

# A Formally Verified Abstract Account of Gödel’s Incompleteness Theorems

Andrei Popescu<sup>1</sup> and Dmitriy Traytel<sup>2</sup>

<sup>1</sup> Department of Computer Science, Middlesex University London, UK

<sup>2</sup> Institute of Information Security, Department of Computer Science, ETH Zürich, Switzerland

**Abstract.** We present an abstract development of Gödel’s incompleteness theorems, performed with the help of the Isabelle/HOL theorem prover. We analyze sufficient conditions for the theorems’ applicability to a partially specified logic. In addition to the usual benefits of generality, our abstract perspective enables a comparison between alternative approaches from the literature. These include Rosser’s variation of the first theorem, Jeroslow’s variation of the second theorem, and the Świerczkowski–Paulson semantics-based approach. As part of our framework’s validation, we upgrade Paulson’s Isabelle proof to produce a mechanization of the second theorem that does not assume soundness in the standard model, and in fact does not rely on any notion of model or semantic interpretation.

## 1 Introduction

Gödel’s incompleteness theorems [8, 11] are landmark results in mathematical logic. Both theorems refer to consistent logical theories that satisfy some assumptions, notably that of “containing enough arithmetic.” The first incompleteness theorem ( $\mathcal{IT}_1$ ) says that there are sentences that the theory cannot decide (i.e., neither prove nor disprove); the second theorem ( $\mathcal{IT}_2$ ) says that the theory cannot prove (an internal formulation of) its own consistency. It is generally accepted that  $\mathcal{IT}_1$  and  $\mathcal{IT}_2$  have a wide scope, covering many logics and logical theories. However, when it comes to rigorous presentation, typically these results are only proved for particular, albeit paradigmatic cases, such as theories of arithmetic or hereditarily finite (HF) sets, within classical first-order logic (FOL); and even in these cases the constructions and proofs tend to be “incomplete and (apparently) irremediably messy” [3, p.16]. Hence, the theorems’ scope remains largely unexplored on a rigorous/formal basis.

The emergence of powerful theorem provers has changed the rules of the game and, we argue, the expectation. Using interactive theorems provers, we can reliably keep track of all the constructions and their properties. Proof automation (often powered by fully automatic provers [16, 26]), makes complete, fully rigorous proofs feasible. And indeed, researchers have successfully met the challenge of mechanizing  $\mathcal{IT}_1$  [13, 23, 25, 31] and recently  $\mathcal{IT}_2$  [25]. Besides reassurance, these verification *tours de force* have brought superior technical insight into the theorems. But they have taken place within the same solitary confinement of scope as the informal proofs.

This paper takes steps towards a more comprehensive prover-backed exploration of the incompleteness theorems, based on a detailed analysis of their underlying assumptions. We use Isabelle/HOL [22] to establish general conditions under which these theorems apply to a partially specified logic. Our formalization is publicly available [28].

We start with a notion of logic (Section 2) whose terms, formulas and provability relation are kept abstract (Section 2.1). In particular, substitution and free variables are not defined, but axiomatized by some general properties. On top of this logic substratum, we consider an arithmetic substratum, consisting of a set of closed terms called *numerals* and an order-like relation (Section 2.2). Also factored in our abstract framework are encodings of formulas and proofs into numerals, the representability of various functions and relations as terms or formulas (Section 2.3), variations of the Hilbert-Bernays-Löb derivability conditions [14, 21] (Section 2.4), and standard models (Section 2.5).

Overall, our assumptions capture the notion of “containing enough arithmetics” in a general and flexible way. It is general because only few assumptions are made about the exact nature of formulas and numerals. It is flexible because different versions of the incompleteness theorems consider their own “amount of arithmetics” that makes it “enough,” as proper subsets of these assumptions. Indeed, our formalization of the theorems (Section 3) proceeds in an austere-buffet style: Every result picks just enough infrastructure needed for it to hold—ranging from diagonalization which requires very little (Section 3.1) to Rosser’s version of  $\mathcal{IT}_1$  which is quite demanding. This approach caters for a sharp comparison between different formulations of the theorems, highlighting their trade-offs: Gödel’s original formulation of  $\mathcal{IT}_1$  versus Rosser’s improvement (Section 3.2), proof-theoretic versus semantic versions of  $\mathcal{IT}_1$  (Section 3.2), and Gödel’s original formulation of the  $\mathcal{IT}_2$  versus Jeroslow’s improvement (Section 3.3).

Abstractness is our development’s main strength, but also a potential weakness: Are our hypotheses reasonable? Are they consistent? These questions particularly concern our axiomatization of free variables and substitution—a notoriously error-prone area. As a remedy, we instantiate our framework to Paulson’s semantics-based  $\mathcal{IT}_1$  and  $\mathcal{IT}_2$  for HF set theory [25], also performing an upgrade of Paulson’s  $\mathcal{IT}_2$  to a more general and more standard formulation: not restricted to sound theories, but applicable to any consistent theory (Section 4).

In the rest of this section, we discuss some guiding principles we followed when developing the formalization, and place our line of work in the context of related work.

**Formal Design Principles** Our long-term goal is a framework that makes it easy to instantiate the incompleteness theorems and related results to different logics. This is a daunting task, especially for  $\mathcal{IT}_2$ , where a lot of seemingly logic-specific technicalities are required to even formulate the theorem. The challenge is to push as much as possible of the technical constructions and lemmas to a largely logic-independent layer.

To this end, we strive to make minimal assumptions in terms of structure and properties when inferring the results—we will call this the *Economy* principle. For example, we do not define, but axiomatize syntax in terms of a minimal amount of structure. We assume a generic single-point substitution, then define simultaneous substitution and infer its properties. This is laborious, but worthwhile: Any logic that provides a single-point substitution satisfying our assumptions gets the simultaneous substitution for free.

As another instance of Economy, when faced with two different ways of formulating a theorem’s conclusion we prefer the one that is *stronger under fewer assumptions*. (And dually, we prefer weakness for a theorem’s assumptions.) For example, we discuss two variants of consistency: (1) “does not prove false” or (2) “there exists no formula such that itself and its negation are provable” (Section 3.3). The two statements are equivalent

at the meta-level, but the situation changes when they are represented as object-logic formulas: (1) implies (2) under mild assumptions, but the converse requires a lot more knowledge about the provability representation. So in our abstract theorems we strive to conclude (1) rather than (2). Even if (2) would also imply (1) in all reasonable instances, why postpone for the instantiation time any facts that we can acquire abstractly?

Applying the Economy principle not only stocks up generality for instantiations, but also accurately outlines trade-offs: How much does it cost (in terms of other added assumptions) to improve the conclusion, or to weaken an assumption of a theorem? For example, an Economy-based proof of Rosser’s variant of  $\mathcal{IT}_1$  reveals how much arithmetic we must factor in for weakening the  $\omega$ -consistency assumption into consistency.

**Related Work** In his seminal publication, Gödel gave a proof of  $\mathcal{IT}_1$  and the rough proof idea of  $\mathcal{IT}_2$  [11]. Hilbert and Bernays gave a first detailed proof of  $\mathcal{IT}_2$  [14]. Subsequently, a vast amount of literature was dedicated to the (re)formulation, proof, and analysis of these results [3,29,32,33]. The now canonical line of reasoning goes through the three derivability conditions devised by Bernays and Hilbert [14] and simplified by Löb [21]: One condition is used for  $\mathcal{IT}_1$ , and all three conditions for  $\mathcal{IT}_2$ . The derivability conditions have inspired a new branch of modal logic called provability logic [3]. Jeroslow has argued that one condition is redundant when proving  $\mathcal{IT}_2$  [15].

Kreisel [18] and Jeroslow [15] were the first to study abstract conditions on logics under which the incompleteness theorems apply. Buldt [4] surveys the state of the art focusing on  $\mathcal{IT}_1$ , also sketching the applicability to non-standard logics. Our approach appears to be more abstract than all previous approaches: Being based on generic syntax and provability and truth predicates, it resembles the style of institution-independent model theory [7, 12] and our previous work on abstract completeness [2] and completeness of ordered resolution [30]. However, there are certainly dimensions of generality and strength of this vast topic that our formalized work does not (yet) explore. These include quantifier-free logics [15] and arithmetical hierarchy refinements [17]. Our syntax axiomatization is inspired by algebraic theories of the  $\lambda$ -calculi syntax [9, 10, 27].

In the realm of mechanical proofs,  $\mathcal{IT}_1$  has been proved by Shankar [31] in the Boyer-Moore prover, Harrison in HOL Light [13] and O’Connor in Coq [23].  $\mathcal{IT}_2$  has only been proved recently—by Paulson in Isabelle/HOL [24,25] (who also proved  $\mathcal{IT}_1$ ). All these mechanizations target theories over a fixed language in classical FOL: that of arithmetic (Harrison and O’Connor) and that of HF sets (Shankar and Paulson). These mechanizations are mostly focused on “getting all the work done” in a particular setting (although Harrison targets a more abstract class of theories in the given language). On their way to  $\mathcal{IT}_1$ , Shankar and O’Connor also prove the representability of all partial, respectively primitive recursive functions—which are important standalone results.

By contrast, we explore conditions that enable different formulations for an abstract logic, where aspects such as recursiveness are below our abstraction level. The two approaches are complementary, and they both contribute to formally taming the complex ramifications of the incompleteness theorems. When instantiating our abstract assumptions to recover and upgrade Paulson’s results, we took advantage of Paulson’s substantial work on proving the many low-level lemmas towards the derivability conditions. More should be done at an abstract level to avoid duplicating some of these laborious lemmas when instantiating the theorems to different logics. This will be future work.

## 2 Abstract Assumptions

Roughly, the incompleteness theorems are considered to hold for logical theories that (1) contain enough arithmetic and (2) are “effective” in that they themselves can be arithmetized. Our goal is to give a general expression of these favorable conditions. To this end, we identify some logic and arithmetic substrata consisting of structure and axioms that express the containment of (various degrees of) arithmetic more abstractly and flexibly than relative interpretations [35]. We also identify abstract notions of encodings and representability that have just what it takes for a working arithmetization.

### 2.1 The logical substratum

We start with some unspecified sets of variables ( $\text{Var}$ , ranged over by  $x, y, z$ ), terms ( $\text{Term}$ , ranged over by  $s, t$ ) and formulas ( $\text{Fmla}$ , ranged over by  $\varphi, \psi, \chi$ ). We assume that variables are particular terms,  $\text{Var} \subseteq \text{Term}$ , and that  $\text{Var}$  is infinite. Free-variables and substitution operators,  $\text{FVars}$  and  $[_/_]$ , are assumed for both terms and formulas. We think of  $\text{FVars}(t)$  as the (finite) set of free variables of the term  $t$ , and similarly for formulas. We call *sentence* any formula with no free variable, and let  $\text{Sen}$  denote the set of sentences. We think of  $s [t/x]$  as the term obtained from  $s$  by the (capture-avoiding) substitution of  $t$  for the free occurrences of variable  $x$ ; we think of  $\varphi [t/x]$  as the formula obtained from  $\varphi$  by the substitution of  $t$  for the free occurrences of variable  $x$ .

In FOL, terms introduce no bindings, so any occurring variable is free. FOL terms fall under our framework, and so do terms with bindings as in  $\lambda$ -calculi and higher-order logic (HOL). To achieve this degree of inclusiveness while also being able to prove interesting results, we work under some well-behavedness assumptions about the free-variables and substitution operators. For example, free-variables distribute over substitution,  $\text{FVars}(\varphi [s/x]) = \text{FVars}(\varphi) - \{x\} \cup \text{FVars}(s)$  if  $x \in \text{FVars}(\varphi)$ , and substitution is compositional,  $\varphi [s_1/x_1] [s_2/x_2] = \varphi [s_2/x_2] [(s_1 [s_2/x_2]) / x_1]$  if  $x_1 \neq x_2$  and  $x_1 \notin \text{FVars}(s_2)$ . The full list of our generic syntax axioms is shown in Appendix A.1.

