  in  $\mathcal{L}_{ar}$  such that

$$(2) \quad (a) \vdash_{PA} \delta_f(\underline{a}, \underline{f\vec{a}}) \text{ for all } \vec{a} \in \mathbb{N}^n; \quad (b) \vdash_{PA} \forall \vec{x} \exists ! y \delta_f(\vec{x}, y).$$

Since PA is  $\Sigma_1$ -complete, 2(a) is equivalent to  $\mathcal{N} \models \delta_f(\underline{a}, \underline{f\vec{a}})$  for all  $\vec{a}$ , which is often more easily verified than 2(a) and could replace 2(a). We will show that *all* p.r. functions are  $\Sigma_1$ -definable in PA, which strengthens their explicit definability in PA. Thereafter we may treat all occurring p.r. functions in PA as if they had been available in the language right from the outset. Essentially this fundamental fact allows a treatment of elementary number theory and combinatorics within the boundaries of PA and hence is particularly interesting for a critical foundation of mathematics.

If  $\delta_f(\vec{x}, y)$  in (2) is  $\Delta_0$  then  $f$  is called  $\Delta_0$ -definable. An example is the  $\beta$ -function (Exercise 1), which from now on may be supposed to be present in PA. Basic for the  $\Sigma_1$ -definability of all p.r. functions is  $\beta$ 's main property, Lemma 6.4.1, of which we need, of course, some provable version in PA. Since Euclid's lemma and the Chinese remainder theorem are involved here, these should be derived first. Clearly, the basic arithmetical laws applied in their proofs in 6.3 should be at our disposal, including those on the order relation and on  $a - b$  for  $a \geq b$ , all provable in  $\mathbb{N}$ .

The proof of Euclid's lemma is straightforward, Exercise 2. As for the Chinese remainder theorem, we avoid the quantification over finite

<sup>2</sup>In [Gö2], Gödel presented a list of 45 definable p.r. functions; the last was  $\chi_{bew}$ . Following [WR], Gödel considered a higher-order arithmetical theory. That Gödel's theorems also hold in first-order arithmetic was probably first noticed in [HB].

sequences for the time being, by stating the theorem as a scheme. Let  $\mathbf{c}, \mathbf{d}$  denote unary provably recursive functions, which may depend on further parameters. Each such  $\mathbf{c}$  determines for given  $n$  the sequence  $c_0, \dots, c_n$ , with  $c_\nu = \mathbf{c}(\nu)$  for  $\nu \leq n$ . For suggestive reasons from now on also letters such as  $n, \nu, \dots$  may denote variables in  $\mathcal{L}_{ar}$ . With the  $\Delta_0$ -definable relation  $\perp$  of coprimeness, the Chinese remainder theorem can provisionally be stated as follows: for arbitrary  $\mathbf{c}, \mathbf{d}$  as arranged above, we get

$$(3) \quad \vdash_{\text{PA}} \forall n [(\forall i, j \leq n)(c_i < d_i \wedge (i \neq j \rightarrow d_i \perp d_j)) \\ \rightarrow \exists a (\forall \nu \leq n) \text{rem}(a, d_\nu) = c_\nu].$$

To convert the original proof of the remainder theorem to one for (3) we require, for given provably recursive  $\mathbf{d}$ , the term  $\text{lcm}\{d_\nu \mid \nu \leq n\}$ , the least common multiple of  $d_0, \dots, d_n$ . **Claim:**  $f : n \mapsto \text{lcm}\{d_\nu \mid \nu \leq n\}$  is defined in PA by the  $\Sigma_1$ -formula

$$\delta_f(x, y) := (\forall \nu \leq x) d_\nu \mid y \wedge (\forall z < y) (\exists \nu \leq x) d_\nu \nmid z.$$

More precisely,  $\delta_f(x, y)$  describes a  $\Sigma_1$ -formula in  $\mathcal{L}_{ar}$  that is even  $\Delta_0$ , provided  $\mathbf{d}$  is  $\Delta_0$ -definable. Clearly  $\mathcal{N} \models \delta(\underline{n}, \text{lcm}\{d_\nu \mid \nu \leq n\})$  for all  $n$ . Thus, 2(a) holds. With the minimum schema (Exercise 4 in 3.3) applied to  $\beta(x, y) := (\forall \nu \leq x) d_\nu \mid y$ , we obtain  $\vdash_{\text{PA}} \exists! y \delta_f(x, y)$ , provided it has been shown that  $\vdash_{\text{PA}} \exists y \beta(x, y)$  (' $c_0, \dots, c_x$  have a common multiple'), which is easily derived by induction on  $x$ ; see Example 1 in 2.5. This proves the claim. After having derived Euclid's lemma in PA (Exercise 2) we confirm (3) by following the proof of the remainder theorem in 6.2, and, writing  $\beta st$  for  $\beta(s, t)$ , a suitable version of Lemma 6.4.1 as follows:

$$(4) \quad \vdash_{\text{PA}} \forall n \exists u (\forall \nu \leq n) c_\nu = \beta u \nu, \text{ for any given provably recursive } \mathbf{c}.$$

**Theorem 1.1.** *Each p.r. function  $f$  is provably recursive. Moreover, the recursion equations for  $f$  are provable in PA whenever  $f = \mathbf{Op}(g, h)$ .*

**Proof.** For the initial functions and  $+, \cdot$  the formulas  $\mathbf{v}_0 = 0$ ,  $\mathbf{v}_1 = \mathbf{Sv}_0$ ,  $\mathbf{v}_n = \mathbf{v}_\nu$  along with  $\mathbf{v}_2 = \mathbf{v}_0 + \mathbf{v}_1$  and  $\mathbf{v}_2 = \mathbf{v}_0 \cdot \mathbf{v}_1$  are obviously defining  $\Sigma_1$ -formulas. For the composition  $f = h[g_1, \dots, g_m]$ , let  $\delta_f(\vec{x}, y)$  be the formula  $y = h(g_1 \vec{x}, \dots, g_m \vec{x})$ . In this case (2) is clear, because we might think of  $h, g_1, \dots, g_m$  as being already introduced in PA, so that  $\delta_f(\vec{x}, y)$  belongs to the expanded language. Only the construction of  $\delta_f$  for the case  $f = \mathbf{Op}(g, h)$  requires some skill. We may assume that besides  $\beta$  also  $g, h$  have already been introduced in the language. Consider

$$(5) \delta_f(\vec{x}, y, z) := \exists u[\underbrace{\beta u 0 = g\vec{x} \wedge (\forall v < y)\beta u S v = h(\vec{x}, v, \beta u v)}_{\gamma(u, \vec{x}, y, z)} \wedge \beta u y = z].$$

$\delta_f$  is similar to  $\delta_{\text{exp}}$  from Remark 1 in 6.4. It is  $\Sigma_1$ , because  $\beta, g, h$  are  $\Sigma_1$ -definable. Lemma 6.4.1 applied with  $c_i = f(\vec{a}, i)$  for  $i \leq b$  shows that  $\mathcal{N} \models \delta_f(\vec{a}, b, f\vec{a})$ , equivalently 2(a). Uniqueness in 2(b), that is,

$$\delta_f(\vec{x}, y, z) \wedge \delta_f(\vec{x}, y, z') \vdash_{\text{PA}} z = z',$$

derives easily from  $\gamma(u, \vec{x}, y, z) \wedge \gamma(u', \vec{x}, y, z') \vdash_{\text{PA}} z = z'$ , which clearly follows from  $\gamma(u, \vec{x}, y, z) \wedge \gamma(u', \vec{x}, y, z') \vdash_{\text{PA}} (\forall v \leq y)\beta u v = \beta u' v$ . This is easily shown by induction on  $y$ . Also,  $\vdash_{\text{PA}} \exists z \delta_f(\vec{x}, y, z)$  will be shown inductively on  $y$ . We get  $\vdash_{\text{PA}} \exists u \beta u 0 = g\vec{x}$  (hence  $\vdash_{\text{PA}} \exists z \delta_f(\vec{x}, 0, z)$ ) from (4), choosing  $\mathbf{c}$  therein such that  $\mathbf{c}_0 = g\vec{x}$  and  $\mathbf{c}_\nu = 0$  for  $\nu \neq 0$ .  $\mathbf{c}$  is provably recursive, for the term  $g\vec{x}$  is  $\Sigma_1$ -definable. The inductive step will be verified informally, that is, we shall prove

$$(*) \quad \exists z \delta_f(\vec{x}, y, z) \vdash_{\text{PA}} \exists z' \delta_f(\vec{x}, S y, z').$$

Suppose  $\gamma(u, \vec{x}, y, z)$ . Consider the provably recursive  $\mathbf{c}: \nu \mapsto \mathbf{c}_\nu$  defined by  $\mathbf{c}_\nu = \beta u \nu$  for  $\nu \leq S y$  and  $\mathbf{c}_{S y} = h(\vec{x}, y, \beta u y)$ . Here  $u, \vec{x}, y$  are parameters in the defining  $\Sigma_1$ -formula for  $\mathbf{c}$ . So by (4) (taking  $S y$  for  $n$ ) there is some  $u'$  with  $\beta u' \nu = \mathbf{c}_\nu = \beta u \nu$  for all  $\nu \leq y$  and  $\beta u' S y = \mathbf{c}_{S y} = h(\vec{x}, y, \beta u y)$ . With this  $u'$  and  $z' = \beta u' S y$  we obtain  $\gamma(u', \vec{x}, S y, z')$ , and so  $\exists z' \delta_f(\vec{x}, S y, z')$ . This confirms  $(*)$  and hence 2(b). Thus,  $f$  is provably recursive and may now be introduced in PA. We finally sketch a proof of the recursion equations for  $f$  in PA, which also in PA may be written as usual, i.e.,

$$(A) \quad \vdash_{\text{PA}} f(\vec{x}, 0) = g\vec{x}, \quad (B) \quad \vdash_{\text{PA}} f(\vec{x}, S y) = h(\vec{x}, y, f(\vec{x}, y)).$$

(A) holds because  $\vdash_{\text{PA}} \delta_f(\vec{x}, 0, f(\vec{x}, 0)) \equiv_{\text{PA}} \exists u(\beta u 0 = g\vec{x} \wedge \beta u 0 = f(\vec{x}, 0))$  and clearly  $\exists u(\beta u 0 = g\vec{x} \wedge \beta u 0 = f(\vec{x}, 0)) \vdash f(\vec{x}, 0) = g\vec{x}$ . (B) follows by  $<$ -induction on  $y$  applied to  $\alpha = \alpha(\vec{x}, y) := f(\vec{x}, S y) = h(\vec{x}, y, f(\vec{x}, y))$ . Assume that  $(\forall v < y)\alpha \frac{v}{y}$ . Choosing  $u$  in (5) such that  $\gamma(u, \vec{x}, S y, f(\vec{x}, S y))$ , we readily obtain  $(\forall v \leq y)f(\vec{x}, v) = \beta u v$ , so that

$$f(\vec{x}, S y) = \beta u S y = h(\vec{x}, y, \beta u y) = h(\vec{x}, y, f(\vec{x}, y)).$$

This confirms  $\forall y((\forall v < y)\alpha \frac{v}{y} \rightarrow \alpha)$ , hence  $\vdash_{\text{PA}} \forall y \alpha$  by  $<$ -induction.  $\square$

We thus have achieved our first goal. Next observe that the properties of  $*$ ,  $\ell, \dots$  from the remark on page 230 along with the basic property (5) stated there are also readily proved *within* PA. This is a little extra program that includes the proof of unique prime factorization, see Exercise 4.

Thus,  $D2^*$  and hence  $D2$  are indeed provable for  $T = \text{PA}$ . In particular, the property (6) from page 230 carries over to  $\text{PA}$ , so that

$$(6) \quad \Box(x \dot{\sim} y) \vdash_{\text{PA}} \Box(x) \rightarrow \Box(y).$$

We mention that  $\Box$  in (3) may even denote the formula  $\text{bwb}_T$  for any axiomatizable (and arithmetizable) theory  $T$ .  $D3$  will be proved in the next section in a somewhat broader context.

**Remark 2.** The formalized equations of Exercise 3 in 6.4 are now also provable in  $\text{PA}$ . For instance, item (b) reads  $\vdash_{\text{PA}} \text{sb}_{\vec{x}}(\ulcorner \varphi \urcorner, \vec{x}) = \text{sb}_{\vec{x}'}(\ulcorner \varphi \urcorner, \vec{x}')$  for  $\varphi = \varphi(\vec{x})$ , where  $\vec{x}' (\subseteq \vec{x})$  enumerates the free variables of  $\varphi$ . As regards (c), consider first a special case. Let  $\varphi$  be  $\text{S}x = y$ . Then  $\text{sb}_{xy}(\ulcorner \varphi \urcorner, x, \text{S}x) = \text{sb}_x((\ulcorner \varphi \frac{\text{S}x}{y} \urcorner), x)$ , formalized  $\text{sb}_{xy}(\ulcorner \varphi \urcorner, x, y) \frac{\text{S}x}{y} = \text{sb}_x(\ulcorner \varphi \frac{\text{S}x}{y} \urcorner, x)$ . For the proof of this equation in  $\text{PA}$ , just  $\vdash_{\text{PA}}$  cf  $\text{S}x = \text{S}cx$  is required, which holds by Theorem 1.1. Whoever wants to write down a detailed proof should follow the example on page 249.

