

## 5. Second incompleteness theorem: research developments and consequences

### 5.1. Intensionality of the consistency statements

In Gödel's 1931 article there is only the sketch of what is been called “second theorem”, with the comment that it is a formalization of the first theorem, accompanied by the announcement of the imminent publication of a second part of the article, which is never happened. The first real proof of the “second theorem” appeared in Hilbert and Bernays (1934), vol. II, and it is in this context that emerges the decisive importance of *the way* in which the sentence that affirms the consistency of a certain theory is formalized. Hilbert and Bernays pointed out some conditions that the predicate numbering theorems must satisfy, called “derivability conditions”, that were further developed by Löb. These conditions are sufficient for obtaining the second incompleteness theorem.

Gödel's second theorem, states that given a consistent extension  $T$  of Robinson's  $Q$  that meet certain conditions, there a un sentence (that we denote  $Con(T)$ ) who claim to represent in a natural way the consistency of the theory, and which behaves as the undecidable sentence of the first theorem; i.e. we can produce, under the same conditions, a sentence undecidable which is the arithmetization of the metatheoretical sentence expressing in the language of the theory, the consistency of the theory itself:

1. If  $T$  is consistent, then  $T \not\vdash Con(T)$
2. If  $T$  is  $\omega$ -consistent, then  $T \not\vdash \neg Con(T)$

Actually, the second theorem implies that no theory for the formal arithmetic that incorporates the “finitary mathematics” is able to demonstrate the consistency of transfinite mathematics with the only means of finitary mathematics and this, as Von Neumann first observed, undermines Hilbert's programme. However, in view of the second theorem, the meaning of the provability predicate deserves a deep reflection. As we have seen, Rosser used a particular predicate. The statement of consistency  $Con^R(T)$  that results does not fulfil the second Gödel's theorem. If we mistakenly thought we could use  $Pr^R(x)$ , instead of the standard proof predicate, we would be faced with the following result, from which “Rosser consistency”  $\neg Pr^R(\overline{\Gamma 1 = 0})$  follows by considering that  $Q \vdash \neg(\overline{1} = \overline{0})$ .

**Lemma 30.** *If  $Q \vdash \neg\phi$ , then  $Q \vdash \neg Pr^R(\overline{\Gamma\phi})$ .*

*Proof.* Suppose that  $\neg\phi$  is provable; we claim that  $\neg Pr^R(\overline{\Gamma\phi})$ , namely:

$$\forall x(\neg Prf_Q(x, \overline{\Gamma\phi}) \vee \exists y < x Prf_Q(y, \overline{\Gamma\neg\phi}))$$

is provable too.

Hence work in  $Q$  and recall that the theory proves  $x \leq \bar{n} \vee x > \bar{n}$ . Now, if  $\neg\phi$  is provable, then a number  $n$  exists that codes a proof of it and therefore by binumerability  $\vdash Prf_Q(\bar{n}, \overline{\Gamma\neg\phi})$ . This  $n$  cannot be a code of a proof of  $\phi$  too. For any  $x$  we have two possibilities:

1. If  $\bar{n} < x$ , then the second disjunct of the above disjunction is true.
2.  $\bar{n} \geq x$  then, since  $x = \bar{0} \vee \dots \vee x = \bar{n}$ , if any of these numbers codes a proof of  $\phi$  we have a contradiction, again consistency. Hence  $x \leq \bar{n} \rightarrow \neg \text{Prf}_{\mathcal{Q}}(x, \overline{\phi})$ , i.e. the first disjunct holds.

We conclude that the above disjunction holds for all  $x$ .

QED

On the other hand, if we maintain the usual proof-predicate, we must pay attention on the definition of proper axioms. The American logician Solomon Feferman was highly insistent on this point:

At a given theory  $\mathsf{T}$  we associate a class of formulas  $\tau(x)$  that number the set of proper axioms of  $\mathsf{T}[\dots]$  Then we can associate with each of these formulas, in a uniform way, a formula  $\text{Prf}_{\tau}(x, y)$  and therefore the sentence  $\text{Con}_{\tau}$ . When the formula  $\tau(x)$  is recognized to be a correct expression of the fact that  $x$  codes an axiom of  $\mathsf{T}$ , then the sentence associated with it  $\text{Con}_{\tau}$  will be recognized as correct expression of the proposition that  $\mathsf{T}$  is consistent (Feferman (1960), p. 38).

The problem of how arithmetize the predicate “ $x$  is an axiom of  $\mathsf{T}$ ” does not have a single answer: different formulas defining the proper axioms correspond to different notions of provability. Some of these actually result in provability predicates of  $\text{Prf}_{\tau}$  for which it is provable that  $\neg \exists y \text{Prf}_{\tau}(y, \overline{\Gamma 1 = \bar{0}})$ , against Gödel’s result. According to Feferman, although extensionally such  $\tau$  correspond to the set of axioms, actually they do not express properly belonging to them. In order to be not only correct extensionally, but also *intensionally* (from the Latin word *intentio*, which descends from medieval logic), i.e. conceptually appropriate, it is necessary that the corresponding provability predicates meet also other conditions, typically, those due to Löb. Following Feferman, many logicians actually favour an intensional approach to arithmetisation in general, which is not satisfied with the mere representability of functions and relations, but demands that the theory be able to prove general properties of them. On a similar approach is based, for example, Buss (1986), where an intensional arithmetisation of metamathematics is developed, already in the weak theory  $S_2^1$ .

**Theorem 82.** (Feferman (1962)) *There is a binumeration  $\tau$  of the PA’s axioms such that  $\text{PA} \vdash \text{Con}_{\tau}$ .*

*Proof.* Since PA is recursively axiomatizable, it will have a  $\Sigma_1$  binumeration  $\sigma$  and a  $\Pi_1$  binumeration  $\pi$ . Let now  $\tau(x)$  the conjunction of the following:

1.  $\sigma(x) \wedge \forall y \leq x (\sigma(y) \leftrightarrow \pi(y))$
2.  $\neg \exists z \text{Prf}_{\pi \upharpoonright x}(z, \overline{\Gamma 1 = \bar{0}})$

where  $\pi \upharpoonright x(y) \leftrightarrow \pi(y) \wedge y \leq x$ . Recalling now that PA is *reflexive*, i.e. demonstrate the consistency of all its finite subtheories, we establish the equivalence between  $\tau(\bar{k})$ ,  $\sigma(\bar{k})$  and  $\pi(\bar{k})$ , for all  $k$ . (*Exercise*).

Reasoning inside PA, it can be proved that:

1. if  $\text{Con}_{\pi}$  holds, since  $\tau(x) \rightarrow \pi(x)$ , we also have  $\text{Con}_{\tau}$ .
2. if  $\neg \text{Con}_{\pi}$  holds, take the minimum  $z$  such that  $\exists y \text{Prf}_{\pi \upharpoonright_{z+1}}(y, \overline{\Gamma 1 = \bar{0}})$ . But then  $\text{Con}_{\pi \upharpoonright z}$  and  $\tau(x) \rightarrow \pi \upharpoonright z(x)$ : suppose  $\neg \pi \upharpoonright z(x)$ . This means that  $\pi(x) \rightarrow (x > z)$  and by the definition of  $z$ ,  $\exists y \text{Prf}_{\pi \upharpoonright x}(y, \overline{\Gamma 1 = \bar{0}})$ , namely  $\neg \tau(x)$ . But from  $\tau(x) \rightarrow \pi \upharpoonright z(x)$  follows that  $\text{Con}_{\pi \upharpoonright z} \rightarrow \text{Con}_{\tau}$ .

In both cases  $\text{Con}_{\tau}$ . Notice that, unlike the standard case, the formula  $\text{Con}_{\tau}$  is provably  $\Pi_2$  in PA, the formula  $\tau$  is provably  $\Delta_2$  in PA (see Hájek and Pudlák (1993), ch.III for further details). QED

A deep investigation of Feferman's interesting, though strange proof predicate and its relations with Gödel's standard one was made in Montagna (1978), Montagna (1987) and Visser (1989). But why should we favor a formalization of the intuitive notion of consistency to another? Thinking in more general terms, in the first theorem we required only a certain formula  $P(x)$  that numerate the theorems of a  $\omega$ -consistent extension of  $\mathsf{T}$  of Robinson arithmetic: to it, we asked just of being *extensionally correct*, i.e. it was required the following:

$$\mathsf{T} \vdash P(\bar{n}) \text{ iff } n \text{ codes a theorem of } \mathsf{T}.$$

This is not sufficient for deriving the second theorem: there are numerations extensionally correct, but that give rise to predicates of consistency derivable in  $\mathsf{T}$ . In *second* Gödel's theorem it is required something more: if  $P(x)$  is an enumeration of  $\mathsf{T}$ 's theorems, then sufficient condition to the sentence  $Con_P$  and its negation are both unprovable, is that are fulfilled certain conditions of provability, introduced by Hilbert and Bernays in 1939 later simplified in this propositional form by Löb in 1955 and Jeroslow (1973). Such conditions are indeed so natural that one may wonder if it really would be a provability predicate, that predicate that does not satisfy it. However, there is no general consensus on this fact.

**Definition 39.** *Let  $P \in \Sigma_1$  a definition of theorems of  $\mathsf{T}$ . Löb's conditions are the following:*

1. *If  $\mathsf{T} \vdash \phi$ , then  $\mathsf{T} \vdash P(\overline{\Gamma\phi})$ .*
2.  *$\mathsf{T} \vdash P(\overline{\Gamma\phi}) \wedge P(\overline{\Gamma\phi \rightarrow \psi}) \rightarrow P(\overline{\Gamma\psi})$ .*
3.  *$\mathsf{T} \vdash P(\overline{\Gamma\phi}) \rightarrow P(\overline{\Gamma P(\overline{\Gamma\phi})})$ .*

The standard proof predicate  $\text{Pr}_{\mathsf{T}}(x) = \exists y \text{Prf}_{\mathsf{T}}(y, x)$  satisfies these conditions in sufficiently strong theories. But Robinson's arithmetic is not strong enough for this purpose. It proves just the first condition. Concerning the third condition, this is an application of the so-called  $\Sigma_1$  formalized completeness. This means that if  $\mathsf{T}$  is a sufficiently strong theory, it proves  $\phi \rightarrow \text{Pr}_{\mathsf{T}}(\overline{\Gamma\phi})$ , for  $\phi \in \Sigma_1$ , noting that  $\text{Pr}_{\mathsf{T}}(y) \in \Sigma_1$ . Once more, Robinson's arithmetic is not strong enough to prove it. The third condition is actually the most problematic and difficult to justify (see Detlefsen (2001) for an in-depth discussion around this condition). In fact the proof can be performed in the fragment of Peano arithmetic denoted  $I\Delta_0 + \text{exp}$ , i.e. with induction restricted to  $\Delta_0$ -formulas, together with an axiom that states the totality of the function  $2^x$ . Falling below these levels, for example in the important theory of *Bounded Arithmetic* named  $I\Delta_0 + \Omega_1$  (see 3) some problems may arise. In Berarducci and Verbrugge (1991) it is shown that if this theory proves the  $\Sigma_1$ -completeness, then it would also be valid the equation  $\text{NP} = \text{co-NP}$ . More exactly, if  $\text{NP} \neq \text{co-NP}$ , then there exist  $\phi$  and  $\psi$  such that the  $\Sigma_1$ -completeness of  $\Sigma_1$ -formula:

$$\exists x (\text{Prf}_{I\Delta_0 + \Omega_1}(x, \overline{\Gamma\phi}) \wedge \forall z \leq x \neg \text{Prf}_{I\Delta_0 + \Omega_1}(z, \overline{\Gamma\psi}))$$

is not provable in  $I\Delta_0 + \Omega_1$ . However, we can carefully obtain versions of the  $\Sigma_1$ -completeness restricted to specific subclasses, but sufficient for deducing the derivability conditions already in Buss's  $\mathsf{S}_2^1$ . As regards the second condition, it is in practice to demonstrate that, for all  $a, b$ , if  $a$  codes a proof of  $\phi$ , if  $b$  codes a proof of  $\phi \rightarrow \psi$ , then  $a * b * \overline{\Gamma\psi}$  codes a proof of  $\psi$ , and therefore it is needed the amount of induction required (induction on  $\Sigma_1$ -formulas is sufficient). Rosser's predicate does not meet at least one of them and Feferman's predicate of theorem on p.122 does not satisfy 3. because it is not  $\Sigma_1$ .

The idea behind the derivability conditions was that all formulas correctly expressing provability would have to satisfy them, although, this generally accepted idea has also found proud opponents (see primarily the works Detlefsen (1986) and Detlefsen (1979)). Regarding the second incompleteness theorem for  $\mathsf{Q}$ , we wonder which relation exists, in  $\mathsf{Q}$ , between the arithmetic sentence  $Con(\mathsf{Q})$  and the metamathematical statement " $\mathsf{Q}$  is consistent". Pivotal to this debate, albeit of controversial interpretation, was a version of the second incompleteness theorem for  $\mathsf{Q}$  in Bezboruah and Shepherdson (1976) using different methods, where nevertheless the two logicians, amazingly, sharing Kreisel's point of view, do not attach importance to this result.

Indeed, according to Georg Kreisel, which had shown similar results well in advance,  $Con(\mathbb{Q})$  should not be considered, in Robinson arithmetic, as expressing its formalized consistency, but a mere algebraic property, that only in stronger systems can be said to constitute a formalization of consistency of  $\mathbb{Q}$ . Nowadays many people think that  $\mathbb{Q}$  can actually prove the unprovability of its own consistency, on the basis of an argument due to Pudlák (1996), which starts from the consideration that the weak theory  $I\Delta_0 + \Omega_1$  is interpretable in  $\mathbb{Q}$  (see on p.192). The argument is the following. A formula  $J(x)$  is called a *cut* of a theory  $\mathbb{T}$ , if this theory proves:

1.  $J(0)$
2.  $J(x) \rightarrow J(x+1)$
3.  $J(x) \wedge y \leq x \rightarrow J(y)$

Let therefore  $Con_{\mathbb{T}}^J = \neg \exists x (J(x) \wedge \overline{Prf_{\mathbb{T}}(x, \overline{\Gamma 1} = \overline{0^1})})$ . Pudlák showed that for every consistent extension  $\mathbb{T}$  of  $\mathbb{Q}$  and every cut  $J(x)$ ,  $\mathbb{T} \not\vdash Con_{\mathbb{T}}^J$ . Actually it is possible to build cuts  $J(x)$  with further properties, in particular, in such a way that satisfy the axioms of  $I\Delta_0 + \Omega_1$  relativized to  $J(x)$ . Take  $\mathbb{T} = \mathbb{Q}$ . Hence it is consistent to assume that a proof of a contradiction from  $\mathbb{Q}$  is encoded in a model of  $I\Delta_0 + \Omega_1$ , a theory in which the meaning of  $Con_{\mathbb{Q}}$  is not ambiguous.

But let us now see the usual argument based on the derivability conditions.

**Lemma 31.** *Let  $\mathbb{T}$  be a consistent extension of Robinson arithmetic and let  $Con_P$  the formula  $\neg P(\overline{\Gamma 0} = \overline{1^1})$ , where  $P$  satisfies Löb's conditions; then  $\mathbb{T} \vdash Con_P \leftrightarrow \neg P(\overline{\Gamma \phi}^1) \vee \neg P(\overline{\Gamma \neg \phi}^1)$ , for any  $\phi$ .*

*Proof.*  $\Leftarrow$  Recall that  $\neg(\overline{0} = \overline{1})$  is an axiom; from the *first derivability condition* in  $\mathbb{T}$  we get  $P(\overline{\Gamma \neg(\overline{0} = \overline{1})}^1)$ , from which, by propositional tautology

$$A \rightarrow (B \rightarrow (A \wedge B))$$

we obtain  $\neg Con_P \rightarrow (P(\overline{\Gamma \neg(\overline{0} = \overline{1})}^1) \wedge P(\overline{\Gamma 0} = \overline{1^1}))$ . From the first condition and tautology  $A \rightarrow (\neg A \rightarrow B)$  we get  $P(\overline{\Gamma(\overline{0} = \overline{1})} \rightarrow (\neg(\overline{0} = \overline{1}) \rightarrow \phi)^1)$  and by the second derivability condition and tautology  $(A \rightarrow (B \rightarrow C)) \leftrightarrow ((A \wedge B) \rightarrow C)$  lastly  $P(\overline{\Gamma(\overline{0} = \overline{1})} \wedge P(\overline{\Gamma \neg(\overline{0} = \overline{1})}^1) \rightarrow P(\overline{\Gamma \phi}^1))$ . By modus ponens we obtain  $\neg Con_P \rightarrow P(\overline{\Gamma \phi}^1)$ ; we can repeat the argument and in specular way to obtain  $\neg Con_P \rightarrow P(\overline{\Gamma \neg \phi}^1)$ . Now  $\Rightarrow$ . In  $\mathbb{T}$ , from a propositional tautology and the first condition we get  $P(\overline{\Gamma \phi} \rightarrow (\neg \phi \rightarrow (\overline{0} = \overline{1})))^1$  and still for simple propositional transformations, using the second condition, we finally obtain  $P(\overline{\Gamma \neg \phi}^1) \wedge P(\overline{\Gamma \phi}^1) \rightarrow \neg Con_P$ . Take the contronominale. QED

**Theorem 83.** (The second incompleteness theorem) *Let  $\mathbb{T}$  a consistent extension of Robinson arithmetic and  $P(x)$  a  $\Sigma_1$ -definition of its theorems, satisfying the derivability conditions. Hence  $\mathbb{T} \not\vdash Con_P$ ; if moreover  $\mathbb{T}$  is also  $\omega$ -consistent (or  $\Sigma_1$ -sound), then  $\mathbb{T} \not\vdash \neg Con_P$ .*

*Proof.* Let  $\nu$  be a fixed point of  $\neg P(x)$  in  $\mathbb{T}$ , namely  $\mathbb{T} \vdash \nu \leftrightarrow \neg P(\overline{\Gamma \nu}^1)$ . By the previous lemma  $\neg P(\overline{\Gamma \neg \nu}^1) \vee \neg P(\overline{\Gamma \nu}^1) \rightarrow Con_P$  and therefore consider the meaning of  $\nu$ , lastly  $\nu \rightarrow Con_P$ . On the other hand  $\mathbb{T}$  we also have  $\neg \nu \rightarrow P(\overline{\Gamma \nu}^1)$ . By the third condition  $P(\overline{\Gamma \nu}^1) \rightarrow P(\overline{\Gamma P(\overline{\Gamma \nu}^1)}^1)$  and by the meaning of  $\nu$ , from this we get  $P(\overline{\Gamma \nu}^1) \rightarrow P(\overline{\Gamma \neg \nu}^1)$  (where we have replaced  $\neg \nu$  in place of  $P(\overline{\Gamma \nu}^1)$  in the consequent); hence we have  $\neg \nu \rightarrow P(\overline{\Gamma \neg \nu}^1) \wedge P(\overline{\Gamma \nu}^1)$  and in a few logical steps, using the previous lemma, finally, we get also  $\neg \nu \rightarrow \neg Con_P$ . QED

What about sentences asserting their own *provability*? This question is known as the *Henkin's question*.