The incompleteness theorems rely heavily on simultaneous substitution, written  $\varphi [t_1/x_1, \dots, t_n/x_n]$ , whose properties are tricky to formalize—for example, Paulson’s formalization paper dedicates them ample space [25, 6.2]. To address this problem once and for all generically, we define simultaneous substitution from the single-point substitution,  $\varphi [t/x]$ , and infer its properties from the single-point substitution axioms. For example, we prove that  $\text{FVars}(\varphi [s_1/x_1, \dots, s_n/x_n]) = \text{FVars}(\varphi) \cup \bigcup \{\text{FVars}(s_i) - \{x_i\} \mid i \in \{1, \dots, n\} \text{ and } x_i \in \text{FVars}(\varphi)\}$ . The technicalities are delicate: To avoid undesired variable replacements,  $\varphi [s_1/x_1, \dots, s_n/x_n]$  must be defined as  $\varphi [y_1/x_1] \dots [y_n/x_n] [s_1/y_1] \dots [s_n/y_n]$  for some fresh  $y_1, \dots, y_n$ , the choice of which we must show to be immaterial. This definition’s complexity is reflected in the properties’ proofs. But again, this one-time effort benefits any “customer” logic: In exchange for a well-behaved single-point substitution, it gets back a well-behaved simultaneous substitution.

We let  $v_1, v_2, \dots$  be fixed mutually distinct variables. We write  $\text{Fmla}_k$  for the set of formulas whose free variables are precisely  $\{v_1, \dots, v_k\}$ , and  $\text{Fmla}_k^{\subseteq}$  for the set of formulas whose variables are among  $\{v_1, \dots, v_k\}$ . Note that  $\text{Fmla}_k \subseteq \text{Fmla}_k^{\subseteq}$  and  $\text{Fmla}_0 = \text{Fmla}_0^{\subseteq} = \text{Sen}$ . Given  $\varphi \in \text{Fmla}_k^{\subseteq}$ , we write  $\varphi(t_1, \dots, t_n)$  instead of  $\varphi [t_1/v_1, \dots, t_n/v_n]$ .

In addition to free variables and substitution, our theorems will require formulas to be equipped with term equality ( $\equiv$ ), Boolean connectives ( $\perp, \top, \rightarrow, \neg, \wedge, \vee$ ), universal and existential quantifiers ( $\forall, \exists$ ). When needed, we will assume them not as constructors (syntax builders), but as operators on terms and formulas, e.g.,  $\equiv : \text{Term} \rightarrow \text{Term} \rightarrow \text{Fmla}$ ,  $\perp \in \text{Fmla}$ ,  $\forall : \text{Var} \times \text{Fmla} \rightarrow \text{Fmla}$ . This caters for logics that do not have them as primitives, but can define them. For example, higher-order logic (HOL) defines all connectives and quantifiers from  $\lambda$ -abstraction and either equality or implication.

We fix a unary relation  $\vdash \subseteq \text{Fmla}$  on formulas, called *provability*. We write  $\vdash \varphi$  instead of  $\varphi \in \vdash$ , and say the formula  $\varphi$  is *provable*. Whenever certain formula connectives or quantifiers are assumed present, we will assume that  $\vdash$  behaves intuitionistically w.r.t. them—namely, we assume the usual (Hilbert-style) intuitionistic FOL axioms with respect to the abstract connectives and quantifiers. Stronger systems, such as those of classical logic, also satisfy these assumptions.

Consistency, denoted  $\text{Con}$ , is defined as the impossibility to prove false, namely  $\not\vdash \perp$ . Another central concept is  $\omega$ -consistency—we carefully choose a formulation that works intuitionistically, with conclusion reminiscent of Gödel’s negative translation [6]:

$\text{OCon}$ : For all  $\varphi \in \text{Fmla}_1^{\subseteq}$ , if  $\vdash \neg \varphi(n)$  for all  $n \in \text{Num}$  then  $\not\vdash \neg \neg (\exists x. \varphi(x))$ .

Assuming classic deduction in  $\vdash$ , this is equivalent to the standard formulation: For all  $\varphi \in \text{Fmla}_1^{\subseteq}$ , it is not the case that  $\vdash \varphi(n)$  for all  $n \in \text{Num}$  and  $\vdash \neg (\forall x. \varphi(x))$ .

Occasionally we will need to consider not only provability, but also explicit proofs. In such cases, we fix a set  $\text{Proof}$  of (entities which we call) *proofs*, ranged over by  $p, q$ , and a binary relation between proofs  $p$  and sentences  $\varphi$ , written  $p \Vdash \varphi$ , of which we think of as stating that  $p$  is a proof of  $\varphi$ . We also assume that  $\vdash$  and  $\Vdash$  are related as expected, in that provability means the existence of a proof:

$\text{Rel}_{\vdash}^{\Vdash}$ : For all  $\varphi \in \text{Sen}$ ,  $\vdash \varphi$  iff there exists  $p \in \text{Proof}$  such that  $p \Vdash \varphi$ .

## 2.2 The arithmetic substratum

We extend the generic syntax assumptions with a subset  $\text{Num} \subseteq \text{Term}$ , of *numerals*, ranged over by  $m, n$ , which are assumed to be closed, i.e., have no free variables.

**Convention 1.** In all the shown results we implicitly assume: (1) the generic syntax (free variable and substitution) axioms, (2) at least  $\rightarrow$  and  $\perp$  plus whatever connectives and quantifiers appear in the statement, (3) closedness of  $\vdash$  under intuitionistic deduction rules, and (4) the existence of numerals. Other assumptions (e.g., order-like relation axioms, consistency, standard models, etc.) will be indicated explicitly.

On one occasion, we will assume an order-like binary relation modeled by a formula  $\prec \in \text{Fmla}_2$ . We write  $t_1 \prec t_2$  instead of  $\prec(t_1, t_2)$  and  $\forall x \prec n. \varphi$  instead of  $\forall x. x \prec n \rightarrow \varphi$ . It turns out that at our level of abstraction it does not matter whether  $\prec$  is a strict or a non-strict order. Indeed, we only require the following two properties, where  $x \in M$  denotes  $\bigvee_{m \in M} x \equiv m$  and  $\bigvee$  expresses the disjunction of a finite set of formulas:

$\text{Ord}_1$ : For all  $\varphi \in \text{Fmla}_1$  and  $n \in \text{Num}$ , if  $\vdash \varphi(m)$  for all  $m \in \text{Num}$ , then  $\vdash \forall x \prec n. \varphi(x)$ .

$\text{Ord}_2$ : For all  $n \in \text{Num}$ , there exists a finite set  $M \subseteq \text{Num}$  such that  $\vdash \forall x. x \in M \vee n \prec x$ .

$\text{Ord}_1$  states that if a property  $\varphi$  is provable for all numerals, then its universal quantification bounded by any given numeral  $n$  is also provable. Having in mind the arith-

metic interpretation of numerals, it would also make sense to assume a stronger version of  $\text{Ord}_1$ , replacing “if  $\vdash \varphi(m)$  for all  $m \in \text{Num}$ ” by the weaker hypothesis “if  $\vdash \varphi(m)$  for all  $m \in \text{Num}$  such that  $\vdash m < n$ ”. But this stronger version will not be needed.

$\text{Ord}_2$  states that, for any numeral  $n$ , any element  $x$  in the domain of discourse is either greater than  $n$  or equal to one of a finite set  $M$  of numerals. If we instantiate our syntax to that of first-order arithmetic, then the natural number model satisfies  $\text{Ord}_1$  and  $\text{Ord}_2$  when interpreting  $<$  as either  $<$  or  $\leq$ . Moreover, these properties are provable in intuitionistic Robinson arithmetic, again for both  $<$  and  $\leq$ .

### 2.3 Encodings and representability

Central in the incompleteness theorems are functions that encode formulas and proofs as numerals,  $\langle \_ \rangle : \text{Fmla} \rightarrow \text{Num}$  and  $\langle \_ \rangle : \text{Proof} \rightarrow \text{Num}$ . For our abstract results, the encodings are not required to be injective or surjective.

Let  $A_1, \dots, A_m$  be sets, and let, for each of them,  $\langle \_ \rangle : A_i \rightarrow \text{Num}$  be an “encoding” function to numerals. Then, an  $m$ -ary relation  $R \subseteq A_1 \times \dots \times A_m$  is said to be *represented* by a formula  $\textcircled{R} \in \text{Fmla}_m$  if the following hold for all  $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$ :

- $(a_1, \dots, a_m) \in R$  implies  $\vdash \textcircled{R}(\langle a_1 \rangle, \dots, \langle a_m \rangle)$
- $(a_1, \dots, a_m) \notin R$  implies  $\vdash \neg \textcircled{R}(\langle a_1 \rangle, \dots, \langle a_m \rangle)$

Let  $A$  be another set with  $\langle \_ \rangle : A \rightarrow \text{Num}$ . An  $m$ -ary function  $f : A_1 \times \dots \times A_m \rightarrow A$  is said to be *represented* by a formula  $\textcircled{f} \in \text{Fmla}_{m+1}$  if for all  $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$ :

- $\vdash \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, \langle f(a_1, \dots, a_m) \rangle)$
- $\vdash \forall x, y. \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, x) \wedge \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, y) \rightarrow x = y$

The notion of a function being represented is stronger than that of its graph being represented (as a relation)—but with enough deductive power they are equivalent [32, §16]. We will need an even stronger notion: A function  $f$  as above is *term-represented* by an operator  $\textcircled{f} : \text{Term}^m \rightarrow \text{Term}$  if  $\vdash \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle) \equiv \langle f(a_1, \dots, a_m) \rangle$  for all  $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$ . When the formula by which a relation/function  $P$  is represented or term-represented is irrelevant, we call  $P$  *representable* or *term-representable*.

We will also need an enhancement of relation representability: Given  $i < m$ , we call the representation of an  $m$ -ary relation  $R$  by  $\textcircled{R}$   *$i$ -clean* if  $\vdash \neg \textcircled{R}(n_1, \dots, n_m)$  for all  $n_1, \dots, n_m$  such that  $n_i$  is outside the image of  $\langle \_ \rangle$  (i.e., there is no  $a \in A_i$  with  $n_i = \langle a \rangle$ ). Cleaness would be trivially satisfied if the encodings were surjective. However, surjectivity is not a reasonable assumption. For example, most of the numeric encodings used in the literature are injective but not surjective.

We let  $S : \text{Fmla}_1 \rightarrow \text{Sen}$  be the *self-substitution* function, which sends any  $\varphi \in \text{Fmla}_1$  to  $\varphi(\langle \varphi \rangle)$ , i.e., to the sentence obtained from  $\varphi$  by substituting the encoding of  $\varphi$  for the unique variable of  $\varphi$ . An alternative is the following “soft” version of  $S$ , which sends any  $\varphi \in \text{Fmla}_1$  to  $\exists v_1. v_1 \equiv \langle \varphi \rangle \wedge \varphi$ , where  $v_1$  is the single free variable of  $\varphi$ . The soft version yields provably equivalent formulas and has the advantage that it is easier to represent inside the logic, since it does not require formalizing the complexities of capture-avoiding substitution. All our results involving  $S$  have been proved for both versions.

We will consider the properties  $\text{Repr}_{\neg}$ ,  $\text{Repr}_S$ , and  $\text{Repr}_{\Vdash}$ , stating the representability of the functions  $\neg$  and  $S$ , and of the relation  $\Vdash$ ; and  $\text{TRepr}_{\neg}$ , stating the term-representability of  $\neg$ . In addition,  $\text{Clean}_{\Vdash}$  will state that the considered representation of  $\Vdash$  is 1-clean, i.e., it is clean on the proof component. For the representing formulas for the above relations and functions we will use their circled names,  $\ominus$ ,  $\oplus$ , etc.; for example,  $\text{Repr}_{\Vdash}$  means that (1)  $p \Vdash \varphi$  implies  $\vdash \oplus(\langle p \rangle, \langle \varphi \rangle)$  and (2)  $p \not\Vdash \varphi$  implies  $\vdash \ominus(\langle p \rangle, \langle \varphi \rangle)$  for all  $p \in \text{Proof}$  and  $\varphi \in \text{Sen}$ .  $\vdash \boxed{\neg} \langle \varphi \rangle \equiv \langle \neg \varphi \rangle$  for all  $\varphi \in \text{Sen}$ .

## 2.4 Derivability conditions

Most of our assumptions refer to representability. An important exception is the provability relation  $\vdash$ , for which only a weakening of representability is reasonable. Let  $\oplus \in \text{Fml}_{a_1}$  be the formula for this task. We consider the following assumptions about  $\oplus$ , known as the Hilbert-Bernays-Löb derivability conditions:

- HBL<sub>1</sub>:  $\vdash \varphi$  implies  $\vdash \oplus \langle \varphi \rangle$  for all  $\varphi \in \text{Sen}$ .
- HBL<sub>2</sub>:  $\vdash \oplus \langle \varphi \rangle \wedge \oplus \langle \varphi \rightarrow \psi \rangle \rightarrow \oplus \langle \psi \rangle$  for all  $\varphi, \psi \in \text{Sen}$ .
- HBL<sub>3</sub>:  $\vdash \oplus \langle \varphi \rangle \rightarrow \oplus \langle \oplus \langle \varphi \rangle \rangle$  for all  $\varphi \in \text{Sen}$ .