## Exercises

1. Prove in  $\text{PA}$  the  $\Delta_0$ -definability of the remainder function  $\text{rem}$ , the pairing function, and the  $\beta$ -function; see 6.4. In particular,  $\text{rem}$  is defined by  $\delta_{\text{rem}}(a, b, r) := (\exists q \leq a)(a = b \cdot q + r \wedge r < b) \vee b = r = 0$ . The laws of arithmetic as given by  $\mathbb{N}$  (page 235) may be used.
2. Prove in  $\text{PA}$  (a)  $(\forall a, b > 0) \exists x \exists y (a \perp b \rightarrow ax + 1 = by)$ , that is, Euclid's lemma. (b)  $(\forall a > 1) \exists p (\text{prim } p \wedge p \mid a)$  ('each number  $\geq 2$  has a prime divisor'), (c)  $\vdash_{\text{PA}} (\forall a, b > 0) \forall p (\text{prim } p \wedge p \mid ab \rightarrow p \mid a \vee p \mid b)$ .
3. Show that  $\vdash_{\text{PA}} \text{prim } p \wedge p \mid \text{lcm}\{\mathbf{d}_\nu \mid \nu \leq n\} \rightarrow (\exists i \leq n) p \mid \mathbf{d}_i$ , required for carrying out the proof of the Chinese remainder theorem in  $\text{PA}$ .
4. One of several possibilities of formalizing the prime factorization in  $\text{PA}$  is  $(\forall n \geq 2) (\exists m \geq 2) n = \prod_{i \leq \ell_m} p_i^{(m)_i}$ , where  $m$  serves as a variable for the sequence of prime exponents.<sup>3</sup> Prove this in  $\text{PA}$ , as well as its uniqueness, which is essentially based on Exercise 2.
5. Let  $T' = T + \alpha$  and  $T$  satisfy  $D1$ – $D3$ . Show that
  - (a)  $\vdash_T \Box_{T'} \varphi \leftrightarrow \Box_T (\alpha \rightarrow \varphi)$  (the formalized deduction theorem),
  - (b)  $D1$ – $D3$  hold also for  $T'$ .

<sup>3</sup> An equivalent formalization of the prime factorization in  $\text{PA}$  using the  $\beta$ -function is  $(\forall k \geq 2) \exists u \exists n (k = \prod_{i \leq n} p_i^{\beta_{ui}} \wedge \beta_{un} \neq 0)$ .



## 7.2 The Provable $\Sigma_1$ -Completeness

$D3$  is a special case of the *provable*  $\Sigma_1$ -completeness. This is essentially the statement  $\vdash_{\text{PA}} \alpha \rightarrow \Box\alpha$  for  $\Sigma_1$ -sentences  $\alpha$ . The proof demands still additional preparation, and even good textbooks do not carry out all proof steps. All steps described in this section and not handled in detail can easily be completed in full by the sufficiently assiduous reader. Life could be made easier through the mutual interpretability of  $\text{PA}$  and  $\text{ZFC}_{\text{fin}}$  mentioned in 6.6. Let  $\Box = \Box(x)$  denote the formula  $\text{bwb}_{\text{PA}}(x)$  till the end of this section. We first introduce an additional notation. Let  $\varphi = \varphi(\vec{x})$ .

**Definition.**  $\Box[\varphi] := \Box(\text{sb}_{\vec{x}}(\ulcorner\varphi\urcorner, \vec{x})) (= \text{bwb}_{\text{PA}} \frac{\text{sb}_{\vec{x}}(\ulcorner\varphi\urcorner, \vec{x})}{x})$ .

By Remark 2 in 7.1,  $\vdash_{\text{PA}} \text{sb}_{\vec{x}}(\ulcorner\varphi\urcorner, \vec{x}) = \text{sb}_{\vec{x}'}(\ulcorner\varphi\urcorner, \vec{x}')$ , where  $\vec{x}'$  enumerates *free*  $\varphi$ . Hence, we may assume w.l.o.g. that *free*  $\Box[\varphi] = \text{free } \varphi$ . Moreover, for  $\alpha \in \mathcal{L}_{\text{ar}}^0$  we have  $\vdash_{\text{PA}} \text{sb}_{\vec{x}}(\ulcorner\alpha\urcorner, \vec{x}) = \text{sb}_{\emptyset}(\ulcorner\alpha\urcorner) = \ulcorner\alpha\urcorner$ , hence  $\Box[\alpha]$  and  $\Box\alpha$  may be identified. ‘ $\vdash_{\text{PA}} \varphi(\vec{a})$  for all  $\vec{a} \in \mathbb{N}^n$ ’ is reflected in  $\text{PA}$  by ‘ $\vdash_{\text{PA}} \forall \vec{x} \Box[\varphi]$ ’. The latter thus reflects in  $\text{PA}$  the existence of a *collection of proofs* which, due to the  $\omega$ -incompleteness of  $\text{PA}$ , may be less than  $\vdash_{\text{PA}} \Box \forall \vec{x} \varphi$ , or what amounts to the same,  $\vdash_{\text{PA}} \Box \varphi$ .

**Example.** Let  $\varphi = \varphi(x, y)$  be  $Sx = y$ . We prove  $\varphi \vdash_{\text{PA}} \Box[\varphi]$ , or equivalently,  $\vdash_{\text{PA}} \Box[\varphi] \frac{Sx}{y}$ , where w.l.o.g.  $x, y$  do not occur bound in  $\Box(x)$ . In order to prove  $\vdash_{\text{PA}} \Box[\varphi] \frac{Sx}{y}$  observe that in view of Remark 2 in 7.1,

$$\Box[\varphi] \frac{Sx}{y} = \Box(\text{sb}_{xy}(\ulcorner\varphi\urcorner, x, Sy)) \equiv_{\text{PA}} \Box(\text{sb}_x(\ulcorner\varphi \frac{Sx}{y}\urcorner, x)) = \Box[\alpha(x)]$$

with  $\alpha(x) := Sx = Sx$ . Thus, it suffices to verify  $\vdash_{\text{PA}} \Box[\alpha(x)]$  (equivalently,  $\vdash_{\text{PA}} \forall x \Box[\alpha(x)]$ ). This reflects in  $\text{PA}$  ‘for arbitrary  $n$ ,  $\vdash_{\text{PA}} \mathbf{S}\underline{n} = \mathbf{S}\underline{n}$ ’. We verify  $\vdash_{\text{PA}} \Box[\alpha(x)]$  in detail. Consider the p.r. function  $\tilde{\alpha}: n \mapsto \text{sb}_x(\hat{\alpha}, n)$  (the Gödel number of  $\alpha(\underline{n})$ ). By axiom  $\Lambda 9$ ,  $\langle \tilde{\alpha}(n) \rangle$  is for each  $n$  a trivial arithmetized proof of length 1. Stated within  $\text{PA}$ ,  $\vdash_{\text{PA}} \text{bew}_{\text{PA}}(\langle \tilde{\alpha}(x) \rangle, \tilde{\alpha}(x))$ . This clearly yields  $\vdash_{\text{PA}} \exists y \text{bew}_{\text{PA}}(y, \tilde{\alpha}(x)) = \Box(\tilde{\alpha}(x)) = \Box[\alpha]$ .

Next we prove some modifications  $D1$ ,  $D2$  for  $\alpha = \alpha(\vec{x})$  and  $\beta = \beta(\vec{x})$ :

- (1) (a)  $\vdash_{\text{PA}} \alpha \Rightarrow \vdash_{\text{PA}} \Box[\alpha]$ ; (b)  $\Box[\alpha \rightarrow \beta] \vdash_{\text{PA}} \Box[\alpha] \rightarrow \Box[\beta]$ .

To see (a) let  $\vdash_{\text{PA}} \alpha$ , hence also  $\vdash_{\text{PA}} \forall \vec{x} \alpha$  and so  $\vdash_{\text{PA}} \Box \forall \vec{x} \alpha$ . Just as in the above example, a proof for  $\forall \vec{x} \alpha$  provides one for  $\alpha_{\vec{x}}(\vec{a})$  in a p.r. way, or stated *within*  $\text{PA}$ :  $\Box \forall \vec{x} \vdash_{\text{PA}} \Box(\text{sb}_{\vec{x}}(\ulcorner\alpha\urcorner, \vec{x})) (= \Box[\alpha]$ ; thus,  $\vdash_{\text{PA}} \Box[\alpha]$ ).

(b) follows from (6) in 7.1 with  $\text{sb}_{\vec{x}}(\ulcorner \alpha \urcorner, \vec{x}), \text{sb}_{\vec{x}}(\ulcorner \beta \urcorner, \vec{x})$  for  $x, y$ , observing that  $\vdash_{\text{PA}} \text{sb}_{\vec{x}}(\ulcorner \alpha \rightarrow \beta \urcorner, \vec{x}) = \text{sb}_{\vec{x}}(\ulcorner \alpha \urcorner, \vec{x}) \dot{\sim} \text{sb}_{\vec{x}}(\ulcorner \beta \urcorner, \vec{x})$ , see Exercise 3 in 6.4. (c) of this exercise yields for all not necessarily distinct  $x, y$

$$(2) \quad \Box[\alpha]_{\vec{x}}^t \equiv_{\text{PA}} \Box[\alpha]_{\vec{x}} \quad (t \in \{0, y, \text{Sy}\} \text{ and } y \notin \text{bnd} \alpha).$$

Now, D3 is only a special case of the *provable*  $\Sigma_1$ -*completeness* of PA, stated not only for sentences, but for arbitrary formulas as follows:

$$(3) \quad \varphi \vdash_{\text{PA}} \Box[\varphi] \text{ (equivalently, } \vdash_{\text{PA}} \varphi \rightarrow \Box[\varphi]), \text{ for all } \Sigma_1\text{-formulas } \varphi.$$

Indeed, choose in (3) for  $\varphi$  the  $\Sigma_1$ -sentence  $\Box\alpha$  for any  $\alpha \in \mathcal{L}_{ar}^0$ . Then  $\Box\alpha \vdash_{\text{PA}} \Box[\Box\alpha] \equiv \Box\Box\alpha$ , and D3 is proved. We obtain (3) from Theorem 2.1 below, since by (1), (2), and since w.l.o.g.  $\text{free} \alpha = \text{free} \Box[\alpha]$ , the operator  $\partial: \alpha \mapsto \Box[\alpha]$  satisfies the conditions of the theorem.

**Theorem 2.1.** *Let  $\partial: \mathcal{L}_{ar} \rightarrow \mathcal{L}_{ar}$  be any operator with  $\text{free} \partial\alpha \subseteq \text{free} \alpha$  satisfying*

$$d1: \quad \vdash_{\text{PA}} \alpha \Rightarrow \vdash_{\text{PA}} \partial\alpha,$$

$$d2: \quad \partial(\alpha \rightarrow \beta) \vdash_{\text{PA}} \partial\alpha \rightarrow \partial\beta,$$

$$ds: \quad \partial\alpha_{\vec{x}}^t \equiv_{\text{PA}} \partial(\alpha_{\vec{x}}^t) \quad (t \in \{0, y, \text{Sy}\}, \quad y \notin \text{bnd} \alpha).$$

Then  $\vdash_{\text{PA}} \varphi \rightarrow \partial\varphi$  holds for all  $\Sigma_1$ -formulas  $\varphi \in \mathcal{L}_{ar}$ .

**Proof.**  $\partial$  satisfies also  $d0, d00$ , and  $d\wedge$  (see Remark 1 in 7.1). Hence, by Theorem 6.7.2 and  $d00$  we need to carry out the proof only for special  $\Sigma_1$ -formulas. First let  $\varphi$  be  $\text{S}x = y$ . Clearly,  $\vdash_{\text{PA}} \varphi \rightarrow \partial\varphi$  is equivalent to  $\vdash_{\text{PA}} \partial\varphi_{\vec{y}}^{\text{S}x}$ , and this to  $\vdash_{\text{PA}} \partial\text{S}x = \text{S}x$  by  $ds$ , which is obvious from  $d1$ . Now let  $\varphi$  be  $x + y = z$ . We shall prove  $\vdash_{\text{PA}} \forall yz(\varphi \rightarrow \partial\varphi)$  by induction on  $x$ . Observing that  $y = z \vdash_{\text{PA}} \partial y = z$  (equivalently  $\vdash_{\text{PA}} \partial z = z$ ), we obtain  $\varphi_{\vec{x}}^0 \vdash_{\text{PA}} y = z \vdash_{\text{PA}} \partial y = z \equiv_{\text{PA}} \partial(\varphi_{\vec{x}}^0) \equiv_{\text{PA}} \partial\varphi_{\vec{x}}^0$ . Thus,  $\vdash_{\text{PA}} \forall yz(\varphi \rightarrow \partial\varphi)_{\vec{x}}^0$ . Now  $\varphi_{\vec{y}}^{\text{S}y} \equiv_{\text{PA}} \varphi_{\vec{x}}^{\text{S}x}$ ; hence  $\partial\varphi_{\vec{y}}^{\text{S}y} \equiv_{\text{PA}} \partial\varphi_{\vec{x}}^{\text{S}x}$ , by  $d00, ds$ . The induction step  $\forall yz(\varphi \rightarrow \partial\varphi) \vdash_{\text{PA}} \forall yz(\varphi \rightarrow \partial\varphi)_{\vec{x}}^{\text{S}x}$  follows then from

$$\forall yz(\varphi \rightarrow \partial\varphi) \vdash_{\text{PA}} \varphi_{\vec{y}}^{\text{S}y} \rightarrow \partial\varphi_{\vec{y}}^{\text{S}y} \vdash_{\text{PA}} \varphi_{\vec{x}}^{\text{S}x} \rightarrow \partial\varphi_{\vec{x}}^{\text{S}x} = (\varphi \rightarrow \partial\varphi)_{\vec{x}}^{\text{S}x}.$$

The formula  $x \cdot y = z$  is left to the reader, who should observe  $d\wedge, d2$ , the induction steps for  $\wedge$  and  $\exists$ , and  $\text{S}x \cdot y = z \equiv_{\text{PA}} \exists u(x \cdot y = u \wedge u + y = z)$ .

