



Truth, disjunction, and induction

Ali Enayat¹ · Fedor Pakhomov²

Received: 25 June 2018 / Accepted: 17 December 2018
© The Author(s) 2019

Abstract

By a well-known result of Kotlarski et al. (1981), first-order Peano arithmetic PA can be conservatively extended to the theory $CT^-[PA]$ of a truth predicate satisfying compositional axioms, i.e., axioms stating that the truth predicate is correct on atomic formulae and commutes with all the propositional connectives and quantifiers. This result motivates the general question of determining natural axioms concerning the truth predicate that can be added to $CT^-[PA]$ while maintaining conservativity over PA. Our main result shows that conservativity fails even for the extension of $CT^-[PA]$ obtained by the seemingly weak axiom of disjunctive correctness DC that asserts that the truth predicate commutes with disjunctions of arbitrary finite size. In particular, $CT^-[PA] + DC$ implies $Con(PA)$. Our main result states that the theory $CT^-[PA] + DC$ coincides with the theory $CT_0[PA]$ obtained by adding Δ_0 -induction in the language with the truth predicate. This result strengthens earlier work by Kotlarski (1986) and Cieśliński (2010). For our proof we develop a new general form of Visser's theorem on non-existence of infinite descending chains of truth definitions and prove it by reduction to (Löb's version of) Gödel's second incompleteness theorem, rather than by using the Visser–Yablo paradox, as in Visser's original proof (1989).

Keywords Axiomatic truth · Compositional theory of truth · Conservativity

Mathematics Subject Classification 03F30

The work of Fedor Pakhomov is supported by the Russian Science Foundation under Grant 16-11-10252 and performed at Steklov Mathematical Institute of Russian Academy of Sciences. F. Pakhomov was selected as one of the Young Russian Mathematics award winners, and he would like to thank its sponsors and jury; his work on this particular paper was funded from another source. Sections 2 and 4 of this paper were contributed by A. Enayat; Section 3 was contributed by F. Pakhomov; Sections 1 and 5 were contributed jointly by the authors. The authors are grateful to the anonymous referees for their meticulous and valuable feedback.

✉ Ali Enayat
ali.enayat@gu.se

Extended author information available on the last page of the article

1 Introduction

By a theorem of Krajewski, Kotlarski, and Lachlan [12], every countable recursively saturated model \mathcal{M} of PA (Peano Arithmetic) carries a ‘full satisfaction class’, i.e., there is a subset S of the universe M of \mathcal{M} that ‘decides’ the truth/falsity of each sentence of arithmetic in the sense of \mathcal{M} —even sentences of nonstandard length—while obeying the usual recursive clauses guiding the behavior of a Tarskian satisfaction predicate. This remarkable theorem implies that theory $\text{CT}^-[\text{PA}]$ (compositional truth over PA with induction only for the language \mathcal{L}_A of arithmetic) is *conservative* over PA, i.e., if an \mathcal{L}_A -sentence φ is provable in $\text{CT}^-[\text{PA}]$, then φ is already provable in PA. New proofs of this conservativity result were given by Visser and Enayat [5] using basic model theoretic ideas, and by Leigh [14] using proof theoretic tools; these new proofs make it clear that in the Krajewski–Kotlarski–Lachlan theorem the theory PA can be replaced by any ‘base’ theory that supports a modicum of coding machinery for handling elementary syntax.

On the other hand, it is well-known [9, Thm. 8.39 and Cor. 8.40] that the consistency of PA (and much more) is readily provable in the stronger theory $\text{CT}[\text{PA}]$, which is the result of strengthening $\text{CT}^-[\text{PA}]$ with the scheme of induction over natural numbers for all \mathcal{L}_{A+T} -formulae, where $\mathcal{L}_{A+T} := \mathcal{L}_A \cup \{T(x)\}$.¹ Indeed, it is straightforward to demonstrate the consistency of PA within the subsystem $\text{CT}_1[\text{PA}]$ of $\text{CT}[\text{PA}]$, where $\text{CT}_n[\text{PA}]$ is the subtheory of $\text{CT}[\text{PA}]$ with the scheme of induction over natural numbers limited to \mathcal{L}_{A+T} -formulae that are at most of complexity Σ_n [16, Thm. 2.8].

The discussion above leaves open whether $\text{CT}_0[\text{PA}]$ is conservative over PA. Kotlarski [11] established that $\text{CT}_0[\text{PA}]$ is a subtheory of $\text{CT}^-[\text{PA}] + \text{Ref}(\text{PA})$, where $\text{Ref}(\text{PA})$ is the \mathcal{L}_{A+T} -sentence stating that “every first order consequence of PA is true”. Recently Łełyk [15] demonstrated that the converse also holds, which immediately implies that $\text{CT}_0[\text{PA}]$ is not conservative over PA since $\text{Con}(\text{PA})$ is readily provable in $\text{CT}^-[\text{PA}] + \text{Ref}(\text{PA})$.² Kotlarski’s aforementioned theorem was refined by Cieśliński [3] who proved that $\text{CT}^-[\text{PA}] + \text{“T is closed under propositional proofs”}$ and $\text{CT}^-[\text{PA}] + \text{Ref}(\text{PA})$ axiomatize the same theory. The main result of this paper, in turn, refines Cieśliński’s result by demonstrating:

Theorem 1 $\text{CT}^-[\Delta_0 + \text{Exp}] + \text{DC}$ and $\text{CT}_0[\text{PA}]$ axiomatize the same first order theory.

In the above theorem, $\text{CT}^-[\Delta_0 + \text{Exp}]$ is the weakening of $\text{CT}^-[\text{PA}]$ obtained by replacing the ‘base theory’ PA with its fragment consisting of Robinson’s arithmetic Q, along with the scheme for Δ_0 -induction and the totality of the exponential function; and DC (disjunctive correctness) is the statement asserting that a disjunction of

¹ $\text{CT}[\text{PA}]$ dwarfs PA in arithmetical strength: By a classical theorem (discovered by a number of researchers, including Feferman and Takeuti, and explained in [9]) the arithmetical consequences of $\text{CT}[\text{PA}]$ are the same as the arithmetical consequences of ACA, the subsystem of second order arithmetic obtained by adding the full scheme of induction over natural numbers (in the language of second order arithmetic) to the well-known subsystem ACA_0 of second order arithmetic.