Above and elsewhere, to lighten notation we omit parentheses when instantiating one-variable formulas with encodings of formulas—e.g., writing  $\oplus \langle \varphi \rangle$  instead of  $\oplus(\langle \varphi \rangle)$ .

HBL<sub>1</sub> states that, if a sentence is provable, then its encoding is also provable inside the representation. HBL<sub>3</sub> is roughly a formulation of HBL<sub>1</sub> “one level up,” inside the proof system  $\vdash$ . Finally, note that the provability relation is closed under *modus ponens*, in that  $\vdash \varphi$  and  $\vdash \varphi \rightarrow \psi$  implies  $\vdash \psi$  for all  $\varphi, \psi \in \text{Sen}$ . Thus, HBL<sub>2</sub> roughly states the same property inside the proof system. In short, the derivability conditions state that the representation of provability acts partly similarly to the provability relation. Note that the representability of “proof of” implies HBL<sub>1</sub>, taking  $\oplus(x)$  to be  $\exists y. \oplus(y, x)$ .

**Convention 2.** In this paper, we focus on the standard provability representation: Whenever we assume explicit proofs and representability of “proof of,” the formula  $\oplus$  will be defined from  $\oplus$  as shown above.

We will also be interested in the following variations of the derivability conditions:

- HBL<sub>4</sub>:  $\vdash \oplus \langle \varphi \rangle \wedge \oplus \langle \psi \rangle \rightarrow \oplus \langle \varphi \wedge \psi \rangle$  for all  $\varphi, \psi \in \text{Sen}$ .
- HBL<sub>1</sub><sup>≠</sup>:  $\vdash \oplus \langle \varphi \rangle$  implies  $\vdash \varphi$  for all  $\varphi \in \text{Sen}$ .
- SHBL<sub>3</sub>:  $\vdash \oplus \langle t \rangle \rightarrow \oplus \langle \oplus \langle t \rangle \rangle$  for all closed terms  $t$ .
- WHBL<sub>2</sub>:  $\vdash \varphi \rightarrow \psi$  implies  $\vdash \oplus \langle \varphi \rangle \rightarrow \oplus \langle \psi \rangle$  for all  $\varphi, \psi \in \text{Sen}$ .

HBL<sub>4</sub> has a similar flavor as HBL<sub>2</sub>, but refers to conjunction: It states that the conjunction introduction rule holds inside the proof system. HBL<sub>1</sub><sup>≠</sup> is the converse of HBL<sub>1</sub>. Finally, SHBL<sub>3</sub> is a strengthening of HBL<sub>3</sub> holding for all closed terms and not only those that encode sentences, and (if we assume HBL<sub>1</sub>) WHBL<sub>2</sub> is a weakening of HBL<sub>2</sub>.

## 2.5 Standard models

We fix a unary relation  $\models \subseteq \text{Sen}$ , representing *truth of a sentence in the standard model*. We write  $\models \varphi$  instead of  $\varphi \in \models$ , and read it as “ $\varphi$  is true.” We consider the assumptions:

- Syn<sub>\models</sub>: Syntactic entities (logical connectives and quantifiers) handle truth as expected:
  - (1)  $\not\models \perp$ ;
  - (2) for all  $\varphi, \psi \in \text{Sen}$ ,  $\models \varphi$  and  $\models \varphi \rightarrow \psi$  imply  $\models \psi$ ;

- (3) for all  $\varphi \in \text{Fmla}_1$ , if  $\models \varphi(n)$  for all  $n \in \text{Num}$  then  $\models \forall x. \varphi(x)$ ;
- (4) for all  $\varphi \in \text{Fmla}_1$ , if  $\models \exists x. \varphi(x)$  then  $\models \varphi(n)$  for some  $n \in \text{Num}$ ;
- (5) for all  $\varphi \in \text{Sen}$ ,  $\models \varphi$  or  $\models \neg \varphi$ .

Soundness (of provability with respect to truth):  $\vdash \varphi$  implies  $\models \varphi$  for all  $\varphi \in \text{Sen}$ .

$\text{Syn}_{\models}(1-4)$  only contains a partial description of the syntactic entities' behavior—corresponding to elimination rules for  $\perp$ ,  $\rightarrow$  and  $\exists$  and introduction rule for  $\forall$ . For our results this suffices.  $\text{Syn}_{\models}(5)$  states that standard models decide every sentence.

On his way to formalizing  $\mathcal{IT}_2$  for extensions of the HF set theory, after proving  $\text{HBL}_1$  Paulson notes [25, p.21]: “The reverse implication [namely  $\text{HBL}_1^{\Leftarrow}$ ], despite its usefulness, is not always proved.” In his abstract account, Buldt also assumes  $\text{HBL}_1^{\Leftarrow}$  in his most general formulation of  $\mathcal{IT}_1$  [4, Theorem 3.1]; that formulation has in mind not necessarily the standard provability representation (our Convention 2), but any formula that weakly represents  $\vdash$ , which is acceptable for  $\mathcal{IT}_1$  but not for  $\mathcal{IT}_2$  [1].

We avoid such  $\mathcal{IT}_1$  versus  $\mathcal{IT}_2$  divergence by remaining focused on the standard provability representation. In this case, for arithmetics and related theories,  $\text{HBL}_1^{\Leftarrow}$  cannot be inferred without assuming soundness in the standard model (which Paulson does), or at least  $\omega$ -consistency. We can depict the situation abstractly, without knowing how standard models look like:

- Lemma 3.** (1) Assume  $\text{Rel}_{\vdash}^{\perp}$ ,  $\text{Repr}_{\vdash}$ ,  $\text{Clean}_{\vdash}$  and  $\text{OCon}$ . Then  $\text{HBL}_1^{\Leftarrow}$  holds.  
(2) Assume Soundness and  $\text{Syn}_{\models}(1,2,3)$ . Then  $\text{OCon}$  holds.  
(3) Assume  $\text{Rel}_{\vdash}^{\perp}$ ,  $\text{Repr}_{\vdash}$ ,  $\text{Clean}_{\vdash}$ , Soundness and  $\text{Syn}_{\models}(1,2,4)$ . Then  $\models \oplus \langle \varphi \rangle$  implies  $\vdash \varphi$  for all  $\varphi \in \text{Sen}$ . In particular,  $\text{HBL}_1^{\Leftarrow}$  holds.

Thus, staying in a proof-theoretic world,  $\omega$ -consistency ensures  $\text{HBL}_1^{\Leftarrow}$  if the “proof of” relation is cleanly represented (1). In turn,  $\omega$ -consistency is ensured by minimal semantic requirements, including the soundness of provability (2). Finally, putting together representability and semantics, we can infer something stronger than  $\text{HBL}_1^{\Leftarrow}$ : That the mere truth (and not just the provability) of a sentence's provability representation implies the provability of the sentence itself (3).

It follows from either points (1,2) or point (3) of the lemma that, in the presence of standard models and soundness, clean representability of the “proof of” relation implies  $\text{HBL}_1^{\Leftarrow}$ ; and recall that it also implies  $\text{HBL}_1$ . So it implies an “iff” version of  $\text{HBL}_1$ :  $\vdash \varphi$  if and only if  $\vdash \oplus \langle \varphi \rangle$ . Interestingly, a converse of this implication also holds. To state it, we initially assume there is no “outer” notion of proof (i.e., no set  $\text{Proof}$  and no relation  $\Vdash$ ), but only an “inner” one, given by a formula  $P \in \text{Fmla}_2$  such that:

- $\text{Rel}_{\oplus}^P: \vdash \oplus \langle \varphi \rangle \leftrightarrow \exists x. P(x, \langle \varphi \rangle)$ .
- $\text{Compl}_P: \models P(n, \langle \varphi \rangle)$  implies  $\vdash P(n, \langle \varphi \rangle)$  for all  $n \in \text{Num}$  and  $\varphi \in \text{Sen}$ .
- $\text{Compl}_{\neg P}: \models \neg P(n, \langle \varphi \rangle)$  implies  $\vdash \neg P(n, \langle \varphi \rangle)$  for all  $n \in \text{Num}$  and  $\varphi \in \text{Sen}$ .

$\text{Rel}_{\oplus}^P$  is the inner version of  $\text{Rel}_{\vdash}^{\perp}$ : It expresses that, *inside the representation*, proofs and provability are connected as expected.  $\text{Compl}_P$  and  $\text{Compl}_{\neg P}$  state that provability is complete on  $P$  statements about formula encodings, as well as their negations; in traditional settings, this is true thanks to  $P$  being a bounded arithmetical formula ( $\mathcal{A}_0$ ). Now the converse result states that, thanks to (standard models and) the “iff” version of  $\text{HBL}_1$ , we can define an outer notion of proof that is represented by the inner notion  $P$ :



**Lemma 4.** Assume  $\text{Rel}_{\oplus}^{\text{P}}$ , Soundness,  $\text{Syn}_{\models}(4,5)$ ,  $\text{Compl}_{\text{P}}$ ,  $\text{Compl}_{\neg\text{P}}$ ,  $\text{HBL}_1$  and  $\text{HBL}_1^{\Leftarrow}$ . Take  $\text{Proof} = \text{Num}$  and define  $\Vdash$  by  $n \Vdash \varphi$  iff  $\vdash \text{P}(n, \langle \varphi \rangle)$ . Then  $\text{Rel}_{\oplus}^{\text{P}}$ ,  $\text{Repr}_{\Vdash}$  and  $\text{Clean}_{\Vdash}$  hold, with  $\Vdash$  being represented by  $\text{P}$ .

### 3 Abstract Incompleteness Theorems

After last section’s preparations, we are now ready to discuss different versions of the incompleteness theorems and their major lemmas, based on alternative assumptions.

#### 3.1 Diagonalization

The formula diagonalization technique (due to Gödel and Carnap [5]) yields “self-referential” sentences. All we need for it to work is the representability of substitution.

**Prop 5.** Assuming  $\text{Repr}_{\text{S}}$ , for all  $\psi \in \text{Fmla}_1$  there exists  $\varphi \in \text{Fmla}_1$  with  $\vdash \varphi \leftrightarrow \psi\langle \varphi \rangle$ .

A sentence  $\varphi \in \text{Sen}$  is called a *Gödel sentence* if  $\vdash \varphi \leftrightarrow \neg \oplus\langle \varphi \rangle$ ; it is called a *Rosser sentence* if  $\vdash \varphi \leftrightarrow \neg (\exists x. \oplus(x, \langle \varphi \rangle) \wedge \text{RosserTwist}(x, \langle \varphi \rangle))$ , where we define  $\text{RosserTwist}(x, y) = \forall x'. x' \prec x \rightarrow \forall y'. \ominus(y, y') \rightarrow \neg \oplus(x', y')$ . The existence of Gödel and Rosser sentences follows immediately from diagonalization.

**Prop 6.** Assuming  $\text{Repr}_{\text{S}}$ , there exist Gödel and Rosser sentences.

Thus, any Gödel sentence is provably equivalent to the negation of its own provability; in Gödel’s words [11], it “says about itself that it is not provable.” A Rosser sentence  $\varphi$  asserts its own unprovability in a weaker fashion: Rather than saying “Myself,  $\varphi$ , am not provable” (i.e., “it is not the case that there exists a proof  $p$  of  $\varphi$ ”), it says “it is not the case that there exists a proof  $p$  of  $\varphi$  such that, for all smaller proofs  $q$ ,  $q$  is not a proof of  $\neg \varphi$ .” Here, “smaller” refers to the order the encoding of proofs as numerals imposes.

#### 3.2 The incompleteness theorems

$\mathcal{IT}_1$  identifies sentences that are neither provable nor disprovable—which often holds for Gödel and Rosser sentences with the help of a provability relation satisfying  $\text{HBL}_1$ .

**Prop 7.** Assume  $\text{Con}$  and  $\text{HBL}_1$ . Then  $\not\vdash G$  for all Gödel sentences  $G$ .

For showing that the Gödel sentences are not disprovable, a standard route is to assume explicit proofs, strengthen the consistency assumption to  $\omega$ -consistency, and strengthen  $\text{HBL}_1$  to representability of the “proof of” relation.