We now treat the logical connectives. The induction steps for  $\wedge, \vee, \exists$  are simple. Indeed, from  $d\wedge$  we obtain

$$\alpha \wedge \beta \vdash_{\text{PA}} \alpha, \beta \vdash_{\text{PA}} \partial\alpha \wedge \partial\beta \vdash_{\text{PA}} \partial(\alpha \wedge \beta).$$

For  $\vee$  note that  $\alpha \vdash_{\text{PA}} \partial\alpha \vdash_{\text{PA}} \partial(\alpha \vee \beta)$ , and similarly for  $\beta$ . Further, since  $\varphi \vdash \exists x\varphi$  we get  $\varphi \vdash_{\text{PA}} \partial\varphi \vdash_{\text{PA}} \partial\exists x\varphi$  by  $d0$ , and from  $x \notin \text{free} \partial\exists x\varphi$

follows  $\exists x\varphi \vdash_{\text{PA}} \partial\exists x\varphi$ . The prime-term substitution step ( $t$  is prime in  $\frac{t}{x}$ ) also runs smoothly:  $\varphi \vdash_{\text{PA}} \partial\varphi$  yields  $\varphi \frac{t}{x} \vdash_{\text{PA}} \partial\varphi \frac{t}{x} \vdash_{\text{PA}} \partial(\varphi \frac{t}{x})$  by *ds*.

It remains to verify the step for bounded quantification. Suppose that  $\alpha \vdash_{\text{PA}} \partial\alpha$  and  $y \notin \text{var}\alpha$ . We prove  $\varphi := (\forall x < y)\alpha \vdash_{\text{PA}} \partial\varphi$  by induction on  $y$ . The initial step is obvious:  $\vdash_{\text{PA}} \varphi \frac{0}{y}$ , and therefore

$$\vdash_{\text{PA}} \partial(\varphi \frac{0}{y}) \vdash_{\text{PA}} \partial\varphi \frac{0}{y} \vdash_{\text{PA}} \varphi \frac{0}{y} \rightarrow \partial\varphi \frac{0}{y}.$$

Clearly,  $\varphi \frac{Sy}{y} \equiv_{\text{PA}} \varphi \wedge \alpha \frac{y}{x}$ . Hence  $\alpha \frac{y}{x} \vdash_{\text{PA}} \partial\alpha \frac{y}{x} \vdash_{\text{PA}} \partial(\alpha \frac{y}{x})$  because of  $\alpha \vdash_{\text{PA}} \partial\alpha$ . That leads to

$$\begin{aligned} \varphi \frac{Sy}{y} \wedge (\varphi \rightarrow \partial\varphi) \vdash_{\text{PA}} \varphi \wedge \alpha \frac{y}{x} \wedge (\varphi \rightarrow \partial\varphi) \vdash_{\text{PA}} \partial\varphi \wedge \partial(\alpha \frac{y}{x}) \\ \vdash_{\text{PA}} \partial(\varphi \wedge \alpha \frac{y}{x}) \vdash_{\text{PA}} \partial(\varphi \frac{Sy}{y}). \end{aligned}$$

Thus,  $\varphi \rightarrow \partial\varphi \vdash_{\text{PA}} \varphi \frac{Sy}{y} \rightarrow \partial(\varphi \frac{Sy}{y})$ , which is obviously equivalent to the inductive step.  $\square$

**Remark 3.**  $D1$ – $D3$  are also provable for much weaker theories than PA, e.g., for the so-called *elementary arithmetic*  $\text{EA} = I\Delta_0 + \forall xy\exists z\delta_{\text{exp}}(x, y, z)$ . Here  $I\Delta_0$  is defined in Remark 1 in 6.3 and  $\delta_{\text{exp}}$  is a defining  $\Delta_0$ -formula for exp, see also [FS]. Also Theorem 1.1 can essentially be strengthened and has many variants. For instance, the provably recursive functions of  $I\Sigma_1$  (like PA but IS restricted to  $\Sigma_1$ -formulas) are precisely the p.r. ones, [Tak]. The same provably recursive functions has EA augmented by the  $\Pi_2$ -induction schema without parameters, [Be4]. It is noteworthy that the provable recursive functions of EA itself are precisely the elementary ones, [Si]. For more material on the metatheory of PA and related theories see [Bar, Part D], and in particular [HP].

## 7.3 The Theorems of Gödel and Löb

We are now in a position to harvest the yields of our efforts. As long as not stated otherwise, let  $T$  denote any arithmetizable axiomatic theory in  $\mathcal{L}$ , that satisfies the derivability conditions  $D1$ – $D3$  of 7.1 along with the fixed point lemma of 6.5. We direct attention straight away to the uniqueness statement of Lemma 3.1(b) below. According to this claim, up to equivalence in  $T$  at most  $\Box\alpha \rightarrow \alpha$  can be the fixed point of the formula  $\Box(x) \rightarrow \alpha$ . The proof of Theorem 3.2 will show that  $\neg\Box(x)$  too has only one fixed point modulo  $T$ . Beneath all this lies, as we shall see from Corollary 5.6, a completely general result.

**Lemma 3.1.** *Let  $T$  be as arranged above, and let  $\alpha, \gamma \in \mathcal{L}^0$  be such that  $\gamma \equiv_T \Box\gamma \rightarrow \alpha$ . Then (a)  $\Box\gamma \equiv_T \Box\alpha$  and (b)  $\gamma \equiv_T \Box\alpha \rightarrow \alpha$ .*

**Proof.** The supposition yields  $\Box\gamma \vdash_T \Box(\Box\gamma \rightarrow \alpha) \vdash_T \Box\Box\gamma \rightarrow \Box\alpha$ , by *D0* and *D2*. Now by *D3*, we clearly obtain  $\Box\gamma \vdash_T \Box\Box\gamma$ , hence  $\Box\gamma \vdash_T \Box\alpha$ . Since  $\alpha \vdash_T \Box\gamma \rightarrow \alpha \equiv_T \gamma$  and so  $\alpha \vdash_T \gamma$ , it follows that  $\Box\alpha \vdash_T \Box\gamma$  by *D0*. Together with the already verified  $\Box\gamma \vdash_T \Box\alpha$  we get (a). Using (a) we may replace  $\Box\gamma$  with  $\Box\alpha$  in  $\gamma \equiv_T \Box\gamma \rightarrow \alpha$ , which results in (b).  $\square$

**Theorem 3.2 (Second incompleteness theorem).** *PA satisfies alongside the fixed point lemma also *D1–D3*. Every theory  $T$  with these properties satisfies the conditions*

$$(1) \not\vdash_T \mathbf{Con}_T \text{ provided } T \text{ is consistent,} \quad (2) \vdash_T \mathbf{Con}_T \rightarrow \neg\Box\mathbf{Con}_T.$$

**Proof.** *D1–D3* were proved for PA in 7.1. (1) follows from (2). Assume  $\vdash_T \mathbf{Con}_T$ . Then  $\vdash_T \Box\mathbf{Con}_T$  by *D1*, as well as  $\vdash_T \neg\Box\mathbf{Con}_T$  by (2). Thus,  $T$  is inconsistent. To verify (2), let  $\gamma$  be a fixed point of  $\neg\Box(x)$ , i.e.,

$$(*) \quad \gamma \equiv_T \neg\Box\gamma \quad (\equiv \Box\gamma \rightarrow \perp).$$

By Lemma 3.1(b) with  $\alpha = \perp$ , we obtain  $\gamma \equiv_T \Box\perp \rightarrow \perp \equiv \neg\Box\perp = \mathbf{Con}_T$ . Replacing  $\gamma$  in (\*) with  $\mathbf{Con}_T$  gives  $\mathbf{Con}_T \equiv_T \neg\Box\mathbf{Con}_T$ . Half of this is the claim (2).  $\square$

Thus, by (1), no sufficiently strong consistent theory can prove its own consistency. In particular,  $\not\vdash_{\text{PA}} \mathbf{Con}_{\text{PA}}$  as long as PA is consistent which is assumed throughout this book and is a minimal assumption for a far-reaching metamathematics. The above proof shows that  $\mathbf{Con}_T$  is the only fixed point of  $\neg\mathbf{bwb}_T$  modulo  $T$ . Actually, it shows a bit more, namely

$$(3) \quad \mathbf{Con}_T \equiv_T \neg\Box\mathbf{Con}_T.$$

This strengthens (2), but only by a little:  $\neg\Box\mathbf{Con}_T \vdash_T \mathbf{Con}_T$  is just a special case of

$$(4) \quad \neg\Box\alpha \vdash_T \mathbf{Con}_T \text{ (equivalently, } \neg\mathbf{Con}_T \vdash_T \Box\alpha), \text{ for every } \alpha \in \mathcal{L}.$$

This follows from  $\perp \vdash_T \alpha$ , since  $\neg\mathbf{Con}_T \equiv \Box\perp \vdash_T \Box\alpha$  by *D0*. (4) reflects in  $T$  ‘If  $T$  is inconsistent then every formula is provable’. From (1) and (3) we get in particular  $\not\vdash_{\text{PA}} \neg\Box_{\text{PA}} \mathbf{Con}_{\text{PA}}$ , although ‘ $\mathbf{Con}_{\text{PA}}$  is unprovable in PA’ is true according to (1) (again we tacitly use the consistence of PA).  $\neg\Box_{\text{PA}} \mathbf{Con}_{\text{PA}}$  reflects ‘ $\mathbf{Con}_{\text{PA}}$  is unprovable in PA’; hence  $\not\vdash_{\text{PA}} \neg\Box_{\text{PA}} \mathbf{Con}_{\text{PA}}$  is just another formulation of the second incompleteness theorem.

The above claims hold independently of the “truth content” of the sentences provable in  $T$ . Namely, a consequence of the second incompleteness theorem is the existence of consistent theories  $T \supseteq \text{PA}$  in which along with claims true in  $\mathcal{N}$  also false ones are provable, i.e., in which truth and untruth live in peaceful coexistence with each other. Such “dream theories” are highly rich in content, for all of them include ordinary number theory. An example is  $\text{PA}^\perp := \text{PA} + \neg \text{Con}_{\text{PA}}$ . This theory is consistent because *the consistency of  $\text{PA}^\perp$  is equivalent to the unprovability of  $\text{Con}_{\text{PA}}$  in  $\text{PA}$* . The italicized sentence is even provable in  $\text{PA}$ , as (5) below will show. By the formalized deduction theorem (Exercise 5 in 7.1),  $\Box_{T+\alpha^\perp} \equiv_T \Box(\alpha \rightarrow \perp) \equiv \Box\neg\alpha$ ; hence  $\neg\Box_{T+\alpha^\perp} \equiv_T \neg\Box\neg\alpha$  ( $\equiv \Diamond\alpha$ ), and consequently,

$$(5) \quad \text{Con}_{T+\alpha^\perp} \equiv_T \neg\Box\neg\alpha \quad (\text{in particular, } \text{Con}_{\text{PA}^\perp} \equiv_{\text{PA}} \neg\Box_{\text{PA}} \text{Con}_{\text{PA}}).$$

The special cases under (5) and (3) for  $T = \text{PA}$  now clearly yield

$$(6) \quad \text{Con}_{\text{PA}} \equiv_{\text{PA}} \text{Con}_{\text{PA}^\perp} \quad (\text{hence also } \text{Con}_{\text{PA}} \equiv_{\text{PA}^\perp} \text{Con}_{\text{PA}^\perp}).$$

Put together,  $\text{PA}^\perp$  contains ordinary number theory as known to us, but also proves the indubitably false sentence  $\text{bwb}_{\text{PA}}(\ulcorner 0 \neq 0 \urcorner)$ . Moreover, because of  $\vdash_{\text{PA}^\perp} \neg \text{Con}_{\text{PA}}$  and hence  $\vdash_{\text{PA}^\perp} \neg \text{Con}_{\text{PA}^\perp}$  by (6),  $\text{PA}^\perp$  proves (the reflection of) its own inconsistency, although along with  $\text{PA}$  also  $\text{PA}^\perp$  is consistent. It claims to have a mysterious proof of  $\perp$ . Thus, consistency of  $T$  can have a different meaning within  $T$  and seen from outside, just as the meanings of *countable* diverge, depending on whether one is situated in  $\text{ZFC}$  or is looking at it from outside. One may even say that  $\text{PA}^\perp$  is lying to us with the claim  $\neg \text{Con}_{\text{PA}^\perp}$ .

We learn from the preceding that the extension  $T + \text{Con}_T$  of a consistent theory  $T$  need not be consistent.  $T = \text{PA}^\perp$  is a concrete example, and in fact only one of arbitrarily many others. More on the meaning of  $\neg \text{Con}_T$  will be said in Theorem 3.4.

We now discuss what is, along with (3), the most famous example of a self-referential sentence. Clearly, a fixed point  $\alpha$  of  $\Box(x)$  claims just its own provability, that is,  $\alpha \equiv_T \Box\alpha$ . A trivial example is  $\alpha = \top$ , because  $\vdash_T \Box\top \rightarrow \top$ , and since  $\vdash_T \top$ , clearly  $\vdash_T \Box\top$ , so that  $\top \equiv_T \Box\top$ . What is surprising here is that  $\top$  turns out to be the only fixed point of  $\Box(x)$  modulo  $T$ . By  $D4^\circ$  below,  $\vdash_T \Box\alpha \rightarrow \alpha$  implies  $\vdash_T \alpha$  and so  $\alpha \equiv_T \top$  (which confirms the uniqueness), although one might perhaps expect  $\vdash_T \Box\alpha \rightarrow \alpha$  for all  $\alpha \in \mathcal{L}^0$  because  $\Box\alpha \rightarrow \alpha$  is intuitively true.