**Theorem 84.** (Löb 1955) *Let  $\phi$  be a sentence and  $Pr_{\mathbb{T}}(x)$  the standard proof predicate for  $\mathbb{T}$ . Then  $\mathbb{T} \vdash Pr_{\mathbb{T}}(\overline{\Gamma \phi}^1) \rightarrow \phi$  iff  $\mathbb{T} \vdash \phi$ .*

*Proof.*  $\Leftarrow$  clear.  $\Rightarrow$  Let us suppose that  $\mathsf{T} \not\vdash \phi$ . Hence  $\mathsf{T} + \neg\phi$  is consistent and therefore  $\mathsf{T} + \neg\phi \not\vdash \text{Con}_{\mathsf{T} + \neg\phi}$ , namely  $\mathsf{T} \not\vdash \neg\phi \rightarrow \neg\text{Pr}_{\mathsf{T}}(\overline{\neg\phi \rightarrow (1=0)})$  and therefore  $\mathsf{T} \not\vdash \neg\phi \rightarrow \neg\text{Pr}_{\mathsf{T}}(\overline{\neg\phi})$ . QED

We get a solution to Henkin's question as a corollary.

**Corollary 19.**  $\mathsf{T} \vdash \text{Pr}_{\mathsf{T}}(\overline{\neg\phi}) \leftrightarrow \phi$  iff  $\mathsf{T} \vdash \phi$ .

This result (sometimes seen as a generalization of Gödel's result, taking  $1=0$  for  $\phi$ ) focuses our attention on so called *reflection shemas*.

*Hilbert's program of elimination of abstract entities.* Principles like the above " $\text{Pr}_{\mathsf{T}}(\overline{\neg\phi}) \rightarrow \phi$ ", represent a formalization of *soundness* called *reflection principles*; the main schematic representations of them are the following:

1. (*Local Reflection*  $Rfn_{\mathsf{T}}$ )  $\text{Pr}_{\mathsf{T}}(\overline{\neg\phi}) \rightarrow \phi$
2. (*Uniform Reflection*  $RFN_{\mathsf{T}}$ )  $\forall x(\text{Pr}_{\mathsf{T}}(\overline{\neg\phi(\dot{x})}) \rightarrow \phi(x))$

(where  $\overline{\neg\phi(\dot{x}_0, \dots, \dot{x}_n)}$  is the function that associates to a given sequence of numbers  $m_0, \dots, m_n$ , the number  $\overline{\neg\phi(\overline{m_0}, \dots, \overline{m_n})}$ ). The second principle is stronger than the first. Actually  $\mathsf{T} + RFN_{\mathsf{T}}$  proves  $\text{Con}(\mathsf{T} + Rfn_{\mathsf{T}})$ .

If we consider *partial scheme*  $RFN_{\mathsf{T}, \Pi_1}$ ,  $Rfn_{\mathsf{T}, \Pi_1}$ , i.e. restricted, in particular, to  $\Pi_1$ -formulas, and if  $\text{Con}_{\mathsf{T}}$  formalize consistency of  $\mathsf{T}$ , then we have that the principles  $\text{Con}_{\mathsf{T}}$ ,  $RFN_{\mathsf{T}, \Pi_1}$  and  $Rfn_{\mathsf{T}, \Pi_1}$  are equivalent. This allows us to better clarify the link between the Hilbert program of consistency and that of  $\Pi_1$ -conservativity, i.e. the statement that is exposed by Hilbert in the 1926 essay on infinity, in these terms:

the extension by the addition of ideal elements is legitimate, if no contradiction is determined in the old, restricted domain, namely if relations that result for old objects, when ideal objects are eliminated are valid in the old domain.

Since we tied the finitary mathematics to Skolem arithmetic, suppose here to fix ideas, that  $\mathsf{T}$  is an extension of PRA. Hence consider the formula  $\psi(x) \in \Pi_1$ : since  $\neg\psi(x)$  will be then  $\Sigma_1$ , if the theory  $\mathsf{T}$  is sufficiently strong (as it is PRA) to show the  $\Sigma_1$ -formalized completeness, then it will show  $\neg\psi(x) \rightarrow \text{Pr}_{\mathsf{T}}(\overline{\neg\psi(\dot{x})})$ , and therefore  $\neg\text{Pr}_{\mathsf{T}}(\overline{\neg\psi(\dot{x})}) \rightarrow \psi(x)$ .

But  $\text{PRA} + \text{Con}_{\mathsf{T}}$  proves  $\text{Pr}_{\mathsf{T}}(\overline{\psi(\dot{x})}) \rightarrow \neg\text{Pr}_{\mathsf{T}}(\overline{\neg\psi(\dot{x})})$  - see the preparatory lemma to the proof of the second Gödel's theorem - and therefore, lastly,  $\text{Pr}_{\mathsf{T}}(\overline{\psi(\dot{x})}) \rightarrow \psi(x)$ , as claimed. On the other hand, if we use the principle of reflection for  $\Pi_1$ -formulas, this will apply in particular to the formula  $\overline{1=0}$ ; however it is an axiom that  $\neg(\overline{1=0})$ , from which  $\neg\text{Pr}_{\mathsf{T}}(\overline{\neg(\overline{1=0})})$ , namely  $\text{Con}_{\mathsf{T}}$ .

Thus we have shown in Skolem's theory PRA the equivalence between  $RFN_{\Pi_1}$  and  $\text{Con}_{\mathsf{T}}$ . But from this follows that if a formula  $\psi(x)$  of the language of Skolem arithmetic, of complexity  $\Pi_1$  is provable in  $\mathsf{T}$ , then it is also in  $\text{PRA} + \text{Con}_{\mathsf{T}}$ : indeed, for a theory strong enough as Skolem arithmetic, the first provability condition can be better refined in this sense: "if  $\mathsf{T} \vdash \psi(x)$ , then  $\text{PRA} \vdash \text{Pr}_{\mathsf{T}}(\overline{\psi(\dot{x})})$ ", from which  $\text{PRA} + RFN_{\Pi_1} \vdash \psi(x)$  and for the equivalences have just demonstrated,  $\text{PRA} + \text{Con}_{\mathsf{T}} \vdash \psi(x)$ .

Suppose, then, that  $\mathsf{T}$  embodies the infinitary mathematics in sense of Hilbert, and PRA, as in Tait thesis, incorporates finitary mathematics. Remember that, according to Hilbert, the "concrete" sentences have the form  $\forall x(f(x) = 0)$ , with  $f(x)$  primitive recursive, and are therefore  $\Pi_1$ -sentences of language PRA. If as expected by Hilbert (and refuted by the second theorem of incompleteness) the consistency of  $\mathsf{T}$  was demonstrated in PRA, then, as now established, we would have that if  $\psi(x)$  be provable in  $\mathsf{T}$ , the would already in PRA: in the proof of "concrete sentences", the "ideal entities" (to quote Hilbert) could be eliminated.

This allows us to make another point. For a time it is approached primarily to *second* Gödel theorem as a refutation of Hilbert's program. Subsequently, between the '70s and late '80s, by logicians like Kreisel, Prawitz, Simpson, or Smorynski, the focus is shifted more on the first theorem.

For example Smorynski drew from the connection that we have now illustrated between the program of *consistency* and the *conservativity* program, the conclusion that already from the first theorem we can derive a refutation of Hilbert program on consistency, in this sense: the sentence which gives us the first incompleteness theorem, has precisely the form of a  $\Pi_1$ -sentence “concrete”  $\forall x R(x)$  which is undecidable (and notice that its particular instances  $R(\bar{0})$ ,  $R(\bar{1})$ ,  $R(\bar{2})$ ... are provable).

If  $T$  is the transfinite system containing ideal elements, it shows all true sentences of finitary mathematics; but the program of conservativity precisely requires all  $\Pi_1$ -sentences “concrete” provable in  $T$  are then demonstrable already in PRA; but the first theorem of incompleteness, rather, shows a true  $\Pi_1$ -sentence of language PRA (and therefore provable in  $T$ ), but unprovable in PRA itself. The program of conservativity is therefore impossible: hence, for what we said above, it is also that of consistency.

## 5.2. Beating incompleteness: Turing’s progressions

The second Gödel’s theorem gives a way to extend a theory to a stronger theory and allows to see the phenomenon of inexhaustibility of mathematics from a particular point of view (see Franzen (2003)): the consistency statement  $Con(T)$  is independent from the consistent theory  $T$ , although true; hence, if  $T$  is a sound theory (i.e. does not prove false things), also the theory  $T'$  obtained from  $T$  adding the sentence  $Con(T)$  as a new axiom will be sound; moreover it will be stronger than  $T$  (for instance it will show that  $Con(T)$ ). But also  $T'$  will be incomplete and therefore  $Con(T')$  will be independent from it, but we can define a stronger theory  $T''$  that decides  $Con(T')$ , by adding to  $T'$  the true sentence  $Con(T')$  and so on.

This leads to the idea of an effective association of formal systems  $S_\alpha$  with ordinals  $\alpha$ , but that can be done only for countable ordinals and to deal with limits in an effective way, it turns out that we must work not with ordinals per se, but with recursive ordinals, or *notations* for ordinals. In his PhD thesis. 1937 in Princeton, Alan Turing formalized this intuition, by introducing the notion of ordinal logic and a suggestive idea to “overcome incompleteness” iterating a transfinite number of times the operation of adding an undecidable sentence to a theory, of the kind of reflection statements or consistency statements, hoping to get to a certain point a complete theory, with respect to a significant class of sentences. This began a series of studies around transfinite recursive hierarchies axiomatic theories. As Turing (1939) says:

The well-known theorem of Gödel (1931) shows that every system of logic is in a certain sense incomplete, but at the same time it indicates means whereby, from a system  $L$  of logic, a more complete system  $L'$  may be obtained. By repeating the process we get a sequence  $L_1 = L', L_2 = L'_1, \dots$  each theory more complete than the preceding... A logic  $L_\omega$  may be constructed in which the provable theorems are the totality of theorems provable with the help of logics  $L, L', L'' \dots$

Recall Feferman’s approach to the proof-predicate which we discussed earlier. He considers primarily the arithmetization of axiomhood: in the proof-predicate  $Prf_\tau$  he clearly distinguishes a proof-predicate  $Prf$  for the first order logic, which is fixed, from the definition  $\tau(x)$  of the specific mathematical axioms, which may change. A major case is when  $\tau(x) \in \Sigma_1^0$ , and *binumerates* the axioms, and therefore the whole expression  $\exists y Prf_\tau(y, x)$  becomes in turn  $\Sigma_1^0$ . The derivability conditions are satisfied.

We now need to recall some set theory concepts. Recall that from a set theoretic point of view an ordinal is just a transitive set (i.e. an  $y$  such that if  $x \in y$ , then  $x \subseteq y$ ) well (and therefore linearly) strictly ordered by the appartenance  $\in$ , i.e. every non-empty subset of it contains a  $\in$ -least element. This notion generalises the definition of natural numbers *à la* Von Neumann  $0 = \emptyset$ ,  $n + 1 = n \cup \{n\} = \{0, 1, 2, \dots, n\}$  by admitting limit ordinals as for example  $\omega = \{0, 1, 2, \dots\}$ , and so on. Let us recall some basic properties (see e.g. Jech (1978)):

1. Every member of an ordinal is an ordinal, however the class of all ordinals is not a set, but a proper class: since  $\in$  well-orders this class, otherwise it would be in turn an ordinal.

2. For every ordinal  $\alpha$ , also  $\alpha \cup \{\alpha\}$  is an ordinal, and this is the successor of it, i.e.,  $\alpha < \alpha \cup \{\alpha\}$  and there is no  $\beta$  such that  $\alpha < \beta < \alpha \cup \{\alpha\}$ . An ordinal is a *successor*, if it has this form: otherwise, either is 0, or it is a *limit*.
3. If  $\alpha$  is 0 or a limit ordinal then  $\sup \alpha = \alpha = \bigcup \alpha$ .
4. If  $\alpha$  is a successor ordinal then  $\sup \alpha$  is the predecessor of  $\alpha$ .
5. If  $\alpha$  is a limit ordinal and  $\gamma < \alpha$ , then there is an ordinal  $\delta$  such that  $\gamma < \delta < \alpha$ , and in particular  $\gamma < \gamma \cup \{\gamma\} < \alpha$ . Say that  $\alpha$  is a finite ordinal, or a natural number, if  $\alpha = 0$ , or  $\alpha$  is a successor and every ordinal  $\beta < \alpha$ , in turn is 0 or a successor.
6. Every well ordered set is isomorphic to a unique ordinal and the isomorphism too is unique.

A *countable* ordinal is an ordinal whose set of predecessors can be matched up with natural numbers. For instance, all these are countable:

$$\left. \begin{array}{l} \omega, \omega + n, \omega + \omega, \omega \times n, \omega \times \omega, \omega^\omega \dots \epsilon_0 = \omega^{\omega^{\omega^{\omega^{\dots}}}} \end{array} \right\} \omega - \text{times}$$

These have all the same cardinality and are essentially different ordering of natural numbers. *Uncountable* ordinals exist. The first is denoted by  $\omega_1$ . Recall that each well ordering is isomorphic to a unique ordinal. A *computable* ordinal is an ordinal isomorphic to a computable well ordering, i.e. if there is a computable relation on a subset of the integers that is well-ordered and isomorphic to it. The first *non-computable* ordinal is denoted by  $\omega_1^{CK}$ . It holds that  $\omega_1^{CK} < \omega_1$ . There are uncountably many countable ordinals, so, if we want “give a name” to each one (i.e. code by numbers), we cannot invent distinct names. We shall see a method, due to Kleene, to assign names to so-called *constructive* ordinals. Spector has proved that the computable ordinals are just Kleene’s constructive ordinals.

Given an effective description of formal systems  $S_0, S_1, S_2, \dots$  in the same language, starting e.g. from  $S_0 = PA$ , we can form the union  $S_\omega = \bigcup_i S_i$  i.e. the formal system made of the set the axioms of all  $S_i$ , but if each  $S_{i+1}$  results by the interaction of one and the same operation, for example  $S_{i+1} = S_i + Con(S_i)$ , we can continue in the transfinite, defining  $S_{\omega+1} = S_\omega + Con(S_\omega)$  and so on. To the limit steps  $\lambda$  we take the union  $S_\lambda = \bigcup_{\beta < \lambda} S_\beta$ . If  $S_i$  is sound (i.e. correct), then (being  $Con(S_i)$  a true sentence, although independent from  $S_i$ ) also  $S_{i+1}$  will be sound. From here the idea to associate in an effective way to a countable ordinal  $\alpha$ , a formal system  $S_\alpha$ , in the framework of an inductive construction of the kind which we have referred.

There are several ways to make rigorous this topic: it is necessary first to give a  $\Sigma_1$ -formula defining arithmetically the axioms of the theories that constitute the sequence, so as to be able to express, for example, that a certain theory is consistent (however, in relation to what we said about consistency predicates, it should be emphasized the character in some respects intensional of this construction, which depends significantly on how the axioms of a theory  $S_i$  are defined) and the problematic case concerns the limit steps. For instance, if  $\tau_0(x)$  is a definition of axioms of PA, then we can define  $S_0 = PA$  and the axioms of  $S_{i+1}$  will be defined by the formula  $\tau_{i+1}$  satisfying  $\tau_{i+1}(x) \leftrightarrow \tau_i(x) \vee x = \overline{Con_{\tau_i}}$ .

The first problem now concerns the definition of the axioms of  $S_\omega$ . For this purpose we can introduce a notation for countable ordinals, i.e. their coding by integers that allows us to describe them arithmetically. The idea is then to think a limit ordinal as a sequence of ordinals converging to it, enumerated from a function  $\phi_e(x)$ , so that for the definition of the axioms of a theory indexed with a limit ordinal, is intuitively fulfilled this condition “ $\tau_{lim(e)}(y)$  iff there is an  $n$  such that  $\tau_{\phi_e(n)}(y)$ ”. With Ordinal Logic we mean a sequence of theories  $S_{a_0}, S_{a_1}, S_{a_2}, \dots$  where each  $a_i$  is the name of an ordinal, i.e. a number in Kleene’s  $\mathcal{O}$  we are going to define.

It should be noted, however, that the same limit ordinal can have different notations (think for example that  $\omega$  is the limit of all strictly increasing computable infinite sequence of natural numbers) and is not sure that if two numbers  $a_j \in a_i$  denote the same ordinal, then the theories

$S_{a_i}, S_{a_j}$  prove the same theorems: if this is the case, then the logic is *invariant*. Turing proved, however, that an ordinal logic can be invariant, or complete for  $\Pi_1$  statements, but not both things together.

**Definition 40.** *An ordinal is said computable iff it is isomorphic to a recursive well order.*

It is well known that there is at least one *countable* ordinal, but *not computable*, and the least of these ordinals is denoted with  $\omega_1^{CK}$ .

1. In 1936 Church and Kleene introduced a satisfactory characterization for this concept, through the notion of *constructive ordinal* and a system  $\mathcal{O}$  of notations, i.e. of codes for countable ordinals.
2. the set  $\mathcal{O}$  of constructive ordinals provably coincide with that of computable ordinals. If we write  $\tilde{a}_i$  for the ordinal denoted by  $a_i$ , then we can assert that  $\omega_1^{CK}$  is the smallest ordinal not in the form  $\tilde{a}$ , for some  $a \in \mathcal{O}$ .