**Prop 8.** Assume  $\text{OCon}$ ,  $\text{Rel}_{\oplus}^{\text{P}}$ ,  $\text{Repr}_{\Vdash}$ ,  $\text{Clean}_{\Vdash}$ . Then  $\not\vdash \neg G$  for all Gödel sentences  $G$ .

*Proof.* Let  $G$  be a Gödel sentence. We prove  $\not\vdash \neg G$  by contradiction. Assume (1)  $\vdash \neg G$ .

- By consistency (which is implied by  $\text{OCon}$ ), we obtain  $\not\vdash G$ .
- From this and  $\text{Rel}_{\oplus}^{\text{P}}$ , we obtain  $p \not\vdash G$  for all  $p \in \text{Proof}$ .
- From this,  $\text{Repr}_{\Vdash}$  and  $\text{Clean}_{\Vdash}$ , we obtain  $\vdash \neg \oplus(n, \langle G \rangle)$  for all  $n \in \text{Num}$ .
- From this and  $\text{OCon}$ , we obtain  $\not\vdash \neg \exists x. \oplus(x, \langle G \rangle)$ , i.e.,  $\not\vdash \neg \oplus\langle G \rangle$ .
- Hence, since  $G$  is a Gödel sentence, we obtain  $\not\vdash \neg G$ , which contradicts (1).  $\square$

While the line of reasoning in the above proof is mostly well-known, it contains two subtle points about which the literature is not explicit (due to the usual focus on classical first-order arithmetic and particular choices of encodings).

First, we must assume the representation of the “proof of” relation to be 1-clean, i.e., clean with respect to the proof component. Indeed, the argument crucially relies on converting the statement “ $p \not\vdash G$  for all  $p \in \text{Proof}$ ” into “ $\vdash \neg \oplus(n, \langle G \rangle)$  for all  $n \in \text{Num}$ ,” which is only possible for 1-clean encodings. This assumption will be repeatedly needed in later results. By contrast, cleanness is never required with respect to the sentence component of “proof of” or for the provability relation (which only involves sentence encodings). In short, cleanness is only needed for proofs, not for sentences.

Second, to reach the desired contradiction for our intuitionistic proof system  $\vdash$ , from “ $\vdash \neg \oplus(n, \langle G \rangle)$  for all  $n \in \text{Num}$ ” it is not sufficient to employ standard  $\omega$ -consistency, which would only give us  $\not\vdash \exists x. \oplus(x, \langle G \rangle)$ , i.e.,  $\not\vdash \oplus(G)$ ; the last together with  $\vdash G \leftrightarrow \neg \oplus(G)$  would be insufficient for obtaining  $\not\vdash \neg G$ . However, our stronger version of  $\omega$ -consistency,  $\text{OCon}$ , does the trick.  $\mathcal{IT}_1$  now follows by putting together Props. 6–8:

**Theorem 9.** ( $\mathcal{IT}_1$ ) Assume  $\text{OCon}$ ,  $\text{Rel}_{\vdash}^{\perp}$ ,  $\text{Repr}_{\perp}$ ,  $\text{Clean}_{\perp}$  and  $\text{Repr}_{\mathcal{S}}$ . Then:

- (1) There exists a Gödel sentence.      (2)  $\not\vdash G$  and  $\not\vdash \neg G$  for all Gödel sentences  $G$ .

Rosser’s contribution to  $\mathcal{IT}_1$  was an ingenious trick for weakening the  $\omega$ -consistency assumption into plain consistency—as such, it is usually seen as a *strict improvement* over Gödel’s version. While this is true for the concrete case of FOL theories extending arithmetic, from an abstract perspective the situation is more nuanced: The improvement is achieved at the cost of asking more from the logic. Our framework makes this trade-off clearly visible. The idea is to use Rosser sentences instead of Gödel sentences to “repair” the  $\omega$ -consistency assumption of Theorem 9 (inherited from Prop. 8):

**Theorem 10.** (Rosser’s  $\mathcal{IT}_1$ ) Assume  $\text{Con}$ ,  $\text{Ord}_1$ ,  $\text{Ord}_2$ ,  $\text{Repr}_{\neg}$ ,  $\text{Rel}_{\vdash}^{\perp}$ ,  $\text{Repr}_{\perp}$ ,  $\text{Clean}_{\perp}$  and  $\text{Repr}_{\mathcal{S}}$ . Then:

- (1) There exists a Rosser sentence.      (2)  $\not\vdash R$  and  $\not\vdash \neg R$  for all Rosser sentences  $R$ .

Highlighted is the assumption trade-off between the two versions: Rosser’s weakening of  $\omega$ -consistency into consistency is paid by additionally assuming representability of negation and an order-like relation satisfying  $\text{Ord}_1$  and  $\text{Ord}_2$ . Certainly, negation representability is not a big price, since for concrete logics this tends to be a lemma that is anyway needed when proving  $\text{HBL}_1$ . On the other hand, the ordering assumptions seem to be a significant generality gap in favor of Gödel’s version. A clear manifestation of this gap is in our inference of a semantic version of  $\mathcal{IT}_1$ —which we obtain from Theorem 9 with the help of Lemmas 3(2) and 4:

**Theorem 11.** (Semantic  $\mathcal{IT}_1$ ) Assume  $\text{Rel}_{\oplus}^{\text{P}}$ ,  $\text{Soundness}$ ,  $\text{Syn}_{\models}$ ,  $\text{Compl}_{\text{P}}$ ,  $\text{HBL}_1$ ,  $\text{HBL}_1^{\leftarrow}$ , and  $\text{Repr}_{\mathcal{S}}$ . Then:

- (1) There exists a Gödel sentence. (2)  $\models G$ ,  $\not\vdash G$  and  $\not\vdash \neg G$  for all Gödel sentences  $G$ .

We have highlighted the assumptions specific to the semantic treatment. They replace  $\text{OCon}$ ,  $\text{Rel}_{\vdash}^{\perp}$ ,  $\text{Repr}_{\perp}$  and  $\text{Clean}_{\perp}$  from the proof-theoretic Theorem 9. Also highlighted is the additional fact concluded: that the Gödel sentences are true.

We have inferred the semantic version from Gödel’s proof-theoretic version (Theorem 9), and not from Rosser’s variation (Theorem 10). This is because in the semantic version  $\omega$ -consistency comes for free (from Lemma 3(2)). By contrast, for deploying Rosser’s version we would need to explicitly consider the order-like relation with its own hypotheses. This would have led to a *strictly less general* abstract result (if we ignore the difference in the way Gödel and Rosser sentences are actually defined).

The semantic  $\mathcal{IT}_1$  relies on  $\text{HBL}_1^{\Leftarrow}$ . If we additionally commit to classical logic (i.e., assume  $\vdash \neg \neg \varphi \rightarrow \varphi$ ), we can more directly show, taking advantage of  $\text{HBL}_1^{\Leftarrow}$ , that the Gödel sentences are not disprovable, which immediately proves  $\mathcal{IT}_1$ :

**Theorem 12.** (Classical  $\mathcal{IT}_1$ ) Assume classical logic, Con,  $\text{HBL}_1$ ,  $\text{HBL}_1^{\Leftarrow}$ ,  $\text{Repr}_S$ . Then:

- (1) There exists a Gödel sentence.      (2)  $\not\vdash G$  and  $\not\vdash \neg G$  for all Gödel sentences  $G$ .

Even though  $\mathcal{IT}_1$  uses a predicate  $\oplus$  that weakly represents provability, its conclusion (the existence of undecided sentences) is meaningful regardless of whether  $\oplus$  *adequately expresses provability*. By contrast, the meaning of  $\mathcal{IT}_2$ ’s conclusion, the theory cannot prove its own consistency, relies on this (non-mathematical) “intensional” assumption [1]. In this case, consistency is adequately expressed by the sentence  $\neg \oplus \langle \perp \rangle$ . The standard formulation (and proof) of  $\mathcal{IT}_2$  uses all three derivability conditions:

**Theorem 13.** ( $\mathcal{IT}_2$ ) Assume Con,  $\text{HBL}_1$ ,  $\text{HBL}_2$ ,  $\text{HBL}_3$  and  $\text{Repr}_S$ . Then  $\not\vdash \neg \oplus \langle \perp \rangle$ .

### 3.3 Jeroslow’s approach

Next we look into an alternative line of reasoning due to Jeroslow [15], often cited as a simplification of the canonical route to prove  $\mathcal{IT}_2$  [29, 32, 33]. To study its features and pitfalls, we need some standard notation used by Jeroslow. A *pseudo-term* is a formula  $\varphi \in \text{Fmla}_m$  expressing a provably functional relation via “exists unique”:  $\vdash \forall x_1, \dots, x_m. \exists! y. \varphi(x_1, \dots, x_m, y)$ . Next we only discuss the case  $m = 2$ ; the general case is similar.

**Notation 14.** Given a pseudo-term  $\varphi \in \text{Fmla}_2$ , we treat it as if it is a one-variable term:

- for any terms  $s$  and  $t$ , we write  $t \equiv \varphi(s)$  instead of  $\varphi(s, t)$ ;
- for any term  $s$  and formula  $\psi \in \text{Fmla}_1$ , we write  $\psi(\varphi(s))$  instead of  $\exists y. \varphi(s, y) \wedge \psi(y)$ .

This notation smoothly integrates pseudo-terms with terms: If  $\vdash t \equiv \varphi(s)$  and  $\vdash \psi(\varphi(s))$  then  $\vdash \psi(t)$ , where  $\psi(t)$  denotes actual substitution of terms in formulas.

Jeroslow relies on an abstract class of  $m$ -ary functions,  $\mathcal{F}_m \subseteq \text{Num}^m \rightarrow \text{Num}$ , for all arities  $m \in \mathbb{N}$ , on which he considers following assumptions:

$\text{Repr}_{\mathcal{F}}$ : Every  $f \in \mathcal{F}_m$  is represented by some pseudo-term  $\langle f \rangle \in \text{Fmla}_{m+1}$  under the identity encoding  $\text{Num} \rightarrow \text{Num}$ .

$\text{CapN}$ : Some  $\mathbb{N} \in \mathcal{F}_1$  correctly captures negation:  $\mathbb{N}\langle \varphi \rangle = \langle \neg \varphi \rangle$  for all  $\varphi \in \text{Sen}$ .

$\text{CapSS}$ : Some  $\text{ssap} : \text{Fmla}_1 \rightarrow \mathcal{F}_1$  correctly captures substituted self-application:

$$\text{ssap } \psi \langle f \rangle = \langle \psi(\langle f \rangle \langle f \rangle) \rangle \text{ for all } \psi \in \text{Fmla}_1 \text{ and } f \in \mathcal{F}_1.$$

In  $\text{CapSS}$ , following Jeroslow we employed Notation 14 taking advantage of the fact that  $\langle f \rangle$  are pseudo-terms: The highlighted text denotes  $\exists y. \langle f \rangle(\langle f \rangle, y) \wedge \psi(y)$ . Moreover, using the same notation, the statement of  $\text{Repr}_{\mathcal{F}}$  for some  $f \in \mathcal{F}_1$  and  $n \in \text{Num}$  would be written as  $\vdash f(n) \equiv \langle f \rangle(n)$ . Similarly, combining  $\text{CapN}$  with the instance of  $\text{Repr}_{\mathcal{F}}$ , we obtain a fact that can be written as  $\vdash \langle \neg \varphi \rangle \equiv \mathbb{N}\langle \varphi \rangle$ .

When our logical theory is a recursive extension of the Robinson arithmetic and  $\text{Num} = \mathbb{N}$ ,  $\mathcal{F}_m$  could be the set of  $m$ -ary computable functions. Then every  $f \in \mathcal{F}_m$  would indeed be represented by a formula  $\langle f \rangle$ . Moreover, assuming a computable and injective encoding of formulas,  $\langle \_ \rangle : \text{Fmla}_1 \rightarrow \mathbb{N}$ , we can take  $N : \mathbb{N} \rightarrow \mathbb{N}$  to be the following computable function: Given input  $n$ , it checks if  $n$  has the form  $\langle \varphi \rangle$ ; if so, it returns  $\langle \neg \varphi \rangle$ ; if not, it returns any value (e.g., 0). And  $\text{ssap } \psi$  can be defined similarly, obtaining the desired property for every  $\varphi \in \text{Fmla}_2$ , not necessarily of the form  $\langle f \rangle$ . In short, Jeroslow’s assumptions cover arithmetic (but also potentially many other systems).