**Theorem 3.3 (Löb’s theorem).** *Take  $T$  to satisfy D1–D3 and the fixed point lemma. Then  $T$  has the properties*

$$D4: \vdash_T \Box(\Box\alpha \rightarrow \alpha) \rightarrow \Box\alpha, \quad D4^\circ: \vdash_T \Box\alpha \rightarrow \alpha \Rightarrow \vdash_T \alpha \quad (\alpha \in \mathcal{L}^0).$$

**Proof.** Let  $\gamma$  be a fixed point of  $\Box(x) \rightarrow \alpha$ , i.e.,  $\gamma \equiv_T \Box\gamma \rightarrow \alpha$ . Then  $\gamma \equiv_T \Box\alpha \rightarrow \alpha$  by Lemma 3.1(b). This and D0 imply  $\Box\gamma \equiv_T \Box(\Box\alpha \rightarrow \alpha)$ . Lemma 3.1(a) states  $\Box\gamma \equiv_T \Box\alpha$ , hence  $\Box\alpha \equiv_T \Box(\Box\alpha \rightarrow \alpha)$ . Half of this is D4. Now suppose  $\vdash_T \Box\alpha \rightarrow \alpha$ . Then by D1,  $\vdash_T \Box(\Box\alpha \rightarrow \alpha)$ . Using D4 results in  $\vdash_T \Box\alpha$ , and  $\vdash_T \Box\alpha \rightarrow \alpha$  yields  $\vdash_T \alpha$ , thus proving  $D4^\circ$ .  $\square$

D4 reflects just  $D4^\circ$  in  $T$ . One application of Löb’s theorem is an extremely easy proof of  $\not\vdash_{\text{PA}} \text{Con}_{\text{PA}}$ . Indeed,  $\vdash_{\text{PA}} \text{Con}_{\text{PA}} (\equiv \Box\perp \rightarrow \perp)$  implies  $\vdash_{\text{PA}} \perp$  by  $D4^\circ$ . That’s all. Similarly, D4 implies (2) for  $\alpha = \perp$  by contraposition. Thus, Löb’s theorem is stronger than Gödel’s second incompleteness theorem, which is not obvious at first glance.

Unlike  $\text{PA}^+$ ,  $\text{PA} + \text{Con}_{\text{PA}}$  conforms to truth (in  $\mathcal{N}$ ). Unfortunately it is not quite clear what  $\text{Con}_{\text{PA}}$  means in number-theoretic terms. This is clear, however, for an arithmetical statement discovered by Paris and Harrington (see [Bar]) that implies  $\text{Con}_{\text{PA}}$ ; this statement is provable in ZFC but not in PA. Since then, many such sentences have been found, mostly of a combinatorial nature. A popular example is

**Goodstein’s theorem.** *Every Goodstein sequence ends in 0.*

A *Goodstein sequence* is a number sequence  $(a_n)_{n \in \mathbb{N}}$ , with arbitrary  $a_0$  given in advance, such that  $a_{n+1}$  is obtained from  $a_n$  as follows: Let  $b_n = n + 2$ , so that  $b_0 = 2$ ,  $b_1 = 3$ , etc. Expand  $a_n$  in  $b$ -adic base for  $b := b_n$ , so that for suitable  $k$ ,

$$(*) \quad a_n = \sum_{i \leq k} b^{k-i} c_i, \quad \text{with } 0 \leq c_i < b.$$

Also the powers  $k - i$  are represented in  $b$ -adic form, so too the powers of powers, and so on. Now replace  $b$  everywhere with  $b + 1 (= b_{n+1})$  and subtract 1 from the output. The result is  $a_{n+1}$ . The table below gives an example beginning with  $a_0 = 11$ ; already  $a_6$  has the value 134 217 727.

$a_0 = 11 = 2^{2+1} + 2 + 1$	$2 \rightsquigarrow 3$	$3^{3+1} + 3 + 1 = 85$
$a_1 = 84 = 3^{3+1} + 3$	$3 \rightsquigarrow 4$	$4^{4+1} + 4 = 1028$
$a_2 = 1027 = 4^{4+1} + 3$	$4 \rightsquigarrow 5$	$5^{5+1} + 3 = 15\,628$
$a_3 = 15\,627 = 5^{5+1} + 2$	$5 \rightsquigarrow 6$	$6^{6+1} + 2 = 279\,938$
$a_4 = 279\,937 = 6^{6+1} + 1$	$6 \rightsquigarrow 7$	$7^{7+1} + 1 = 5\,764\,802$

As one sees from this example,  $a_n$  initially increases enormously, and it is hardly believable that the sequence ever starts to decrease and ends in 0. But the proof of the theorem is not particularly difficult; one estimates  $a_n$  from above by the ordinal number  $\lambda_n$ , which, crudely put, results from  $a_n$  on replacing the basis  $b$  in  $(*)$  by  $\omega$ . With some ordinal arithmetic it can readily be shown that  $\lambda_{n+1} < \lambda_n$  as long as  $\lambda_n \neq 0$ . Since there is no properly decreasing infinite sequence of ordinal numbers (these are well-ordered), the sequence  $(a_n)_{n \in \mathbb{N}}$  must eventually end in 0. For more detailed information see for instance [HP].

Many metatheoretic properties can be expressed using the provability operator  $\Box$  in  $T$ , often using sentence schemata. The following ones turn out to be equivalent and facilitate a better understanding of the meaning of  $\neg \text{Con}_T$  within  $T$ . None of these properties hold for a consistent  $T$  from the outside (Theorem 6.5.1'), but all of them are provable in  $T = \text{PA}^\perp$ .

- (i)  $\neg \text{Con}_T : \Box \perp$  (provable inconsistency),
- (ii) SyComp :  $\Box \alpha \vee \Box \neg \alpha$  (syntactic completeness),
- (iii) SeComp :  $\alpha \rightarrow \Box \alpha$  (semantic completeness),
- (iv)  $\omega$ -Comp :  $\forall x \Box [\varphi(x)] \rightarrow \Box \forall x \varphi(x)$  ( $\omega$ -completeness).

**Theorem 3.4.** *The properties (i)–(iv) are all equivalent in a theory  $T$  satisfying the properties named at the beginning of this section.*

**Proof.** By (4) (i) $\Rightarrow$ (ii),(iii),(iv) are clear. (ii) $\Rightarrow$ (i): By Rosser's theorem formulated in  $T$  (see 7.5),  $\text{Con}_T \vdash_T \neg \Box \alpha \wedge \neg \Box \neg \alpha$  for some  $\alpha$ . Thus,  $\Box \alpha \vee \Box \neg \alpha \vdash_T \neg \text{Con}_T$ . (iii) $\Rightarrow$ (i): For  $\alpha := \text{Con}_T$ , SeComp and (2) yield  $\alpha \vdash_T \Box \alpha, \neg \Box \alpha$  and so  $\vdash_T \neg \alpha$ . (iv) $\Rightarrow$ (i): By (3) in 7.2, we obtain  $\neg \text{bew}_T(x, \ulcorner \perp \urcorner) \vdash_T \Box [\neg \text{bew}_T(x, \ulcorner \perp \urcorner)]$ , for  $\neg \text{bew}_T(x, \ulcorner \perp \urcorner)$  is  $\Sigma_1$ . Hence,

$$\text{Con}_T = \forall x \neg \text{bew}_T(x, \ulcorner \perp \urcorner) \vdash_T \forall x \Box [\neg \text{bew}_T(x, \ulcorner \perp \urcorner)].$$

$\omega$ -Comp and (2) yield  $\text{Con}_T \vdash_T \Box \forall x \neg \text{bew}(x, \ulcorner \perp \urcorner) = \Box \text{Con}_T \vdash_T \neg \text{Con}_T$ . Therefore,  $\vdash_T \neg \text{Con}_T$ .  $\square$

**Remark.**  $\text{Con}_T$  is also equivalent in  $T$  to other properties, for example to the schema  $\Box \alpha \rightarrow \alpha$  for  $\Pi_1$ -formulas  $\alpha$  (the *local*  $\Pi_1$ -reflection principle) as well as the *uniform*  $\Pi_1$ -reflection principle  $\forall x \Box [\alpha(x)] \rightarrow \forall x \alpha(x)$  for  $\Pi_1$ -formulas  $\alpha$ . Both the theorems of Paris–Harrington and of Goodstein are equivalent in PA to the uniform  $\Sigma_1$ -reflection, or equivalently, to the consistency of PA plus all true  $\Pi_1$ -sentences; see e.g. [Bar, D8].

Define inductively  $T^0 = T$  and  $T^{n+1} = T^n + \text{Con}_{T^n}$ . This *n-times-iterated consistency extension*  $T^n$  can be written as  $T^n = T + \neg \Box^n \perp$  with  $\Box = \text{bwb}_T$ ,  $\Box^0 \alpha = \alpha$  and  $\Box^{n+1} \alpha = \Box \Box^n \alpha$  (Exercise 3). Thus, the consistency of  $T^n$  can be expressed by an iterated consistency statement on  $T$ . Let  $T^\omega := \bigcup_{n \in \omega} T^n$ . Since  $T^n \subseteq T^{n+1}$  and  $T^n = T + \neg \Box^n \perp$  (hence  $T^\omega = T \cup \{\neg \Box^n \perp \mid n \in \omega\}$ ), the following three items are equivalent:

(i)  $T^\omega$  is consistent, (ii)  $T^n$  is consistent for all  $n$ , (iii)  $\not\vdash_T \Box^n \perp$  for all  $n$ .

Like  $\text{PA}^1 = \text{PA} + \text{Con}_{\text{PA}}$ , also  $\text{PA}^\omega$  conforms to truth looking at PA from outside. When considered more closely, this means only that  $\text{PA}^\omega$  is relatively consistent with respect to ZFC. In other terms,  $\vdash_{\text{ZFC}} \text{Con}_{\text{PA}^\omega}$ . The argument (to be formalized in ZFC) runs as follows:  $\vdash_{\text{PA}^\omega} \perp$  implies  $\vdash_{\text{PA}^n} \perp$  for some  $n$ , as was noticed above, hence  $\vdash_{\text{PA}} \Box^n \perp$ . But this is impossible, as is seen by a repeated application of  $D1^*$  (p. 271) on PA.

## Exercises

1. Prove  $D4^\circ$  for  $T$  by applying Theorem 3.2 to  $T' = T + \neg \alpha$ .
2. Show by means of Löb's theorem that  $\text{Con}_{\text{PA}} \rightarrow \neg \Box \neg \text{Con}_{\text{PA}}$  is unprovable in PA, although this formula is true if seen from outside.
3. Let  $T^n$  recursively be defined as in the text above. Prove that  $T^n = T + \neg \Box^n \perp$  and  $\text{Con}_{T^n} \equiv_T \neg \Box^{n+1} \perp$ , where  $\Box$  is  $\text{bwb}_T$ .
4. Show that  $\vdash_{\text{ZFC}} \Box_{\text{PA}} \alpha \rightarrow \alpha$  for all arithmetical sentences  $\alpha$  from  $\mathcal{L}_\epsilon$  (the  $\mathcal{L}_\epsilon$ -sentences relativized to  $\omega$ ).

## 7.4 The Provability Logic G

In 7.3 first-order logic was hardly required. It comes then as no surprise that many of the results there can be obtained propositionally, more precisely, in a certain modal propositional calculus. This calculus contains alongside  $\wedge, \neg$  the falsum symbol  $\perp$ , and a further unary connective  $\Box$  to be interpreted as the proof operator in  $\mathcal{L}_{ar}$ , denoted by  $\Box$  as well. First we define a propositional language  $\mathcal{F}_\Box$ , whose formulas are denoted by  $H, G, F$ : (a) the variables  $p_1, p_2, \dots$  from PV (page 4) and  $\perp$  belong to  $\mathcal{F}_\Box$ ; (b) if  $H, G$  belong to  $\mathcal{F}_\Box$  then so too  $(H \wedge G)$ ,  $\neg H$ , and  $\Box H$ .



No other strings belong to  $\mathcal{F}_\square$  in this context.  $H \vee G$ ,  $H \rightarrow G$ , and  $H \leftrightarrow G$  are defined as in 1.4,  $\top := \neg\perp$ . Further, set  $\diamond H := \neg\square\neg H$  and define recursively  $\square^0 H = H$ ,  $\square^{n+1} H = \square\square^n H$ . Let  $\mathbf{G}$  be the set of those formulas in  $\mathcal{F}_\square$  derivable using substitution in  $\mathcal{F}_\square$ , modus ponens MP, and the rule MN:  $H/\square H$  from the tautologies of two-valued propositional logic, augmented by the axioms (called also the G-axioms)

$$\square(p \rightarrow q) \rightarrow \square p \rightarrow \square q, \quad \square p \rightarrow \square\square p,^4 \quad \square(\square p \rightarrow p) \rightarrow \square p.$$

For  $H \in \mathbf{G}$  we mostly write  $\vdash_{\mathbf{G}} H$  (read “ $H$  is derivable in  $\mathbf{G}$ ”). Rule MN corresponds to *D1*. The first G-axiom reflects *D2*, the middle *D3*, and the last (called *Löb’s formula*) *D4*, hence the name *provability logic*. The connection between  $\mathbf{G}$  and PA is described in 7.5. Here we are concerned with the modal logic  $\mathbf{G}$  and its *Kripke semantics*. For simplicity, we restrict ourselves to finite *Kripke frames*, which are just finite directed graphs. We can do so, since all modal logics considered here have the finite model property. We begin without further ado with the following

**Definition.** A *G-frame* or *Kripke frame for G* is a finite poset  $(g, <)$ . A *valuation* is a mapping  $w$  that assigns to every variable  $p$  a subset  $wp$  of  $g$ . The relation  $P \Vdash H$ , dependent on  $w$ , between points  $P \in g$  and formulas  $H \in \mathcal{F}_\square$  (read “ $P$  accepts  $H$ ”) is defined inductively by

$$\begin{aligned} P \Vdash p &\text{ iff } P \in wp, & P \not\Vdash \perp, & P \Vdash H \wedge G &\text{ iff } P \Vdash H \ \&\ \& \ P \Vdash G, \\ P \Vdash \neg H &\text{ iff } P \not\Vdash H, & P \Vdash \square H &\text{ iff } P' \Vdash H \text{ for all } P' > P. \end{aligned}$$

These conditions easily imply  $P \Vdash \diamond H$  iff  $P' \Vdash H$  for some  $P' > P$ , and  $P \Vdash H \rightarrow G$  iff  $P \Vdash H \Rightarrow P \Vdash G$ . If  $P \Vdash H$  for all  $w$  and all  $P \in g$ , we write  $g \models H$  and say  $H$  holds in  $g$ . If  $g \models H$  for all  $\mathbf{G}$ -frames  $g$ , we write  $\models_{\mathbf{G}} H$  and say  $H$  is *G-valid*. The  $\mathbf{G}$ -frame on the right, consisting of two points  $P, P'$  with  $P < P'$ , shows that  $\not\models_{\mathbf{G}} p \rightarrow \square p$ .