² This result was first claimed by Kotlarski [11], but his proof outline of $\text{Ref}(\text{PA})$ within $\text{CT}_0[\text{PA}]$ was found to contain a serious gap in 2011 by Heck and Visser; this gap cast doubt over the veracity of Kotlarski’s claim until the issue was resolved by Łełyk in his doctoral dissertation [15]. Łełyk’s work was preceded by the discovery of an elegant proof of the nonconservativity of $\text{CT}_0[\text{PA}]$ over PA by Wcisło and Łełyk [16].

arithmetical sentences of arbitrary finite length is true (in the sense of T) iff one of the disjuncts is true. Coupled with Łełyk’s aforementioned result [15], Theorem 1 shows that $CT^-[\Delta_0 + \text{Exp}] + DC$, $CT^-[\text{PA}] + DC$, and $CT_0[\text{PA}] + \text{Ref}(\text{PA})$ are axiomatizations of the same theory.

The plan of the paper is as follows: in Sect. 2 we review preliminary definitions and results, including more precise versions of those definitions and results mentioned in this introduction. In Sect. 3 we establish the veracity of the principle IC (Inductive Correctness, often referred to in the literature as “internal induction”) within $CT^-[\Delta_0 + \text{Exp}] + DC$. This is demonstrated by first establishing a new general form of Visser’s theorem [20] on nonexistence of infinite descending chains of truth definitions with the help of (Löb’s version of) Gödel’s second incompleteness theorem instead of the Visser–Yablo paradox. In Sect. 4 we show that $CT_0[\text{PA}]$ is a subtheory of $CT^-[\text{PA}] + DC + IC$; thus Sects. 3 and 4 together constitute the proof of the hard direction of Theorem 1 since it is routine to verify that both DC and IC are theorems of $CT_0[\text{PA}]$. We close the paper with some open problems in Sect. 5.

Historical note The concept of disjunctive correctness first appeared in the work of Krajewski [13, p.133], who called it “ \vee -completeness”; the current terminology was coined in a working paper of Enayat and Visser that was privately circulated in 2011, only a fragment [5] of which has been published so far. The working paper included the claim that $CT^-[\text{PA}] + DC$ is conservative over PA, but the proof outline presented in the paper was found in 2013 to contain a significant gap by Cieśliński and his (then) doctoral students Łełyk and Weisło. On the other hand, in 2012 Enayat found a proof of $CT_0[\text{PA}]$ within $CT^-[\Delta_0 + \text{Exp}] + DC + IC$; his proof was only privately circulated, and later was presented in the doctoral dissertation of Łełyk [15]. This proof forms the content of Sect. 4 of this paper. In light of these developments, and the well-known conservativity of $CT^-[\text{PA}] + IC$ over PA (see Theorem 2.3), the question of conservativity of $CT^-[\text{PA}] + DC$ over PA came to prominence amongst truth theory experts [4, p.226], and had been unsuccessfully attacked by a number of researchers since 2013, until Pakhomov established IC within $CT^-[\Delta_0 + \text{Exp}] + DC$ as in Sect. 3 of this paper, which, coupled with Enayat’s aforementioned result, yields Theorem 1 and exhibits the unexpected arithmetical strength of DC .

2 Preliminaries

Definition 2.1 (a) \mathcal{L}_A is the usual language of first order arithmetic $\{+, \cdot, S(x), <, 0\}$.

To simplify matters, we will assume that the logical constants of first order logic consist only of \neg (negation), \vee (disjunction), and \exists (existential quantification); in particular \forall (universal quantification) as well as \wedge (conjunction) and other propositional connectives are treated here as derived notions.

(b) Given a language $\mathcal{L} \supseteq \mathcal{L}_A$, an \mathcal{L} -formula φ is said to be a $\Delta_0(\mathcal{L})$ -formula if all the quantifiers of φ are bounded by \mathcal{L} -terms, i.e., they are of the form $\exists x \leq t$, or of the form $\forall x \leq t$, where t is a term of \mathcal{L} not involving x . Given a predicate $U(x)$, \mathcal{L}_{A+U} is the language $\mathcal{L}_A \cup \{U(x)\}$.

- (c) Given a language $\mathcal{L} \supseteq \mathcal{L}_A$, $\text{I}\Delta_0(\mathcal{L})$ is the scheme of induction over natural numbers for $\Delta_0(\mathcal{L})$ -formulae. We shall omit the reference to \mathcal{L} if $\mathcal{L} = \mathcal{L}_A$, e.g., a Δ_0 -formula is a $\Delta_0(\mathcal{L}_A)$ -formula; and we shall use $\text{I}\Delta_0(\text{U})$ to abbreviate $\text{I}\Delta_0(\mathcal{L}_{A+\text{U}})$.
- (d) $\text{I}\Delta_0 + \text{Exp}$ is the fragment of Peano arithmetic obtained by strengthening Robinson's arithmetic Q with $\text{I}\Delta_0$ and with the sentence Exp that expresses the totality of the exponential function $y = 2^x$. It is well-known that Exp can be written as $\forall x \exists y \text{Exp}(x, y)$, where $\text{Exp}(x, y)$ is a Δ_0 -predicate which, provably in $\text{I}\Delta_0$, satisfies the familiar algebraic laws governing the graph of the exponential function, cf. [8, Sec. V3(c)].
- (e) $\text{Sent}_A(x)$ is the \mathcal{L}_A -formula that expresses “ x is the Gödel-number of a formula of \mathcal{L}_A with no free variables”, and $\text{Form}_A^n(x)$ is the \mathcal{L}_A -formula that expresses “ x is the Gödel-number of a formula of \mathcal{L}_A with precisely n free variables”. We use Sent_A and Form_A^n to refer to the corresponding definable classes of Gödel-numbers of \mathcal{L}_A -formulae.
- (f) Given a (base) theory B whose language is \mathcal{L}_A and which extends $\text{I}\Delta_0 + \text{Exp}$, $\text{CT}^-[\text{B}]$ is the theory obtained by strengthening B by adding the sentences tarski_0 through tarski_4 described below, where we use the following conventions: τ_1 and τ_2 vary over Gödel-numbers of closed \mathcal{L}_A -terms, τ_i° is the value of the term coded by τ_i , φ and ψ range over Gödel-numbers of \mathcal{L}_A -sentences, v ranges over variables, $\gamma(v)$ ranges over Form_A^1 , and \underline{x} is the numeral representing the value of x .