At the heart of Jeroslow’s approach lies an alternative diagonalization technique, producing *term* fixpoints, not just formula fixpoints:

**Lemma 15.** Assume  $\text{CapSS}$  and  $\text{Repr}_{\mathcal{F}}$  and let  $\psi \in \text{Fmla}_1$ . Then there exists a closed pseudo-term  $t$  such that  $\vdash t \equiv \langle \psi(t) \rangle$ . Moreover, taking  $\varphi = \psi(t)$ , we have  $\vdash \varphi \leftrightarrow \psi\langle \varphi \rangle$ .

*Proof.* Let  $f = \text{ssap } \psi$  and  $t = \langle f \rangle \langle f \rangle$ . From  $\text{CapSS}$ , we obtain  $f\langle f \rangle = \langle \psi(\langle f \rangle \langle f \rangle) \rangle$ . From this and  $\text{Repr}_{\mathcal{F}}$ , we obtain  $\vdash \langle f \rangle \langle f \rangle \equiv \langle \psi(\langle f \rangle \langle f \rangle) \rangle$ , i.e.,  $\vdash t \equiv \langle \psi(t) \rangle$ . With the equality rules, we obtain  $\vdash \psi(t) \leftrightarrow \psi(\langle \psi(t) \rangle)$ , i.e.,  $\vdash \varphi \leftrightarrow \psi\langle \varphi \rangle$ .  $\square$

This lemma offers us Gödel and Rosser sentences, which can be used like in Sections 3.1 and 3.2, leading to corresponding variants of  $\mathcal{IT}_1$ . But Jeroslow’s main innovation affects  $\mathcal{IT}_2$ : While traditionally  $\mathcal{IT}_2$  requires all three derivability conditions, Jeroslow’s version does not make use of the second,  $\text{HBL}_2$ :

**Theorem 16.** ( $\mathcal{IT}_2$  à la Jeroslow) Assume  $\text{Con}$ ,  $\text{HBL}_1$ ,  $\text{SHBL}_3$ ,  $\text{Repr}_{\mathcal{F}}$ ,  $\text{CapN}$ ,  $\text{CapSS}$ . Then  $\not\vdash \text{jcon}$ , where  $\text{jcon}$  denotes  $\forall x. \neg (\oplus(x) \wedge \ominus(\mathbb{N}(x)))$ .

Like with Rosser’s trick, we analyze this innovation’s trade-offs from an abstract perspective. A first trade-off is in the employment of a stronger version of the third condition,  $\text{SHBL}_3$  (extended to affect all closed pseudo-terms via Notation 14). Another is in the way consistency is expressed in the logic. Jeroslow does not conclude  $\not\vdash \neg \oplus(\perp)$ , but something more elaborate, namely  $\not\vdash \text{jcon}$ . While the formula  $\neg \oplus(\perp)$  internalizes the statement  $\not\vdash \perp$ ,  $\text{jcon}$  internalizes the equivalent statement “for all  $\varphi$ , it is not the case that  $\vdash \varphi$  and  $\vdash \neg \varphi$ .” But are the internalizations themselves equivalent, i.e., is it the case that  $\vdash \neg \oplus(\perp)$  iff  $\vdash \text{jcon}$ ? This surely holds for many concrete logics, but it is one direction that we can infer logic-independently: Assuming  $\text{HBL}_1$ ,  $\text{Repr}_{\mathcal{F}}$  and  $\text{CapN}$ ,  $\vdash \text{jcon}$  implies  $\vdash \neg \oplus(\perp)$ . And it seems we cannot infer the other direction without knowing how  $\oplus$  looks like more concretely. Therefore,  $\not\vdash \neg \oplus(\perp)$ , the conclusion of the original  $\mathcal{IT}_2$ , is *abstractly stronger than*, hence *preferable to*  $\not\vdash \text{jcon}$ . In short, Jeroslow somewhat weakens the theorem’s conclusion.

Let us now look at (a slight rephrasing of) Jeroslow's proof:

*Proof of Theorem 16.* We assume  $(*) \vdash \text{jcon}$  and aim to reach a contradiction.

- Applying Lemma 15 to  $\oplus(\mathbb{N}(x))$ , we obtain the closed term  $t$  such that  $(**) \vdash t \equiv \langle \oplus(\mathbb{N}(t)) \rangle$ .
- By SHBL<sub>3</sub> applied to  $\mathbb{N}(t)$ , we obtain  $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\oplus(\mathbb{N}(t)))$ .
- From  $(**)$  and the equality rules, we obtain  $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\mathbb{N}(\oplus(\mathbb{N}(t))))$ .
- The last two facts give us  $\vdash \varphi \rightarrow \oplus(\varphi) \wedge \oplus(\mathbb{N}(\varphi))$ , where  $\varphi$  denotes  $\oplus(\mathbb{N}(t))$ .
- On the other hand,  $(*)$  instantiated with  $\langle \varphi \rangle$  gives us  $\vdash \neg(\oplus(\varphi) \wedge \oplus(\mathbb{N}(\varphi)))$ .
- From the last two facts, we obtain  $(***) \vdash \neg \varphi$ .
- With HBL<sub>1</sub>, we obtain  $\vdash \oplus(\neg \varphi)$ .
- With CapN and Repr <sub>$\mathcal{F}$</sub> , we obtain  $\vdash \oplus(\mathbb{N}(\varphi))$ .
- From  $(**)$  and the equality rules, we obtain  $\vdash \oplus(\mathbb{N}(\oplus(\mathbb{N}(t)))) \rightarrow \oplus(\mathbb{N}(t))$ , i.e.,  $\vdash \oplus(\mathbb{N}(\varphi)) \rightarrow \varphi$ .
- From the last two facts, we obtain  $\vdash \varphi$ .
- This and  $(***)$  contradict the consistency assumption. □

A first major observation is that, under the stated assumptions, the above proof is *incorrect*. It uses an implicit assumption, hidden under Notation 14: When we disam-

biguate the notation, we see that Lemma 15 gives us a pseudo-term  $t$  that does not exactly satisfy  $(1) \vdash t \equiv \langle \psi(t) \rangle$  (which is what the theorem’s proof needs), but something weaker, namely  $(2) \vdash t \equiv \langle \chi \rangle$ , where  $\chi$  is  $\vdash \exists x. \mathcal{D}(\langle \mathcal{D} \rangle, x) \wedge \psi(x)$ . And although  $\vdash \chi \leftrightarrow \psi(t)$ , we still cannot infer (1) from (2), unless *the encodings of provably equivalent formulas are assumed provably equal*. But this assumption is unreasonable: Usually formula equivalence is undecidable, so no computable encoding can achieve that. (Incidentally, this problem is also the reason why we need SHBL<sub>3</sub> instead of HBL<sub>3</sub>: In the proof’s application of SHBL<sub>3</sub> to obtain  $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\oplus(\mathbb{N}(t)))$ , we cannot work with  $\langle \neg \varphi \rangle$  instead of  $\mathbb{N}(t)$ , even though  $\vdash \langle \neg \varphi \rangle \equiv \mathbb{N}(t)$ .)

To repair that, we can replace representation by pseudo-terms with actual term-representation. Namely, if we change  $\text{Repr}_{\mathcal{F}}$  into “Every  $f \in \mathcal{F}_m$  is term-represented by some  $\mathcal{J} : \text{Term}^m \rightarrow \text{Term}$  under the identity encoding  $\text{Num} \rightarrow \text{Num}$ ,” then CapSS, Lemma 15, and all proofs can operate with terms rather than pseudo-terms and everything will be formally correct. In summary, it seems that Jeroslow’s approach to  $\mathcal{IT}_2$  fails for pseudo-terms representing computable functions, but requires actual terms. This usually means that the logic must have built-in Skolem symbols and axioms.

Finally, let us see what it takes to alleviate the second trade-off: from  $\not\vdash \text{jcon}$  to the more desirable  $\not\vdash \neg \oplus(\perp)$ . We see that Theorem 16’s proof uses  $\vdash \text{jcon}$  not at  $\text{jcon}$ ’s full generality but only instantiated with formula encodings, which thanks to  $\text{Repr}_{\mathcal{F}}$  and CapN would follow from  $(*) \vdash \neg (\oplus\langle \varphi \rangle \wedge \oplus\langle \neg \varphi \rangle)$ . And it only takes WHBL<sub>2</sub> (a weaker version of HBL<sub>2</sub>) and HBL<sub>4</sub> to prove  $\vdash (\oplus\langle \varphi \rangle \wedge \oplus\langle \neg \varphi \rangle) \rightarrow \oplus(\perp)$ , allowing us to infer  $(*)$  from  $\vdash \neg \oplus(\perp)$ ; meaning that the latter could have been used. We obtain:

**Theorem 17.** If in the (corrected) Theorem 16 we additionally assume WHBL<sub>2</sub> and HBL<sub>4</sub>, its conclusion can be upgraded to  $\not\vdash \neg \oplus(\perp)$ .

Whether WHBL<sub>2</sub> and HBL<sub>4</sub> are a good trade-off for HBL<sub>2</sub> will of course depend on the logic’s specificity, in particular, on its primitive rules of inference.

Jeroslow presented his approach for an abstract logical theory over a FOL language, which is not necessarily a FOL theory—so it found a natural fit in our generic framework. To our knowledge, very few subsequent authors present Jeroslow’s approach rigorously, and none at its original level of generality. Smith’s monograph gives a rigorous account for arithmetic [32, §33], silently performing the correction we have shown here, but failing to detect the need for SHBL<sub>3</sub> instead of HBL<sub>3</sub> (which Jeroslow did not). A mechanical prover is of invaluable help for detecting such nuances and pitfalls.

We conclude with an anecdote involving our prover and Jeroslow’s notations. Given the simplicity of Lemma 15 (version with terms, not pseudo-terms), it came as no surprise that Isabelle’s Sledgehammer [26] was able to prove it automatically. Sledgehammer reported to have used  $\vdash$ -reflexivity in the proof. And indeed, it had found a term  $t$  for which it had proved not just  $\vdash t \equiv \langle \psi(t) \rangle$ , but actual equality,  $t = \langle \psi(t) \rangle$ ; in particular, the term was a numeral. All this was too good to be true. It took us some time to realize why that happened: Due to another one of Jeroslow’s notations, who wrote  $f$  instead of  $\mathcal{J}$  (thus identifying a function with its representing pseudo-term), we had at first misstated CapSS, writing  $\langle \psi(f\langle \mathcal{J} \rangle) \rangle$  instead of  $\langle \psi(\mathcal{J}\langle \mathcal{J} \rangle) \rangle$ ; the former is still a valid expression, since  $f$  is a function between numerals which are particular terms. Embarrassingly, it took us even longer to realize why this variation discovered by chance was not an improvement of Jeroslow’s diagonalization lemma: because the

assumption CapSS becomes unreasonable. Indeed, no concrete computable function can act like the intended ssap  $\psi$ : Given an input  $n$ , (1) decode it into a unique formula  $\varphi$  such that  $n = \langle \varphi \rangle$ , (2) decode  $\varphi$  into a unique function  $f$  such that  $\varphi = \langle f \rangle$  and (3) proceed to apply  $f$  as part of producing  $\langle \psi(f(\langle \mathcal{I} \rangle)) \rangle$ . The second step requires an injective and computable encoding of computable functions into formulas, which is impossible.

**Summary** Using our generic infrastructure (Section 2), we have formally proved several abstract incompleteness results. They include four versions of  $\mathcal{IT}_1$ :

- Gödel’s original  $\mathcal{IT}_1$  (Theorem 9) and an  $\mathcal{IT}_1$  based on classical logic (Theorem 12) required the formalization of some well-known arguments without change.
- Rosser’s  $\mathcal{IT}_1$  (Theorem 10) involved the generalization of a well-known argument: distilling two abstract conditions,  $\text{Ord}_1$  and  $\text{Ord}_2$ .
- A novel semantic  $\mathcal{IT}_1$  (Theorem 11) was born from analyzing the relationship between standard models, the “iff” version of  $\text{HBL}_1$ , and proof representability.

They also include two versions of  $\mathcal{IT}_2$ :

- The standard  $\mathcal{IT}_2$  based on the three derivability conditions (Theorem 13) again only required formalizing a well-known argument.
- The alternative, Jeroslow-style  $\mathcal{IT}_2$  (Theorems 16 and 17) involved a detailed analysis and correction of an existing abstract result.

## 4 Instances of the Abstract Results

We first validate the assumptions about our abstract logic and arithmetic:

**Prop 18.** (1) Any FOL theory that extends the Robinson arithmetic or the HF set theory satisfies all the axioms in our logical and arithmetical substrata (in Sections 2.1 and 2.2). (2) If, in addition, the theory is sound, then, together with its corresponding standard model, it also satisfies all our model-theoretic axioms (in Section 2.5).