A significant part of the theory of ordinal numbers can be formulated as a theory of *ordinal notations*.

*Notation system.* It is a function  $\nu$  having domain a subset  $D$  of the naturals and codomain  $X \subset On$ , for which there exists a recursive partial function  $k(x)$  such that:

$$k(x) = \begin{cases} 0 & \text{if } \nu(x) = 0 \\ 1 & \text{if } \nu(x) = \beta + 1 \\ 2 & \text{if } \nu(x) = \lambda \text{ limit} \end{cases}$$

1. There is a partial recursive function  $p(x)$  such that  $p(x)$  converges and  $\nu(x) = \nu(p(x)) + 1$ , if  $\nu(x) = \beta + 1$ .
2. there is a function  $q(x)$  such that  $\phi_{q(x)}$  is total and the sequence:

$$\nu(\phi_{q(x)}(0)) < \nu(\phi_{q(x)}(1)) < \nu(\phi_{q(x)}(2)) < \dots$$

has limit  $\nu(x)$ , if  $\nu(x) = \lambda$  limit.

In 1936, Church and Kleene had introduced a system of constructive ordinal notations, given by certain expressions in the lambda-calculus. A variant of this uses numerical codes, and associates with each number a countable ordinal.

*Kleene's system  $\mathcal{O}$ .* The intuitive idea is to code  $0, 1, 2, 3, \dots$  with the powers  $1, 2, 2^2, 2^{2^2}, \dots$  and if  $\alpha$  is a limit ordinal, then its notations are all numbers  $3 \cdot 5^e$  such that  $\phi_e$  is a total function such that  $\phi_e(n)$  is a number in  $\mathcal{O}$  that denotes  $\alpha_n$  and  $\alpha_0, \alpha_1, \alpha_2, \dots$  is an increasing sequence converging to  $\alpha$ .

*Kleene's system  $\mathcal{O}$ .* For the sake of simplicity let us write  $\tilde{b}$  in place of  $\nu(b)$ .

1. 0 receives notation 1.
2. Suppose we have assigned a notation to all ordinals less than  $\alpha$  and having defined  $<_o$  on it:
  - (a) if  $\alpha = \beta + 1$  and  $\tilde{b} = \beta$ , then  $\tilde{2^b} = \alpha$  and add the pairs  $\langle z, 2^b \rangle$  to the relation  $<_o$ , for all  $z \leq_o b$ .
  - (b) If  $\alpha$  is a limit ordinal, it can be understood as a sequence of ordinals converging to it. Suppose that this sequence is enumerated by a total function  $\phi_e$  with values in  $\mathcal{O}$ , such that for all  $n$ ,  $\phi_e(n) <_o \phi_e(n+1)$ , where  $\phi_e(n) = a_n$  and the increasing sequence  $\tilde{a}_0, \tilde{a}_1, \tilde{a}_2, \dots$  has limit  $\alpha$ ; then  $\tilde{3 \cdot 5^e} = \alpha$ ; add moreover each pair  $\langle z, 3 \cdot 5^e \rangle$  such that  $z <_o \phi_e(n)$  for some  $n$ , to the relation  $<_o$ .

$$(c) \quad k(1) = 0, k(2^x) = 1, k(3 \cdot 5^y) = 2, p(2^x) = x, q(3 \cdot 5^y) = y.$$

Note that the natural numbers receive a fixed notation, but not limit ordinals. If the sequence  $\phi_e(0), \phi_e(1), \phi_e(2), \dots$  denotes effective enumeration of ordinals converging to  $\alpha$ , then there are infinite choices for  $e$ , both because each function has infinitely many index and because there are infinitely many increasing sequences of ordinals converging to  $\alpha$ . In summary, the system begins with  $1, 2, 2^2, 2^{2^2}, 2^{2^{2^2}} \dots$ , as notations for  $0, 1, 2, 3, \dots$ . To the limit ordinal  $\alpha$  we assign a notation  $3 \cdot 5^e$ , where  $\phi_e(x)$  is total recursive such that if  $\phi_e(n) = b_n$ , then  $b_0 <_o b_1 <_o b_2 <_o \dots$  and  $\alpha = \lim_n \alpha_m$ , where  $b_n$  denotes  $\alpha_n$ .

Hence, for instance,  $\omega$  is denoted by  $3 \cdot 5^e$ , for infinite  $e$ . The order  $<_o$  is therefore a tree, that in correspondence with a node labeled by a limit ordinal, branches in infinite branches.

$$1 <_o 2 <_o 2^2 <_o 2^{2^2} <_o \dots \left\{ \begin{array}{l} 3 \cdot 5^{e_0} <_o 2^{3 \cdot 5^{e_0}} <_o 2^{2^{3 \cdot 5^{e_0}}} <_o \dots \\ 3 \cdot 5^{e_1} <_o 2^{3 \cdot 5^{e_1}} <_o 2^{2^{3 \cdot 5^{e_1}}} <_o \dots \\ \dots \\ 3 \cdot 5^{e_n} <_o 2^{3 \cdot 5^{e_n}} <_o 2^{2^{3 \cdot 5^{e_n}}} <_o \dots \\ \dots \end{array} \right.$$

The following facts are well known (see Sacks (2017), ch.1):

1. There exist a recursive function  $+_o$  such that for all  $x, y \in \mathcal{O}$ :
  - (a)  $x +_o y \in \mathcal{O}$
  - (b)  $\widetilde{x +_o y} = \widetilde{x} + \widetilde{y}$
  - (c)  $y \neq 1 \rightarrow x <_o x +_o y$
2. For every other notation system  $S$  exists  $\phi : S \rightarrow \mathcal{O}$  partial recursive, such that:
3. if  $x \in S$ , then  $\nu_S(x) \leq \widetilde{\phi(x)}$ .
4.  $\omega_1^{CK} = \bigcup_{a \in \mathcal{O}} \widetilde{a}$  = set of recursive ordinals.
5. (Kleene)  $\mathcal{O}$  is  $\Pi_1^1$ -complete, i.e. is  $\Pi_1^1$ -definable and if  $Y \in \Pi_1^1$ , then  $Y \leq_m \mathcal{O}$ .

The version of the second incompleteness theorem presented by Feferman (1960) has this general form.

**Theorem 85.** *Let  $\mathbb{T}$  be a consistent extension of PA, and let  $\tau(x)$  be a  $\Sigma_1$ -formula which numerates the axioms of  $\mathbb{T}$ , and  $\text{Cons}(\mathbb{T})$  be a consistency statement constructed from  $\tau(x)$  and  $\text{Pr}f_\tau(x, y)$ . Then  $\text{Cons}(\mathbb{T})$  is not provable in  $\mathbb{T}$ .*

Hence we start with a  $\Sigma_1$  numeration of the axioms of such a theory  $\mathbb{T}$ . Given a numeration  $\tau(x)$ , we naturally construct the formula  $\text{Pr}f_\tau(y, x)$  that expresses the predicate is the Gödel number of the proof, in  $\mathbb{T}$ , of the formula with the number  $x$ ” and the formula of provability in  $\mathbb{T}$ ,  $\exists y \text{Pr}f_\tau(y, x)$ . Following Beklemishev (1992) we can now make more precise what we have said in the introduction, through this definition. First of all we introduce a provability predicate with a parameter  $n$  for the “level”  $\text{Pr}f_\tau(n, y, x)$ , which is the provability predicate for the theory axiomatized by  $\tau(n, x)$  defined below (hence  $\text{Con}_\tau(n)$  will be  $\neg \exists y \text{Pr}f_\tau(n, y, \overline{\ulcorner 1 = 0 \urcorner})$ ). Moreover, recall that “ $\phi_e(x) \simeq y$ ” can be formalized by a  $\Sigma_1$  formula  $\sigma(e, x, y)$ .

**Definition 41.** *A  $\Sigma_1$  formula  $\tau(n, x)$  is a verifiable enumeration for  $\tau(x)$  if PA proves the following:*

1.  $a = 1 \vee (\forall u(a \neq 2^u) \wedge \forall u(a \neq 3 \cdot 5^u)) \rightarrow (\tau(a, x) \leftrightarrow \tau(x))$
2.  $a \geq 1 \rightarrow (\tau(2^a, x) \leftrightarrow \tau(a, x) \vee \overline{\text{Con}_\tau(\dot{a})})$
3.  $\tau(3 \cdot 5^a, x) \leftrightarrow \tau(x) \vee \exists u \exists w (\sigma(e, u, w) \wedge \tau(w, x))$

**Theorem 86.** *For all  $\Sigma_1$  binumeration  $\tau(x)$ , there exists a  $\Sigma_1$  formula  $\tau(z, x)$  provably equivalent to the disjunction of the following:*

1.  $(z = \bar{1} \vee \forall u \leq z(z \neq 2^u \wedge z \neq s \cdot 5^u)) \wedge \tau(x)$
2.  $\exists u \leq z(u \neq \bar{0} \wedge z = 2^u \wedge (\tau(u, x) \vee x = \overline{\text{Con}_\tau(\dot{u})})$
3.  $\exists u \leq z(z = 3 \cdot 5^u \wedge (\tau(x) \vee \exists v \exists w (\sigma(u, v, w) \wedge (w \neq z \wedge \tau(w, x)) \vee (w = z \wedge x = v)))$

*Proof.* By using the fixed point theorem and partial truth predicates (recall that for all  $e \in \mathcal{O}$  and  $n$ ,  $\phi_e(n) \neq 3 \cdot 5^e$ ). QED

We have the following:

1. If  $a \in \mathcal{O}$ , the formula  $\tau(\bar{a}, x)$  gives the axioms of the theory  $S_a$  in the progression.
2.  $S_1 = \text{PA}$ .
3.  $S_{2^a} = S_a \cup \{\text{Con}_\tau(\bar{a})\}$ .
4. If  $\lambda$  denotes a limit ordinal, then  $S_\lambda = \bigcup_{d < \mathcal{O} \lambda} S_d$ .
5. Lastly, for all  $a, b \in \mathcal{O}$ , PA proves that if  $a <_{\mathcal{O}} b$ , then  $\tau(a, x) \rightarrow \tau(b, x)$ .

The *Turing progression* with enumeration  $\tau(e, x)$ , is the sequence  $\{S_e\}_{e \in \mathcal{O}}$  of theories enumerated by the formula  $\tau(e, x)$ .

**Theorem 87.** (Turing 1937) *For all progressions and all true  $\Pi_1^0$ -sentence  $\forall x \psi(x)$ , there exist a notation  $a \in \mathcal{O}$  such that  $\bar{a} = \omega + 1$  and  $\forall x \psi(x)$  is provable in  $S_a$ .*

*Proof.* We give an informal sketch of the proof, following Feferman (2006). Turing proceeded as follows: we denote for simplicity  $S(a)$  the code of the successor  $2^a$  and with  $\text{lim}(a)$  the code of the limit  $3 \cdot 5^a$ . Let  $n_{\mathcal{O}}$  the notation for the natural number  $n$ . We begin by defining (using the recursion theorem) function:

$$\phi_e(n) = \begin{cases} n_{\mathcal{O}} & \text{if for all } k \leq n, \psi(\bar{k}) \text{ is true} \\ S(\text{lim}(e)) & \text{if there is } k \leq n \text{ such that } \psi(\bar{k}) \text{ is false} \end{cases}$$

where  $\psi(x)$  is decidable. Note that if by hypothesis  $\forall x \psi(x)$  is a true  $\Pi_1$  sentence then for all  $n$ , we have  $\phi_e(n) = n_{\mathcal{O}}$  and therefore the sequence of values  $\phi_e(0), \phi_e(1), \phi_e(2), \dots$  is just the sequence  $0_{\mathcal{O}}, 1_{\mathcal{O}}, 2_{\mathcal{O}}, \dots$  and  $\text{lim}(e)$  is therefore an element of  $\mathcal{O}$  that denotes  $\omega$ . Reason inside  $S_{S(\text{lim}(e))}$ , checking in this theory that if  $S(\text{lim}(e))$  is consistent, then the sentence  $\forall x \psi(x)$  is true: indeed, suppose by contradiction that the sentence  $\forall x \psi(x)$  is false; hence we have that for some  $n$ , the sentence  $\psi(\bar{n})$  is false. But then the theory  $S(\text{lim}(e))$ , i.e. the union of the sequence  $S_{\phi_e(0)}, S_{\phi_e(1)}, S_{\phi_e(2)} \dots$  will be such that for some  $n$  and for all  $k \geq n$  (i.e. from a certain point onwards) we will have  $\phi_e(k) = S(\text{lim}(e))$ ; hence from a certain point onwards  $S_{S(\text{lim}(e))}$  and  $S(\text{lim}(e))$  will coincide, and therefore  $S(\text{lim}(e))$  will prove its own consistency (being  $\text{Con}(S(\text{lim}(e)))$  contained in  $S_{S(\text{lim}(e))}$ ); it follows from the second Gödel's result that  $S(\text{lim}(e))$  is inconsistent. But  $S_{S(\text{lim}(e))}$  actually proves the consistency of  $S(\text{lim}(e))$ . Ergo  $S_{S(\text{lim}(e))}$  proves  $\forall x \psi(x)$ . Observe that  $S(\text{lim}(e))$  denotes  $\omega + 1$ . QED

Turing, however, was not satisfied with this result, which in his view shifted the problem to determine whether a  $\Pi_1^0$ -sentence is true, to the problem considerably more complex to determine whether a number belongs to  $\mathcal{O}$ :

My completeness theorem [...] is completely useless for the purpose of producing proofs (Turing (1940)).

Turing also obtained a completeness result for  $\Pi_1$  statements via the transfinite iteration of the local reflection principle, in place of the consistency statement. Further results were obtained in Feferman (1962) that strengthen Turing's theorems, where progressions obtained starting from  $PA = S_0$  were studied, iterating the principle of *universal* reflection principle (see on p.125). It was proved, for instance, that for all true arithmetical sentences  $\theta$  there exists an  $a \in \mathcal{O}$ , denoting an ordinal less or equal to  $\omega^{\omega^{\omega+1}}$  such that  $S_a \vdash \theta$  (where on the contrary Turing's progression based on iteration of consistency statements is not complete for true  $\Pi_2$ -sentences). Feferman also proved that Turing's progression based on iteration of consistency statements is not complete for true  $\Pi_2$  statements.

However, it was disappointing the fact that although two notations  $a$  and  $b$  denote the same ordinal, this does not imply that  $S_a$  proves the same sentence that  $S_b$ . We say that this ordinal logic is *not invariant* and Turing actually showed that an ordinal logic cannot be both complete for  $\Pi_1$  sentences and invariant (see Feferman (2006)). In Franzen (2003) it is explained that the reason why a  $\Pi_1$  sentence  $\psi$  can only be proved at stage  $\omega + 1$  in a Turing consistency sequence is that at stage  $\omega$  of the construction a non-standard definition of the axioms can be introduced in such a sequence, depending on the definition of  $\phi_e$ . In Feferman's theorem a path (i.e. a subset  $P \subseteq \mathcal{O}$  linearly ordered by  $<_{\mathcal{O}}$  and such that if  $b \in P$  and  $c <_{\mathcal{O}} b$ , then  $c \in P$ ) of length  $\omega^{\omega^{\omega+1}}$  was generated, such that  $\bigcup_{a \in P} S_a$  is complete for all arithmetical sentences, but that completeness result essentially depends on the choice of the path and the result does not hold for other paths. Hence the invariance property still fails. In other words, these result are sensitive to the choice of notation. As a consequence of this, the crucial problem then became to discover natural conditions to impose on the choice of ordinal notations used to index theories in a progression. In an attempt to resolve this issue, in Kreisel (1960) it was required that progressions should satisfy a kind of "autonomy" requirement and the purpose was to generate the hierarchy of theories via a kind of boot-strapping process, through so-called *autonomous progressions*, that are in some sense self-justifying, being characterised by the fact that we are allowed to advance to the step  $S_a$ , only if we have obtained a proof that  $a \in \mathcal{O}$  in some previously accepted theory  $S_b$  where  $b <_{\mathcal{O}} a$ .

This line of research has had a strong impact in Proof Theory. The notion of autonomy was used by Kreisel in his proposals to characterize constructive philosophies of mathematics as finitism and predicativity by suitable autonomous progressions of theories, say  $\{F_a\}$  and  $\{R_a\}$ , respectively. While Kreisel (1960) established that a least upper bound for  $\tilde{\alpha}$  appearing in the first sequence is  $\varepsilon_0$ , Feferman (1964) and Schütte (1964, 1965) considered the so-called *Veblen hierarchy* of functions  $\phi_\alpha$ :

1.  $\phi_\alpha(\beta) = \omega^\beta$ , if  $\alpha = 0$ ,
2.  $\phi_\alpha$  enumerates the set  $\{\xi \mid \phi_\gamma(\xi) = \xi, \text{ for all } \gamma < \alpha\}$  if  $\alpha > 0$

and came to the conclusion that an ordinal  $\Gamma_0$  which is the least  $\alpha$  such that  $\phi_\alpha(0) = \alpha$  is the analogous limiting ordinal in the second progression. The received view, even if not agreed unanimously (see for instance Weaver (2005), that strongly disagree with this interpretation), is that this represents the smallest predicatively non-provable ordinal. See Beklemishev (1995), Franzen (2003), Franzen (2004) for a broader overview and Feferman's appendix to Takeuti (1987) to explore these developments further. On Feferman's predicativism see Crosilla (2017).