$$\begin{aligned} \text{tarski}_0 &:= \forall x (\text{T}(x) \rightarrow \text{Sent}_A(x)). \\ \text{tarski}_1 &:= \forall \tau_1, \tau_2 (\text{T}(\tau_1 = \tau_2) \leftrightarrow \tau_1^\circ = \tau_2^\circ). \\ \text{tarski}_2 &:= \forall \varphi (\text{T}(\neg \varphi) \leftrightarrow \neg \text{T}(\varphi)). \\ \text{tarski}_3 &:= \forall \varphi, \psi (\text{T}(\varphi \vee \psi) \leftrightarrow (\text{T}(\varphi) \vee \text{T}(\psi))). \\ \text{tarski}_4 &:= \forall v \forall \gamma(v) (\text{T}(\exists v \gamma(v)) \leftrightarrow \exists x \text{T}(\gamma(\underline{x}))). \end{aligned}$$

- (g) $\text{CT}_0[\text{B}] := \text{CT}^-[\text{B}] \cup \text{I}\Delta_0(\text{T})$.
- (h) DC (disjunctive correctness) is the $\mathcal{L}_{A+\text{T}}$ -sentence asserting that T commutes with disjunctions of arbitrary length, i.e., DC asserts that for all numbers s and for all sequences $\langle \varphi_i : i < s \rangle$ from Sent_A , the following equivalence holds:

$$\text{T} \left(\bigvee_{i < s} \varphi_i \right) \leftrightarrow \exists i < s \text{T}(\varphi_i),$$

where for definiteness $\bigvee_{i < s} \varphi_i$ is defined³ by recursion: $\bigvee_{i < 0} \varphi_i := \varphi_0$ and $\bigvee_{i < t+1} \varphi_i := (\bigvee_{i < t} \varphi_i) \vee \varphi_t$.

- (i) We will employ the abbreviation $\bigwedge_{i < s} \varphi_i$ for $\neg \bigvee_{i < s} \neg \varphi_i$, and CC (conjunctive correctness) for the $\mathcal{L}_{A+\text{T}}$ -sentence

$$\text{T} \left(\bigwedge_{i < s} \varphi_i \right) \leftrightarrow \forall i < s \text{T}(\varphi_i).$$

³ One can also formulate disjunctive correctness in a stronger way by asserting that all disjunctions (and not just the ones that are grouped to the left) are well-behaved with respect to T . But the current frugal form, as shown by Theorem 1, ends up implying the seemingly stronger form since it is easy to show in CT_0 that the two forms are equivalent.

- Note that the commutativity of T with negation implies that DC and CC are equivalent.
- (j) IC (inductive correctness⁴) is the \mathcal{L}_{A+T} -sentence asserting that each \mathcal{L}_A -instance of induction over natural numbers is true, i.e., IC asserts that for all unary \mathcal{L}_A -formulae $\psi = \psi(x)$, $T(\ulcorner \text{Ind}_{\psi} \urcorner)$ holds, where Ind_{ψ} is the following \mathcal{L}_A -sentence that asserts that ψ is inductive:

$$\psi(0) \rightarrow (\forall x (\psi(x) \rightarrow \psi(x + 1))) \rightarrow \forall x \psi(x).$$

The $B = \text{PA}$ case of Theorem 2.2 below, and its elaboration Theorem 2.3, were first established in the work of Krajewski et al. [12] for $B = \text{PA}$, where PA is formulated in a relational language, and ‘domain constants’ play the role of numerals. As mentioned in the introduction to this paper, their result was couched in model theoretic terms involving the notions of recursive saturation and satisfaction classes, but it is well-known that their formulation is equivalent to appropriately formulated conservativity assertions (the key ingredients of this equivalence are the following facts: Every consistent theory in a countable language has a recursively saturated model, and countable recursively saturated models are resplendent). Later Kaye [10] developed the theory of satisfaction classes over models of PA in languages incorporating function symbols; his work was extended by Engström [6] to truth classes over models of PA in functional languages.⁵ More recently, newer and more informative proofs of Theorems 2.2 and 2.3 have been found in the joint work of Visser and Enayat [5] (with base theories that support a modicum of coding, and which are formulated in a relational language), and by Leigh [14] (for functional base theories that support a modicum of coding). As verified by Cieśliński [4, Ch. 7], the Visser-Enayat model theoretic methodology can be extended so as to accommodate functional languages.

Theorem 2.2 $CT^-[B]$ is conservative over B for every arithmetical base theory B extending $\Delta_0 + \text{Exp}$.

Theorem 2.3 $CT^-[PA] + IC$ is conservative over PA.

The direction (a) \Rightarrow (b) of Theorem 2.4 below is due to Kotlarski [11]; the other direction is due to Łełyk [15].

Theorem 2.4 (Kotlarski–Łełyk) *The following theories are deductively equivalent:*

- (a) $CT^-[PA] + \text{Ref}(PA)$.
- (b) $CT_0[PA]$.

The direction (a) \Rightarrow (b) of Theorem 2.5 below is due to Cieśliński [3], who refined Kotlarski’s proof of the direction (a) \Rightarrow (b) of Theorem 2.4; the other direction involves a routine induction.

Theorem 2.5 (Cieśliński) *The following theories are deductively equivalent:*

- (a) $CT^-[PA] + \text{“}T \text{ is closed under propositional proofs”}$.
- (b) $CT_0[PA]$.

⁴ This condition has been referred to as Int (internal induction) in the literature.

⁵ The subtle distinction between satisfaction classes and truth classes, and their close relationship, is explained in [5] and [4, Ch. 7].