Indeed, let  $wp = \{P\}$ . Then  $P \Vdash p$ , but  $P \not\Vdash \square p$  because  $P' \not\Vdash p$ . Note also  $\not\models_{\mathbf{G}} \square p \rightarrow p$ , for  $P' \not\Vdash p$  but  $P' \Vdash \square p$  because there is no  $P'' > P'$ .

We may tacitly assume that  $\mathbf{G}$ -frames are *initial* (have a smallest point), for  $g \models H$  is verified pointwise. We write  $H \equiv_{\mathbf{G}} H'$  for  $\models_{\mathbf{G}} H \leftrightarrow H'$ . It is readily seen that  $\equiv_{\mathbf{G}}$  is a congruence in  $\mathcal{F}_\square$  that extends the usual logical equivalence conservatively. For instance,  $\neg\square H \equiv_{\mathbf{G}} \neg\square\neg\neg H \equiv_{\mathbf{G}} \diamond\neg H$ . Many more equivalences are presented in the following examples. These will later be translated into statements about self-reference.

<sup>4</sup>This axiom is dispensable; it is provable from the remaining, see e.g. [Boo] or [Ral].

**Examples.** (a) Let  $g$  be an arbitrary  $\mathbf{G}$ -frame. Although always  $P \not\Vdash \perp$ , we have  $P \Vdash \Box\perp$ , provided  $P$  is maximal in  $g$ , that is, no  $Q > P$  exists. Likewise,  $\Box\neg\Box\perp$  is accepted precisely at the maximal points of  $g$ . Thus,  $\Box\perp \equiv_{\mathbf{G}} \Box\neg\Box\perp$ , or equivalently,  $\neg\Box\perp \equiv_{\mathbf{G}} \Diamond\Box\perp (= \neg\Box\neg\Box\perp)$ . This reflects in  $\mathbf{G}$  the second incompleteness theorem, as will be seen in 7.5.

(b) Let  $\{P_0, \dots, P_n\}$  be the ordered  $\mathbf{G}$ -frame with  $P_n < \dots < P_0$ . Clearly,  $P_0 \Vdash \Box^m\perp$  for each  $m > 0$ . Induction on  $n$  shows that  $P_n \Vdash \Box^m\perp$  for all  $m > n$ , but  $P_n \not\Vdash \Box^n\perp$ , and therefore  $P_n \not\Vdash \Box^{n+1}\perp \rightarrow \Box^n\perp$ . Hence,  $\not\equiv_{\mathbf{G}} \Box^{n+1}\perp \rightarrow \Box^n\perp$ , and a fortiori  $\not\equiv_{\mathbf{G}} \Box^n\perp$  and  $\not\equiv_{\mathbf{G}} \neg\Box^{n+1}\perp$ , for all  $n$ .

(c)  $\vDash_{\mathbf{G}} \Box(\Box p \rightarrow p) \rightarrow \Box p$ . For take an arbitrary  $g$  and  $P \in g$ . If  $P \not\Vdash \Box p$  then there is, since  $g$  is finite, some  $Q > P$  with  $Q \Vdash \neg p$  and  $Q' \Vdash p$  for all  $Q' > Q$ . Thus  $Q \Vdash \Box p$ ; hence  $Q \not\Vdash \Box p \rightarrow p$  and so  $P \not\Vdash \Box(\Box p \rightarrow p)$ . Consequently,  $P \Vdash \Box(\Box p \rightarrow p) \rightarrow \Box p$ , which proves our claim. Note also that  $\vDash_{\mathbf{G}} \Box p \rightarrow \Box\Box p$ . Only the transitivity of  $<$  is relevant for the proof.

(d)  $\vDash_{\mathbf{G}} \neg\Box^{n+1}\perp \rightarrow \Diamond R_n$ , where  $R_n := \bigwedge_{i=1}^n (\Box p_i \rightarrow p_i)$ . For let  $P \in g$ ,  $P \Vdash \neg\Box^{n+1}\perp$ . Then there must be a chain  $P = P_0 < P_1 < \dots < P_{n+1}$  in  $g$ . Now, it is a nice separate exercise to verify that each conjunct of  $R_n$  fails to be accepted by at most one of the  $n+1$  points  $P_1, \dots, P_{n+1}$ . Thus, at least one of these accepts all conjuncts. In other words,  $P_i \Vdash R_n$  for some  $i > 0$ ; hence  $P \Vdash \Diamond R_n$ . This nontrivial example will essentially be employed in the proof of Theorem 7.1.

By induction on  $\vdash_{\mathbf{G}} H$  one easily proves  $\vdash_{\mathbf{G}} H \Rightarrow \vDash_{\mathbf{G}} H$  (soundness of Kripke semantics for  $\vdash_{\mathbf{G}}$ ). Example (c) is a part of the initial step. The induction steps over the rules are easy. For instance,  $g \vDash H$  clearly implies  $g \vDash \Box H$ . The converse,  $\vDash_{\mathbf{G}} H \Rightarrow \vdash_{\mathbf{G}} H$ , holds as well. Thus,  $\vdash_{\mathbf{G}} H$  can be confirmed by proving  $\vDash_{\mathbf{G}} H$ , and vice versa. This is the content of

**Theorem 4.1 (Completeness of Kripke semantics for  $\mathbf{G}$ ).** *For each formula  $H$  from  $\mathcal{F}_{\Box}$  it holds that  $\vdash_{\mathbf{G}} H \Leftrightarrow \vDash_{\mathbf{G}} H$ .*

The nontrivial direction  $\Leftarrow$  follows directly from the finite model property of  $\mathbf{G}$ , i.e., each  $H \notin \mathbf{G}$  is falsified or refuted by some finite  $\mathbf{G}$ -frame, proved, for example, in [Boo], [Ra1], and [CZ]. For the relatively simple formulas considered here,  $\vDash_{\mathbf{G}} H$  is in general more easily checked than  $\vdash_{\mathbf{G}} H$ .

Both the formulas provable in  $\mathbf{G}$  and those refutable are clearly recursively enumerable, thanks to the finite model property of  $\mathbf{G}$ . Thus, in analogy to Exercise 2 in 3.6, we obtain

**Theorem 4.2.**  $\mathbf{G}$  is decidable.

**Remark.** The finite model property, decidability, and some other properties such as interpolation can all be proved in one move, see e.g. [Ra2]. An important fragment of  $\mathbf{G}$  is  $\mathbf{G}^0 := \mathbf{G} \cap \mathcal{F}_\square^0$ , where  $\mathcal{F}_\square^0$  denotes the set of variable-free formulas of  $\mathcal{F}_\square$ . The formulas  $\neg \square^n \perp$  ( $\equiv_{\mathbf{G}} \diamond^n \top$ ) form a Boolean base in  $\mathbf{G}^0$ . One proves this most easily by showing that  $\mathbf{G}^0$  is complete with respect to all (totally) ordered  $\mathbf{G}$ -frames, including the infinite ones, and applying Theorem 5.2.3 accordingly.

### Exercises

1. Let  $g$  be any finite Kripke frame (a graph) that satisfies the axioms of  $\mathbf{G}$ . Show that  $g$  is necessarily a poset. Only this fact justifies the identification of  $\mathbf{G}$ -frames with posets.
2. Prove  $\vdash_{\mathbf{G}} \square p \rightarrow \square(\square p \rightarrow p)$ , the inverse of Löb's formula. (Only the first of the three  $\mathbf{G}$ -axioms is needed in the proof.)

## 7.5 The Modal Treatment of Self-Reference

Let  $T$  be a theory as in 7.3. A mapping  $\iota$  from  $PV$  to  $\mathcal{L}^0$  with  $p_i^\iota = \alpha_i$  is called an *insertion*.  $\iota$  can be extended to the whole of  $\mathcal{F}_\square$  by the clauses  $\perp^\iota = \perp$ ,  $(\neg H)^\iota = \neg H^\iota$ ,  $(H \wedge G)^\iota = H^\iota \wedge G^\iota$ , and  $(\square H)^\iota = \square H^\iota$  ( $= \mathbf{bwb}_T(\ulcorner H^\iota \urcorner)$ ). Briefly speaking,  $H^\iota$  results from  $H(p_1, \dots, p_n)$  by replacing the  $p_\nu$  by the sentences  $\alpha_\nu$  from  $\mathcal{L}$ . For instance, if  $p^\iota = \alpha$  then  $(\square p \wedge \neg \square \perp)^\iota = \square \alpha \wedge \neg \square \perp$ , and  $(\neg \square \perp)^\iota = \neg \square \perp = \mathbf{Con}_T$ . The following lemma shows that  $\vdash_{\mathbf{G}}$  is “sound” for  $\vdash_T$ . Already this simple fact considerably simplifies proofs about self-referential statements.

**Lemma 5.1.** For each  $H$  with  $\vdash_{\mathbf{G}} H$  and each insertion  $\iota$ ,  $\vdash_T H^\iota$ .

**Proof** by induction on  $\vdash_{\mathbf{G}} H$ . If  $H$  is a propositional tautology then  $H^\iota \in \mathbf{Taut}_{\mathcal{L}} \subseteq T$ . If  $H$  is one of the modal axioms of  $\mathbf{G}$ , then  $\vdash_T H^\iota$  by D2, D3, or D4. If  $\vdash_{\mathbf{G}} H$  and  $\sigma: \mathcal{F}_\square \rightarrow \mathcal{F}_\square$  is a substitution, then  $\vdash_T H^{\sigma^\iota}$ , since  $H^{\sigma^\iota} = H^{\iota'}$  with  $\iota': p \mapsto p^{\sigma^\iota}$ , and  $\vdash_T H^{\iota'}$  holds by the induction hypothesis. As regards the induction step over MP, consider  $(F \rightarrow G)^\iota = F^\iota \rightarrow G^\iota$ . Finally, if MN is applied, and  $\vdash_T H^\iota$  by the induction hypothesis, then  $\vdash_T \square H^\iota = (\square H)^\iota$ , due to D1.  $\square$

**Example 1.** We prove (3) of Theorem 3.2 with the calculus  $\vdash_{\mathbf{G}}$ . By Lemma 5.1 and Theorem 4.1 it suffices to show that  $\vDash_{\mathbf{G}} \neg \square \perp \leftrightarrow \neg \square \neg \square \perp$ .

This holds by Example (a) in 7.4. Next example:  $\vDash_G \Box(p \leftrightarrow \Diamond p) \rightarrow \neg \Diamond p$  is easily confirmed. Thus,  $\vdash_T \Box(\alpha \leftrightarrow \Diamond \alpha) \rightarrow \neg \Diamond \alpha$ . This tells us (if everything is related to  $T = \text{PA}$ ) that a sentence claiming its own consistency with PA is incompatible with PA, which hardly seems plausible. Even the converse is provable in PA since  $\vDash_G \neg \Diamond p \rightarrow \Box(p \leftrightarrow \Diamond p)$ .

We now explain certain facts that expand upon the reasoning of above. For PA and related theories, the converse of Lemma 5.1 holds as well. That is to say, the derivability conditions and Löb's theorem already contain everything worth knowing about self-referential formulas or schemes. This is essentially the content of Theorem 5.2. For the subtle proofs of Theorems 5.2, 5.4, and 5.5, the reader is referred to [Boo].

**Theorem 5.2 (Solovay's completeness theorem).** *For all  $H \in \mathcal{F}_\Box$ :  $\vdash_G H$  (equivalently  $\vDash_G H$ ) if and only if  $\vdash_{\text{PA}} H^i$  for all insertions  $i$ .*

**Example 2 (applications).** (a)  $\not\vdash_{\text{PA}} \Box^{n+1} \perp \rightarrow \Box^n \perp$  because by Example (b) in 7.4,  $\not\vdash_G \Box^{n+1} \perp \rightarrow \Box^n \perp$ . In particular,  $\not\vdash_{\text{PA}} \text{Con}_{\text{PA}} (\equiv \Box \perp \rightarrow \perp)$ . (b)  $\not\vdash_{\text{PA}} \neg \Box^{n+1} \perp$ , since  $\not\vdash_G \neg \Box^{n+1} \perp$ . (c) It is easily verified with the 2-point frame on page 285 that  $\not\vdash_G \neg \Box p \rightarrow \Box \neg \Box p$ , in particular  $\not\vdash_G \neg \Box \perp \rightarrow \Box \neg \Box \perp$ . Therefore,  $\not\vdash_{\text{PA}} \text{Con}_{\text{PA}} \rightarrow \Box \text{Con}_{\text{PA}}$ . (d)  $\text{PA}_n := \text{PA} + \Box^n \perp$  is consistent for  $n > 0$  by (b), but is  $\omega$ -inconsistent. Otherwise, by  $D1^*$  (page 271),  $\vdash_{\text{PA}_n} \Box^n \perp \Rightarrow \vdash_{\text{PA}_n} \Box^{n-1} \perp \Rightarrow \dots \Rightarrow \vdash_{\text{PA}_n} \perp$ , contradicting  $\not\vdash_{\text{PA}_n} \perp$ . Since  $\vdash_{\text{PA}} \Box^n \perp \rightarrow \Box^{n+1} \perp$  by  $D3$ , we get  $\text{PA}_n \supseteq \text{PA}_{n+1}$ , and since  $\text{PA}_n \neq \text{PA}_{n+1}$  by (a), we have  $\text{PA}_0 \supset \text{PA}_1 \supset \dots \supset \text{PA}$ . Observe that  $\text{PA}_1$  is just  $\text{PA}^\perp$ .

Note also the following: Since  $\not\vdash_G \Box p \rightarrow p$ , there must be some  $\alpha \in \mathcal{L}_{ar}^0$  such that  $\not\vdash_{\text{PA}} \Box \alpha \rightarrow \alpha$ . Indeed, choose  $\alpha = \perp$ . The above examples point out that Theorem 5.2 and the decidability of G are very efficient tools in deciding the provability of self-referential statements.