### 5.3. Propositional Provability Logic: classical vs. intuitionistic arithmetic

*Provability Logic* is a modal logic, in which the boxed formula  $\Box\phi$  ("it is necessary that  $\phi$ ") is interpreted as "it is provable (in some formal theory) that  $\phi$ ". Directly related to the second incompleteness theorem, it was a vast research programme that involved many Italian logicians and characterised the Sienese school of logic in particular from the 1970s to the 1990s. The literature too is vast, but here, in keeping with the slant of the book, we would like to draw attention to a particular chapter: the relations with intuitionist constructive logic. Why consider

intuitionist theories? In fact, this interest dates back to the origins of studies in this field, and not only on the part of the Dutch school. Later on, we will explain in more detail that the first, partial examples of “arithmetical completeness” results come from this area. However, the generalisation of Solovay’s seminal scientific achievement (see below), moving from classical mathematical theories to intuitionistic ones, although it has been a constant interest, especially on the part of the Dutch school (de Jongh, Visser, Iemhoff), has proved fraught with difficulties and has only made significant progress in very recent years. From the early years of this research programme, the issue also affected Italian school. What has been made clear is that the provability logic of HA is not a sublogic of that of PA. Let’s start in this regard with a little history. At the three seminars of the Siena’s annual meeting named “Incontri di logica matematica” in 1982, many talks were focused on the topic of *constructive logic and mathematics*. At the first seminar Franco Montagna gave a talk on the subject: *Solovay’s theorem and Heyting Arithmetic*. Solovay’s theorem is the main result in the field of so-called *Provability Logic*, and essentially states that the modal logic GL captures everything Peano arithmetic PA (or other fragments of arithmetic considered) can truthfully say about their own provability predicate. Heyting arithmetic HA is the intuitionistic analogous of classical Peano arithmetic (only the logic is different). In the subsequent years, research on these two topics intertwined. Montagna concluded his talk by posing these problems:

1. Give an axiomatization of the provability logic of Heyting’s intuitionistic arithmetic.
2. Determine whether this logic is decidable.

The whole problem remained unsolved for a long time (see Visser and Beklemishev (2006)) and important results have been achieved only recently. Research on provability logic was carried out mainly between the ’70s and ’90s in various universities: Prague, Moscow, Amsterdam, Utrecht, Siena, Oxford, Manchester. In particular, in Siena, in the early 70s, R. Magari and other people proposed an algebraic approach to the formal provability which led to the concept of diagonalizable algebra. Sambin and Ursini studied an intuitionistic version. The research in this field continued after the 90s at other universities, and especially in the Netherlands addressed to the provability logic of *intuitionistic arithmetic*. Those who know the modal logic will recognize in the Löb conditions, discussed in the framework of the second incompleteness theorem, modal axioms and rules which are reflected in the propositional Provability Logic GL (“Gödel-Löb logic”). This is a (classical) modal logic in which the box operator  $\Box$  is interpreted as the provability predicate  $\text{Pr}_T$ . GL can be given as an extension of the basic normal modal system K, where the distribution axiom  $\Box(\alpha \rightarrow \beta) \rightarrow (\Box\alpha \rightarrow \Box\beta)$  interprets the second derivability condition on p.123, the Necessitation rule:

$$\frac{\vdash \alpha}{\vdash \Box\alpha}$$

interprets the first of these conditions and Löb’s axiom  $\Box(\Box\phi \rightarrow \phi) \rightarrow \Box\phi$  (which implies the third condition) internalizes Löb’s theorem. This theory is incompatible with  $\Box\phi \rightarrow \phi$  (note that  $\Box(\Box\perp \rightarrow \perp)$  implies the provability of consistency  $\Box\neg\Box\perp$ ). It is (weakly) complete w.r.t. Kripke finite treelike models (hence is decidable), we are interested also in another kind of completeness: the arithmetical completeness, proved in Solovay (1976), of which Montagna (1979) gave a uniform version. Indeed, the most important results in this area are the so-called *fixed point theorem* (early results in the 70s as Bernardi (1975) or Sambin (1976)) and the fact that this logic was shown to be the logic of the formal proof predicate for many classical arithmetical theories: this is Solovay’s arithmetical completeness theorem (1976). What is the problem with intuitionistic arithmetic? Recall that HA and PA are formulated in the same language, with the same mathematical axioms and share many properties: they have the same provably recursive functions and are proof-theoretically equivalent and are equi-consistent, i.e. if HA is consistent, then PA is also consistent. Moreover, each is capable of expressing its own proof predicate and by Gödel’s results, if HA is consistent then neither HA nor PA can prove its own consistency (see e.g. Van Dalen (2001) and Avigad and Feferman (1998)). However, the provability logic of intuitionistic arithmetic goes beyond the intuitionistic version IGL of Löb’s logic: Solovay’s first

theorem does not hold in its immediate transposition, by replacing PA with HA and GL with IGL, this emerged clearly in the context of early algebraic researchs.

For reasons of space, we will not address here the thorny issue of Gentzen’s proof systems for Provability Logic. In general, a Gentzen approach to modal systems is itself problematic, except for a few systems. Therefore, more complex formal systems, such as hypersequents, are often preferred. The cut elimination theorem for GL was first proved in Valentini (1983). There was a long controversy surrounding the validity of this theorem, where sequents were formalized as sets rather than multisets or sequences, however, improvements and corrections were then made in Sasaki (2001), in Goré and Ramanayake (2012), and, applying the hypersequent method, in Poggiolesi (2009). Many early investigations on provability logic, especially in Siena, where instead of algebraic character. We just quickly mention them, before reformulating these results in perhaps more understandable logical and proof-theoretical terms. Magari in the 70s introduced the equational class of *Diagonalizable algebras* DA, boolean algebras equipped with an additional unary operator  $\tau(x)$  intended to “mimic” the provability predicate “ $x$  is provable”, and fulfills conditions that reflect the Löb’s conditions. Sambin and Ursini studied the *intuitionistic* version of these algebras, based on *Heyting algebras*, rather than boolean algebras. A primary example of diagonalizable algebra is the following. Let  $\mathcal{M}_{\text{PA}}$  be the *Lindenbaum algebra of Peano Arithmetic*, i.e. the well-known boolean algebra whose universe is the set of equivalence classes  $[\phi]$  of sentences of this theory, modulo provable equivalence in PA:

$$\phi \sim_{\text{PA}} \psi \text{ if and only if } \text{PA} \vdash \phi \leftrightarrow \psi$$

equipped with a Boolean algebra structure, where the operations join  $\sqcup$ , meet  $\sqcap$  and complementation  $-$  are defined as usual by:

1.  $[\phi] \sqcup [\psi] = [\phi \vee \psi]$
2.  $[\phi] \sqcap [\psi] = [\phi \wedge \psi]$
3.  $-[\phi] = [\neg\phi]$

where the square brackets denote the equivalence classes. This *is* a diagonalizable algebra, when enriched with an operator  $\tau(x)$  defined as  $\tau[\phi] = [\exists y \text{Prf}_{\text{PA}}(y, \ulcorner \phi \urcorner)]$ . Solovay’s celebrated theorem says that  $\mathcal{M}_{\text{PA}}$  is *functionally free* in the equational class DA of diagonalizable algebras. In other words, any identity true in  $\mathcal{M}_{\text{PA}}$  is true in *every* Diagonalizable algebras.

Since the beginning was raised the problem of extending this result to the intuitionistic case.

*Problem.* What, if in the above statement we replace PA with HA and diagonalizable algebras with *intuitionistic* diagonalizable algebra? Many problems arose immediately. Ursini in 1977 proved that  $\mathcal{M}_{\text{HA}}$  (the Lindenbaum algebra of HA) *is not functionally free* in the class of diagonalizable intuitionistic algebras (with Heyting’s algebra in place of Boolean algebras): there are principles, as the law

$$\tau(x + y) \leq \tau(\tau(x) + \tau(y))$$

that holds in the diagonalizable Lindenbaum algebra  $\mathcal{M}_{\text{HA}}$ , but not in all intuitionistic diagonalizable algebras. This is the problem on which Montagna’s talk focused. The reasons why this happens were clearly highlighted by Montagna in the above mentioned conference:

... one of the reasons why this happens is that there are classically invalid rules that are valid in HA. By arithmetizing such rules, we obtain properties of the provability predicate for HA that are not deducible from the identities of intuitionist DA’s.

We will better understand this problem after having reformulated Solovay’s results in a general proof theoretical framework. From now on  $\mathbb{T}$  will denote a “reasonable” (i.e. we assume that it is at least recursively enumerable and  $\Sigma_1^0$  – *sound*) extension of  $\text{ID}_0 + \text{exp}$  (a formalization of the so-called *Elementary Arithmetic*, see at p. 192).

**Definition 42.** An arithmetical interpretation in  $\mathbb{T}$  is a function  $*$  from the modal language to the arithmetical language that:

1. associates to each propositional variable  $p_i$  a formula  $p_i^*$  of the language of arithmetical theory  $\mathsf{T}$ .
2. commutes with connectives,
3.  $\perp^* = (1 = 0)$ .
4.  $(\Box\psi)^* = \exists y \text{Prf}_{\mathsf{T}}(y, \overline{\ulcorner \psi^* \urcorner})$ .

A propositional formula  $\phi$  is arithmetically valid in  $\mathsf{T}$  iff for all arithmetical interpretations in  $\mathsf{T}$  as above, we have that  $\phi^*$  is provable in  $\mathsf{T}$ .

The Provability Logic of  $\mathsf{T}$ , that we denote  $\text{PL}_{\mathsf{T}}$ , is the set of propositional modal formulas arithmetically valid in  $\mathsf{T}$ . The most remarkable results in this area are the mentioned Solovay's arithmetical completeness theorems:

1. *Solovay's First theorem* (1976).  $\text{PL}_{\mathsf{T}} = \text{GL}$
2. *Solovay's Second theorem* (1976). The set of modal propositional  $\phi$  such that  $\phi^*$  is true in the standard model of arithmetic, for all interpretations  $*$  in  $\mathsf{T}$  (here, for any *sound* extension  $\mathsf{T}$  of  $\text{I}\Delta_0 + \text{exp}$ ) are those provable in  $\text{GL}^- + \Box\phi \rightarrow \phi$ , where  $\text{GL}^-$  is  $\text{GL}$  minus the necessitation rule.
3. *Montagna's uniform theorem* (1979). There is an interpretation  $*$  in  $\mathsf{T}$  such that for all  $\phi$ ,

$$\mathsf{T} \vdash \phi^* \text{ iff } \text{GL} \vdash \phi$$

Other proofs of the uniform theorem were obtained independently using different methods in the same year or in the years immediately following by logicians as Artemov, Visser, Boolos, and Avron. See, for example, Boolos (2008) 132–136. However, Montagna's proof is the most original and has an algebraic flavour, as does Solovay's original proof. What is actually proved is the following equivalent statement: if  $\mathsf{T}$  is a theory as above, and has infinite characteristic (see below for a definition of this concept), then the free Magari algebra of countably many generators is embeddable in the provability algebra of  $\mathsf{T}$ . See Beklemishev and Flaminio (2016) for a detailed discussion.

The logic of provability  $\text{GL}$  actually has several semantics, including Kripke semantic and a fundamental arithmetic interpretation, which we will focus on in particular. With regard to Kripke semantics, we recall that a Kripke structure is a pair  $\langle W, R \rangle$ , where  $W$  is a nonempty set and  $R$  is a binary relation on it; an evaluation in a Kripkean structure is a function that associates each propositional variable  $p$  with a set  $V(p) \subseteq W$  that can be seen as the set of states in which  $p$  is considered true. By a *model* we mean the triple  $\mathcal{M} = \langle W, R, V \rangle$  on which we define a relation  $x \Vdash_{\mathcal{M}} \phi$  to be read as: “ $\phi$  is true in state  $x$  in model  $\mathcal{M}$ ”, inductively, as follows:

- (a)  $x \Vdash_{\mathcal{M}} p$  if and only if  $x \in V(p)$
- (b)  $x \Vdash_{\mathcal{M}} \neg\phi$  if and only if it is not true that  $x \Vdash_{\mathcal{M}} \phi$
- (c)  $x \Vdash_{\mathcal{M}} \phi \wedge \psi$  if and only if  $x \Vdash_{\mathcal{M}} \phi$  and  $x \Vdash_{\mathcal{M}} \psi$
- (d)  $x \Vdash_{\mathcal{M}} \phi \vee \psi$  if and only if  $x \Vdash_{\mathcal{M}} \phi$  or  $x \Vdash_{\mathcal{M}} \psi$
- (e)  $x \Vdash_{\mathcal{M}} \phi \rightarrow \psi$  if and only if, if  $x \Vdash_{\mathcal{M}} \phi$ , then  $x \Vdash_{\mathcal{M}} \psi$
- (f)  $x \Vdash_{\mathcal{M}} \Box\phi$  if and only if for all  $y \in W$ , if  $xRy$ , then  $y \Vdash_{\mathcal{M}} \phi$ .

Actually, a Kripke model for GL is such that  $R$  is a converse well-founded (hence irreflexive) strict partial ordering on  $W$ . We assume that every model has a *root*. A formula is valid in a model, if it is forced at the root. In Segerberg (1971) was demonstrated that GL is weakly complete with respect to the class of structures that are finite trees, namely that a formula is provable in GL if and only if it is valid in all finite treelike models, a property that we will use in Solovay's theorem<sup>1</sup>.

*Solovay's proof for classical theories.* Solovay's technique consists in defining a function  $F$ , which constitutes an immersion of a Kripke model into arithmetic: the behaviour of this function is often described, at a heuristic level, as that of a refugee who moves from country to country, obtaining permission to cross the border on condition that he promises not to settle permanently in the country of arrival: if the refugee is not allowed to return twice to the same country, there must nevertheless be a country where he or she can settle permanently. Therefore, if the refugee is honest, he or she will never have to leave the country of origin. The story is told in Artemov and Beklemishev (2005), which we also follow for its version of the proof.

Solovay's function  $F$  (containing an obvious circularity) can be defined by appealing to the formalised principle of recursion (Kleene), and it is introduced in Solovay's original work, even if a series of simplifications can be applied to the original definition.

Suppose that  $\mathcal{M} = \langle W, R, r, \Vdash \rangle$  is a finite model, where without loss of generality we assume that  $W = \{1, 2, 3, \dots, n\}$  and that  $r = 1$ , where  $R$  does not necessarily coincide with the natural order of  $\mathbb{N}$ ; for purely technical reasons, another node 0 is generally added so that  $0Ri$ , for every  $i \leq n$ , without assuming anything about the forcing in it.

Solovay's (partially recursive) function  $F : \mathbb{N} \rightarrow W \cup \{0\}$  is informally defined as follows:

- (a)  $F(0) = 0$
- (b) At step  $x + 1$ , suppose that  $F(x)$  has already been defined: check whether  $x + 1$  is the code for a proof of the fact that  $\lim_{k \rightarrow \infty} F(k) \neq z$ , for some  $z$  accessible to  $F(x)$ : if YES, then  $F(x + 1) = z$ ; if NO  $F(x + 1) = F(x)$ .

where  $\lim_{k \rightarrow \infty} F(k) = z$  is an abbreviation for  $\exists x \forall y > x \psi_F(x, y)$ , and  $\psi_F(x, y)$  is the  $\Sigma_1^0$ -graph of  $F$ .

Let us remember as we defined the hierarchy of theories  $S_0, S_1, S_3 \dots$  obtained by iterated addition of consistency statements of the previous chapter. The *characteristic* of  $S = S_0$  is the last number  $n$  such that  $S_n$  is inconsistent. If such a number does not exist, we say that it has an *infinite characteristic*: for instance, if  $S$  is  $\Sigma_1$ -sound, then it has an *infinite characteristic*. We establish Solovay's result for axiomatizable extensions  $S$  of *Elementary Arithmetic*  $\text{I}\Delta_0 + \text{exp}$  with infinite characteristic.

Let us abbreviate the expression " $\exists m \forall n > m (h(n) = z)$ " as " $L_z$ ".

**Lemma 32.** *The following properties of the function  $F$  are provable in the Elementary Arithmetic  $\text{I}\Delta_0 + \text{exp}$ :*

- (a)  $\bigvee_{z \in W \cup \{0\}} L_{\bar{z}}$
- (b)  $\forall u \forall v (u \neq v \rightarrow \neg(L_u \wedge L_v))$
- (c)  $L_{\bar{m}} \wedge \bar{m} R \bar{s} \rightarrow \neg Pr(\overline{\neg L_{\bar{s}}})$
- (d)  $L_{\bar{m}} \wedge \bar{m} \neq \bar{0} \rightarrow Pr(\overline{\bigvee_{m R w} L_{\bar{w}}})$

<sup>1</sup> Recall that *strong* completeness does not hold, because this logic is not compact.

*Proof.* (see Artemov and Beklemishev (2005) 481-482).

QED

The idea behind the proof in Solovay (1976) is to now provide an arithmetic simulation of Kripke's model; let us therefore define a *Solovay interpretation*  $*$ , for a model  $\langle \{0, 1, 2, 3 \dots n\}, 1, \Vdash \rangle$ , the one in which, in particular, propositional variables are interpreted as follows:

$$p^* = \bigvee \{L_x | x \Vdash p \text{ and } 0 \leq x \leq n\}$$

where the statements  $L_x$  assume the role of the nodes  $x$  of the model, and where  $*$ , in our case, is an *arithmetical interpretation in the theory S*. The empty disjunction is  $\bar{0} = \bar{1}$ . The fundamental step in proving the arithmetical completeness theorem consists in demonstrating the *faithfulness* of the interpretation, i.e. this result:

**Lemma 33.** *For  $1 \leq x \leq n$  and  $*$  as above, the following apply:*

- (a) *If  $x \Vdash \psi$  then  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{x}} \rightarrow \psi^*$*
- (b) *If  $x \not\Vdash \psi$  then  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{x}} \rightarrow \neg\psi^*$*

Before proving this lemma, note that arithmetic completeness follows directly from it: if  $\phi$  is not a theorem of  $\mathbf{GL}$ , hence it is false at the root of a finite treelike model  $W = \{1, 2, 3, \dots, n\}$ . We add a new root 0 such that  $0Rm$  for all  $m \in W$  (forcing on 0 is arbitrarily defined). From the above lemma, if  $1 \not\Vdash \phi$ , then  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{1}} \rightarrow \neg\phi^*$ . Now, if  $\mathbf{S} \vdash \phi^*$  we have  $\mathbf{S} \vdash \neg L_{\bar{1}}$  and by the point c) of Lemma 32 we also obtain  $\neg L_{\bar{0}}$ . It follows that  $L_{\bar{m}}$  must be true, for some  $m \in W$ . Observe that  $m \Vdash \Box^{d(m)+1} \perp$ , where  $d(m) = \sup\{d(s) + 1 | mRs\}$  and  $\Box^{d(m)+1}$  denotes a block of  $d(m) + 1$  boxes. Hence  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{m}} \rightarrow (\Box^{d(m)+1} \perp)^*$  and therefore that  $(\Box^{d(m)+1} \perp)^*$  is true, against the hypothesis that  $\mathbf{S}$  had *infinite* characteristic.