3 Disjunctive correctness implies inductive correctness

Definition 3.1 ITB (iterated truth biconditionals) is a theory formulated in two-sorted first order logic. The first sort $x, y, z \dots$ of ITB is for the ‘natural numbers’. The second sort $\alpha, \beta, \gamma, \dots$ is for the *indices* of truth definition. The language \mathcal{L}_{ITB} of ITB is obtained by augmenting the language \mathcal{L}_A of arithmetic with two binary predicates: $\alpha < \beta$ and $\text{T}(\alpha, x)$, but we shall write $\text{T}(\alpha, x)$ as $\text{T}_\alpha(x)$ to display the indexicality of α . The axioms of ITB come in three groups. The first group consists of the axioms of Q (Robinson arithmetic); the second group consists of a single axiom asserting that $<$ is a transitive relation; and the third group consists of the following biconditionals B_φ :

$$\text{B}_\varphi := \forall \alpha (\text{T}_\alpha(\ulcorner \varphi \urcorner) \leftrightarrow \varphi^{<\alpha}),$$

where φ ranges over all \mathcal{L}_{ITB} -sentences, and for each index variable α , $\varphi^{<\alpha}$ denotes the relativization of φ to the cone of indices below α , i.e. the formula obtained by replacing all the quantifiers of the form $\forall \beta (\exists \beta)$ with $\forall \beta < \alpha (\exists \beta < \alpha)$, and if there is a bounded instance of α we make the appropriate renaming. Clearly $\varphi^{<\alpha} = \varphi$ if φ is a purely arithmetical formula.

- Note that we take the theory ITB over the variant of many-sorted logic that allows domains of some sorts to be empty.
- Although we haven’t required $<$ to be irreflexive, we treat it as an irreflexive relation; in particular we define a minimal element α to be an element such that $\forall \beta \neg(\beta < \alpha)$. The reason is that Theorem 3.2 below implies that ITB proves the irreflexivity of $<$. Alternatively one could show that ITB proves irreflexivity of $<$ by observing that the existence of a model of ITB with a reflexive point contradicts Tarski’s undefinability of truth theorem.
- We will use the following convention to lighten the notation: The notation $\ulcorner \varphi \urcorner$ for the Gödel number of a formula φ will be generally used, but the corner-notation will be omitted when φ appears inside of a truth predicate T , or inside an indexed version of T .

The proof of the following theorem was inspired by the recent James Walsh proof [18] of nonexistence of infinite recursive provably descending chains of sentences with respect to $<_{\text{Con}}$ -order. We note that Theorem 3.2 is similar to a result by Flumini and Sato [7], which states that the second order principle asserting the existence of iteration of Π_1^0 -comprehension over a preorder $<$ implies that $<$ is well-founded.

Theorem 3.2 *ITB + $\exists \alpha (\alpha = \alpha)$ proves the existence of a $<$ -minimal element. Equivalently, the following theory DTB (descending truth biconditionals) is inconsistent:*

$$\text{DTB} := \text{ITB} + \forall \alpha \exists \beta (\beta < \alpha) + \exists \alpha (\alpha = \alpha).$$

Proof We prove the inconsistency of DTB by Löb’s version of Gödel’s second incompleteness theorem⁶: We exhibit a formula $\theta(x)$ that satisfies the HBL (Hilbert–Bernays–Löb) conditions for a provability predicate over the theory DTB. This enables

⁶ Löb’s paper [17], in which the venerable ‘Löb’s Theorem’ was proved, is responsible for the now common standard textbook framework for the presentation of ‘abstract’ form of Gödel’s second incompleteness

us to justify the inconsistency of DTB by showing that DTB proves the “consistency” sentence $\neg\theta(0 = 1)$.

Consider the formula $\theta(x)$:

$$\theta(x) := \forall\alpha(\mathsf{T}_\alpha(x)).$$

We will verify that $\theta(x)$ satisfies the HBL conditions listed below provably in DTB; in what follows φ and ψ range over all sentences of the language of DTB:

- HBL-1. $\text{DTB} \vdash \varphi \implies \text{DTB} \vdash \theta(\ulcorner\varphi\urcorner)$.
- HBL-2. $\text{DTB} \vdash \theta(\ulcorner\varphi \rightarrow \psi\urcorner) \rightarrow (\theta(\ulcorner\varphi\urcorner) \rightarrow \theta(\ulcorner\psi\urcorner))$.
- HBL-3. $\text{DTB} \vdash \theta(\ulcorner\varphi\urcorner) \rightarrow \theta(\ulcorner\theta(\ulcorner\varphi\urcorner)\urcorner)$.

Since we have biconditionals, in order to prove HBL-1 it is enough to show that for each sentence φ , if $\text{DTB} \vdash \varphi$ then $\text{DTB} \vdash \forall\alpha \varphi^{<\alpha}$. The latter is the case since for any model \mathcal{M} of DTB and index a in \mathcal{M} , the theory DTB holds in the model $\mathcal{M}^{<a}$ that is the restriction of \mathcal{M} to all indices $< a$.

For a given φ and ψ , HBL-2 follows directly from the biconditional axioms $\mathsf{B}_{\varphi \rightarrow \psi}$, B_φ , and B_ψ of ITB.

Finally, HBL-3 holds since:

$$\text{ITB} \vdash \theta(\ulcorner\theta(\ulcorner\varphi\urcorner)\urcorner) \iff (\forall\alpha\forall\beta((\beta < \alpha) \rightarrow \mathsf{T}_\beta(\varphi))).$$

On the other hand, the formula $\forall\alpha \neg\mathsf{T}_\alpha(\ulcorner 0 = 1\urcorner)$ is provable in ITB, hence the formula $\neg\theta(\ulcorner 0 = 1\urcorner)$ is provable in ITB, and therefore in $\text{ITB} + \exists\alpha(\alpha = \alpha)$. So by Löb’s version of Gödel’s second incompleteness theorem, DTB is inconsistent. \square

Lemma 3.3 $\text{CT}^-[\Delta_0 + \text{Exp}] + \text{DC}$ proves IC.

Proof By Theorem 3.2 we can fix an inconsistent finite subtheory DTB^- of DTB. Suppose DTB^- contains only biconditionals for the formulae $\varphi_0, \dots, \varphi_{k-1}$. We will use ITB^- to denote the subtheory of ITB whose only biconditional axioms are $\{\mathsf{B}_{\varphi_i} : i < k\}$.