In particular, point (2) shows that our discussion of standard models applies equally well to  $\mathbb{N}$  and the datatype of HF sets. (In the latter case, Num becomes the entire set of closed terms, so that numerals can denote arbitrary HF sets. This shows the versatility of our abstract concept of numeral.) Then we instantiate three of our main theorems:

**Theorem 19.** (1) Any FOL theory that extends the HF set theory with a finite set of axioms and **is sound in the standard HF set model** satisfies the hypotheses of Theorems 11 and 13. Hence  $\mathcal{IT}_1$  (semantic version) and  $\mathcal{IT}_2$  hold for it. (2) Any FOL theory that extends the HF set theory with a finite set of axioms and **is consistent** satisfies the hypotheses of Theorem 13. Hence  $\mathcal{IT}_2$  holds for it.

These instances are heavily based on the lemmas proved by Paulson in his formalization of  $\mathcal{IT}_1$  and  $\mathcal{IT}_2$  [24, 25], who follows and corrects Świerczkowski’s detailed informal account [34]. Point (1) is a restatement of Paulson’s formalized results: theorems *Goedel\_I* and *Goedel\_II* in [25]. (His theorems also assume consistency, but that is redundant: Consistency follows from his underlying soundness assumption.)

By contrast, point (2) is an upgrade of Paulson’s *Goedel\_II*, applicable to any consistent, though possibly unsound theory. This stronger version is in fact  $\mathcal{IT}_2$ ’s standard

form, free from any model-theoretic considerations.<sup>3</sup> Paulson had proved both  $\text{HBL}_1$  and  $\text{HBL}_1^{\Leftarrow}$  taking advantage of soundness, so we needed to discard  $\text{HBL}_1^{\Leftarrow}$  and re-prove  $\text{HBL}_1$  by replacing any semantic arguments with proofs within the HF calculus. We also removed all invocations of a convenient “truth implies provability for  $\Sigma$ -sentences” lemma, which depended on soundness due to Paulson’s choice of  $\Sigma$ -sentence definition.

This instantiation process has offered important feedback into the abstract results. A formal development such as ours is (largely) immune to reasoning errors, but not to missing out on useful pieces of generality. We experienced this firsthand with our assumptions about substitution. An *a priori* natural choice was to assume representability of the numeral substitution  $\text{Sb} : \text{Fmla}_1 \times \text{Num} \rightarrow \text{Sen}$  (defined as  $\text{Sb}(\varphi, n) = \varphi(n)$ ), part of which means (1)  $\vdash \text{SB}(\langle \varphi \rangle, n, \text{Sb}(\varphi, n))$ . But Paulson had instead proved (2)  $\vdash \text{SB}(\langle \varphi \rangle, \langle n \rangle, \text{Sb}(\varphi, n))$ . The key difference from (1) is that (2) applies the term encoding function  $\langle \_ \rangle : \text{Term} \rightarrow \text{Num}$  to numerals as well (as particular terms); and since his  $\langle \_ \rangle$  function is injective, it is far from the case that  $\langle n \rangle = n$  for all numerals  $n$ . Paulson’s version makes more sense than ours when building the results bottom-up: Representability should not discriminate numerals, but filter them through the encodings like other terms. However, top-down our version also made sense: It yielded the incompleteness theorems under reasonable assumptions, which do hold, by the way, for the HF set theory—even though in a bottom-up development one is unlikely to prove them. We resolved this discrepancy through a common denominator: the representability of self-substitution  $S : \text{Fmla}_1 \rightarrow \text{Sen}$  (Section 2.3), which made our results more general.

Paulson’s formalization has also inspired our abstract treatment of standard models (Section 2.5). Since Paulson proves  $\text{HBL}_1^{\Leftarrow}$  and uses classical logic, an obvious “port of entry” of his  $\mathcal{IT}_2$  into our framework is Theorem 12. But this theorem tells us nothing about the truth of the Gödel sentences. Delving deeper into Paulson’s proof, we noted that he (unconventionally) completely avoids  $\text{Repr}_{\Vdash}$ , and does not even define  $\Vdash$ . This raised the question of whether  $\text{HBL}_1^{\Leftarrow}$  and  $\text{Repr}_{\Vdash}$  are somehow interchangeable in the presence of standard models—and we found that they indeed are, under mild assumptions about truth. Incidentally, these assumptions were also sufficient for establishing the Gödel sentences’ truth, leading to our semantic  $\mathcal{IT}_1$  (Theorem 11), which we believe expresses the essence of the Świerczkowski–Paulson approach.

Many other logics and logical theories satisfy our theorems’ assumptions. We do *not* require the logic to be reducible to a single syntactic category of formulas,  $\text{Fmla}$ , a single syntactic judgment,  $\vdash$ , etc.; but only that such (well-behaved) formulas, provability relation, etc. are identifiable as part of that logic, e.g., localized to a given type and/or relativised by a given predicate. This allows our framework to capture most variants of higher-order logic and type theory, and also, we believe, many of the logics surveyed by Buldt [4], including non-classical and fuzzy. But enabling “mass instantiation” that is both formal and painless requires more progress on the agenda we started here: recognizing reusable construction and proof patterns and formalizing them as abstract results.

**Acknowledgments.** We are grateful to Bernd Buldt for his patient explanations on some of the material in his monograph.

<sup>3</sup> The value of freeing  $\mathcal{IT}_2$  from its reliance on models becomes even greater (1) for expressive systems such as ZF set theory and (2) for systems that lack obvious notions of (standard) models—such as the logic of Isabelle/HOL itself [19, 20].



## References

1. Auerbach, D.: Intensionality and the Gödel theorems. *Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition* 48(3), 337–351 (1985)
2. Blanchette, J.C., Popescu, A., Traytel, D.: Unified classical logic completeness—A coinductive pearl. In: Demri, S., Kapur, D., Weidenbach, C. (eds.) *IJCAR 2014*. LNCS, vol. 8562, pp. 46–60. Springer (2014)
3. Boolos, G.: *The Logic of Provability*. Cambridge University Press (1993)
4. Buldt, B.: The scope of Gödel’s first incompleteness theorem. *Logica Universalis* 8(3), 499–552 (2014)
5. Carnap, R.: *Logische syntax der sprache*. *Philosophical Review* 44(4), 394–397 (1935)
6. Davis, M.: *The Undecidable: Basic Papers on Undecidable Propositions, Unsolvability Problems, and Computable Functions*. Dover Publication (1965)
7. Diaconescu, R.: *Institution-independent Model Theory*. Birkhäuser, 1st edn. (2008)
8. Feferman, S., Dawson, Jr., J.W., Kleene, S.C., Moore, G.H., Solovay, R.M., van Heijenoort, J. (eds.): *Kurt Gödel: Collected Works. Vol. 1: Publications 1929–1936*. Oxford University Press (1986)
9. Fiore, M.P., Plotkin, G.D., Turi, D.: Abstract syntax and variable binding. In: *Logic in Computer Science (LICS) 1999*, pp. 193–202. IEEE Computer Society (1999)
10. Gabbay, M.J., Mathijssen, A.: Nominal (universal) algebra: Equational logic with names and binding. *J. Log. Comput.* 19(6), 1455–1508 (2009)
11. Gödel, K.: Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik* 38(1), 173–198 (1931)
12. Goguen, J.A., Burstall, R.M.: *Institutions: Abstract model theory for specification and programming*. *J. ACM* 39(1), 95–146 (1992)
13. Harrison, J.: HOL Light proof of Gödel’s first incompleteness theorem, located at <http://code.google.com/p/hol-light/>, directory Arithmetic
14. Hilbert, D., Bernays, P.: *Grundlagen der Mathematik, Vol. II* (1939)
15. Jeroslow, R.G.: Redundancies in the Hilbert-Bernays derivability conditions for Gödel’s second incompleteness theorem. *J. Symb. Log.* 38(3), 359–367 (1973)
16. Kaliszky, C., Urban, J.: HOL(y)Hammer: Online ATP service for HOL light. *Mathematics in Computer Science* 9(1), 5–22 (2015)
17. Kikuchi, M., Kurahashi, T.: Generalizations of Gödel’s incompleteness theorems for  $\Sigma_n$ -definable theories of arithmetic. *Rev. Symb. Logic* 10(4), 603–616 (2017)
18. Kreisel, G.: *Mathematical logic*. In: Saaty, T.L. (ed.) *Lectures on modern mathematics*, vol. 3. Wiley (1963)
19. Kunčar, O., Popescu, A.: A Consistent Foundation for Isabelle/HOL. In: *ITP*. pp. 234–252 (2015)
20. Kunčar, O., Popescu, A.: Comprehending Isabelle/HOL’s consistency. In: *ESOP*. pp. 724–749 (2017)
21. Löb, M.: Solution of a Problem of Leon Henkin. *The Journal of Symbolic Logic* 20(2), 115–118 (1955)
22. Nipkow, T., Paulson, L., Wenzel, M.: *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*, LNCS, vol. 2283. Springer (2002)
23. O’Connor, R.: Essential incompleteness of arithmetic verified by Coq. In: *TPHOLS*. pp. 245–260 (2005)
24. Paulson, L.C.: A machine-assisted proof of Gödel’s incompleteness theorems for the theory of hereditarily finite sets. *Rev. Symb. Logic* 7(3), 484–498 (2014)
25. Paulson, L.C.: A mechanised proof of Gödel’s incompleteness theorems using Nominal Isabelle. *J. Autom. Reasoning* 55(1), 1–37 (2015)

26. Paulson, L.C., Blanchette, J.C.: Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In: The 8th International Workshop on the Implementation of Logics, IWIL 2010, Yogyakarta, Indonesia, October 9, 2011. pp. 1–11 (2010)
27. Popescu, A., Rosu, G.: Term-generic logic. *Theor. Comput. Sci.* 577, 1–24 (2015)
28. Popescu, A., Traytel, D.: Formalization associated with this paper. [https://bitbucket.org/traytel/abstract\\_incompleteness/](https://bitbucket.org/traytel/abstract_incompleteness/) (2019)
29. Raatikainen, P.: Gödel’s incompleteness theorems. In: The Stanford Encyclopedia of Philosophy. Metaphysics Research Lab, Stanford University (2018)
30. Schlichtkrull, A., Blanchette, J.C., Traytel, D., Waldmann, U.: Formalizing Bachmair and Ganzinger’s ordered resolution prover. In: Galmiche, D., Schulz, S., Sebastiani, R. (eds.) IJCAR 2018. LNCS, vol. 10900, pp. 89–107. Springer (2018)
31. Shankar, N.: *Metamathematics, Machines, and Gödel’s Proof*. Cambridge University Press (1994)
32. Smith, P.: *An introduction to Gödel’s incompleteness theorems*. Cambridge University Press (2007)
33. Smoryński, C.: The incompleteness theorems. In: Barwise, J. (ed.) *Handbook of Mathematical Logic*, pp. 821–865. North-Holland (1977)
34. Świerczkowski, S.: Finite sets and Gödel’s incompleteness theorems. *Dissertationes Mathematicae* 422, 1–58 (2003)
35. Traski, A., Mostowski, A., Robinson, R.: *Undecidable Theories*. Studies in Logic and the Foundations of Mathematics. North-Holland (1953), 3rd edition, 1971
36. Troelstra, A.S., Schwichtenberg, H.: *Basic Proof Theory*. Cambridge University Press (1996)

# APPENDIX

This appendix contains more details on the concepts discussed in the main paper (Section A) and proof sketches for all the abstract results (Section B). The proof sketches have been provided for readers who are interested in the main line of reasoning but do not wish to inspect the formal Isabelle scripts.

## A More details about various concepts

### A.1 Free variables and substitution

The free-variable and substitution operators on terms and formulas have the following types, where  $\mathcal{P}_{\text{fin}}(\text{Var})$  denotes the set of finite subsets of  $\text{Var}$ .