Let us now prove the lemma:

*Proof.* Induction on the complexity of  $\phi$ . If  $\phi = p$ , the result follows from the definition; we leave the Boolean cases as an exercise and come to the fundamental case, where  $\phi = \Box\psi$ ; note that  $x \Vdash \Box\psi$  implies that for every  $y$ , if  $xRy$ , then  $y \Vdash \psi$ , from which for every  $y$ , if  $xRy$  then  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{y}} \rightarrow \psi^*$  by inductive hypothesis and therefore:

$$\mathbf{I}\Delta_0 + \text{exp} \vdash \bigvee \{L_{\bar{y}} \rightarrow \psi^* | xRy\}$$

Hence  $\mathbf{I}\Delta_0 + \text{exp} \vdash Pr(\overline{\bigvee \{L_{\bar{y}} | xRy\}}) \rightarrow Pr(\overline{\psi^*})$  by logic and by the derivability conditions. Lastly  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{x}} \rightarrow Pr(\overline{\psi^*})$ , by the point d) of the previous lemma. If on the contrary  $x \not\Vdash \Box\psi$ , then there exists  $y$  such that  $xRy$ , but  $y \not\Vdash \psi$ . It follows that there exists  $y$  such that  $xRy$  and  $\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{y}} \rightarrow \neg\psi^*$ , or, equivalently, by contranomial  $\mathbf{I}\Delta_0 + \text{exp} \vdash \psi^* \rightarrow \neg L_{\bar{y}}$ ; hence, for some  $y$ ,  $xRy$  and  $\mathbf{I}\Delta_0 + \text{exp} \vdash Pr(\overline{\psi^*}) \rightarrow Pr(\overline{\neg L_{\bar{y}}})$ . Taking the contrapositive of this implication, from point c) of the previous lemma, it follows that:

$$\mathbf{I}\Delta_0 + \text{exp} \vdash L_{\bar{x}} \rightarrow \neg Pr(\overline{\psi^*})$$

QED

This concludes the proof of Solovay's *first* theorem, and from it we draw the conclusion that  $\mathbf{GL}$  axiomatises the modal axiom schemes *provable* in  $\mathbf{S}$ : but what about the *true* schemes in the standard model? The answer is provided by the modal system  $\mathbf{GL}^-$ , obtained by removing the necessitation rule from  $\mathbf{GL}$  and adding the schema  $\Box\alpha \rightarrow \alpha$ . Note incidentally that if we did not remove the necessitation rule, from  $\Box\perp \rightarrow \perp$ , we could derive  $\Box(\Box\perp \rightarrow \perp)$  and by Löb  $\Box\perp$  and finally  $\perp$ . The *second Solovay theorem* follows from the fact that these propositions are equivalent:

- (a)  $GL^- \vdash \alpha$
- (b)  $GL \vdash \bigwedge_{\square\beta \in Sub(\alpha)} (\square\beta \rightarrow \beta) \rightarrow \alpha$ , where  $Sub(\alpha)$  is the set of subformulas of  $\alpha$ .
- (c)  $\alpha$  is true at the root of every  $\alpha$ -sound model (i.e., whose root forces  $\square\beta \rightarrow \beta$ , for every  $\square\beta \in Sub(\alpha)$ ).
- (d)  $\mathbb{N} \models \alpha^*$ , for every arithmetic interpretation  $*$ .

For further details and information, see Boolos (2008), ch.9, Artemov and Beklemishev (2005), ch.3 and De Jongh and Japaridze (1998), ch.3. Among the improvements that have been made we like to remember this: although the above proof employed the recursion theorem, in De Jongh, Jumelet and Montagna (1991) it is shown that the use of this theorem actually is not necessary. Using the recursion theorem makes the procedure more intuitive, but, according to these logicians, it “adds to the mystery of the proof”. The alternative proof proposed is closer to the spirit of modal logic, replaces the recursion theorem by Gödel’s diagonalization lemma and is also applicable to another system of modal logic, yielding the arithmetical completeness of the so-called Rosser logic of Gauspari-Solovay with respect to extensions of the *Elementary Arithmetic*. It is still not known what is the provability logic of the important weak fragment  $S_2^1$ , discussed in section 7.3.

We conclude the general introduction here and ask ourselves: what happens if we consider *intuitionistic* theories of arithmetic rather than classical theories? We like to point up that already at the end of the 60s, a particular kind of Solovay style theorem had been shown in De Jongh (1969). It was the first theorem of this kind and concerned the intuitionistic version HA of Peano Arithmetic. De Jongh’s theorem provides an answer concerning the *ordinary* propositional logic of the intuitionistic arithmetical theory HA, i.e. the box-free part of the provability logic of HA, and it is the first kind of “arithmetical completeness” theorem discovered. We state here this theorem for IPC (the intuitionistic propositional calculus).

**Theorem 88.** (De Jongh 1970) *The following statements are equivalent, for all propositional formulas  $\phi(p_0, \dots, p_n)$ :*

- (a)  $IPC \vdash \phi(p_0, \dots, p_n)$
- (b)  $HA \vdash \phi(B_0, \dots, B_n)$ , for all arithmetical sentences  $B_0, \dots, B_n$ .

A proof can be found in Smorynski (1973). In particular since an arithmetical interpretation is in fact a substitution of the above kind, this means that  $\phi(p_0, \dots, p_n)$  is in the provability logic of HA. Refinements were given by several logicians. In particular Friedman (1973) proved that the choice of  $B_0, \dots, B_n$  can be actually made uniformly.

**Theorem 89.** (Friedman 1973) *There exists a computable sequence  $B_0, \dots, B_n$  of arithmetical sentences, such that for all  $\phi(p_0, \dots, p_n)$ :*

$$IPC \vdash \phi(p_0, \dots, p_n) \text{ if and only if } HA \vdash \phi(B_0, \dots, B_n)$$

Further extensions of these results are obtained by adding to HA other principles e.g. accepted by Russian constructivists (see Trolestra and Van Dalen (1988)), but considered problematic by other constructivists. Typically, the *Extended Church Thesis*  $ECT_0$  and *Markov’s principle* MP. Smorynski proved that the propositional logic of  $HA + MP$  is still IPC, where MP is Markov’s principle:

$$\forall x(A(x) \vee \neg A(x)) \wedge \neg \forall x \neg A(x) \rightarrow \exists x A(x)$$

Gavrilenko proved that if we replace MP with  $ECT_0$  we still get IPC, where  $ECT_0$  is the extended Church thesis:

$$\forall x(A(x) \rightarrow \exists y B(x, y)) \rightarrow \exists z \forall x(A(x) \rightarrow \exists u(T(z, x, u) \wedge B(x, U(u))))$$

where  $T(z, x, u)$  is Kleene's predicate of recursion theory and  $U(u)$  the result of computation  $u$ , and  $A(x)$  does not contain  $\vee$ , and  $\exists$  only in front. But it is not known what is the propositional logic of *Markov's Arithmetic*, i.e.  $\text{HA} + \text{MP} + \text{ECT}_0$ . However it is known that it is *not* IPC. For instance, if  $\alpha = \neg p \vee \neg q$ , then the following:

$$((\neg\neg\alpha \rightarrow \alpha) \rightarrow (\neg\neg\alpha \vee \neg\alpha)) \rightarrow (\neg\neg\alpha \vee \neg\alpha)$$

is in the propositional provability logic of Markov's arithmetic, but is not a theorem of IPC. Therefore we come to the intermediate logics. The interesting purpose of De Jongh, Verbrugge and Visser (2011) is instead that of strengthen the propositional logic, rather than the arithmetical theory, and then to consider a large class of intermediate logics. The intermediate (or *superintuitionistic*) logics are logics between IPC and CPC. The authors conjecture that, if  $L$  is an intermediate logic and  $\text{HA}_L$  is obtained by adding this logic to the intuitionistic arithmetic, then  $L$  is the propositional logic (in the sense of de Jongh's theorem) of  $\text{HA}_L$ . Actually they proved the conjecture only for logics satisfying the so called *finite frame property*.

Another line of research investigates the consequences of strengthening the *propositional logic* (while remaining in the constructive field), rather than the arithmetical theory. Coming back to the problem of *Provability Logic for Intuitionistic Arithmetic*, the problem is to find an intuitionistic modal logic  $I$  such that  $I = \text{PL}_{\text{HA}}$ , namely, a modal formula  $\phi$  is derivable in  $I$  if and only if  $\phi^*$  is derivable in  $\text{HA}$ , for all arithmetical interpretation  $*$ . We remark that Provability logics in general are *not monotone*, i.e. if a theory  $T$  extends a theory  $V$ , this is *not generally true* for the relative provability logics. For this reason Solovay's results in the classical case (e.g. the first theorem still hold for all  $\Sigma_1^0$  and r.e. theories extending, or interpreting, *Elementary Arithmetic*) are very surprising: Solovay's theorems concerning the classical logic are very *stable*. On the contrary the situation for constructive arithmetical theories is different. Different constructive theories may have different logic. What principles belong to  $\text{PL}_{\text{HA}}$ ? Since the standard proof predicate for  $\text{HA}$  and the relative Löb derivability conditions are derivable already in the intuitionistic version of *Elementary Arithmetic*, we conclude that the provability logic for  $\text{HA}$ , must at least contain the modal axioms of Löb's logic  $\text{GL}$  to IPC. The provability logic of Heyting arithmetic contains in general *principles that the provability logic of Peano Arithmetic does not share* and therefore is not a sublogic of that of Peano Arithmetic (non-monotonicity). For example, Leivant's principle (see below). On the other hand, there are classical provability principles that cannot be accepted by the intuitionists, e.g.  $\Box(p \vee \neg p)$ . What do we know until today? Some examples:

- (a) The formalization of the so-called *Markov's rule* (an intuitionistically *admissible* rule) is a principle of the Provability logic of  $\text{HA}$ :

$$\Box\neg\neg\Box\phi \rightarrow \Box\Box\phi$$

- (b) Since, on the contrary, the formalization of the intuitionistic disjunction property:

$$\Box(\phi \vee \psi) \rightarrow (\Box\phi \vee \Box\psi)$$

is *not* provable in  $\text{HA}$  (H. Friedman), one can add in its place a (provable in  $\text{HA}$ ) weak form of this principle, due to Leivant:

$$\Box(\phi \vee \psi) \rightarrow \Box(\Box\phi \vee \psi)$$

- (c) However Leivant's axiom is *not* a theorem of  $\text{GL}$ : in classical framework it allows to derive the iterated inconsistency statement  $\Box\Box\perp$ .

After Visser had found partial results concerning the letterless fragment (i.e. based only on constants  $\perp, \top$ ), in particular Ardeshtir and Mojtabehi (2014) provided an answer to the problem for the  $\Sigma_1^0$  provability logic of  $\text{HA}$ , while Zoethout and Visser (2019) provided an alternative route to this problem.

**Definition 43.** *Let us call an arithmetical interpretation  $*$  a  $\Sigma_1^0$ -interpretation, if and only if for all propositional variables  $p_i$ , we have that  $p_i^*$  is  $\Sigma_1^0$ , i.e. has the form  $\exists x\theta$ , where  $\theta$  is atomic, or contains only connectives and bounded quantifiers.*

In 1990 Visser had previously introduced an algorithm to associate to propositional  $\phi$  a particular form  $\phi^+$ , called NNIL-form (No Nested Implications to the Left, as for instance  $(p \rightarrow (q \rightarrow \perp)) \vee (q \rightarrow p)$ ). In some sense  $\phi^+$  is the “best approximation from below” of  $\phi$ , in the sense of this theorem.

**Theorem 90.**  $\text{IPC} \vdash \phi^+ \rightarrow \phi$ . *Moreover, if  $\text{IPC} \vdash \alpha \rightarrow \phi$ , then also  $\text{IPC} \vdash \alpha \rightarrow \phi^+$ .*

This class has been studied independently of our problem, due to its interesting properties in Visser, de Jongh, Van Benthem and Renardel de Lavalette (1995). This algorithm has been extended to modal formulas in Ardeshir and Mojtaehedi (2014), that introduced the class TNNIL: “no  $\rightarrow$  in the left side of an implication, except those in the scope of a  $\Box$ ” (example,  $\Box(p \rightarrow q) \rightarrow q$  is in this class, but  $(p \rightarrow q) \rightarrow q$  is not). The result is the following:

**Theorem 91.** (Ardeshir and Mojtaehedi 2014) *A modal formula  $\phi$  is in the  $\Sigma_1^0$ -provability logic of HA, if and only if  $\text{IGLC} \vdash \phi^+$ , where IGLC is the intuitionist version of the Löb logic GL, plus the completeness principle  $\psi \rightarrow \Box\psi$ .*

Other directions of the research involve the notion of relative interpretability; actually, all known principles of  $\text{PL}_{\text{HA}}$  are derivable in a bi-modal system of interpretability logic introduced by Albert Visser and studied in particular in Iemhoff (2003). This logic is based on a binary modal operator  $\alpha \triangleright \beta$  (where  $\Box\alpha = \top \triangleright \alpha$ ) where an arithmetical interpretation  $*$  of it is the formalization in the language of arithmetic of the fact that if  $\text{HA} \vdash \sigma \rightarrow \alpha^*$ , then  $\text{HA} \vdash \sigma \rightarrow \beta^*$ , for all  $\sigma$  belonging to  $\Sigma_1$ . Subsequent work by logicians such as Visser and de Jongh led to an axiomatisation of this modal logic named IPH and the demonstration of the arithmetical soundness of this axiomatisation. Iemhoff conjectured that it was also arithmetically complete, i.e. that  $\text{IPH} = \text{PL}_{\text{HA}}$ , but this is still open.

As far as we know, a definitive solution to the problem raised by Franco Montagna in the mentioned conference has finally been achieved recently, largely relying on the concepts and results we have briefly summarised. Mojtaehedi’s theorem is contained in Mojtaehedi (2022). It is a not yet published work of frightening complexity. Axioms in particular are rather complex. Roughly speaking,  $\text{PL}_{\text{HA}}$  consists in IGL, i.e. Löb’s logic of provability on an intuitionistic basis formulated in a non standard way, plus all formulas  $\Box\alpha \rightarrow \Box\beta$  for  $\alpha$  and  $\beta$  that satisfy a condition inspired by Iemhoff’s definition of intuitionistic interpretability:  $\text{IGL} \vdash \sigma \rightarrow \alpha$ , then  $\text{IGL} \vdash \sigma \rightarrow \beta$ , for all  $\sigma$  belonging to a set of formulas that can be *projected to a NNIL formula*, a rather technical notion coming from the theory of intuitionistic unification studied in Ghilardi (1999).

#### 5.4. First-order Provability Logic for classical arithmetic

After the results of arithmetic completeness of the logic of provability at the classical propositional level that we have briefly illustrated, it seemed natural to raise the question of the characterisation of the logic of *predicative* provability. We thought it is interesting to provide some information on a more advanced but not very well known topic, namely the first-order quantified version of the Provability Logic, which is in fact a first-order modal logic, subject around which there have been several philosophical controversies. The famous logician Ruth Barcan Marcus (1921-2012) published her first paper on Quantified Modal Logic (QML) in 1946; in 1947 she extended QML to the second order. In these works for the first time is investigated the relationship among modal operators  $\Box, \Diamond$  (“is necessary”, “is possible”) and quantifiers  $\forall, \exists$ , and between *de re* modalities  $\forall x\Box\phi$  (“for all  $x$ , necessarily  $\phi$ ”) and *de dicto* modalities  $\Box\forall x\phi$  (“is necessary that for all  $x$ ,  $\phi$ ”). In she placed among the axioms of her system, the schema  $\forall x\Box\psi(x) \rightarrow \Box\forall x\psi(x)$  that connect those modalities.

The main opponent of QML was Quine, that rejected it since, in his opinion, combination of quantifiers and modalities produced “unintelligible” results. Quine’s objections are logical arguments based on failure of substitution: the sentence ‘8 is necessarily greater than 7’ is true, while the sentence “the number of planets is necessarily greater than 7” is false: but the latter was obtained from the first by substitutipn of the coreferential term “the number of planets” in place of term ‘8’. It follows that such occurrences of singular terms are not “purely referential”, and therefore quantification into modal context is unintelligible. Quine had also a metaphysical objection: quantification in modal contexts commits us to accepting essentialism, that he refused as indefensible. Leaving this debate in the background (and not addressing the issue of propositional quantifiers in Provability Logic), we simply highlight the role in this context of the debated Barcan Marcus axioms:

The converse Barcan formula says that, as we move to an alternative situation, nothing passes out of existence. The Barcan formula says that, under the same circumstances, nothing comes into existence. The two together say the same things exist no matter what situation (Fitting and Mendelsohn (1998) 114).

The language of QGL, the first-order quantified version of Gödel-Löb logic, adds the symbol  $\Box$  to first order logic without identity, constants or functional symbols. Rules and axioms are those of GL plus those of predicate calculus, for all QGL formulas.

**Definition 44.** *An arithmetical interpretation for the first-order case is given as follows:*

- (a) *For all atomic formula  $P(x_0, \dots, x_n)$  of the language of QGL, the formula  $(P(x_0, \dots, x_n))^*$  is a formula of the language of arithmetic with the same free variable.*
- (b)  $(\phi \rightarrow \psi)^* = (\phi^* \rightarrow \psi^*)$  (analogously for  $\vee$  and  $\wedge$ ).
- (c)  $(\neg\phi)^* = \neg(\phi^*)$ .
- (d)  $(\exists x\psi)^* = \exists x(\psi^*)$  (analogous for  $\forall$ ).
- (e)  $(\Box\psi)^* = \exists y \text{Prf}_{PA}(y, \overline{\Box\psi^*})$ .

We say that a formula  $\psi$  of the language of QGL is *always provable*, iff for all arithmetical interpretations  $*$ ,  $\psi^*$  is provable in PA. We say that a formula  $\psi$  of the language of QGL is *always true*, iff for all arithmetical interpretations  $*$ ,  $\psi^*$  is true in the standard model  $\mathbb{N}$ . As far as propositional theories GL and GLS (the theory axiomatized by all theorems of GL and that replace the necessitation rule with the schema  $\Box A \rightarrow A$ ) the situation is troublefree:

- (a) GL axiomatizes the class of *always provable* sentences.
- (b) GLS axiomatizes the class of *always true* sentences.