For the rest of the proof we will reason in $\text{CT}^-[\Delta_0 + \text{Exp}] + \text{DC}$. In order to prove IC we assume for a contradiction that some arithmetical $\psi(x)$ is not inductive in the sense of T , i.e., we have:

$$\neg\mathsf{T}(\psi(0) \rightarrow (\forall x(\psi(x) \rightarrow \psi(x + 1)) \rightarrow \forall x \psi(x))).$$

Within $\text{CT}^-[\Delta_0 + \text{Exp}] + \text{DC}$ we use induction on n to define translations ι_n from the language of ITB to the language of first-order arithmetic such that from the point of view of T all of them will be interpretations of ITB^- , i.e., we will have $\mathsf{T}(\iota_n(\varphi))$, for all axioms φ of ITB^- . We will arrive at a contradiction by showing that $\neg\mathsf{T}(\psi(\underline{n}))$

Footnote 6 continued

theorem: If T is a consistent theory extending Q that supports a unary predicate $\theta(x)$ satisfying conditions HBL-1, HBL-2, and HBL-3, then T doesn’t prove $\theta(\ulcorner 0 = 1\urcorner)$, i.e., the consistency sentence corresponding to θ (the intended meaning of $\theta(x)$ is “the sentence with Gödel number x is provable in T ”). See, e.g., [1, Ch. 18], for the presentation of such a general form of Gödel’s second incompleteness theorem.

implies that t_n is an interpretation of DTB^- ; thus it will be necessary to consider n that are non-standard from external point of view.

Note that each translation t_n consists of finitely many formulae, giving the interpretation of the domains and all symbols of the signature of ITB , and thus could be easily represented by a number. The interpretation of arithmetic in each t_n is the identity interpretation, but the domain of indices of truth definition t_n is given by the following formula $D^{(n)}(x)$:

$$x < \underline{n} \wedge \neg\psi(x).$$

For all n the relation $<$ is interpreted by $<$. The formula $T_\alpha(x)$ is interpreted by the formula $\text{IT}^{(n)}(y, x)$, where y corresponds to α , and x corresponds to itself:

$$\bigwedge_{i < k} (x = \varphi_i \rightarrow \bigwedge_{m < n} ((y = \underline{m} \wedge \neg\psi(\underline{m})) \rightarrow t_m(\varphi_i))),$$

where $t_m(\varphi_i)$ is the t_m -translation of the sentence φ_i . It is easy to see that this definition could be carried out in $\text{I}\Delta_0 + \text{Exp}$.

Let us now prove that the translations given by t_n are indeed the desired interpretations inside T , i.e., we need to prove that for all n and axioms A of ITB^- we have $T(t_n(A))$. Clearly it is the case for all the axioms of Q and the axioms of partial order for $<$. Now let us show that for any $s < k$:

$$T(t_n(\forall\alpha(T_\alpha(\varphi_s) \leftrightarrow \varphi_s^{<\alpha}))). \tag{\blacktriangle}$$

By compositional axioms, we just need to show that for all u such that $u < n$ and $T(\neg\psi(\underline{u}))$ we have:

$$T\left(\bigwedge_{i < k} \left(\varphi_s = \varphi_i \rightarrow \bigwedge_{m < n} ((\underline{u} = \underline{m} \wedge \neg\psi(\underline{m})) \rightarrow t_m(\varphi_i))\right)\right) \leftrightarrow T(t_n(\varphi_s^{<\underline{u}})).$$

Now by compositional axioms and DC (in the form of CC, as explained in part (i) of Definition 2.1) our task can be reduced to proving the equivalence:

$$T(t_u(\varphi_s)) \leftrightarrow T(t_n(\varphi_s^{<\underline{u}})).$$

In order to prove this we will show by induction on subformulae θ of φ_s that for the universal closure $\bar{\theta}$ of θ :

$$T(t_u(\bar{\theta})) \leftrightarrow T(t_n(\bar{\theta}^{<\underline{u}})).$$

Note that since φ_s is a fixed formula with finitely many subformulae, actually this external induction will be formalizable in $\text{CT}^-[\text{I}\Delta_0 + \text{Exp}] + \text{DC}$ despite the fact that it lacks the induction axiom for the appropriate class of formulae. The only non-trivial case here is the case when θ is $T_\alpha(x)$:

$$T(t_u(\forall\alpha\forall x T_\alpha(x))) \leftrightarrow T(t_n(\forall\alpha < \underline{u} \forall x(T_\alpha(x)))).$$

Hence we just need to show that for all $p < u$ such that $\text{T}(\neg\psi(\underline{p}))$, and for all o , the following pair of formulae (whose only formal difference is in the bound for indices of the second conjunction) are equivalent:

$$\begin{aligned} & \text{T} \left(\bigwedge_{i < k} \left(\underline{o} = \varphi_i \rightarrow \bigwedge_{m < u} \left((\underline{p} = \underline{m} \wedge \neg\psi(\underline{m})) \rightarrow \iota_m(\varphi_i) \right) \right) \right), \\ & \text{T} \left(\bigwedge_{i < k} \left(\underline{o} = \varphi_i \rightarrow \bigwedge_{m < n} \left((\underline{p} = \underline{m} \wedge \neg\psi(\underline{m})) \rightarrow \iota_m(\varphi_i) \right) \right) \right). \end{aligned}$$

But since $p < u < n$, we trivially use DC (in the form of CC) to show that the formulae are indeed equivalent. Thus we conclude that (\blacktriangle) holds.

Choose some n such that $\text{T}(\neg\psi(\underline{n}))$; this is possible since we assumed that induction fails for $\psi(x)$ in the sense of T . It is easy to see that ι_n actually is an interpretation of DTB^- inside T . We externally fix some proof of inconsistency from axioms of DTB^- and follow it inside T to derive a contradiction, thereby completing the proof of IC. \square

Corollary 3.4 $\text{CT}^-[\text{I}\Delta_0 + \text{Exp}] + \text{DC}$ proves PA, and therefore $\text{CT}^-[\text{I}\Delta_0 + \text{Exp}] + \text{DC}$ and $\text{CT}^-[\text{PA}] + \text{DC}$ axiomatize the same theory.