Many other theories have the same provability logic as PA, where in general a modal propositional logic H is the *provability logic for T* when the analogue of Theorem 5.2 holds with respect to T and H. For some theories, the provability logic may be a proper extension of G. For example, the  $\omega$ -inconsistent theory  $\text{PA}_n$  from Example 2(d) has the provability logic  $\text{G}_n := \text{G} + \Box^n \perp$ , the smallest extension of G closed under all rules of G with the additional axiom  $\Box^n \perp$  (Exercise 1; note that  $\text{G}_0$  is inconsistent). By the following theorem, which will be proved in 7.7, other extensions of G to be considered as provability logics are out of the question.

**Theorem 5.3 ([Vil]).** *Let  $T$  be at least as strong as PA. Then*

- (a) *If  $T^\omega$  (page 284) is consistent, then  $\mathsf{G}$  is the provability logic of  $T$ ;*
- (b) *if  $\vdash_{T^\omega} \perp$  and  $n$  is minimal such that  $\vdash_{T^n} \perp$ , then  $T$ 's provability logic is  $\mathsf{G}_n$ .*

The formulas  $H \in \mathcal{F}_\square$  such that  $\mathcal{N} \vDash H^i$  for all insertions  $i$  in  $\mathcal{L}_{ar}$  can also be surprisingly easily characterized. All  $H \in \mathsf{G}$  are obviously included; but in addition also  $\Box p \rightarrow p$  belongs to this sort of formula, because  $\mathcal{N} \vDash \Box \alpha \rightarrow \alpha$  for  $\alpha \in \mathcal{L}_{ar}^0$ . Indeed, if  $\mathcal{N} \vDash \Box \alpha$  then *there is* some  $n$  that codes a proof of  $\alpha$  in PA, hence  $\mathcal{N} \vDash \alpha$ .

Let  $\mathsf{GS} (\supseteq \mathsf{G})$  be the set of all formulas in  $\mathcal{F}_\square$  that can be obtained from those in  $\mathsf{G} \cup \{\Box p \rightarrow p\}$  using substitution and modus ponens only. Induction in  $\mathsf{GS}$  readily yields  $H \in \mathsf{GS} \Rightarrow \mathcal{N} \vDash H^i$  for all  $i$ . Again, the converse holds as well:

**Theorem 5.4 ([So]).**  *$H \in \mathsf{GS}$  if and only if  $\mathcal{N} \vDash H^i$  for all insertions  $i$ .*

$\mathsf{GS}$  is decidable as well, because it can be shown that  $H \in \mathsf{GS} \Leftrightarrow H^* \in \mathsf{G}$ , where  $H^* := [\bigwedge_{\Box G \in \mathsf{Sf}^\square} (\Box G \rightarrow G)] \rightarrow H$ . Here  $\mathsf{Sf}^\square H$  is the set of subformulas of  $H$  of the form  $\Box G$ . Thus, Theorem 5.4 reduces the decidability of  $\mathsf{GS}$  to that of  $\mathsf{G}$ . Using this theorem, many questions concerning the relations between *provable* and *true* are effectively decidable. For instance,

$$H(p) := \neg \Box (\neg \Box \perp \rightarrow \neg \Box p \wedge \neg \Box \neg p) \notin \mathsf{GS}$$

is readily verified. Hence  $\mathcal{N} \vDash \neg H(\alpha) \equiv \Box (\neg \Box \perp \rightarrow \neg \Box \alpha \wedge \neg \Box \neg \alpha)$  for some  $\alpha \in \mathcal{L}_{ar}^0$  by Theorem 5.4. Translated into English: *It is provable in PA that the consistency of PA implies the independence of  $\alpha$  for some sentence  $\alpha$ .* This is exactly Rosser's theorem, which in this way turns out to be provable in PA. As was shown in [Be1], the box in the formulas  $H \in \mathsf{GS}$  in Theorem 5.4 may denote  $\mathsf{bwb}_T$  for *any* axiomatizable  $T \supseteq \text{PA}$ , provided  $T \subseteq \text{Th}\mathcal{N}$ . However, if  $T$  proves false sentences (as does e.g.  $\text{PA}^\perp$ ) then  $\mathsf{GS}$  has to be redefined in a feasible manner and is always decidable.

A variable  $p$  in  $H$  is called *modalized in  $H$*  if every occurrence of  $p$  is contained within the scope of a  $\Box$ , as is the case in  $\neg \Box p$ ,  $\neg \Box \neg p$ , and  $\Box(p \rightarrow q)$ . By contrast,  $p$  is not modalized in  $\Box p \rightarrow p$ . Another particularly interesting theorem is

**Theorem 5.5 (DeJongh–Sambin fixed point theorem).** *Let  $p$  be modalized in  $H(p, q_1, \dots, q_n)$ ,  $n \geq 0$ . Then a formula  $F = F(\vec{q})$  from  $\mathcal{F}_\square$  can effectively be constructed such that*

- (a)  $F \equiv_{\mathcal{G}} H(F, \vec{q})$ ,
- (b)  $\vdash_{\mathcal{G}} \bigwedge_{i=1}^2 [(p_i \leftrightarrow H(p_i, \vec{q})) \wedge \square(p_i \leftrightarrow H(p_i, \vec{q}))] \rightarrow (p_1 \leftrightarrow p_2)$ .

This theorem easily yields a corresponding result for theories  $T$ :

**Corollary 5.6.** *Let  $p$  be modalized in  $H = H(p, \vec{q})$  and suppose  $T$  satisfies D1–D4. Then there is an  $F = F(\vec{q}) \in \mathcal{F}_\square$  with  $F(\vec{\alpha}) \equiv_T H(F(\vec{\alpha}), \vec{\alpha})$  for all  $\vec{\alpha} = (\alpha_1, \dots, \alpha_n)$ ,  $\alpha_i \in \mathcal{L}^0$ . For each  $\vec{\alpha}$  there is only one  $\beta \in \mathcal{L}^0$  modulo  $T$  such that  $\beta \equiv_T H(\beta, \vec{\alpha})$ .*

**Proof.** Choose  $F$  as in (a) of the theorem. Then  $F(\vec{\alpha}) \equiv_T H(F(\vec{\alpha}), \vec{\alpha})$  by Lemma 5.1 ( $\vec{q}^n = \vec{\alpha}$ ). To prove uniqueness let  $\beta_i \equiv_T H(\beta_i, \vec{\alpha})$  for  $i = 1, 2$ . By D1,  $\vdash_T (\beta_i \leftrightarrow H(\beta_i, \vec{\alpha})) \wedge \square(\beta_i \leftrightarrow H(\beta_i, \vec{\alpha}))$ . Inserting  $\beta_i$  for  $p_i$  and  $\alpha_i$  for  $q_i$  in the formula under (b) in the theorem then yields  $\vdash_T \beta_1 \leftrightarrow \beta_2$  by Lemma 5.1.  $\square$

**Example 3.** For  $H = \neg\square p$  ( $n = 0$ ),  $F = \neg\square\perp$  is a “solution” of (a) in Theorem 5.5 because  $\neg\square\perp \equiv_{\mathcal{G}} \neg\square(\neg\square\perp)$ . According to Corollary 5.6,  $\text{Con}_T (= \neg\square\perp)$  is modulo  $T$  the only fixed point of  $\neg\text{bwb}_T$ . This is just the claim of (3) from 7.3.

Many special cases of the corollary represent older self-reference results from Gödel, Löb, Rogers, Jeroslow, and Kreisel, which, stated in terms of modal logic, concern fixed points of  $\neg\square p$ ,  $\square p$ ,  $\neg\square\neg p$ ,  $\square\neg p$ , and  $\square(p \rightarrow q)$  in PA. Incidentally, one gets the fixed points of these formulas—namely  $\neg\square\perp$ ,  $\top$ ,  $\perp$ ,  $\square\perp$ , and  $\square q$ —according to a simple recipe. All first listed formulas are of the form  $H = G \frac{\square H'}{p}$ , where  $p$  is not modalized in  $G(p, \vec{q})$  and  $H'(p, \vec{q})$  is chosen appropriately. In this case,  $F = H \frac{G(\top, \vec{q})}{p}$  is the fixed point of  $H$ , as is seen after some calculation. For  $H = \neg\square p$  from Example 3 is  $G = \neg p$ . Thus, according to the recipe the fixed point is

$$F = \neg\square p \frac{\neg\top}{p} = \neg\square\neg\top \equiv_{\mathcal{G}} \neg\square\perp.$$

For Kreisel’s formula  $\square(p \rightarrow q)$  is  $G = p$ . Hence, it has the fixed point

$$F = \square(p \rightarrow q) \frac{\top}{p} = \square(\top \rightarrow q) \equiv_{\mathcal{G}} \square q.$$

The recipe also works for  $H = \square p \rightarrow q$ , by choosing  $G = p \rightarrow q$ . Hence  $F = (\square p \rightarrow q) \frac{\top \rightarrow q}{p} = \square(\top \rightarrow q) \rightarrow q \equiv_{\mathcal{G}} \square q \rightarrow q$  is the only fixed point of

$H$  modulo  $T$ . Exactly this is the claim of Lemma 3.1(b), used in Gödel's second incompleteness theorem.

### Exercises

1. Prove that the theory  $\text{PA}_n$  from Example 2(d) has the provability logic  $\mathbf{G}_n$ .
2. Show that  $\text{PA}_\perp^n := \text{PA}^n + \neg \text{Con}_{\text{PA}^n}$  equals  $\text{PA} + \Box^{n+1}\perp \wedge \neg\Box^n\perp$  and that it has the provability logic  $\mathbf{G}_1 = \mathbf{G} + \Box\perp$ . Here  $\Box$  means  $\Box_{\text{PA}}$ .
3. Prove that  $\top$ ,  $\perp$ , and  $\Box\perp$  are the fixed points of  $\Box p$ ,  $\neg\Box\neg p$ , and  $\Box\neg p$ .
4. (Mostowski). Let  $T \supseteq \text{PA}$  be axiomatizable and suppose  $\mathcal{N} \models T$ . Show that there are two mutually independent  $\Sigma_1$ -sentences  $\alpha, \beta$  in  $T$ , that is,  $\alpha \rightarrow \beta$ ,  $\alpha \rightarrow \neg\beta$ ,  $\beta \rightarrow \alpha$ ,  $\beta \rightarrow \neg\alpha$  (hence also  $\alpha$ ,  $\beta$ ,  $\neg\alpha$ , and  $\neg\beta$ ) are unprovable in  $T$ .

## 7.6 A Bimodal Provability Logic for PA

Hilbert remarked jokingly that the incompleteness phenomenon can be forcefully removed from the world by use of the so-called  $\omega$ -rule

$$\rho_\omega : \frac{X \vdash \varphi(n) \text{ for all } n}{X \vdash \forall x \varphi}.$$

$\rho_\omega$  has infinitely many premises. It is an easy exercise to derive with the aid of  $\rho_\omega$  every sentence  $\alpha$  valid in  $\mathcal{N}$  from the axioms of PA, even from those of Q. Indeed, all sentences can (up to equivalence) be obtained from variable-free literals with  $\wedge, \vee, \forall, \exists$ , bypassing formulas with free variables. Due to the  $\Sigma_1$ -completeness of Q, all valid variable-free literals are derivable. The inductive steps for  $\wedge, \vee, \exists$  are simple, applying  $\Sigma_1$ -completeness in the  $\exists$ -step once again. Only in the  $\forall$ -step is  $\rho_\omega$  used.

Clearly, an unrestricted use of the infinitistic rule  $\rho_\omega$  (in spite of its relevance for higher order arithmetic) contradicts Hilbert's own intention of giving mathematics a finitistic foundation. However, things look different if we restrict  $\rho_\omega$  each time to a *single* application. In view of Remark 1 in 6.2, we no longer distinguish between  $\varphi$  and  $\dot{\varphi}$ , so that  $\varphi$  itself is a number and  $\ulcorner \varphi \urcorner = \underline{\varphi}$  is the corresponding Gödel term. Let us define

$$1bwb_{\text{PA}}(\alpha) := (\exists \varphi \in \mathcal{L}_{ar}^1)[bwb_{\text{PA}}(\forall x \varphi \rightarrow \alpha) \ \& \ \forall n \ bwb_{\text{PA}}(\varphi(\underline{n}))].$$

$1bwb_{\text{PA}}$  is arithmetical, in fact it is  $\Sigma_3$ , for  $bwb_{\text{PA}}$  is  $\Sigma_1$  and  $\forall n \ bwb_{\text{PA}}(\varphi(\underline{n}))$  is  $\Pi_1$ . We read  $1bwb_{\text{PA}}(\alpha)$  as “ $\alpha$  is 1-provable.” Let  $\mathbf{1bwb}(x)$  be the  $\Sigma_3$ -formula in  $\mathcal{L}_{ar}$  defining  $1bwb_{\text{PA}}$ . Here let  $x$  be  $\mathbf{v}_0$ . Write  $\Box\alpha$  for  $\mathbf{1bwb}(\ulcorner \alpha \urcorner)$  and  $\Diamond\alpha$  for  $\neg\Box\neg\alpha$ . Clearly,  $\Box\alpha$  for  $\alpha \in \mathcal{L}_{ar}^0$  ( $\Box = \Box_{\text{PA}}$ ) can be read ‘ $\text{PA} + \neg\alpha$  is inconsistent’, while  $\Box\alpha$ , by Lemma 6.1, formalizes ‘ $\text{PA} + \neg\alpha$  is  $\omega$ -inconsistent’. Thus,  $\Diamond\top$  ( $\equiv \neg\Box\perp$ ) means ‘ $\text{PA}$  ( $= \text{PA} + \neg\perp$ ) is  $\omega$ -consistent’. This explains the interest in the operator  $\Box$ .