Both theories are decidable and therefore so are the above-mentioned classes. What is the status of the Barcan Marcus schema BS,  $\forall x\Box\psi(x) \rightarrow \Box\forall x\psi(x)$ ? Actually this formula is not *always true* (neither *always provable*), but its converse CBS:

$$\exists x\Diamond\psi \rightarrow \Diamond\exists x\psi$$

is *always provable*. The following is a correct argument (see Smorynski (1987)). Recall that  $RFN_{\Sigma_n^0}(\mathbb{T})$  is the following reflection schema:

$$\forall x_0, \dots, \forall x_n (\text{Pr}_{\mathbb{T}}(\overline{\Box\phi(x_0, \dots, x_n)}) \rightarrow \phi(x_0, \dots, x_n))$$

where  $\phi \in \Sigma_n^0$  (analogously for  $\phi \in \Pi_n^0$ ). For  $n \geq 1$ , we show that  $RFN_{\Sigma_n^0}(\mathbb{T})$  is equivalent to  $RFN_{\Pi_{n+1}^0}(\mathbb{T})$  by means of CBS: let  $\phi(x, y) \in \Sigma_n^0$  and suppose that  $\text{Pr}_{\mathbb{T}}(\overline{\Gamma \forall y \phi(\dot{x}, y)})$ ; hence, since we have CBS, it follows  $\forall y \text{Pr}_{\mathbb{T}}(\overline{\Gamma \phi(\dot{x}, \dot{y})})$ . For  $Rfn_{\Sigma_n^0}(\mathbb{T})$  and the predicate calculus we conclude  $\forall y \phi(x, y)$ .

$$\begin{aligned}
 \text{QGL} \quad & \vdash \Box(\forall x \alpha(x) \rightarrow \alpha(u)) \\
 & \vdash \forall u \Box(\forall x \alpha(x) \rightarrow \alpha(u)) \\
 & \vdash \forall u (\Box \forall x \alpha(x) \rightarrow \Box \alpha(u)) \\
 & \vdash (\Box \forall x \alpha(x) \rightarrow \forall u \Box \alpha(u))
 \end{aligned} \tag{1}$$

Analogously  $\text{QGL} \vdash \exists u \Box \alpha(u) \rightarrow \Box \exists u \alpha(u)$ . From these result it follows that a kind of formalized completeness  $\alpha \rightarrow \Box \alpha$  holds in  $\text{QGL} + \text{BS}$  for *any* quantified formula  $\alpha$ .

**Theorem 92.** *The following are equivalent in PA:*

- (a)  $\neg \text{Con}(\text{PA})$
- (b)  $\phi \rightarrow \text{Pr}_{\text{PA}}(\overline{\Gamma \phi})$ , for closed  $\phi$ .
- (c)  $\text{Pr}_{\text{PA}}(\overline{\Gamma \phi}) \vee \neg \text{Pr}_{\text{PA}}(\overline{\Gamma \phi})$
- (d)  $\forall x \text{Pr}_{\text{PA}}(\overline{\Gamma \phi(\dot{x})}) \rightarrow \text{Pr}_{\text{PA}}(\overline{\Gamma \forall x \phi(x)})$

where  $\text{Con}(\text{PA})$  stands as usual for  $\neg \exists y \text{Prf}_{\text{PA}}(y, \overline{\Gamma 1 = 0})$ . The (BS) schema is in this context known as “ $\omega$ -completeness schema”. Recall that according to the Second Gödel’s theorem, under the assumption of consistency of the theory, both,  $\text{Con}(\text{PA})$  and  $\neg \text{Con}(\text{PA})$  are unprovable. Ergo: none of the above principles is provable! For example, let us see that  $\text{PA} + (\text{BS})^* \vdash \neg \text{Con}(\text{PA})$ . But we have also, for all  $\psi$ ,  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \text{Pr}_{\text{PA}}(\dot{y}, \overline{\Gamma \psi(\dot{x})})} \rightarrow \psi(\dot{x}))$ , from which follows in particular:

$$\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \neg(1 = 0) \rightarrow \neg \text{Prf}_{\text{PA}}(\dot{y}, \overline{\Gamma 1 = 0})})$$

Since  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \neg(1 = 0)})$ , by Löb conditions and manipulation of quantifiers we get  $\text{PA} \vdash \forall y \text{Pr}_{\text{PA}}(\overline{\Gamma \neg \text{Prf}_{\text{PA}}(\dot{y}, \overline{\Gamma 1 = 0})})$ . If (BS)\* holds, it follows:

$$\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \forall y \neg \text{Prf}_{\text{PA}}(y, \overline{\Gamma \neg(1 = 0)})})$$

Hence  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \exists y \text{Prf}_{\text{PA}}(y, \overline{\Gamma (1 = 0)})} \rightarrow 1 = 0)$ , in a few steps, from which  $\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma 1 = 0})$ , namely  $\neg \text{Con}(\text{PA})$ , by the formalized Löb theorem, according to which, for all  $\psi$ :

$$\text{PA} \vdash \text{Pr}_{\text{PA}}(\overline{\Gamma \text{Pr}_{\text{PA}}(\overline{\Gamma \psi})} \rightarrow \psi)$$

It is now interesting to understand if the collection of always true or the always provable sentences can be characterized axiomatically. If we denote  $\text{QPL}_{\text{PA}}$  the first-order provability logic of PA, we will see that actually  $\text{QGL} \subset \text{QPL}_{\text{PA}}$ . From Vardanyan’s theorem it follows that cannot have not even a recursive axiomatization.

**Theorem 93.** *The following hold:*

- (a) (Vardanyan 1986) *The class of all predicate modal formulas always provable is  $\Pi_2^0$ -complete.*

(b) ( Boolos and McGee 1987) *The class of all sentences of predicate modal formulas always true is  $\Pi_1^0$  – complete in  $Th(\mathbb{N})$ .*

*Proof.* (See Boolos (2008) pp. 233-41)

QED

Contrast with the propositional case: GL axiomatizes the class of always provable sentences of propositional Provability Logic; GLS axiomatizes the class of always true sentences. Such classes are *decidable*. At the opposite, in respect to QGL, these classes are *not even recursively enumerable*; since the collection of theorems of an axiomatizable theory is r.e. (namely  $\Sigma_1^0$ ), it follows that the class of always provable formulas of QGL *is not axiomatizable*.

Now observe that if  $X$  is an arithmetically definable set, then  $X \leq_T Th(\mathbb{N})$ : indeed, if some  $\phi \in \Sigma_n^0$  defines  $X$ , then  $n \in X$  iff  $\ulcorner \phi(\bar{n}) \urcorner \in Th(\mathbb{N})$ ; hence  $X \leq_m Th(\mathbb{N})$  (and a fortiori Turing reducible) via the function  $n \xrightarrow{\phi(\bar{x})} \ulcorner \phi(\bar{n}) \urcorner$ . But if  $X$  is  $\Pi_1^0$ -complete in  $Th(\mathbb{N})$ , then  $X$  is not  $\Sigma_1^0$  in  $Th(\mathbb{N})$ , namely is not r.e. in  $Th(\mathbb{N})$  and a fortiori, not recursive in it. Since by Boolos and McGee's result the class of *always true* sentences of QGL is  $\Pi_1^0$ -complete in  $Th(\mathbb{N})$ , it follows that such a class is not arithmetically definable.

### 5.5. Semantic for first-order Modal Logic

We start by introducing a Kripke-style semantic for first-order Modal Logic.

**Definition 45.** *A constant domain Kripke frame is a structure  $\mathcal{M} = \langle W, R, D, \{V_w\}_{w \in W} \rangle$  where for all  $w \in W$ ,  $V_w(P(x_1, \dots, x_n))$  is an  $n$ -ary relation on  $D$  and if  $\sigma$  is an assignment in  $D$  to the variables and we denote  $\sigma(a/x)$  its  $x$ -th variant, the forcing relation is given by the following:*

- (a)  $w \Vdash_{\sigma}^{\mathcal{M}} P(x_1, \dots, x_n)$  iff  $\langle \sigma(x_1), \dots, \sigma(x_n) \rangle \in V_w(P(x_1, \dots, x_n))$ .
- (b)  $w \Vdash_{\sigma}^{\mathcal{M}} x = y$  iff  $\sigma(x) = \sigma(y)$
- (c) *analogous to the propositional case, for connectives.*
- (d)  $w \Vdash_{\sigma}^{\mathcal{M}} \Box \phi$  for all  $v$ , if  $wRv$ , then  $v \Vdash_{\sigma}^{\mathcal{M}} \phi$ .
- (e)  $w \Vdash_{\sigma}^{\mathcal{M}} \exists x \phi$  iff for some  $x$ -variant  $\sigma(a/x)$ ,  $w \Vdash_{\sigma(a/x)}^{\mathcal{M}} \phi$ .

A frame with *monotone variable domain* is obtained by adding to the frame  $\langle W, R \rangle$  a collection of sets  $\{D_w\}_{w \in W}$  where for all  $w \in W$ , it holds that if  $wRv$ , then  $D_w \subseteq D_v$  (we call it *antimonotone* if at the opposite  $D_w \supseteq D_v$ ).

A counterexample to the validity of (BS) in frames with variable monotone domain is easily given, taking a model where  $W = \{w, v\}$ ,  $R = \{\langle w, v \rangle\}$ ,  $D_w = \{a\}$ ,  $D_v = \{a, b\}$ , defining  $V_v(P(x)) = \{b\}$  and by considering an assignment  $\sigma$  and its  $x$ -variant  $\sigma(b/x)$ . Actually (BS) and (CBS) are both valid in a semantic of *constant domains*, but as long as variable domains are considered, (BS) holds in a variable domain frame, iff such a frame satisfies *antimonotonicity*, and (CBS) holds in a variable domain frame, iff it satisfies *monotonicity*:

Let us call QK the system given adding to the First Order classical logic with identity these axioms:

- (a)  $\Box(\phi \rightarrow \psi) \rightarrow (\Box \phi \rightarrow \Box \psi)$
- (b) Necessitation rule  $\phi / \Box \phi$
- (c)  $x \neq y \rightarrow \Box(x \neq y)$

(d) Barcan Marcus schema.

A completeness result holds for this semantic: Every valid formula in constant domains frames is a theorem of QK (and *viceversa*). The system obtained removing (BS) is complete with respect to the increasing domain semantic: note that in this system the converse of Barcan Marcus schema is derivable.

Forcing  $w \Vdash A$  is defined on closed formulas  $\alpha$  with parameters in  $D_w$ . A Kripke model  $\mathfrak{R} = \langle W, R, \{D_x\}_{x \in W}, \Vdash \rangle$  for QGL satisfies the following conditions:

- (a)  $W \neq \emptyset$
- (b) If  $xRy$ , then  $D_x \subseteq D_y$
- (c) The forcing  $\Vdash$  is defined as in the propositional case, but with a further clause:

$$x \Vdash \exists u \alpha(u) \text{ iff there exists } a \in D_x, x \Vdash \alpha(a)$$

**Theorem 94.** *QGL is complete with respect to the class of models transitive and such that for any closed formula  $\alpha$ , if  $x \in W$  and  $x \Vdash \alpha$ , then there exists  $y \in W$  such that  $y \Vdash \alpha$ , and if  $yRz$ , then  $z \Vdash \neg\alpha$ .*

Recall that a theory  $S$  is *valid within a class of frames*  $C$ , if for all frames  $\mathfrak{R} \in C$ ,  $\mathfrak{R} \models \phi$  (namely if for all models  $\mathcal{M}$  based on that structure,  $\mathcal{M} \models \phi$ ), for all theorems  $\phi$  of  $S$ ; a theory is (weakly) complete with respect to a class of frames  $C$ , iff when  $\mathfrak{R} \models \phi$ , for all  $\mathfrak{R} \in C$ , then  $\phi$  is a theorem of  $S$ .

$S$  is *strongly complete* with respect to  $C$ , if  $\Gamma \models_C \phi$  implies  $\Gamma \vdash_S \phi$ . Strong completeness does not hold for GL because compactness fails: one can device an infinite set of formulas  $\Delta$  that is not satisfiable, but such that each of its finite subset is satisfiable. Hence, if  $\Gamma$  is a finite subset of  $\Delta$ , then  $\Gamma \not\vdash \perp$ , and a fortiori  $\Delta \not\vdash \perp$ , since otherwise  $\perp$  would be provable from some finite subset of  $\Delta$ . However, being  $\Delta$  unsatisfiable, we have also  $\Delta \models \perp$ . But QGL is incomplete with respect to *any class of frame*. Indeed, this is the situation:

- (a) QGL is *valid in a frame*, iff  $R$  is transitive and  $R^{-1}$  is well-founded, but:
- (b) QGL is *not complete* with respect to this class of frames. It follows that:
- (c) QGL is not complete with respect to *any class of frames*

For instance, the sentence:

$$\neg(\exists u \diamond P(u) \wedge \forall v \exists w \Box(P(v) \rightarrow \diamond P(w)))$$

is valid in all transitive and conversely well founded frames, but is not a theorem of QGL.

We will now illustrate the proof of two negative results namely that due to Artemov and that due to Montagna and mentioned above, respectively. This proof of the arithmetical incompleteness of QGL was provided by Montagna (1984).

**Theorem 95.** *QGL is not arithmetically complete w.r.t. PA.*

*Proof.* Let  $T$  a finitely axiomatizable and consistent theory (e.g. the theory NBG) such that  $T \vdash \text{Con}(T) \rightarrow \text{Con}(\text{PA} + \text{Con}(\text{PA}))$ . Let  $\bigwedge T$  be the conjunction of all axioms of  $T$  and let us define  $\alpha = \diamond \bigwedge T \rightarrow \diamond \diamond T$ .

We claim that  $\alpha$  is valid in PA, but not provable in QGL. Indeed, let us suppose that  $*$  is an arithmetical interpretation and let us consider  $\bigwedge T^*$  (assume w.l.o.g. that QGL contains the language of  $T$ ). Hence  $\alpha^*$  is provably equivalent to the sentence  $\text{Con}(\text{PA} + \bigwedge T^*) \rightarrow$

$Con(\text{PA} + Con(\text{PA}))$  and is provable in PA: let  $B_1, \dots, B_n$  be a proof of  $B_n$  in  $\mathsf{T}$ . Hence  $B_1^*, \dots, B_n^*$  is proof of  $B_n^*$  in  $\mathsf{T}^*$  and by a formalization of this, for all  $C$  not containing  $\Box$ , we have  $\text{PA} \vdash \text{Pr}_{\mathsf{T}}(\ulcorner C \urcorner) \rightarrow \text{Pr}_{\mathsf{T}^*}(\ulcorner C^* \urcorner)$  from which follows  $\text{PA} \vdash Con(\mathsf{T}^*) \rightarrow Con(\mathsf{T})$  and finally (since  $\text{PA} \vdash Con(\text{PA} + \bigwedge \mathsf{T}^*) \rightarrow Con(\mathsf{T}^*)$ ), we obtain  $Con(\text{PA} + Con(\text{PA}))$ . Hence  $\alpha$  is PA-valid. However it is not provable in QGL: let  $\mathfrak{R}$  be a model of  $\mathsf{T}$  and let us consider the model  $\mathfrak{S} = \langle W, R, \{D_x\}_{x \in W}, \Vdash \rangle$ , where:

- (a)  $W = \{0, 1\}$
- (b)  $xRy$  iff  $x = 0$  and  $y = 1$
- (c)  $D_0 = D_1 = \mathbb{N}$
- (d)  $i \Vdash B(a_0, \dots, a_n)$  iff  $\mathfrak{R} \models B(a_0, \dots, a_n)$ , for  $a_0, \dots, a_n \in \mathbb{N}$  and  $B(x_0, \dots, x_n)$  atomic formula of the language of  $\mathsf{T}$ .

Since  $1 \Vdash \bigwedge \mathsf{T}$  and  $0R1$ , we have  $0 \Vdash \Diamond \bigwedge \mathsf{T}$ ; but  $0 \Vdash \Box \Box \perp$ , and therefore  $0 \Vdash \neg \alpha$ . Clearly  $R$  is transitive and conversely well founded, so that  $\mathfrak{S} \models \text{QGL}$ .

QED

Another strong negative results is Artemov's theorem. Recall that according to *Tarski-Post* theorem, no single formula  $\Sigma_n^0$  of the Arithmetical Hierarchy can define  $Th(\mathbb{N})$ . We have seen that from a computational point of view,  $Th(\mathbb{N}) \equiv_{\mathsf{T}} \emptyset^\omega$ . Moreover, according to *Tennenbaum's theorem*  $+, \times, <$  are not recursive in (countable) non-standard model of PA, and many subtheories of it.

**Theorem 96.** (Artemov 1985) *Let  $\mathsf{T}$  be an r.e. theory. Then, for all choice of a proof predicate  $\text{Pr}_{\mathsf{T}}(x)$ , the set of predicate modal formulas always true is not arithmetical.*

*Proof.* (see De Jongh and Japaridze (1998)) Let us consider a purely relational arithmetical language with three predicates  $E(x, y)$ ,  $A(x, y, z)$  and  $M(x, y, z)$ , to be interpreted in standard model respectively as  $x = y, x + y = z, x \times y = z$ . Let us consider the following version of Tennenbaum's theorem: there exists an arithmetical sentence  $\theta$  such that:

- (a)  $\theta$  is true in the standard model.
- (b) Every countable model of  $\theta$  where  $E(x, y)$  is identity and  $A(x, y, z)$ ,  $M(x, y, z)$  are recursive, is isomorphic to the standard model.

The second point actually implies that every countable model where  $E(x, y)$  is identity and  $A(x, y, z)$ ,  $M(x, y, z)$  are recursive, is *elementary equivalent* to the standard model. Let now  $C$  be the conjunction of the following sentences:

- (a)  $\forall x \forall y (\Box E(x, y) \vee \Box \neg E(x, y))$
- (b)  $\forall x \forall y \forall z (\Box A(x, y, z) \vee \Box \neg A(x, y, z))$
- (c)  $\forall x \forall y \forall z (\Box M(x, y, z) \vee \Box \neg M(x, y, z))$

*Claim* For any arithmetical formula  $\phi$ ,  $\phi$  is true iff for all arithmetical interpretations  $*$ ,  $((\theta \wedge C) \rightarrow \phi)^*$  is true.