Proof This is an immediate consequence of Lemma 3.3, and the provability of Tarski bi-conditionals in $\text{CT}^-[\text{I}\Delta_0 + \text{Exp}]$. \square

Remark 3.5 Note that Theorem 3.2 could be regarded as a strengthening of Tarski’s undefinability of truth theorem: Tarski’s theorem essentially states that there could be no hierarchy of truth definitions whose set of indices contains a reflexive point. To the best of the authors’ knowledge there is no known proof of Tarski’s theorem that avoids the use of any kind of diagonalization constructions.⁷ In a preprint of this paper we raised a question of whether Lemma 3.3, which we proved using Theorem 3.2, could be proved more directly without the use of diagonalization. Later we noticed that it is possible to replace the use of Theorem 3.2 in the proof of Lemma 3.3 with a result by Flumini and Sato [7, Thm. 1], which has a rather simple proof that in our opinion could be regarded as diagonalization-free. Since this other proof is of methodological interest, we will provide its general outline.

The results of Flumini and Sato in [7] are fairly general and are applicable to various second-order systems. To keep our presentation compact we will just formulate a direct corollary of [7, Thm. 1] that will be relevant to us. The only axioms of our base system of second-order arithmetic B will be those of $\text{I}\Delta_0 + \text{Exp}$. We denote as Ind the usual second-order induction principle:

$$\forall X (0 \in X \wedge \forall x (x \in X \rightarrow x + 1 \in X) \rightarrow \forall x x \in X).$$

For a set X and number a we denote as $(X)_a$ the set $\{b \mid \langle a, b \rangle \in X\}$ and for binary relation $<$ we denote as $X^{<a}$ the set $\{\langle a', b \rangle \mid \langle a', b \rangle \in X \text{ and } a' < a\}$. For a second

⁷ We refer to the introduction of Visser’s paper [21] for a discussion of the role of diagonalization in the proofs of Gödel’s second incompleteness theorem and Tarski’s truth undefinability theorem.

order formula $\varphi(X, y)$ and binary relation $<$ we could naturally write the formula $\text{Hier}(\varphi, <, H)$ that expresses that H is a hierarchy along $<$ produced by φ , i.e., H is such that $\forall x ((H)_x = \{y \mid \varphi(H^{<x}, y)\})$. The usual arithmetic transfinite recursion principle states that hierarchies exist for any well-ordering $<$ and formula φ without second-order quantifiers. The corollary of [7, Thm. 1] that will be relevant to us is that over \mathbf{B} there is an arithmetic (moreover Π_1^0) formula $\varphi(X, y)$ (with additional variables) such that the axiom Ind is implied by the existence of φ -hierarchies along arbitrarily large proper initial segments of natural numbers. Formally this principle is $\forall z \exists H \text{ Hier}(\varphi, < \upharpoonright z, H)$, where $< \upharpoonright z$ is the restriction of the usual order $<$ on naturals to the numbers below z .

Now consider the ω -interpretation of the language of second-order arithmetic in $\text{CT}^-[\text{I}\Delta_0 + \text{Exp}]$, where the range of sets consists of all the sets $\{n \mid \text{T}(\varphi(n))\}$. It is easy to see that IC is equivalent to the validity of Ind in this interpretation.⁸ So in view of the above in order to show that DC implies Ind it will be enough to prove that DC implies that the universal closure of $\forall z \exists H \text{ Hier}(\varphi, < \upharpoonright z, H)$ holds in this interpretation. To achieve the latter we could just directly construct the formula defining the relevant hierarchy and then use DC to verify that it indeed has the desired property.

4 Disjunctive correctness + inductive correctness implies $\Delta_0(\text{T})$ -induction

In this section we shall prove that $\text{I}\Delta_0(\text{T})$ is provable in $\text{CT}^-[\text{PA}] + \text{DC} + \text{IC}$, which, coupled with Lemma 3.3 completes the proof of the nontrivial direction of Theorem 1. We begin with a key definition:

Definition 4.1 In what follows \in_{Ack} is “Ackermann’s epsilon”, i.e., $x \in_{\text{Ack}} y$ is the arithmetical formula that expresses “the x -th bit of the binary expansion of y is a 1”.

- (a) For a unary predicate $\text{U}(x)$, the $\mathcal{L}_{\text{A+U}}$ sentence PC_{U} (read as “U is piece-wise coded”) is the following sentence:

$$\forall u \exists y \forall x [(U(x) \wedge x < u) \leftrightarrow x \in_{\text{Ack}} y].$$

- (b) More generally, given an n -ary $\mathcal{L}_{\text{A+U}}$ -formula $\varphi(\text{U}, x_0, \dots, x_{n-1})$, PC_{φ} is the following $\mathcal{L}_{\text{A+U}}$ -sentence:

$$\forall u \exists y \forall x_0, \dots, \forall x_{n-1} [(\varphi(\text{U}, x_0, \dots, x_{n-1}) \wedge (x_0 < u \wedge \dots \wedge x_{n-1} < u) \leftrightarrow \langle x_i : i < n \rangle \in_{\text{Ack}} y],$$

where $\langle x_i : i < n \rangle$ is a canonical code for the ordered n -tuple (x_0, \dots, x_{n-1}) .

The following lemma shows that over $\text{I}\Delta_0 + \text{Exp}$ the scheme $\text{I}\Delta_0(\text{U})$ is equivalent to the single sentence ‘U is piecewise coded’. The lemma is folklore; we present the proof for the sake of completeness.

⁸ See Remark 4.3.1 for a sharper formulation of this equivalence.

Lemma 4.2 *The following are equivalent over $\text{I}\Delta_0 + \text{Exp}$:*

- (a) $\text{I}\Delta_0(\text{U})$.
- (b) PC_{U} .

Proof We will reason in $\text{I}\Delta_0 + \text{Exp}$. Recall that both $x = 2^y$ and $x \in_{\text{Ack}} y$ have Δ_0 -definitions within $\text{I}\Delta_0$ [8, Ch. V].