If  $bwb_{\text{PA}}(\alpha)$  then certainly  $1bwb_{\text{PA}}(\alpha)$  (choose  $\alpha$  for  $\varphi$ ). The italicized statement is reflected in  $\text{PA}$  as ‘ $\vdash_{\text{PA}} \Box\alpha \rightarrow \ulcorner \alpha \urcorner$  for every  $\alpha \in \mathcal{L}_{ar}^0$ ’. The converse fails, since  $\not\vdash_{\text{PA}} \text{Con}_{\text{PA}}$ , while  $\text{Con}_{\text{PA}}$  is easily 1-provable:  $\vdash_{\text{PA}} \varphi(\underline{n})$  for all  $n$ , with  $\varphi(x) := \neg \text{bew}_{\text{PA}}(x, \perp)$ , and trivially  $\vdash_{\text{PA}} \forall x \varphi(x) \rightarrow \text{Con}_{\text{PA}}$ . In what follows, some claims will not be proved in detail.

Define  $\Omega := \{\varphi \in \mathcal{L}_{ar}^1 \mid \vdash_{\text{PA}} \varphi(\underline{n}) \text{ for all } n\}$ . By its definition,  $\Omega$  and hence also  $\text{PA}^\Omega := \text{PA} + \Omega$  are formally  $\Sigma_3$ . As Theorem 6.2 will show,  $\text{PA}^\Omega$  is properly  $\Sigma_3$  and hence is no longer recursively axiomatizable.

**Lemma 6.1.** *The following properties are equivalent for  $\alpha \in \mathcal{L}_{ar}^0$ :*

- (i)  $1bwb_{\text{PA}}(\alpha)$ , (ii)  $\vdash_{\text{PA}^\Omega} \alpha$ , (iii)  $\text{PA} + \neg\alpha$  is  $\omega$ -inconsistent.

**Proof.** (i) $\Rightarrow$ (ii) follows with a glance at the definitions (read (i) naively). (ii) $\Rightarrow$ (iii): Let  $\vdash_{\text{PA}^\Omega} \alpha$ . Since  $\Omega$  is closed under conjunctions, there is some  $\forall x \varphi(x) \in \Omega$  with  $\forall x \varphi \vdash_{\text{PA}} \alpha$ , hence  $\vdash_{\text{PA}} \neg\alpha \rightarrow \exists x \neg\varphi$  and so  $\vdash_{\text{PA} + \neg\alpha} \exists x \neg\varphi$ . Now,  $\forall x \varphi \in \Omega$ , therefore  $\vdash_{\text{PA}} \varphi(\underline{n})$  and a fortiori  $\vdash_{\text{PA} + \neg\alpha} \varphi(\underline{n})$ , for all  $n$ . Thus,  $\text{PA} + \neg\alpha$  is  $\omega$ -inconsistent. (iii) $\Rightarrow$ (i): Let  $\vdash_{\text{PA} + \neg\alpha} \beta(\underline{n})$  for all  $n$ , but  $\vdash_{\text{PA} + \neg\alpha} \exists x \neg\beta$ . Then  $\vdash_{\text{PA}} \forall x \beta \rightarrow \alpha$ . With  $\varphi(x) := \neg\alpha \rightarrow \beta(x)$  clearly  $\vdash_{\text{PA}} \varphi(\underline{n})$  for all  $n$ . Now,  $\forall x \varphi \equiv \alpha \vee \forall x \beta \vdash_{\text{PA}} \alpha$ . Hence  $\vdash_{\text{PA}} \forall x \varphi \rightarrow \alpha$ . Thus, altogether  $1bwb_{\text{PA}}(\alpha)$ .  $\square$

**Theorem 6.2 (the 1-provable  $\Sigma_3$ -completeness of  $\text{PA}$ ).** *All true  $\Sigma_3$ -sentences are 1-provable. Moreover, for every  $\beta$  of this kind,  $\vdash_{\text{PA}} \beta \rightarrow \Box\beta$ .*

**Proof.** Let  $\mathcal{N} \models \beta := \exists y \forall x \gamma(y, x)$  where  $\gamma(y, x)$  is  $\Sigma_1$ . Then there is some  $m$  such that  $\mathcal{N} \models \gamma(\underline{m}, \underline{n})$  for all  $n$ . Therefore,  $\vdash_{\text{PA}} \gamma(\underline{m}, \underline{n})$  for all  $n$ , because  $\text{PA}$  is  $\Sigma_1$ -complete. Hence,  $\forall x \gamma(\underline{m}, x) \in \Omega$  and so  $\vdash_{\text{PA}^\Omega} \exists z \forall x \gamma$ , or equivalently,  $1bwb_{\text{PA}}(\beta)$  by Lemma 6.1. Because of the provable  $\Sigma_1$ -completeness of  $\text{PA}$ , this argumentation is comprehensible in  $\text{PA}$ , so that also  $\vdash_{\text{PA}} \beta \rightarrow \Box\beta$ .  $\square$



$D1$ – $D4$  are also valid for the operator  $\boxplus$ :  $\mathcal{L}_{ar}^0 \rightarrow \mathcal{L}_{ar}^0$ . Indeed,  $D1$  holds because  $\vdash_{\text{PA}} \alpha \Rightarrow \vdash_{\text{PA}} \Box\alpha \Rightarrow \vdash_{\text{PA}} \boxplus\alpha$ , and  $D2$  formalizes (or reflects) ‘ $\vdash_{\text{PA}\Omega} \alpha, \alpha \rightarrow \beta \Rightarrow \vdash_{\text{PA}\Omega} \beta$ ’ in PA (observe Lemma 6.1).  $D3$  follows from Theorem 6.2 with  $\beta = \boxplus\alpha$ . The proof of  $D4$  in 7.3 uses, along with the fixed point lemma, only  $D1$ – $D3$ ; so  $D4$  holds as well. Therefore, nearly everything said in 7.3 on  $\Box$  applies also to  $\boxplus$ , including Theorem 3.2, which now reads  $\not\vdash_{\text{PA}} \neg\boxplus\perp$  ( $\equiv \Diamond\top$ ). To put it more concisely, although the consistency of PA is provable with the extended means,  $\omega$ -consistency is not. Hence, this property, which is  $\Pi_3$ -definable according to Exercise 3 in 6.7, cannot be  $\Sigma_3$  by Theorem 6.2, and must therefore be properly  $\Pi_3$ . Equivalently,  $\omega$ -inconsistency is properly  $\Sigma_3$ .

Alongside  $\Box\alpha \rightarrow \boxplus\alpha$ , there are other noteworthy interactions between  $\Box$  and  $\boxplus$ , in particular  $\vdash_{\text{PA}} \neg\Box\alpha \rightarrow \boxplus\neg\Box\alpha$ . This formalizes ‘If  $\not\vdash_{\text{PA}} \alpha$  then  $\neg\Box\alpha$  is 1-provable’. To verify the latter notice that  $\not\vdash_{\text{PA}} \alpha$  implies  $\vdash_{\text{PA}} \varphi(\underline{n})$  for all  $n$ , where  $\varphi(x)$  is  $\neg\text{bew}_{\text{PA}}(x, \ulcorner\alpha\urcorner)$ , and since  $\vdash_{\text{PA}} \forall x\varphi \rightarrow \neg\Box\alpha$ , we get  $\vdash_{\text{PA}} \boxplus\neg\Box\alpha$ . On the other hand,  $\vdash_{\text{PA}} \neg\Box\alpha \rightarrow \Box\neg\Box\alpha$  fails in general; Example 2(c) in 7.5 yields a counterexample.

The language of the bimodal propositional logic GD now to be defined results from  $\mathcal{F}_{\Box}$  by adding a further connective  $\boxplus$  to  $\mathcal{F}_{\Box}$ , which is treated syntactically just as  $\Box$ . The axioms of GD are those of G stated both for  $\Box$  and  $\boxplus$ , augmented by the axioms

$$\Box p \rightarrow \boxplus p \quad \text{and} \quad \neg\Box p \rightarrow \boxplus\neg\Box p.$$

The rules of GD are the same as those for G. Insertions  $\iota$  to  $\mathcal{L}_{ar}^0$  are defined as in 7.5, but with the additional clause  $(\boxplus H)^\iota = \boxplus H^\iota$ , that is,  $(\boxplus H)^\iota = \text{1wb}(\ulcorner H^\iota \urcorner)$ . By the reasoning above, all axioms and rules of GD are sound. This proves (the easier) half of the following remarkable theorem from Dzhaparidze (1985):

**Theorem 6.3.**  $\vdash_{\text{GD}} H \Leftrightarrow \vdash_{\text{PA}} H^\iota$  for all insertions  $\iota$  as defined above. Furthermore, GD is decidable.

Thus, the modal system GD completely captures the interaction between  $\text{bwb}_{\text{PA}}$  and  $\text{1wb}_{\text{PA}}$ ; also Theorem 5.5 carries over. However, GD no longer has an adequate Kripke semantics, which complicates the decision procedure. For further references see [Boo] or [Be3].

As an exercise, the reader should derive  $\boxplus(\Box p \rightarrow p)$  from the axioms of GD. Thus,  $\vdash_{\text{PA}} \boxplus(\Box\alpha \rightarrow \alpha)$  for every  $\alpha \in \mathcal{L}_{ar}^0$ , while  $\vdash_{\text{PA}} \Box(\Box\alpha \rightarrow \alpha)$  is the

case only provided  $\vdash_{\text{PA}} \alpha$ . In other words, the *local reflection principle*  $\{\Box\alpha \rightarrow \alpha \mid \alpha \in \mathcal{L}_{ar}^0\}$  is 1-provable in PA. **Be careful:** GD expands G conservatively, so that  $\not\vdash_{\text{GD}} \Box p \rightarrow p$ .

## 7.7 Modal Operators in ZFC

Considerations regarding self-reference in ZFC are technically sometimes easier, but from the foundational point of view more involved because there is no superordinate theory. If ZFC is consistent, as we assume it is, then  $\text{Con}_{\text{ZFC}}$  is a true arithmetical statement that is unprovable in ZFC. Thus, true arithmetical statements may even be unprovable in ZFC, not only in PA or similarly strong arithmetical theories. It makes sense, therefore, to consider  $\text{ZFC}^+ := \text{ZFC} + \text{Con}_{\text{ZFC}}$ , because after all, we want set theory to embrace as many facts about numbers and sets as possible from which interesting consequences may result.

As 7.3 shows, the consistency of ZFC alone does not guarantee that  $\text{ZFC}^+$  is consistent. The second incompleteness theorem clearly excludes  $\vdash_{\text{ZFC}} \text{Con}_{\text{ZFC}}$  but does not preclude  $\vdash_{\text{ZFC}} \text{Con}_{\text{ZFC}} \rightarrow \text{Con}_{\text{ZFC}^+}$ . In this case  $\vdash_{\text{ZFC}^+} \text{Con}_{\text{ZFC}^+}$ , and so  $\vdash_{\text{ZFC}^+} \perp$  by the same theorem. On the other hand, from certain assumptions about the existence of large cardinals, the consistency of  $\text{ZFC}^+$  readily follows. These assumptions would have to be jettisoned in case  $\vdash_{\text{ZFC}^+} \perp$ , i.e.  $\vdash_{\text{ZFC}} \neg \text{Con}_{\text{ZFC}}$ . Moreover, the consistency of ZFC would then not correctly be reflected in ZFC, and ZFC proves along with true arithmetical facts also false ones. This sounds strange, but there is hardly a convincing argument that this cannot be so.

Even if  $\text{ZFC}^+$  is consistent, i.e.  $\not\vdash_{\text{ZFC}} \neg \text{Con}_{\text{ZFC}}$ , it may still be that one of the sentences from the sequence  $\Box \neg \text{Con}_{\text{ZFC}}, \Box \Box \neg \text{Con}_{\text{ZFC}}, \dots$  is provable in ZFC (where  $\Box$  denotes  $\Box_{\text{ZFC}}$  as long as it is not redefined). The latter is excluded only if we assume that the  $\omega$ -iterated consistency extension  $\text{ZFC}^\omega$  is consistent, hence  $\not\vdash_{\text{ZFC}} \Box^n \perp$ , for all  $n$  (see page 284), so that by Theorem 5.3, G would be the provability logic of ZFC.

In fact, the assumption  $(\forall n \in \mathbb{N}) \not\vdash_{\text{ZFC}} \Box^n \perp$  is equivalent to G's being the provability logic of ZFC, by the general Theorem 7.1 below. Therein  $Rf_T := \{\Box\alpha \rightarrow \alpha \mid \alpha \in \mathcal{L}^0\}$  denotes the already encountered reflection principle. Also Theorem 5.3 is a corollary of the theorem, simply because  $(\forall n \in \mathbb{N}) \not\vdash_T \Box^{n+1} \perp$  is equivalent to the consistency of  $T^\omega$ .

**Theorem 7.1.** *For a sufficiently expressive theory  $T$ <sup>5</sup> the following conditions are equivalent:*

- (i)  $T^\omega$  is consistent,
- (ii)  $T + Rf_T$  is consistent,
- (iii)  $G$  is the provability logic of  $T$ .