- $\text{FVars} : \text{Term} \rightarrow \mathcal{P}_{\text{fin}}(\text{Var})$
- $\text{FVars} : \text{Fmla} \rightarrow \mathcal{P}_{\text{fin}}(\text{Var})$
- $\_[-/_] : \text{Term} \times \text{Term} \times \text{Var} \rightarrow \text{Term}$
- $\_[-/_] : \text{Fmla} \times \text{Term} \times \text{Var} \rightarrow \text{Fmla}$

We do not assume a particular syntax for terms or formulas. All we assume is that the above operators on the abstract sets  $\text{Term}$  and  $\text{Fmla}$  satisfy some properties:

**FVV:** Free-variables and substitution on variables behave as expected:

- $\text{FVars}(x) = \{x\}$ ,
- $x [s/x] = s$ , and  $y [s/x] = y$  if  $x \neq y$

**SVac:** Substitution on terms is vacuous outside the free variables:  $x \notin \text{FVars}(t)$  implies  $t [s/x] = t$ , and similarly for substitution on formulas

In addition, for the formula operators we assume the following:

**FVS:** Free-variables distribute over substitution:

$$\text{FVars}(\varphi [s/x]) = \text{FVars}(\varphi) - \{x\} \cup \text{FVars}(s) \text{ if } x \in \text{FVars}(\varphi)$$

**SVacV:** Substitution of a variable for itself is vacuous:  $\varphi [x/x] = \varphi$

**SC:** Substitution composes with itself under certain freshness assumptions:

- $\varphi [s_1/x] [s_2/x] = \varphi [(s_1 [s_2/x]) / x]$
- $\varphi [s_1/x_1] [s_2/x_2] = \varphi [(s_1 [s_2/x_2]) / x_1]$  if  $x_2 \notin \text{FVars}(\varphi)$
- $\varphi [s_1/x_1] [s_2/x_2] = \varphi [s_2/x_2] [(s_1 [s_2/x_2]) / x_1]$  if  $x_1 \neq x_2$  and  $x_1 \notin \text{FVars}(s_2)$

Of the above assumptions, FVV only applies to, and only makes sense for, substitution on terms. By contrast, we assume SVac for both terms and formulas. The last group, FVS, SVacV and SC, would make sense for terms too, but we only assume them for formulas—this is in line with our Economy principle, since we did not need them for terms.

### A.2 Numerals

Numerals are assumed to have no free variables:

$$\text{FVNum}: \text{FVars}(n) = \emptyset \text{ for all } n \in \text{Num}.$$

By SVac, this implies that substitution on numerals is vacuous.

### A.3 Abstract connectives and quantifiers

The main theorems depend on the existence of (subsets of) the following connectives and quantifiers, which we represent as operators on terms:

- Equality,  $\equiv : \text{Term} \rightarrow \text{Term} \rightarrow \text{Fmla}$
- False,  $\perp \in \text{Fmla}$
- Conjunction, implication and disjunction,  $\vee, \wedge, \rightarrow : \text{Fmla} \times \text{Fmla} \rightarrow \text{Fmla}$
- Universal and existential quantification,  $\forall, \exists : \text{Var} \times \text{Fmla} \rightarrow \text{Fmla}$

When we need negation,  $\neg$ , we define it from implication as False, taking  $\neg \varphi$  to be  $\varphi \rightarrow \perp$ . On the other hand, even in the presence of negation, we do not assume that  $\vee$  and  $\exists$  are definable from  $\wedge$  and  $\forall$  or vice versa—this is because, in line with the Economy principle, we use intuitionistic, not classical logic.

### A.4 Provability

In our formalization, we prefer to start with a Hilbert-style system because this is the most economic (requires the least amount of infrastructure). Namely, we consider the assumptions corresponding to the standard implication-based Hilbert system for intuitionistic logic [36]:

$$\begin{aligned} &\vdash \varphi \rightarrow \psi \text{ and } \vdash \varphi \text{ implies } \vdash \psi \\ &\vdash \varphi \rightarrow (\psi \rightarrow \varphi) \\ &\vdash (\varphi \rightarrow (\psi \rightarrow \chi)) \rightarrow (\varphi \rightarrow \psi) \rightarrow (\varphi \rightarrow \chi) \end{aligned}$$

Then each connective is integrated into this system on a need basis—if it is used in the theorems. For example, here are the assumptions for conjunction

$$\begin{aligned} &\vdash \varphi \wedge \psi \rightarrow \varphi \\ &\vdash \varphi \wedge \psi \rightarrow \psi \\ &\vdash \varphi \rightarrow (\psi \rightarrow \varphi \wedge \psi) \end{aligned}$$

and here are those for universal quantification:

$$\begin{aligned} &x \notin \text{FVars}(\varphi) \text{ and } \vdash \varphi \rightarrow \psi \text{ implies } \vdash \varphi \rightarrow (\forall x. \psi) \\ &\vdash (\forall x. \varphi) \rightarrow \varphi[t/x] \end{aligned}$$

And similarly for existential, equality, etc. Notice again that we do take the above to be an inductive definition of  $\vdash$ , but simply assume that  $\vdash$  is closed under those rules. In particular, the system for classical logic satisfies our assumptions.

Since a natural deduction system is more convenient for performing formal deductions, in our formalization we extract one from  $\vdash$ . Namely, in the presence of conjunction and implication, we extend  $\vdash$  to a binary relation  $\_ \vdash \_$  between finite sets  $F$  and formulas and single formulas  $\varphi$ , defined as  $F \vdash \varphi$  iff  $\vdash (\bigwedge F) \rightarrow \varphi$ . We think of  $F \vdash \varphi$  as “ $\varphi$  follows from  $F$ .” For this relation, we prove the standard introduction and elimination rules for each connective and quantifier.



### A.10 More on Rosser's variation of Gödel's First

The following propositions cover the two halves of Rosser's  $\mathcal{IT}_1$  (Theorem 10). The next one works with Rosser sentences instead of Gödel sentences, in order to avoid Prop. 8's  $\omega$ -consistency assumption:

**Prop 25.** Assume  $\text{Con}$ ,  $\text{Ord}_2$ ,  $\text{Rel}_{\perp}^{\perp}$ ,  $\text{Repr}_{\perp}$ ,  $\text{Repr}_{\neg}$  and  $\text{Clean}_{\perp}$ . Then  $\not\vdash \neg R$  for all Rosser sentences  $R$ .

*Proof.* To prove  $\not\vdash \neg R$ , we assume  $(1) \vdash \neg R$  and aim to reach a contradiction.

- By  $\text{Rel}_{\perp}^{\perp}$ , we obtain  $p \Vdash \neg R$  for some  $p \in \text{Proof}$ .
- With  $\text{Repr}_{\perp}$ , we obtain  $(2) \vdash \oplus(\langle p \rangle, \langle \neg R \rangle)$ .
- From (1) and consistency, we obtain  $\not\vdash R$ .
- With  $\text{Rel}_{\perp}^{\perp}$ , we obtain  $q \not\vdash R$  for all  $q \in \text{Proof}$ .
- With  $\text{Repr}_{\perp}$ ,  $\text{Clean}_{\perp}$  and Lemma 22, we obtain  $(3) \vdash \neg \oplus(n, \langle R \rangle)$  for all  $n \in \text{Num}$ .
- By  $\text{Ord}_2$ , we obtain a finite  $M \subseteq \text{Num}$  such that  $(4) \forall x. \vdash (\bigvee_{m \in M} x \equiv m) \vee \langle p \rangle \prec x$
- We prove  $(5) \vdash \forall x. \neg (\oplus(x, \langle R \rangle) \wedge \text{RosserTwist}(x, \langle R \rangle))$ . The proof is performed in our intuitionistic proof system, but we describe it informally: Fixing  $x$ , we assume  $\oplus(x, \langle R \rangle) \wedge \text{RosserTwist}(x, \langle R \rangle)$  and must reach a contradiction. We perform a case distinction according to (4):
  - If  $x$  equals some  $m \in M$ , then  $\oplus(m, \langle R \rangle)$ , which together with (3) leads to contradiction.
  - If  $\langle p \rangle \prec x$ , then from  $\text{RosserTwist}(x, \langle R \rangle)$  and  $\oplus(\langle R \rangle, \langle \neg R \rangle)$  (which holds thanks to  $\text{Repr}_{\neg}$ ), we obtain  $\neg \oplus(\langle p \rangle, \langle \neg R \rangle)$ , which contradicts (2).
- From (5), by (intuitionistic) logic we obtain  $\vdash \neg \exists x. (\oplus(x, \langle R \rangle) \wedge \text{RosserTwist}(x, \langle R \rangle))$ .
- Thanks to  $R$  being a Rosser formula, we obtain  $\vdash R$ .
- Together with (1), this contradicts consistency.  $\square$

Thus,  $\omega$ -consistency (assumption  $\text{OCon}$ ) has been indeed weakened to consistency (assumption  $\text{Con}$ ), but in exchange we needed to additionally assume a special formula  $\prec$  satisfying  $\text{Ord}_2$ . This represents a quite strong commitment to the arithmetical ordering.

Even worse, this fix on the assumptions needed to show the unprovability of the negated formula ( $\neg R$ ) complicates the proof of the unprovability of the *direct* formula ( $R$ ), which was trivial in Gödel's version (Prop. 7). Unlike in Gödel's version, we again need a cleanly representable "proof of" relation, representable negation, and well-behavedness of the order-like relation  $\prec$ :

**Prop 26.** Assume  $\text{Con}$ ,  $\text{Ord}_1$ ,  $\text{Rel}_{\perp}^{\perp}$ ,  $\text{Repr}_{\perp}$ ,  $\text{Repr}_{\neg}$  and  $\text{Clean}_{\perp}$ . Then  $\not\vdash R$  for all Rosser sentences  $R$ .

*Proof.* To prove  $\not\vdash R$ , we assume  $(1) \vdash R$  and aim to reach a contradiction.

- With  $\text{Rel}_{\perp}^{\perp}$ , we obtain  $p \Vdash R$  for some  $p \in \text{Proof}$ .
- With  $\text{Repr}_{\perp}$ , we obtain  $(2) \vdash \oplus(\langle p \rangle, \langle R \rangle)$ .
- With (1) and consistency, we obtain  $\not\vdash \neg R$ .
- With  $\text{Rel}_{\perp}^{\perp}$ , we obtain  $q \not\vdash \neg R$  for all  $q \in \text{Proof}$ .
- With  $\text{Repr}_{\perp}$ ,  $\text{Clean}_{\perp}$  and Lemma 22, we have  $\vdash \neg \oplus(n, \langle \neg R \rangle)$  for all  $n \in \text{Num}$ .
- With (2) and  $\text{Ord}_1$ , we obtain  $\vdash \forall y \prec \langle p \rangle. \neg \oplus(y, \langle \neg R \rangle)$ .

- Since, by  $\text{Repr}_{\neg}$ , the only  $z$  such that  $\ominus(z, \langle \neg R \rangle)$  is  $\langle \neg R \rangle$ , we obtain  $\vdash \text{RosserTwist}(\langle p \rangle, \langle R \rangle)$ .
- From this and (2), we obtain  $\vdash \exists x. \oplus(x, \langle R \rangle) \wedge \text{RosserTwist}(x, \langle R \rangle)$ .
- But since  $R$  is a Rosser sentence, from (1) we obtain  $\vdash \neg \exists x. \oplus(x, \langle R \rangle) \wedge \text{RosserTwist}(x, \langle R \rangle)$ .
- The last two facts contradict consistency.  $\square$

### A.11 Jeroslow's approach

Here is Jeroslow's version of Gödel's First:

**Theorem 27.** Theorems 9–12 still hold if we replace the  $\text{Repr}_S$  assumption with  $\text{CapSS}$  and  $\text{Repr}_{\mathcal{F}}$ .

*Proof.* Lemma 15 allows to construct and use Gödel and Rosser sentences just like in Sections 3.1 and 3.2.  $\square$

The following is claimed in the main paper:

**Lemma 28.** Assume  $\text{HBL}_1$ ,  $\text{Repr}_{\mathcal{F}}$ ,  $\text{CapN}$ . Then  $\vdash \text{jcon}$  implies  $\vdash \neg \oplus(\perp)$ .

*Proof.* Assume  $\vdash \text{jcon}$ .

- From this, we obtain (1)  $\vdash \mathbb{N}(\langle \perp \rangle, \langle \neg \perp \rangle) \rightarrow \neg (\oplus(\perp) \wedge \oplus(\neg \perp))$ .
- From  $\text{Repr}_{\mathcal{F}}$  and  $\text{CapN}$ , we obtain  $\vdash \mathbb{N}(\langle \perp \rangle, \langle \neg \perp \rangle)$ .
- With (1), we obtain (2)  $\vdash \neg (\oplus(\perp) \wedge \oplus(\neg \perp))$ .
- From  $\vdash \neg \perp$  and  $\text{HBL}_1$ , we obtain  $\vdash \oplus(\neg \perp)$ .
- With (2), we obtain  $\vdash \neg (\oplus(\perp))$ , as desired.  $\square$

## B Other proof sketches

### Lemma 3

*Proof.* (1) Assume  $\vdash \oplus(\varphi)$ .