(a \rightarrow b): Assume $\text{I}\Delta_0(\text{U})$. Given u , let w be the Ackermann-code for the set of predecessors of u (i.e., $\forall x (x < u \leftrightarrow x \in_{\text{Ack}} w)$). Clearly $w = \sum_{i < u} 2^i = 2^u - 1$, and w is an upper bound for any w' that codes a subset of the predecessors of u . Let $\delta(u, w)$ be the Δ_0 -formula below:

$$(w = 2^u) \rightarrow [\exists y < w \forall x < u ((\text{U}(x) \wedge x < u) \leftrightarrow x \in_{\text{Ack}} y)].$$

A simple induction on u (where w is treated as parameter) using $\text{I}\Delta_0(\text{U})$ shows that $\forall u \forall w \delta(u, w)$ holds, which completes the proof that PC_{U} holds.

(b \rightarrow a): A straightforward induction on the complexity of $\Delta_0(\text{U})$ -formulae shows that:

(*) If U is piecewise coded and $\delta(\text{U}, x_0, \dots, x_{n-1})$ is a $\Delta_0(\text{U})$ -formula, then PC_{δ} holds.

The base case of the induction is clearly equivalent to the assumption that U is piecewise coded. What allows the inductive steps to be smoothly carried out is that, provably in $\text{I}\Pi_0 + \text{Exp}$, \in_{Ack} obeys many familiar axioms of set theory, as verified in [8, Ch. I, Thm. 1.39]. As an example, in the existential case of the induction, we suppose that $\delta(\text{U}, x_0, \dots, x_{n-1})$ is a $\Delta_0(\text{U})$ -formula such that PC_{δ} holds, and then establish $\text{PC}_{\delta'}$, where $\delta' = \exists x_0 < t(x_0, \dots, x_{n-1}) \delta(\text{U}, x_0, \dots, x_{n-1})$ for some term t . To do so, let us fix any number u and demonstrate that there is an \in_{Ack} -set s' , where:

$$s' = \{\{x_i : 0 < i < n\} \mid x_1, \dots, x_{n-1} < u \text{ and } \delta'(U, x_1, \dots, x_{n-1})\}.$$

Let v be a number such that the value $t(x_0, \dots, x_{n-1}) \leq v$ for $x_0, \dots, x_{n-1} < u$. By PC_{δ} we have the following \in_{Ack} -set s :

$$s = \{\{x_i : i < n\} \mid x_0, \dots, x_{n-1} < v \text{ and } \delta(U, x_0, \dots, x_{n-1})\}.$$

Using Δ_0 -Separation we construct the set s' as:

$$\{\{x_i : 0 < i < n\} \mid x_1, \dots, x_{n-1} < u \text{ and } \exists x_0 < t(x_0, \dots, x_{n-1}) \langle x_i : i < n \rangle \in_{\text{Ack}} s\}.$$

With (*) at our disposal, we could trivially deduce the least number principle for $\Delta_0(\text{U})$ -formula, which is of course equivalent to $\text{I}\Delta_0(\text{U})$. □

- In the next lemma and its proof, $\text{Code}(c, \varphi, u)$ denotes the ternary $\mathcal{L}_{\text{A}+\text{T}}$ -formula $\forall x [(x < u \wedge \text{T}(\varphi(x))) \leftrightarrow x \in_{\text{Ack}} c]$, and $\text{PC}(\varphi)$ denotes the $\mathcal{L}_{\text{A}+\text{T}}$ -formula $\forall u \exists c \text{Code}(c, \varphi, u)$.

Lemma 4.3 $CT^- [I\Delta_0 + \text{Exp}] + IC$ proves $\forall\varphi (\text{Form}_A^1(\varphi) \rightarrow PC(\varphi))$.

Proof We reason in $CT^- [I\Delta_0 + \text{Exp}] + IC$. Given $\varphi(x)$ in Form_A^1 , we need to show:

(1) $\forall u \exists c \text{Code}(c, \varphi, u)$.

By the compositional properties of T , (1) is equivalent to:

(2) $T(\forall u \psi(u))$, where $\psi(u) := \exists c(\forall x < u \varphi(x) \leftrightarrow x \in_{\text{Ack}} c)$.

On the other hand, $\forall u \psi(u)$ is the conclusion of the formula Ind_ψ (asserting the inductive property of ψ) given by IC. So (2) follows directly from IC and the easily verified facts $T(\psi(\underline{0}))$ and $T(\forall u (\psi(u) \rightarrow \psi(u + 1)))$. \square

Remark 4.3.1 Lemma 4.3 can be readily strengthened to a more general result whose proof we leave to the reader: $CT^- [I\Delta_0 + \text{Exp}] + IC$ verifies ACA_0 for the ω -interpretation of the language of second-order arithmetic, where the range of sets consists of all the sets $\{n \mid T(\varphi(\underline{n}))\}$, where φ is any first-order formula. Moreover, it is a theorem of $CT^- [I\Delta_0 + \text{Exp}]$ that IC is equivalent to the veracity of ACA_0 within this interpretation.

Lemma 4.4 $CT^- [I\Delta_0 + \text{Exp}] + DC + IC \vdash I\Delta_0(T)$.

Proof Reason in $CT^- [I\Delta_0 + \text{Exp}] + DC + IC$. By Lemma 4.2, it suffices to show that T is piecewise coded. Let $\langle \varphi_i : i < u \rangle$ be the sequence of arithmetical sentences such that φ_i is the sentence with Gödel-number i if there is such a sentence, and otherwise φ_i is the sentence $\underline{0} = \underline{1}$. We wish to show that $\{i < u : T(\varphi_i)\}$ is coded. Towards this goal, consider the unary formula $\theta(x) \in \text{Form}_A$ given by:

$$\theta(x) := \bigvee_{i < u} ((x = \underline{i}) \wedge \varphi_i).$$

Claim (*) $\forall i < u [T(\varphi_i) \leftrightarrow T(\theta(\underline{i}))]$.

(\rightarrow) Suppose $T(\varphi_i)$ for some $i < u$. Then $T((\underline{i} = \underline{i}) \wedge \varphi_i)$, and hence by DC we have $T(\theta(\underline{i}))$.