$\Rightarrow$  Actually if  $\phi$  is true,  $*$  is an interpretation and  $\theta^* \wedge C^*$  is true and  $\mathsf{T}$  is r.e., the truth of  $C^*$  implies that in the standard model  $E^*, A^*, M^*$  are recursive. Hence we define a countable model  $\mathfrak{R}$  such that:

- (a)  $\mathfrak{R} \models E(k, m)$  iff  $E^*(k, m)$  is true.

(b)  $\mathfrak{R} \models A(k, m, n)$  iff  $A^*(k, m, n)$  is true.

(c)  $\mathfrak{R} \models M(k, m, n)$  iff  $M^*(k, m, n)$  is true.

Hence we have in general that  $\mathfrak{R} \models \phi$  iff  $\phi^*$  is true, and in particular this holds for  $\theta$ . Since by hyp.  $\theta^*$  is true, also  $\mathfrak{R} \models \theta$ , and by our version of Tennenbaum's theorem we conclude that  $\mathfrak{R} \equiv \mathbb{N}$ . But by hyp.  $\phi$  is true; it follows that  $\mathfrak{R} \models \phi$  and once more  $\phi^*$ .  $\Leftarrow$  Let  $\phi$  be false and let  $*$  the trivial interpretation  $E^* = E, A^* = A, M^* = M$ . Hence  $\theta^* = \theta$  and  $\phi^* = \phi$ . We have to check that  $\theta \wedge C^* \rightarrow \phi$  is false, namely that  $\theta \wedge C^*$  is true (since  $\phi$  is false by hypothesis). But  $\theta$  is true by hyp. and from decidability in  $T$  of  $x = y, x + y = z, x \times y = z$  it follows that  $C^*$  is true. It follows that the set  $V$  of *always true* formulas of QGL is not arithmetical. Indeed, we know that  $Th(\mathbb{N})$  is not arithmetical and the previous lemma provide an m-reduction of  $Th(\mathbb{N})$  to  $V$ .

Let  $f(\ulcorner \phi \urcorner) = \ulcorner (\theta \wedge C) \rightarrow \phi \urcorner$ ; therefore  $\ulcorner \phi \urcorner \in Th(\mathbb{N})$  iff  $f(\ulcorner \phi \urcorner) \in V$ , namely  $Th(\mathbb{N}) \leq_m V$ . But if  $X \leq_m V$  and  $V \in \Sigma_n^0$ , also  $X \in \Sigma_n^0$  (contradiction). QED

As already mentioned in the former section, a basic result of propositional Provability Logic GL is the Sambin-de Jongh *fixed point theorem* and we wonder whether an analogous result holds in the predicative case. In other words, we ask whether is possible to prove in QGL the "Diagonalization theorem" of Gödel-Carnap. Let us consider the language of QGL extended with a countable amount of variable for formulas  $p_0, p_1, p_2, \dots$  and the axioms extended to this language. The question is the following:

Let  $\alpha(p_i)$  be a modalized formula of the above language. It is asked: is there a formula  $\beta$  of the original language, with individual variables identical to those of  $\alpha(p_i)$ , such that  $QGL \vdash \beta \leftrightarrow \alpha(\beta)$ ?

The answer is *no*.

**Theorem 97.** *No provable fixed point exists in QGL to the formula:*

$$\forall u \exists v \Box (p_i \rightarrow A(u, v))$$

*Proof.* Let us consider the Kripke model  $\mathfrak{R} = \langle \mathbb{N}, R, \{D_w\}_{w \in \mathbb{N}}, \Vdash \rangle$  where:

(a)  $xRy$  iff  $y < x$

(b)  $D_w = \{y \in \mathbb{N} \mid w \leq y\}$

(c) If  $B$  is a predicate letter:

i. If  $B$  is not  $A$ , then  $w \Vdash B(a_0, \dots, a_n)$ , for all  $a_0, \dots, a_n \in D_w$  and  $w \in \mathbb{N}$ .

ii. If  $B$  is  $A$ , then  $w \Vdash A(a, b)$  iff either  $b = w + 1$  and  $a \neq w + 1$ , or  $a < b$  and  $a, b \neq w + 1$  ( $a, b \in D_w$ ).

The order induced by  $A(a, b)$  is the following:

$$\begin{array}{rcl}
 \vdots & \vdots & \vdots \\
 2 & D_2 & 2, 3, 4 \dots 3 \\
 \downarrow & \cap & \\
 1 & D_1 & 1, 2, 3 \dots 2 \\
 \downarrow & \cap & \\
 0 & D_0 & 0, 1, 2 \dots 1
 \end{array} \tag{2}$$

Since  $R$  is transitive and  $R^{-1}$  well founded, this is a QGL model. Now let us observe that the forcing  $\Vdash$  is definable in the structure  $\langle \mathbb{N}, <, = \rangle$ , namely for any formula  $B(v_0, v_2, v_4 \dots v_{2n})$ , there exists a formula  $B^*(v_1, v_0, v_2, v_4 \dots v_{2n})$  such that for all  $a_0, \dots, a_n \in D_w$  and  $w \in \mathbb{N}$ :

$$w \Vdash B(a_0, a_2, a_4 \dots a_{2n}) \text{ iff } \langle \mathbb{N}, <, = \rangle \models B^*(a_1, a_0, a_2, a_4 \dots a_{2n})$$

Let us define  $B^*$  by induction:

- (a) For  $B(v_0, v_2, v_4 \dots v_{2n})$  atomic, different to  $A$ , let:

$$B^*(v_1, v_0, v_2, v_4 \dots v_{2n}) = \bigwedge_{i \leq n} v_1 \leq v_{2i}$$

- (b) For  $B(v_0, v_2, v_4 \dots v_{2n}) = A(v_0, v_2)$ , let:

$$B^*(v_1, v_0, v_2, v_4 \dots v_{2n}) = A^*(v_1, v_0, v_2) =$$

$$v_1 \leq v_0 \wedge v_1 \leq v_2 \wedge (v_2 = v_1 + 1 \wedge \neg(v_0 = v_1 + 1)) \vee \\ \vee (\neg(v_2 = v_1 + 1) \wedge \neg(v_0 = v_1 + 1) \wedge v_0 < v_2)$$

- (c) for  $B(v_0, v_2, v_4 \dots v_{2n}) = \neg C(v_0, v_2, \dots, v_{2n})$ , let:

$$B^*(v_1, v_0, v_2, v_4 \dots v_{2n}) = \bigwedge_{i \leq n} v_1 \leq v_{2i} \wedge \neg(C^*)(v_1, v_0, v_2, \dots, v_{2n})$$

- (d)  $*$  commutes with  $\exists$  and  $\vee$ .

- (e) for  $B(v_0, v_2, v_4 \dots v_{2n}) = \Box C(v_0, v_2, \dots, v_{2n})$ , let:

$$B^*(v_1, v_0, v_2, v_4 \dots v_{2n}) = \forall w < v_1 (C^*)(v_1, v_0, v_2, \dots, v_{2n})$$

where  $w$  is an odd variable not occurring in  $C^*$ .

It follows that  $\{x \in \mathbb{N} \mid x \Vdash B\}$  is definable in  $\langle \mathbb{N}, <, = \rangle$ ; but it is well known that every set definable in such a structure is *finite or cofinite* (hence, for instance, neither  $2\mathbb{N}$ , nor  $(2\mathbb{N} + 1)$  is definable). Let us suppose now that:

$$\text{QGL} \vdash B \leftrightarrow \forall u \exists v \Box (B \rightarrow A(u, v))$$

hence, since  $\mathfrak{K}$  is a model of QGL we will have that for all  $x \in \mathbb{N}$ ,  $x \Vdash B \leftrightarrow \forall u \exists v \Box (B \rightarrow A(u, v))$ . Since  $0 \Vdash \Box \perp$ , also  $0 \Vdash B$ . Moreover, if  $v \in D_1$ ,  $0 \Vdash \neg A(1, v)$  and therefore  $0 \Vdash B \wedge \neg A(1, v)$ , from which follows  $1 \Vdash \exists u \forall v \Diamond (B \wedge \neg A(u, v))$  and  $1 \Vdash \neg B$ . Note that  $0 \Vdash B$  and  $1 \Vdash \neg B$ . We want generalize this and we make an inductive argument by considering this result as the basis of an induction and we claim that  $B$  is forced at *even* nodes, whereas  $\neg B$  is forced at *odd* nodes. Now let us suppose that for all  $i \leq n$ ,  $2i \Vdash B$  and  $2i + 1 \Vdash \neg B$ . Hence  $2n + 1 \Vdash \neg B$ ; moreover if  $j \leq 2n$ ,  $j \Vdash A(u, u + 1)$ , for  $u \in D_{2n+2} = 2n + 2, 2n + 3, \dots$ ; so, for every  $u \in D_{2n+2}$  there is a  $v \in D_{2n+2}$  such that if  $(2n + 2)Rj$ , then  $j \Vdash \neg B$  or  $j \Vdash A(u, v)$ . It follows that  $2n + 2 \Vdash \forall u \exists v \Box (B \rightarrow A(u, v))$ . Hence  $2n + 2 \Vdash B$ ; moreover, if  $v \in D_{2n+3}$  then  $2n + 2 \Vdash B \wedge \neg A(2n + 3, v)$ , from which follows  $2n + 3 \Vdash \neg \forall u \exists v \Box (B \rightarrow A(u, v))$  that implies  $2n + 3 \Vdash \neg B$ . This concludes the induction step and proves the claim, namely that  $\{x \in \mathbb{N} \mid x \Vdash B\}$  coincides with the set of even numbers: *but this is not finite, neither cofinite* (contradiction).

QED

Which steps of the Sambin-de Jongh theorem fails in predicative case? The fixed point theorem for GL is based on the following substitution lemma:

Let  $A(p), B, C, D$  formulas of GL. If  $\text{GL} \vdash D \rightarrow (B \leftrightarrow C)$ , then:

- (a)  $\text{GL} \vdash D \wedge \Box D \rightarrow (A(B) \leftrightarrow A(C))$
- (b) Moreover, if  $p$  is modalized on  $A$ , then it holds that  $\text{GL} \vdash \Box D \rightarrow (A(B) \leftrightarrow A(C))$

This *does not extend to QGL*. However a weaker version is provable, sufficient to prove *uniqueness* of fixed point: for all modalized  $A(p)$  and all  $B, C$ , if  $\text{QGL} \vdash A(B) \leftrightarrow B$  and  $\text{QGL} \vdash A(C) \leftrightarrow C$ , then  $\text{QGL} \vdash B \leftrightarrow C$ .

*Question* What happens if we add the Barcan Marcus schema? If we add the Barcan Marcus schema, actually some counterexamples to the fixed point theorem in QGL, fail to be counterexamples. For instance, if  $A(p)$  is  $\forall x \exists y \Box(p \rightarrow A(x, y))$  and  $D$  is  $\forall u \exists v \Box A(u, v)$ , then  $\text{QGL} + \text{BS} \vdash D \leftrightarrow A(D)$ . On the other hand, modifying Montagna's proof (and considering Kripke models with fixed domain), it is possible to show there are other formulas, as for example  $\forall u(\Box \Box P(u) \rightarrow \Box(p \rightarrow P(u)))$ , that do not have a fixed point, not even in  $\text{QGL} + \text{BS}$ . Moreover this is also an arithmetical counterexample, in the sense that there is no sentence of the language of QGL, such that for all arithmetical interpretations  $*$ ,  $\text{PA} \vdash D^* \leftrightarrow A(D)^*$ .

We conclude by reporting some important recent developments and pointing to recent lines of research. The above mentioned Vardanyan's Theorem, i.e. the result which establishes the  $\Pi_2$ -completeness of the quantified modal logic of PA, is a benchmark and a barrier in this field of research, excluding the possibility of recursive axiomatisation. However, to paraphrase Visser and De Jonge (2006), we can say that we can glimpse an "escape" from it, although precisely this work dramatically generalises this result, leading to the conclusion that it is impossible to recursively axiomatise the quantified modal logic of a vast class of arithmetic theories. In fact, despite this negative result, some *positive* results still hold: on the one hand the result according to which the existence of a Kripke countermodel implies arithmetic nonvalidity, can be actually extended to the predicate level (see De Jongh and Japaridze (1998), 533-38); on the other hand Artemov and Japaridze (1990) proved that arithmetical completeness still holds, if we restrict ourselves to formulas with just one individual variable. Opening up a further glimmer of light in this gloomy outlook, in Yavorski (2002), Hao and Tournakis (2021) and De Almeida Borges and Joosten (2013), some positive results are achieved. In the last of these works, in particular, a strictly positive first order modal calculus, named  $\text{QRC}_1$  (*Quantified Reflection Calculus with one modality*), is introduced, whose signature is composed by relational symbols and constants, among which  $\top$ , and logical operators are restricted to  $\wedge, \forall$  and  $\Diamond$ . This logic is based on very basic rules that manipulates judgements of the form  $\phi \vdash \psi$ :

- |  |   |  |
|--|---|--|
| (a) $\phi \vdash \top$ and $\phi \vdash \phi$                                      | (e) $\frac{\phi \vdash \psi}{\Diamond \phi \vdash \Diamond \psi}$ | (h) $\frac{\phi[c/x] \vdash \psi[c/x]}{\phi \vdash \psi}$                            |
| (b) $\phi \wedge \psi \vdash \phi$ and $\phi \wedge \psi \vdash \psi$              | (f) $\Diamond \Diamond \phi \vdash \Diamond \phi$                 | (i) $\frac{\phi \vdash \psi}{\phi \vdash \forall x \psi}$ ( $x$ not free in $\phi$ ) |
| (c) $\frac{\phi \vdash \psi \quad \phi \vdash \chi}{\phi \vdash \psi \wedge \chi}$ | (g) $\frac{\phi \vdash \psi}{\phi[t/x] \vdash \psi[t/x]}$         | (j) $\frac{\phi[t/x] \vdash \psi}{\phi \vdash \forall x \psi}$ ( $t$ free for $x$ )  |
| (d) $\frac{\phi \vdash \psi \quad \psi \vdash \chi}{\phi \vdash \chi}$             |   |  |

A modal completeness (and decidability) was shown with respect to finite irreflexive and constant domain Kripke models. A theorem of arithmetical completeness is also demonstrated, providing a more complex notion of arithmetical interpretation and we would like to highlight the peculiar aspects of this interpretation. First a provability predicate *à la* Feferman  $\text{Prf}_\tau$  is considered (see on p.122) for a  $\Sigma_1$ -definition of axioms  $\tau$  of the theory  $\top$ . The basic idea

is that relational symbols  $S(x, c)$  of the modal language are interpreted as  $\Sigma_1$  arithmetical formulas  $S(x, c)^* = \sigma(u, v_x, v_c)$  that define set of axioms of theories and starting from such an interpretation  $*$ , then this is extended to an interpretation  $\tau^*$  of all formulas as follows, where  $\tau(u)$  is added to the definition of the axioms, to ensure that we are dealing with extensions of  $\mathsf{T}$ :

- (a)  $\top^{\tau^*} = \tau(u)$ .
- (b)  $S(x, c)^{\tau^*} = S(x, c)^* \vee \tau(u)$ .
- (c)  $(\psi \wedge \delta)^{\tau^*} = \psi^{\tau^*} \wedge \delta^{\tau^*}$ .
- (d)  $\diamond\psi^{\tau^*} = \tau(u) \vee (u = \ulcorner \text{Cons}_{\psi^{\tau^*}} \urcorner)$
- (e)  $(\forall x\psi)^{\tau^*} = \exists y\psi^{\tau^*}$

where conjunction is interpreted as the union of corresponding sets of axioms and the universal quantifier as an infinite union (see Beklemishev (2014)). For example:

$$(\diamond S(s, c))^{\tau^*} = \tau(u) \vee (u = \ulcorner \text{Cons}_{\sigma(u, v_x, v_c) \vee \tau(u)} \urcorner)$$

A highly sophisticated proof in Solovay's style concludes that  $\text{QRC}_1$  coincides with the set of judgements  $\phi(x, c) \vdash \psi(x, c)$  such that for all arithmetical interpretations  $*$  the theory  $\mathsf{T}$  proves that for all sentences  $\theta$ , the formula  $\forall x\forall y(\text{Pr}_{\psi^{\tau^*}}(\ulcorner \theta \urcorner) \rightarrow \text{Pr}_{\phi^{\tau^*}}(\ulcorner \theta \urcorner))$  holds true, where  $\psi^{\tau^*}$  and  $\phi^{\tau^*}$  generally depend on  $y, z$ , for all recursively enumerable and sound  $\mathsf{T}$  extending  $I\Sigma_1$ . To close the circle with the previous section, the logic  $\text{QRC}_1$  turns out to be arithmetically sound even with respect to  $\text{HA}$ , but the arithmetical completeness problem is still open.