**Proof.** (i) $\Rightarrow$ (ii) indirect: Suppose that  $T + Rf_T$  is inconsistent. Then there are formulas  $\alpha_0, \dots, \alpha_n$  such that  $\vdash_T \neg\varphi$ ,  $\varphi := \bigwedge_{i=1}^n (\Box\alpha_i \rightarrow \alpha_i)$ . Hence  $\vdash_T \Box\neg\varphi \equiv_T \neg\Diamond\varphi$ . Now, because  $\vdash_{T^\omega} \neg\Box^{n+1}\perp$ , by Example (d) in 7.4 and Lemma 5.1, we get  $\vdash_{T^\omega} \Diamond R_n^i (p_i^i = \alpha_i)$ . Clearly,  $R_n^i = \varphi$  and so  $\vdash_{T^\omega} \Diamond\varphi$ . Since also  $\vdash_{T^\omega} \neg\Diamond\varphi$ ,  $T^\omega$  is inconsistent. (ii) $\Rightarrow$ (iii): The proof of Theorem 5.2 for PA, as presented in [Boo], runs nearly in the same way for  $T$ , because PA is transgressed in one place only: one uses the fact that  $\mathcal{N} \models Rf_{PA}$ . However, the existence of a corresponding  $T$ -model is ensured by (ii). (iii) $\Rightarrow$ (i):  $\not\vdash_G \Box^{n+1}\perp$ , hence  $\not\vdash_T \Box^{n+1}\perp \equiv_T \neg\text{Con}_{T^n}$  for all  $n$ , and so  $T^\omega$  is consistent.  $\square$

The equivalence (i) $\Leftrightarrow$ (ii) is a purely proof-theoretic one. It is called *Goryachev's theorem*; see [Gor] or [Be2]. We obtained it using essentially some elementary modal logic. For  $T = \text{ZFC}$ , perhaps a bit more interesting than (i) or (ii) is the assumption of the  $\omega$ -consistency of ZFC, that is,

$$(*) \vdash_{\text{ZFC}} (\exists x \in \omega)\varphi(x) \Rightarrow \not\vdash_{\text{ZFC}} \neg\varphi(\underline{n}) \text{ for some } n \quad (\varphi(x) \in \mathcal{L}_\epsilon).$$

This assumption implies  $D1^*$ , which in turn ensures  $\not\vdash_{\text{ZFC}} \Box^{n+1}\perp$  for all  $n$ , that is, (i), and hence all other conditions in Theorem 7.1 hold for  $T = \text{ZFC}$ . It is worthwhile to observe that the consistency of  $\text{ZFC} + Rf_{\text{ZFC}}$  and thereby the proof of Solovay's completeness theorem for ZFC follow directly from (\*), without appealing to Goryachev's theorem. What is needed to see that the latter is the case is the following

**Lemma.** *Suppose that ZFC is  $\omega$ -consistent. Then there exists a model  $\mathcal{V} \models \text{ZFC}$  such that  $\mathcal{V} \models Rf_{\text{ZFC}}$ .*

**Proof.** Let  $\Omega := \{(\forall x \in \omega)\alpha \mid \alpha = \alpha(x) \in \mathcal{L}_\epsilon, \vdash_{\text{ZFC}} \alpha(\underline{n}) \text{ for all } n\}$ . Then  $\text{ZFC} + \Omega$  is consistent. Indeed, otherwise  $\vdash_{\text{ZFC}} \neg(\forall x \in \omega)\alpha \equiv (\exists x \in \omega)\neg\alpha$  for some  $(\forall x \in \omega)\alpha \in \Omega$  (since  $\Omega$  is closed under conjunction), in contradiction to (\*). Any  $\mathcal{V} \models \text{ZFC} + \Omega$  satisfies the reflection principle  $Rf_{\text{ZFC}}$ , for if

<sup>5</sup>By such a  $T$  we mean that the proof steps of Solovay's Theorem 5.2 not transgressing PA can be carried out in  $T$ . This does not yet imply the provability of the theorem itself. Which steps are transgressing PA is described in the following proof.

$\mathcal{V} \not\models \alpha$  then  $\not\models_{\text{ZFC}} \alpha$  and therefore  $\vdash_{\text{ZFC}} \neg \text{bew}_{\text{ZFC}}(n, \ulcorner \alpha \urcorner)$  for all  $n$ . Hence  $(\forall y \in \omega) \neg \text{bew}_{\text{ZFC}}(y, \ulcorner \alpha \urcorner) \in \Omega$ , which clearly implies  $\mathcal{V} \not\models \Box \alpha$ .  $\square$

Now we interpret the modal operator  $\Box$  no longer as *provable in ZFC*, which is equivalent to *valid in all ZFC-models*, but rather as *valid in particular classes of ZFC-models*. For undefined notions used in the sequel we refer to [Ku]. A ‘model’ is to mean throughout a ZFC-model.

Particularly interesting are *transitive* models, i.e. models  $\mathcal{V} = (V, \in^{\mathcal{V}})$ , where the set  $V$  is *transitive*. This is to mean  $a \in b \in V \Rightarrow a \in V$ . In these models,  $\in^{\mathcal{V}}$  coincides with the ordinary  $\in$ -relation restricted to  $V$  (a set in our metatheory that itself is ZFC). We write  $V$  for  $\mathcal{V}$ . Let  $\rho a$  denote the *ordinal rank* of the set  $a$ , i.e., the smallest ordinal  $\rho$  with  $a \in V_{\rho+1}$ . To prove the soundness half of Theorem 7.3 we will need

**Lemma 7.2.** ([JK]) *Let  $V, W$  be transitive models such that  $\rho V < \rho W$  and let  $V \models \alpha$ . Then  $W \models$  ‘there is a transitive model  $U$  with  $U \models \alpha$ ’.*<sup>6</sup>

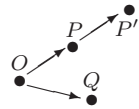
Let the modal logic Gi result from augmenting G by the axiom

$$(i) \quad \Box(\Box p \rightarrow \Box q) \vee \Box(\Box q \rightarrow p).$$

Gi is complete with respect to all *preference orders*  $g$ , i.e.,  $g$  is a finite poset together with some function  $\pi: g \rightarrow n (= \{0, \dots, n-1\})$  such that  $P < Q \Leftrightarrow \pi P < \pi Q$ , for all  $P, Q \in g$ . This implies the finite model property of Gi, which, as for G, ensures the decidability of Gi. More suitable for our aims is the characterization of preference orders  $g$  by the property

$$(p) \quad P < P' \text{ implies } P < Q \text{ or } Q < P', \text{ for all } P, P', Q \in g,$$

which at once follows from the definition: Let  $P < P'$ , hence  $\pi P < \pi P'$ . If  $P \not< Q$ , i.e.  $\pi P \not< \pi Q$ , then  $\pi Q \leq \pi P < \pi P'$ , so that  $Q < P'$ . The proof of the converse is Exercise 1. The figure shows a poset  $g$  that is *not* a preference order (for neither  $P < Q$  nor  $Q < P'$ ).



Axiom (i) is easily refuted in  $g$  choosing  $w_p = \{P'\}$ ,  $w_q = \emptyset$ , and verifying that  $O \Vdash \Diamond(\Box p \wedge \neg \Box q)$  and  $O \Vdash \Diamond(\Box q \wedge \neg p)$  (for notice that  $P \Vdash \Box p \wedge \neg \Diamond q$  and  $Q \Vdash \Box q \wedge \neg p$ ). Hence, (i) does not belong to G, so that Gi is a proper extension of G. We mention that in [So] and in [Boo] a somewhat more complex axiomatization of Gi has been considered.

<sup>6</sup>In transitive models  $W$  the sentence in ‘ ’ (which with some encoding can be formulated in  $\mathcal{L}_\in$ ) is absolute, and therefore equivalent to the existence of a transitive model  $U \in W$  with  $U \models \alpha$ .

**Remark on splittings in modal logic.** The completeness of  $\text{Gi}$  with respect to all preference orders follows also from the fact that  $\text{Gi}$  is the split logic arising from splitting the lattice of all extensions of  $\text{G}$  (see e.g. [Kra]) by the subdirect irreducible  $\text{G}$ -algebra belonging to the frame from the previous page.

We define insertions  $\iota: \mathcal{F}_\square \rightarrow \mathcal{L}_\epsilon^0$  as in 7.5 as usual by  $(\square H)^\iota = \square H^\iota$ , where  $\square\alpha$  for the set-theoretic sentence  $\alpha = H^\iota \in \mathcal{L}_\epsilon^0$  is now to mean ‘ $\alpha$  is valid in all transitive models’. Accordingly,  $\diamond\alpha = \neg\square\neg\alpha$  states ‘ $\alpha$  holds in at least one transitive model’.

**Theorem 7.3.**  $\vdash_{\text{Gi}} H$  iff  $\vdash_{\text{ZFC}} H^\iota$  for all insertions  $\iota$  as defined above.

We prove only the direction  $\Rightarrow$ , that is, soundness. The converse is much more difficult, see [Boo]. As regards the axioms of  $\text{Gi}$ , since  $\square p \rightarrow \square\square p$  is provable from the other axioms of  $\text{G}$  (see 7.4), it suffices to prove

- (A)  $\square(\alpha \rightarrow \beta) \wedge \square\alpha \vdash_{\text{ZFC}} \square\beta$ , (B)  $\square(\square\alpha \rightarrow \alpha) \vdash_{\text{ZFC}} \square\alpha$ ,  
 (C)  $\vdash_{\text{ZFC}} \square(\square\alpha \rightarrow \square\beta) \vee \square(\square\beta \rightarrow \alpha)$ , for all  $\alpha, \beta \in \mathcal{L}_\epsilon^0$ .

(A) is trivial, because the sentences valid in any class of models are closed under MP. (B) is equivalent to (B’):  $\diamond\neg\alpha \vdash_{\text{ZFC}} \diamond(\square\alpha \wedge \neg\alpha)$ . Here is the proof: Suppose  $\diamond\neg\alpha$ , i.e. there is a transitive model in which  $\neg\alpha$  holds. Then there is also one with minimal rank,  $V$  say. We claim  $V \models \square\alpha$ . Otherwise  $V \models \diamond\neg\alpha$ , and hence there would be a transitive model  $U \in V$  with  $U \models \neg\alpha$  and  $\rho U < \rho V$ , contradicting our choice of  $V$ . Therefore,  $V \models \square\alpha \wedge \neg\alpha$ . Thus, there is a transitive model in which  $\square\alpha \wedge \neg\alpha$  holds, which confirms (B’). Finally, (C) is verified by contraposition: suppose there are transitive models  $V, W$  and sentences  $\alpha, \beta$  such that

- (a)  $V \models$  ‘ $\alpha$  holds in all transitive models and there is a transitive model in which  $\neg\beta$  holds’,  
 (b)  $W \models$  ‘ $\beta$  holds in all transitive models’, (c)  $W \models \neg\alpha$ .

From these assumptions it follows first of all that  $\rho W < \rho V$ . Indeed, suppose by (a) that  $U \in V$  is a transitive model for  $\neg\beta$ . If  $\rho V \leq \rho W$  then  $\rho U < \rho W$ . Hence, by Lemma 7.2,  $W \models$  ‘there is a transitive model for  $\neg\beta$ ’, contradicting (b). Now, since  $W \models \neg\alpha$  by (c) and because of  $\rho W < \rho V$ , in  $V$  holds ‘there is some transitive model for  $\neg\alpha$ ’ by Lemma 7.2, in contradiction to (a). This proves (C). Soundness of the substitution rule follows as for  $\text{G}$  in 7.5. MN is trivially sound, because if  $\alpha$  is provable in ZFC then, of course,  $\alpha$  is valid in all transitive models. Also MP is obvious: If  $\alpha$  and  $\alpha \rightarrow \beta$  hold in any class of models, then also  $\beta$ .

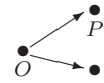
Another interesting model-theoretic interpretation of  $\Box\alpha$  is ‘ $\alpha$  is valid in all  $V_\kappa$ ’. Here  $\kappa$  runs through all inaccessible cardinal numbers. According to [So], the adequate modal logic for this interpretation of  $\Box$  is

$$\mathbf{Gj} := \mathbf{G} + \Box(\Box p \wedge p \rightarrow q) \vee \Box(\Box q \rightarrow p).$$

More precisely, if there are infinitely many inaccessibles then we have

**Theorem 7.4.**  $\vdash_{\mathbf{Gj}} H$  iff  $\vdash_{\mathbf{ZFC}} H^\iota$  for all insertions  $\iota$ , where  $\Box\alpha$  is to mean ‘ $\alpha$  is valid in all  $V_\kappa$ ’,  $\kappa$  running through all inaccessible cardinals.

$\mathbf{Gj}$  is also denoted by  $\mathbf{G.3}$ . This logic is complete with respect to all finite strict linear orders. These, of course, are also frames for  $\mathbf{Gi}$ , so that  $\mathbf{Gi} \subseteq \mathbf{Gj}$ . The figure shows a  $\mathbf{Gi}$ -frame, also called “the fork,” on which the additional axiom is easily refuted at its initial point  $O$  with  $wp = \{P\}$  and  $wq = \emptyset$ . Hence the fork is not a  $\mathbf{Gj}$ -frame, and so  $\mathbf{Gi} \subset \mathbf{Gj}$ . The completeness of  $\mathbf{Gj}$  with respect to finite orders entails the finite model property of  $\mathbf{Gj}$  and hence its decidability.



We recommend that the reader carry out the proof of the soundness part of Theorem 7.4, without consulting the hints to the solutions (Exercise 4). It is easier than the soundness part of Theorem 7.3 proved above. All one needs to know besides Lemma 7.2 is that each  $V_\kappa$  is a transitive ZFC-model and that  $V_\kappa \in V_\lambda$  or  $V_\lambda \in V_\kappa$ , for arbitrary inaccessible cardinals  $\kappa \neq \lambda$ . Maybe the reader can also find a new and lucid proof of the hard direction of Theorem 7.4: If  $\vdash_{\mathbf{ZFC}} H^\iota$  for all  $\iota$  then  $H$  holds in all finite strict linear orders, or equivalently,  $\mathbf{Gj} \vdash H$ .

## Exercises

1. Let  $g$  be a  $\mathbf{G}$ -frame with property **(p)**, page 296. Show by induction on the length of a maximal path in  $g$  that  $g$  is a preference order.
2. Show (using Exercise 1) that axiom **(i)** for  $\mathbf{Gi}$  holds in a  $\mathbf{G}$ -frame  $g$  iff  $g$  is a preference order. This is an essential step in proving the completeness of  $\mathbf{Gi}$  with respect to all preference orders.
3. This exercise is a crucial step in the completeness proof of  $\mathbf{Gj}$ . Show that a  $\mathbf{G}$ -frame  $g$  is a frame for  $\mathbf{Gj}$ , i.e.,  $\Box(\Box p \wedge p \rightarrow q) \vee \Box(\Box q \rightarrow p)$  holds in  $g$  if and only if  $g$  is (totally) ordered.
4. Verify the soundness part of Theorem 7.4, i.e.,  $\vdash_{\mathbf{Gj}} H \Rightarrow \vdash_{\mathbf{ZFC}} H^\iota$  for all insertions  $\iota$ .