(\leftarrow) Suppose $T(\theta(\underline{i}))$ for some $i < u$. Then by DC, there is some $j < u$ such that $T((\underline{i} = \underline{j}) \wedge \varphi_j)$. So $T(\varphi_i)$ holds since T commutes with conjunction and $T(\underline{i} = \underline{j})$ holds iff $i = j$.

By coupling Claim (*) together with Lemma 4.3, we can conclude that $\{i < u : T(\theta(\underline{i}))\}$ is coded. \square

5 Closing remarks and open questions

Question 5.1 *Is the generalization of Theorem 1 in which CT^- is weakened to CS^- (where S stands for satisfaction) true?*

- The notion $CS^- [B]$ is defined in [5] for base theories B formulated in relational languages, using the notation B^{FS} (FS for “full satisfaction”); and in [4, Ch. 7] for functional languages. We expect that an examination of the proofs in Sects. 3 and 4 would show that this question has a positive answer.

Question 5.2 *Is IC provable in $CT^- [S_2^1] + DC$?*

- In the above, S_2^1 is Buss’s well-known arithmetical theory whose provable recursive functions are precisely the functions computable in polynomial time, as in [2]. For the above question to make sense, part (f) of Definition 1.1 should be adjusted so as to accommodate the fact that the language of S_2^1 extends \mathcal{L}_A . In the proof of Lemma 3.3, most likely it is possible to use some tricks with effective formulae (see [19, Section 3]) in order to modify the definition of ι_n in such a way that their sizes will be polynomial. But in order for the construction to work we will also need to ensure that DC is still enough to show that ι_n are indeed interpretations inside the truth predicate.

Question 5.3 *Let DC_{Elim} be the ‘half’ of DC that asserts: if a finite disjunction is true, then at least one of the disjuncts is also true. Is $CT^- [PA + DC_{Elim}]$ conservative over PA?*

- DC can be written as the conjunction of two implications DC_{Elim} and DC_{Intro} , where DC_{Intro} is the converse of DC_{Elim} (i.e., DC_{Intro} asserts: a disjunction is true, whenever at least one of its disjuncts is true). Recent joint work of Wcisło, Łełyk and Enayat (to appear) shows that DC_{Intro} can be conservatively added to $CT^- [PA] + IC + \{\forall x (\text{True}_{\Sigma_n}(x) \rightarrow T(x)) : n < \omega\}$.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

1. Boolos, G.S., Burgess, J.P., Jeffrey, R.C.: Computability and Logic, 5th edn. Cambridge University Press, Cambridge (2007)
2. Buss, S.: Proof theory of arithmetic. In: Buss, S. (ed.) Handbook of Proof Theory, pp. 79–147. Elsevier, Amsterdam (1998)
3. Cieśliński, C.: Deflationary truth and pathologies. J. Philos. Logic **39**, 325–337 (2010)
4. Cieśliński, C.: The Epistemic Lightness of Truth. Deflationism and its Logic. Cambridge University Press, Cambridge (2017)
5. Enayat, A., Visser, A.: New constructions of full satisfaction classes. In: Achourioti, T., Galinon, H., Fujimoto, K., Martínez-Fernández, J. (eds.) Unifying the Philosophy of Truth, pp. 321–325. Springer, Berlin (2015)
6. Engström, F.: Satisfaction classes in nonstandard models of first-order arithmetic (2002). [arXiv:math/0209408](https://arxiv.org/abs/math/0209408)
7. Flumini, D., Sato, K.: From hierarchies to well-foundedness. Arch. Math. Log. **54**, 855–863 (2014)
8. Hájek, P., Pudlák, P.: Metamathematics of First-Order Arithmetic. Springer, Berlin (1993)
9. Halbach, V.: Axiomatic Theories of Truth, 2nd edn. Cambridge University Press, Cambridge (2015)
10. Kaye, R.: Models of Peano Arithmetic. Oxford University Press, Oxford (1991)
11. Kotlarski, H.: Bounded induction and satisfaction classes. Zeitschrift für mathematische Logik und Grundlagen der Mathematik **32**, 531–544 (1986)
12. Kotlarski, H., Krajewski, S., Lachlan, A.: Construction of satisfaction classes for nonstandard models. Can. Math. Bull. **24**, 283–293 (1981)

13. Krajewski, S.: Nonstandard satisfaction classes. In: Marek, W., et al. *Set Theory and Hierarchy Theory: A Memorial Tribute to Andrzej Mostowski*. Lecture Notes in Mathematics, vol. 537, pp. 121–144. Springer, Berlin (1976)
14. Leigh, G.: Conservativity for theories of compositional truth via cut elimination. *J. Symb. Log.* **80**, 845–865 (2015)
15. Łełyk, M.: *Axiomatic theories of truth, bounded induction and reflection principles*. Ph.D. dissertation, University of Warsaw (2017)
16. Łełyk, M., Wcisło, B.: Notes on bounded induction for the compositional truth predicate. *Rev. Symb. Logic* **10**, 455–480 (2017)
17. Löb, M.H.: Solution of a problem of Leon Henkin. *J. Symb. Logic* **20**, 115–118 (1955)
18. Pakhomov, F., Walsh, J.: Reflection ranks and ordinal analysis (2018). [arXiv:1805.02095](https://arxiv.org/abs/1805.02095)
19. Pudlák, P.: The lengths of proofs. In: Buss, S. (ed.) *Handbook of Proof Theory*, pp. 547–637. Elsevier, Amsterdam (1998)
20. Visser, A.: Semantics and the liar paradox. In: Gabbay, D., Günthner, F. (eds.) *Handbook of Philosophical Logic*, vol. 4, pp. 149–240. Reidel, Dordrecht (1989)
21. Visser, A.: From Tarski to Gödel. Or, how to derive the second incompleteness theorem from the undefinability of truth without self-reference (2018). [arXiv:1803.03937](https://arxiv.org/abs/1803.03937)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Affiliations

Ali Enayat¹  · Fedor Pakhomov² 

Fedor Pakhomov
pakhfn@mi.ras.ru

¹ Department of Philosophy, Linguistics, and the Theory of Science, University of Gothenburg, Gothenburg, Sweden

² Steklov Mathematical Institute of Russian Academy of Sciences, Moscow, Russia