## 5.6. The quest for consistency proofs: proposal for further study

By Gödel's theorem, a proof of consistency of arithmetic must necessarily go beyond the means of arithmetic itself, since no consistency proof of a sufficiently strong consistent arithmetical theory can use methods that can be formalized in the theory itself. A proof of the consistency of arithmetic can be given simply by showing that the standard model verifies all its axioms. This in fact is a proof in some set theory like  $\text{ZFC}$ , or in second-order arithmetic, quite far from the idea of a constructive proof. Or it can be given through Gödel's "Dialectica" interpretation using higher type functionals, which we briefly mentioned when discussing the extensions of simply typed lambda calculus. Having to transcend Hilbert's finitistic level with the aim of not straying too far from it, we emphasised how Gödel distinguished between finitistic reasoning and constructive reasoning. Or, finally, such a consistency proof can be given using transfinite induction up to  $\varepsilon_0 = \sup\{\omega, \omega^\omega, \omega^{\omega^\omega}, \dots\}$ . Recall from set theory that a set is called well-ordered if every non-empty subset of it has a least element. Ordinals (see section 5.2) are transitive sets, well ordered by the appartenance. The transfinite induction principle is a generalization of the complete induction principle on natural numbers, to more general well ordered sets. So, for instance, what we call the principle of transfinite induction *up the ordinal*  $\alpha$  can be expressed as:

$$\forall x \in \alpha ((\forall y < x P(y)) \rightarrow P(x)) \rightarrow \forall x \in \alpha P(x)$$

Well orderings are, in particular, linear ordering and actually holds true that any linearly ordered set enjoys the principle of complete induction if and only if it is well-ordered. Recall, by the way, that if we restrict ourselves to countable ordinal numbers, such as  $\varepsilon_0$ , these ordinals are in fact just different orderings of the natural numbers. This brings us to Gentzen's method, based on the well-foundedness of ordinal notations up to  $\varepsilon_0$ . The aim, in this case too,

is to remain as close as possible to finitistically acceptable reasoning. Introducing the sequent calculus we will see what difficulties arise in proving cut-elimination for theories, rather than for pure logic. In Chapter 7, we will demonstrate that the consistency of first-order logic follows from the cut-elimination theorem for its formalization in the sequent calculus because a consequence is that the empty sequent (the formalisation, in this calculus, of the contradiction) cannot be proven. We will also examine the difficulties and some partial solutions for fragments to the difficulties that arise regarding cut elimination for theories (i.e., in the presence of specific axioms or additional rules, such as induction, restricted to some classes of formulas). However cut elimination fails dramatically for *Peano Arithmetic* and therefore this method is not applicable. Gentzen, in alternative, devised a method for assigning an ordinal number smaller than  $\varepsilon_0$  to each finite proof and showed how to effectively transform a proof in *Peano Arithmetic* of the empty sequent, into another proof of the empty sequent such that the latter receives a smaller ordinal than the former. Each reduction step is assigned an ordinal number, and then it is shown that the ordinal number decreases with every step. Actually it would be better to talk about *ordinal notations* and well ordering of these notations, since ordinals are given in the so-called *Cantor normal form*  $\alpha = \omega^{\beta_0} + \dots + \omega^{\beta_n}$  with exponents  $\beta_0 \geq \dots \geq \beta_n$ . These “ordinals” indeed, are syntactic objects, rather than transfinite ordinals, strings of symbols introduced in a purely combinatorial manner, which can be equipped with a well-order and an algebra of elementary operations. In elaborating Gentzen’s proof, Takeuti (1987) 92-100 proposes an effective method for demonstrating that these ordinal notations are well-ordered, based on so-called “eliminators”, which are methods for showing that each strictly decreasing sequence starting from any ordinal notation is finite. This well-ordering property is called *accessibility* and constitutes an attempt, to a certain extent, to fill the gap with finitary mathematics. Starting from the assumption that, by contradiction, the empty sequent is provable, since the ordinal  $\varepsilon_0$  is *accessible*, the proofs of the empty sequent can only be reduced a finite number of times, and this leads to a contradiction, by using the induction on transfinite ordinal numbers up to  $\varepsilon_0$  in the form of an “infinite descent” argument (see Kleene (1952) 12-13). According to Takeuti (1987), the method based on “eliminators” is still acceptable from a finitist point of view, although modified in a more liberal sense to admit “*Gedankenexperimente* on (concrete) operations” (Takeuti (1987) 100-101).

In Gentzen’s proof every step except the well-ordering of the ordering of type  $\varepsilon_0$  can be performed in Skolem’s *Primitive Recursive Arithmetic* PRA. If we agree, following Tait (1981), that what Takeuti calls Hilbert’s “purely finitist standpoint” coincides with Skolem’s primitive recursive arithmetic PRA, we can formally express this argument as follows:

$$\text{PRA} + \text{TI}(\varepsilon_0) \vdash \text{Con}(\text{PA})$$

where  $\text{TI}(\varepsilon_0)$  is the formalization of the transfinite induction principle up to  $\varepsilon_0$ . We can also say that, in a certain sense, the ordinal  $\varepsilon_0$  *measures the strength* of PA. However, to formalise the principle of transfinite induction in the language of arithmetic, we must find a way to represent the ordinals up to  $\varepsilon_0$  in this language and the representation system chosen is not irrelevant to the result. What constitutes a good “natural” order is a much-debated question. Pathological, although *ad hoc*, examples are available. However in Rathjen (1999) (work that we also recommend for a more extensive discussion and bibliography on the subject, primarily by the author himself) is described how “natural” systems typically arise: the ordinals  $\alpha$  smaller than  $\varepsilon_0$  are represented in *Cantor normal form*, whose exponents themselves have Cantor normal forms with even smaller exponents. Since the process must end, this ensures that they can be encoded with natural numbers. It follows that we can devise a coding system  $\ulcorner x \urcorner$  such that for every  $\alpha < \varepsilon_0$  the code  $\ulcorner \alpha \urcorner$  is a natural number that denotes the ordinal  $\alpha$  and the two structures,  $\varepsilon_0$  with the operations  $+$ ,  $\cdot$ ,  $\omega^\alpha$  and the relation  $<$ , on the one hand, and on the other hand the set of these codes, the primitive recursive functions  $\ulcorner \alpha \urcorner + \ulcorner \beta \urcorner = \ulcorner \alpha + \beta \urcorner$ ,  $\ulcorner \alpha \urcorner \cdot \ulcorner \beta \urcorner = \ulcorner \alpha \cdot \beta \urcorner$ ,  $\hat{\omega}^{\ulcorner \alpha \urcorner} = \ulcorner \omega^\alpha \urcorner$  and the primitive recursive relation  $\ulcorner \alpha \urcorner < \ulcorner \beta \urcorner$  if and only if  $\alpha < \beta$  turn out to be isomorphic. Therefore, the principle of transfinite induction can actually be expressed with a formula in the language of primitive

recursive arithmetic:

$$\forall x(\forall y \prec x P(x)) \rightarrow \forall x P(x)$$

where  $P(x)$  is a primitive recursive predicate. In Gentzen (1943) it is showed that PA proves the transfinite induction up to  $\alpha$ , for each  $\alpha < \varepsilon_0$ , so we can say that his consistency result is optimal.

Between 1934 and 1943, Gentzen gave four proofs of consistency (three of which were published, according to the the historical reconstruction in Von Plato (2014)). Nowadays, when talking about ‘‘Gentzen’s consistency proof for arithmetic’’, one usually refers to Gentzen (1938). Some *Proof Theory* textbooks report variations of this method (see for instance Takeuti (1987) and Mancosu, Galvan and Zach (2021)), trying to shed light on its darkest aspects.

The proof in Takeuti (1987) is long and complex and here we will limit ourselves to providing a taste, illustrating a crucial point and showing how induction is replaced by cut, that highlights the method. Again, we refer to Chapter 7 for notation regarding the sequent calculus and for the formalization in this calculus of the arithmetical theories (in particular, of the induction rule). On the contrary, the complexity measures are peculiar to this proof. In the mechanism for assigning ordinals less than  $\varepsilon_0$  to sequents  $S$  within a proof  $\pi$  (notation  $o(S; \pi)$ , or simply  $o(S)$ ), we proceed inductively, looking at the rule by which  $S$  was obtained. Hence in particular, if  $S$  is the lower sequent of a cut inference, where the upper sequents were assigned respectively the ordinals  $\mu$  and  $\nu$ , then  $o(S) = \omega_{k-l}(\mu \# \nu)$  (where  $\omega_0(x) = x$  and  $\omega_{n+1}(x) = \omega^{\omega_n(x)}$ ), where  $k$  and  $l$  are the *heights* (where in this proof, the height of a sequent is the maximum complexity, according to a measure of complexity called *grade*, of the cut-formulas and formulas to which induction under that sequent applies) of the upper sequents and of  $S$  respectively and  $\#$  is the so-called ‘‘natural (or Hessenberg) sum’’:  $\omega^{\mu_0} + \dots + \omega^{\mu_m} \# \omega^{\nu_0} + \dots + \omega^{\nu_n} = \omega^{\xi_0} + \dots + \omega^{\xi_{m+n}}$ , where  $\xi_0, \dots, \xi_{m+n}$  are the exponents  $\mu_0, \dots, \mu_m, \nu_0, \dots, \nu_n$  sorted in nonincreasing order. If  $S$  is the lower sequent of an induction inference, where the upper sequent was assigned the ordinal (in Cantor normal form)  $\mu = \omega^{\mu_0} + \dots + \omega^{\mu_n}$ , then  $o(S) = \omega_{l-k+1}(\mu_0 + 1)$ , where  $l$  and  $k$  are the heights of the upper sequent and of the lower sequent respectively.

The lemma according to which an hypothetical proof of the empty sequent  $\Longrightarrow$  in a sequent calculus formalization of PA can be transformed into another proof of the same sequent, but with a lower ordinal associated with it, is obtained by analysing the so-called *end-piece* of a hypothetical proof of the this sequent, i.e., if the end-sequent is the empty sequent, that part of the proof selected ascending each thread starting from the final sequent, until a logical rule is encountered, and stopping in each branch at the lower sequent of such a logical rule. A crucial step is that of replacing all inductions appearing in such *end-piece* with a sequence of cuts, obtaining a bound to the proof, strictly smaller than the previous one. Hence consider the end-piece of a proof  $\pi$  of the empty sequent and take the lowermost induction inference, i.e. an inference (see Chapter 7.) of this form:

$$\frac{\psi(\bar{x}), \Gamma \Longrightarrow \Delta, \psi(S(x))}{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(s)}$$

Let respectively  $l$  and  $k$  the heights of (both) upper sequents and the height of the lower sequent. By definition, the ordinal assigned to the latter is  $o(S) = \omega_{l-k+1}(\mu_0 + 1)$ , where  $\mu = \omega^{\mu_0} + \dots + \omega^{\mu_n}$  is the one assigned to the upper sequent. It can be shown that in the end-piece any variable that is not an *Eigenvariable* can be replaced by a constant, so we can assume that  $s$  is a closed term and that therefore there is a proof for a certain number  $m$  (its value). that in the final part any variable that is not an *Eigenvariable* can be replaced by a constant, so we can assume that  $s$  is a closed term and that therefore, as a consequence, there is a proof of  $\psi(\bar{m}) \Longrightarrow \psi(s)$ , for some number  $m$  (the value of  $s$ ). Inductive inference will be replaced by a sequence of cuts:

$$\begin{array}{c}
 \frac{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(\bar{1}) \quad \psi(\bar{1}), \Gamma \Longrightarrow \Delta, \psi(\bar{2})}{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(\bar{2})} \quad \psi(\bar{2}), \Gamma \Longrightarrow \Delta, \psi(\bar{3}) \\
 \frac{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(\bar{3})}{\vdots} \\
 \frac{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(\bar{m}) \quad \psi(\bar{m}), \Gamma \Longrightarrow \Delta, \psi(s)}{\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(s)}
 \end{array}$$

Having all  $\psi(\bar{n})$  the same complexity (the *grade* count the number of logical symbols), to all sequents  $\psi(\bar{n}), \Gamma \Longrightarrow \Delta, \psi(S(\bar{n}))$  is assigned the same ordinal  $\mu$  and by definition each  $\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(n)$  is assigned  $\mu \# \dots \# \mu$  ( $n$ -times). But in the proof of  $\psi(\bar{m}), \Gamma \Longrightarrow \Delta, \psi(s)$  no induction or cut on complex formulas is required and in this case the assignment system assigns to it a finite number, say  $q$ . So, actually, the ordinal assigned to the final sequent  $\psi(\bar{0}), \Gamma \Longrightarrow \Delta, \psi(s)$  of this subderivation is  $\omega_{l-k}(\mu \cdot m + q)$ , which is less than the original  $\omega_{l-k+1}(\mu_1 + 1)$ .

The mechanism for assigning ordinals appeared to some to be too ad hoc. Since it is largely the induction rule that causes problems, it seemed considerably preferable to admit infinitary proof systems as a formalization of PA, with infinitary rules replacing the induction rule, and obtaining a more transparent proof. In the wake of Schütte (1960), it is today quite common to overcome both the obstacle of eliminating the cuts and the alleged lack of transparency in the assignment of ordinal numbers to proofs by following a different method by admitting rules with *infinite premises*. This logic is frequently proposed in a certain version of the one-side sequents calculus (i.e. where the antecedent of all sequents is empty), introduced in Tait (1968). The so-called  $\omega$  rule, in the two-side calculus, that we use in these lectures, consists instead of two types of infinitary inference, a possibility already explored by Hilbert, which Gentzen did not admit, however:

$$\frac{\Gamma \Longrightarrow \Delta, \phi(\bar{n}) \quad (\text{for all } n \in \mathbb{N})}{\Gamma \Longrightarrow \Delta, \forall x \phi(x)} \quad \frac{\phi(\bar{n}), \Gamma \Longrightarrow \Delta \quad (\text{for all } n \in \mathbb{N})}{\exists x \psi(x), \Gamma \Longrightarrow \Delta}$$

(see ch.7. for notations) replacing the right universal quantifier rule and the left existential quantifier rule in sequent calculus. The version of PA with these rules instead of the induction rule, named  $\text{PA}_\omega$ , enjoys cut elimination. It is proved that if the empty sequent has a proof in PA, then it has a cut-free proof in  $\text{PA}_\omega$  of height less than  $\varepsilon_0$ , but this is impossible, because, as we will see, there can be no cut-free derivation of the empty sequent. A proof of the consistency theorem along these lines is given, for example, in Schwichtenberg (1977), in Girard (1987) pp. 347-416, and in Troelstra and Schwichtenberg (2000) pp. 259-79.

Let us give at least an idea of how a consistency proof using the above infinitary rules can work. First, note that the derivations now are well-founded trees, which are generally *infinite*. Still arguing informally, we define the *cut-rank*  $k$  of a derivation be the length of the longest cut-formula, and the *height*  $\alpha$  of a derivation the supremum of the heights plus one of all its subderivations (where axioms have height 0). We write  $\text{PA}_\omega \vdash_k^\alpha \Gamma \Longrightarrow \Delta$  to mean that there exists a derivation of that sequent of height  $\alpha$  and cut-rank  $k$ . For example, a  $\text{PA}_\omega$ -derivation of the principle of induction now has this form (see Girard (1987) p. 355). Let  $\Gamma = \psi(\bar{0}), \forall x(\psi(x) \rightarrow \psi(S(x)))$ . We can derive the sequents  $\Gamma \Longrightarrow \psi(\bar{n})$  as follows:

(**h**) for  $n = 0$ , by weakening from  $\psi(\bar{0}) \Longrightarrow \psi(\bar{0})$ .

(b) Suppose that each  $\Gamma \implies \psi(\bar{n})$  has been proved. Hence build this derivation:

$$\frac{\frac{\frac{\Gamma \implies \psi(\bar{n}) \quad \psi(\overline{n+1}) \implies \psi(\overline{n+1})}{\Gamma, \psi(\bar{n}) \rightarrow \psi(\overline{n+1}) \implies \psi(\overline{n+1})}}{\Gamma, \forall x(\psi(x) \rightarrow \psi(S(x))) \implies \psi(\overline{n+1})}}{\Gamma \implies \psi(\overline{n+1})}$$

(c) Now conclude as follows:

$$\frac{\frac{\frac{\Gamma \implies \psi(\bar{0}), \Gamma \implies \psi(\bar{1}), \Gamma \implies \psi(\bar{2}), \dots}{\Gamma \implies \forall x \psi(x)}}{\psi(\bar{0}) \implies (\forall x(\psi(x) \rightarrow \psi(S(x))) \rightarrow \forall x \psi(x))}}{\implies \psi(\bar{0}) \rightarrow (\forall x(\psi(x) \rightarrow \psi(S(x))) \rightarrow \forall x \psi(x))}$$

This is a proof of the induction axiom which is a well founded tree of height  $\omega + 2$ . The crucial result (see Rathjen (2006)) of cut-elimination for this formalization of Peano Arithmetic with infinitary rules establishes that if  $\text{PA}_\omega \vdash_{k+1}^\alpha \Gamma \implies \Delta$ , then  $\text{PA}_\omega \vdash_k^{\omega^\alpha} \Gamma \implies \Delta$  and therefore, by iterating this result we can obtain a cut-free derivation ( $k = 0$ ) at the price of increasing the length of the derivation as:

$$\omega^\omega \dots^\alpha$$

for  $k$ -iterations, and therefore the height of such a derivation is bounded by  $\varepsilon_0$ . But it is provable that if  $\text{PA} \vdash \Gamma \implies \Delta$ , then it is also provable  $\text{PA}_\omega \vdash_k^{\omega+s} \Gamma \implies \Delta$  for some  $s, k < \omega$  and from the above, there exists a *cut-free* derivation of this sequent in  $\text{PA}_\omega$ . Now, if we apply this argument to the empty sequent  $\implies$ , we would obtain a cut-free derivation of it, which is impossible (indeed, it is an immediate application of the subformula property, see ch.7).

Based on these methods, the analysis of the strenght of mathematical theories by means of transfinite induction, now called *ordinal analysis*, has been extended to other theories. For example, sticking to the theories mentioned in this volume, the ordinal associated with  $\text{Q}$  is  $\omega$ ; the one associated with  $\text{ID}_0$  is  $\omega^2$ , while the one associated with  $\text{PRA}$  is  $\omega^\omega$  (see Sommer (1990) and Sommer (1995)). Finally, to cite a fragment of the second order, to which ordinal analysis has been most applied, the proof-theoretic ordinal associated with  $\text{ATR}_0$ , a theory which we mentioned in section 6.4., and of Feferman's *Predicative Analysis* is  $\Gamma_0$ , the so-called "Feferman-Schütte ordinal", a countable ordinal that is the least ordinal "unreachable" by predicative means that we mentioned on p.131. However, analysing stronger and second-order theories goes beyond the purpose of these lecture notes. Rather, remaining within the realm of weak theories, we conclude by pointing to a line of research into the provability of well-foundedness of ordinal notations in weak theories of *Bounded Arithmetic* developed in Sommer (1995) and Beckmann, Pollett and Buss (2003). For instance, the latter authors define a notion of *well-foundedness on bounded domains* and show that the theories  $T_2^1$  and  $S_2^2$ , that we introduce in section 7.3, can prove the well-foundedness on bounded domains of the ordinal notations below  $\varepsilon_0$  and  $\Gamma_0$ .