- Hence  $\vdash \exists x. \oplus(x, \langle \varphi \rangle)$ .
- By (intuitionistic) logic, we obtain  $\vdash \neg \neg (\exists x. \oplus(x, \langle \varphi \rangle))$ .
- With  $\text{OCon}$ , we obtain  $n \in \text{Num}$  such that  $\nVdash \neg \oplus(n, \langle \varphi \rangle)$ .
- With  $\text{Clean}_{\Vdash}$ , we obtain  $p \in \text{Proof}$  such that  $n = \langle p \rangle$ . Hence  $\nVdash \neg \oplus(\langle p \rangle, \langle \varphi \rangle)$ .
- Since, by  $\text{Repr}_{\Vdash}$ , we have that  $p \Vdash \varphi$  implies  $\vdash \neg \oplus(\langle p \rangle, \langle \varphi \rangle)$ , we obtain  $p \Vdash \varphi$ .
- With  $\text{Rel}_{\vdash}^{\Vdash}$ , we obtain  $\vdash \varphi$ , as desired.

(2): Assume  $\vdash \neg \varphi(n)$  for all  $n \in \text{Num}$ .

- With Soundness, we obtain  $\models \neg \varphi(n)$  for all  $n \in \text{Num}$ .
- With  $\text{Syn}_{\models}(3)$ , we obtain  $\models \forall x. \neg \varphi(x)$ .
- With the intuitionistic logic built in  $\vdash$  and Soundness, we obtain  $\models \neg (\exists x. \varphi(x))$ .
- With  $\text{Syn}_{\models}(1,2)$ , we obtain  $\nVdash \neg \neg (\exists x. \varphi(x))$ .
- With Soundness, we obtain  $\nVdash \neg \neg (\exists x. \varphi(x))$ , as desired.

(3): Assume  $\models \oplus(\varphi)$ .

- Then  $\models \exists x. \oplus(x, \langle \varphi \rangle)$ .
- With  $\text{Syn}_{\models}(4)$ , we obtain  $n \in \text{Num}$  such that (i)  $\models \oplus(n, \langle \varphi \rangle)$ .

- With Soundness and  $\text{Syn}_{\models}(1,2)$ , we obtain  $\not\vdash \neg \oplus(n, \langle \varphi \rangle)$ .
- From here, the proof of  $\vdash \varphi$  proceeds just like at point (1): using  $\text{Rel}_{\vdash}^{\oplus}$ ,  $\text{Repr}_{\models}$  and  $\text{Clean}_{\models}$ .  $\square$

#### Lemma 4

*Proof.* To show  $\text{Rel}_{\vdash}^{\oplus}$  in this context (that is, for this particular definitions of Proof and relation  $\models$ ), we must show the equivalence between (i)  $\vdash \varphi$  and (ii) the existence of  $n \in \text{Num}$  such that  $\vdash P(n, \langle \varphi \rangle)$ .

First assume (i).

- With  $\text{HBL}_1$ , we obtain  $\vdash \oplus \langle \varphi \rangle$ .
- With  $\text{Rel}_{\oplus}^P$ , we obtain  $\vdash \exists x. P(x, \langle \varphi \rangle)$ .
- With Soundness, we obtain  $\models \exists x. P(x, \langle \varphi \rangle)$ .
- With  $\text{Syn}_{\models}(4)$ , we obtain  $n \in \text{Num}$  such that  $\models P(n, \langle \varphi \rangle)$ .
- With  $\text{Compl}_P$ , we obtain (ii), as desired.

Now assume (ii).

- By logic, we obtain  $\vdash \exists x. P(x, \langle \varphi \rangle)$ .
- With  $\text{Rel}_{\oplus}^P$ , we obtain  $\vdash \oplus \langle \varphi \rangle$ .
- With  $\text{HBL}_1^{\leftarrow}$ , we obtain (i), as desired.

Showing half of  $\text{Repr}_{\models}$  in this context is trivial, as it amounts to showing that “ $\vdash P(n, \langle \varphi \rangle)$  implies  $\vdash P(n, \langle \varphi \rangle)$ ”. For the other half, assume  $\not\vdash P(n, \langle \varphi \rangle)$ .

- With  $\text{Compl}_P$ , we obtain  $\not\models P(n, \langle \varphi \rangle)$ .
- With  $\text{Syn}_{\models}(5)$ , we obtain  $\models \neg P(n, \langle \varphi \rangle)$ .
- With  $\text{Compl}_{\neg P}$ , we obtain  $\vdash \neg P(n, \langle \varphi \rangle)$ , as desired.

Finally  $\text{Clean}_{\models}$  is trivial in this context, since the encoding of proofs is the identity.  $\square$

#### Prop 5

*Proof.* Assume  $\text{Repr}_S$ , where  $S$  is the “hard” self-substitution function. Let  $\chi \in \text{Fmla}_1$  be  $\exists y. \odot(x, y) \wedge \psi(y)$ . We take  $\varphi$  to be  $\chi \langle \chi \rangle$ .

- We must prove (1)  $\vdash \varphi \leftrightarrow \psi \langle \varphi \rangle$ .
- Using the definition of  $\chi$ , this further means (2)  $\vdash (\exists y. \odot(\langle \chi \rangle, y) \wedge \psi(y)) \leftrightarrow \psi \langle \varphi \rangle$ .
- By the fact that  $S$  is represented by  $\odot$  we obtain (provably in the formal system  $\vdash$ ) that  $\langle \varphi \rangle$  is the unique  $y$  for which  $\odot(\langle \chi \rangle, y)$  holds.
- From this last fact and (2), we obtain (1).

A similar argument works for soft self-substitution.  $\square$

#### Prop 6

*Proof.* Immediately from Prop. 5, taking  $\psi(x)$  to be  $\neg \oplus(x)$  and  $\neg (\exists y. \oplus(y, x) \wedge \text{RosserTwist}(y, x))$ , respectively.  $\square$



### Prop 7

*Proof.* Let  $G$  be a Gödel sentence. To prove  $\not\vdash G$ , we assume  $(1) \vdash G$  and aim to reach a contradiction.

- Since  $G$  is a Gödel sentence, we obtain  $\vdash \neg \oplus \langle G \rangle$ .
- From (1) and  $HBL_1$ , we obtain  $\vdash \oplus \langle G \rangle$ .
- The last two facts contradict the consistency assumption.  $\square$

### Theorem 9

*Proof.* (1): Immediate from Prop. 6.

(2):  $\not\vdash \neg G$  follows by applying Prop. 8 to the assumptions, so it remains to show  $\not\vdash G$ .

- From  $OCon$ , we obtain  $Con$ .
- Applying Lemma 24 to  $Rel_{\vdash}^{\perp}$  and  $Repr_{\Vdash}$ , we obtain  $HBL_1$ .
- Applying Prop. 7 to the last two facts, we obtain  $\not\vdash G$ , as desired.  $\square$

### Theorem 10

*Proof.* (1): Immediate from Prop. 6.

(2):  $\not\vdash R$  follows by applying Prop. 26 to the assumptions, and  $\not\vdash \neg R$  follows by applying Prop. 25 to the assumptions.  $\square$

### Theorem 11

*Proof.*

- Applying Lemma 3(2) to the assumptions, we obtain  $OCon$ .
- Applying Lemma 4 to the assumptions, we obtain  $Rel_{\vdash}^{\perp}$ ,  $Repr_{\Vdash}$  and  $Clean_{\Vdash}$ , where:
  - Proof and  $\Vdash$  are defined as in Lemma 4 and
  - $\Vdash$  is represented by  $P$ .
- Applying Theorem 9 to the last facts,  $OCon$  and the assumptions, we obtain that:
  - there exists a Gödel sentence and
  - letting  $G$  be a Gödel sentence  $\not\vdash G$  and  $\not\vdash \neg G$ .
- It remains to show  $\models G$ .
  - From  $\not\vdash G$ ,  $Rel_{\vdash}^{\perp}$  and  $Repr_{\Vdash}$ , we obtain  $\vdash \neg P(n, \langle G \rangle)$  for all  $n \in Num$ .
  - With Soundness, we obtain  $\models \neg P(n, \langle G \rangle)$  for all  $n \in Num$ .
  - With  $Syn_{\models}(3)$ , we obtain  $(*) \models \forall x. \neg P(x, \langle G \rangle)$ .
  - From  $Rel_{\vdash}^{\perp}$  and intuitionistic logic, we obtain  $\vdash (\forall x. \neg P(x, \langle G \rangle)) \rightarrow \neg \vdash \langle G \rangle$ .
  - With Soundness, we obtain  $\models (\forall x. \neg P(x, \langle G \rangle)) \rightarrow \neg \oplus \langle G \rangle$ .
  - With  $Syn_{\models}(2)$  and  $(*)$ , we obtain  $\models \neg \oplus \langle G \rangle$ .
  - With the definition of Gödel sentence and Soundness, we obtain  $\models G$ , as desired.  $\square$

### Theorem 12

*Proof.* (1): Immediate from Prop. 6.

(2): Let  $G$  be a Gödel sentence. We know that  $\not\vdash G$  from Prop. 7, so we are left to prove  $\not\vdash \neg G$ . To this end, we assume  $(*) \vdash \neg G$  and aim to reach a contradiction.

- Since  $G$  is a Gödel sentence, we obtain  $\vdash \neg \neg \oplus(G)$ .
- With classical logic, we obtain  $\vdash \oplus(G)$ .
- With  $\text{HBL}_1^{\Leftarrow}$ , we obtain  $\vdash G$ .
- With  $(*)$ , this contradicts the consistency assumption. □

### Theorem 13

*Proof.* Let  $G$  be a Gödel sentence, whose existence is ensured by Prop. 6 and Repr<sub>S</sub>; by Prop. 7, Con and  $\text{HBL}_1$ , we also have  $\not\vdash G$ , i.e. (by the Gödel sentence definition), (1)  $\not\vdash \neg \oplus(G)$ .

To prove (2)  $\not\vdash \neg \oplus(\perp)$ , it suffices to prove (3)  $\vdash \oplus(G) \rightarrow \oplus(\perp)$ . Indeed, the last would imply  $\vdash \neg \oplus(\perp) \rightarrow \neg \oplus(G)$ , which together with (1) would imply (2).

So it remains to prove (3).

- Since  $G$  is a Gödel sentence, we obtain  $\vdash G \rightarrow \neg \oplus(G)$ .
- With  $\text{HBL}_1$ , we obtain  $\vdash \oplus(G \rightarrow \neg \oplus(G))$ .
- With  $\text{HBL}_2$  we obtain  $\vdash \oplus(G) \rightarrow \oplus(\neg \oplus(G))$ .
- From  $\text{HBL}_3$ , we obtain  $\vdash \oplus(G) \rightarrow \oplus(\oplus(G))$ .
- From  $\text{HBL}_2$  we obtain  $\vdash \oplus(\oplus(G)) \wedge \oplus(\neg \oplus(G)) \rightarrow \oplus(\perp)$ .
- From the last three facts, we obtain (3), as desired. □

### Theorem 17

*Proof.* The only time when (1)  $\vdash \text{jcon}$  is used in the proof is via its consequence  $\vdash \neg (\oplus(\varphi) \wedge \oplus(\mathbb{N}(\varphi)))$ , which by Repr<sub>F</sub> and CapN would follow from  $(**) \vdash \neg (\oplus(\varphi) \wedge \oplus(\neg \varphi))$ . So it suffices to show that the last follows from  $\vdash \neg \oplus(\perp)$ ,  $\text{WHBL}_2$  and  $\text{HBL}_4$ :

- From  $\text{HBL}_4$ , we obtain  $\vdash \oplus(\varphi) \wedge \oplus(\neg \varphi) \rightarrow \oplus(\varphi \wedge \neg \varphi)$ .
- From  $\text{WHBL}_2$  and  $\vdash \varphi \wedge \neg \varphi \rightarrow \perp$ , we obtain  $\vdash \oplus(\varphi \wedge \neg \varphi) \rightarrow \oplus(\perp)$ .
- From the last two facts, we obtain  $\vdash \oplus(\varphi) \wedge \oplus(\neg \varphi) \rightarrow \oplus(\perp)$ .
- Hence  $\vdash \neg \oplus(\perp) \rightarrow \neg (\oplus(\varphi) \wedge \oplus(\neg \varphi))$ .
- With  $\vdash \neg \oplus(\perp)$ , we obtain  $(**)$ , as desired. □