

# MORE ON SYSTEMS OF TRUTH AND PREDICATIVE COMPREHENSION

CARLO NICOLAI

**ABSTRACT.** In the paper we survey the known connections between theories that extend a common base theory with typed truth axioms on the one hand and predicative set-existence assumptions on the other. How general can the mutual reductions between truth and comprehension be taken to be? In trying to address this question, we consider (typed) classical, positive truth and predicative comprehension as operations on theories.

**Keywords:** Predicative Comprehension. Axiomatic Theories of Truth. Relative Interpretability.

## 1. INTRODUCTION

The proof-theoretic analysis of axiomatic truth theories has uncovered several connections holding between systems extending a base theory with set existence axioms and systems extending the same base theory with axioms characterizing a primitive truth predicate. Taken at face value, these results suggest that assumptions on the existence of certain sets can be replaced by suitable semantic assumptions and vice versa. As we shall see shortly, for instance, the extension of Peano Arithmetic with specific assumptions on arithmetically definable sets is inter-reducible with a full, Tarskian truth theory over PA.

These connections have played an important role in the foundation of mathematics (e.g. in the analysis of the limits of predicativity Feferman (1991)), in the analysis of different solutions to the semantic paradoxes (proving more truth-theoretic iterations has been generally considered a virtue of a theory of truth), in the debate concerning the nature of truth (is truth a light or a substantial property?). They have been also suggested as a tool to carry out ontological reductions.<sup>1</sup>

At any rate, these mutual reductions between truth and set existence axioms involve theories built on a fixed base theory, usually Peano Arithmetic. In this work we set the basis for a different approach: truth and predicative comprehension will be taken as functors applying to arbitrary object theories satisfying some minimal requirements. Crucially, the results of applying these functors to the base theory will turn out to be equivalent tools for uncovering the first step of our commitment implicit in the acceptance of the base theory. Furthermore, the generality given by the proposed approach will also be reflected in a more uniform correspondence between set-theoretic and semantic assumptions, as it will be exemplified in §5.

The paper has two main parts. In the first, we survey several results connecting extensions Peano Arithmetic with truth axioms or predicative comprehension axioms. In the second, we introduce the more general setting just announced.

As this is more congenial to the original, second part of the work we will mostly focus, in the first part, on systems of *typed truth* — that is, theories formulated in languages in

---

<sup>1</sup>Cf. Halbach (2014, Ch. 23) and §5 of the present paper.

which the truth predicate applies to sentences not containing the truth predicate itself — and subsystems of ACA. The core argument, which gives rise to several variations, is the folklore mutual reduction between the system of Tarskian, compositional truth CT and ACA itself. We will also take into account the system PT of positive, typed truth, which is a notational variant of CT: its subsystems are not so easily comparable with their Tarskian counterparts but still enjoy several reductions to subsystems of ACA. All of this will be surveyed in §3.

The second part, starting with §4, will investigate a generalization of the systems considered in §3. We define three main operations on recursively enumerable (RE) theories: (i)  $T[\cdot]$  results in a Tarskian truth theory; (ii)  $\text{Tp}[\cdot]$  in a typed theory of truth simulating positive inductive definitions; (iii)  $\text{PC}[\cdot]$  adds predicative comprehension to the object theory. In order to study  $\text{PC}[\cdot]$  and relate it in full generality to the truth theories also studied, a variant of it,  $\text{PCS}[\cdot]$ , has to be taken into account: it applies to theories axiomatized by schemata in which schematic variables are replaced by ‘second-order’ variables. The upshot of §4 will be that, modulo mutual interpretability, the extension of arbitrary RE base theories via these three operations yields equivalent results. In §5 we reflect on the significance of the technical work carried out in the previous sections and consider possible extensions.

## 2. PRELIMINARIES

**2.1. Theories and Arithmetization.** Unless otherwise specified, arbitrary theories will be formulated in many-sorted, first-order logic. We assume they have a  $\Delta_1^b$  specification<sup>2</sup> — provably in a weak syntax theory to be defined — and that they are formulated in a Hilbert-style calculus in with Modus Ponens as its only rule of inference.<sup>3</sup>

On occasion, we consider *sequential* theories, namely theories that have a nice coding of sequences, which is in fact a safe tool to formalize syntax.<sup>4</sup> We employ  $T, U, V, W, \dots$  to range over arbitrary theories;  $A, B, C, \dots$  are taken to range over finitely axiomatized theories.

Whereas in §3 we work with PA as base theory, so all the usual representations of syntactic notions and operations can be assumed, in §4 we work in the theory  $S_2^1$ .  $S_2^1$  is a theory introduced by Sam Buss in Buss (1986) to study polynomial time computability; the functions that are provably recursive in  $S_2^1$  are exactly the p-time computable functions. Remarkably,  $S_2^1$  is finitely axiomatizable and mutually interpretable with the induction free fragment of first-order arithmetic Q. For details concerning coding in  $S_2^1$ , we refer to Buss (1986), Buss (1998). In particular, for any given language  $\mathcal{L}_W$ , we assume  $\Delta_1^b$ -definitions  $\text{Term}_{\mathcal{L}_W}(x)$ ,  $\text{CTerm}_{\mathcal{L}_W}(x)$ ,  $\text{Fml}_{\mathcal{L}_W}(x)$ ,  $\text{Sent}_{\mathcal{L}_W}(x)$ ,  $\text{Prf}_{\mathcal{L}_W}(x)$ ,  $\text{Proof}_{\mathcal{L}_W}(x, y)$  of the sets of terms, closed terms, formulas, sentences of the language of  $W$ , of proofs in the theory  $W$  and of the relation of being a proof in  $W$  of the  $\mathcal{L}_W$ -formula  $y$ . Also, we take the set of theorems of  $W$  to be

<sup>2</sup>That is, their axiom set is specified by a formula provably equivalent to formulas containing sharply bounded quantifiers (Cf. Hájek & Pudlák (1993, V.4. Def. 4.2)) and one bounded existential or universal quantifier only.

<sup>3</sup>Cf. Enderton (2001) for an example of such axioms system for first-order logic.

<sup>4</sup>More precisely,  $T$  is sequential iff it directly interprets (i.e. identity is mapped into identity and quantifiers are not relativized) *Adjunctive Set Theory*, that is a theory in a (first-order) language with  $\in$  and  $=$  whose axioms are:

$$(AS1) \quad \exists x \forall y y \notin x$$

$$(AS2) \quad \forall u, v \exists x \forall y (y \in x \leftrightarrow (y \in u \vee y = v))$$

Remarkably, Q is not sequential, but many sequential theories are interpretable in it.

defined by the  $\exists\Delta_1^b$ -formula  $\text{Pr}_W(x) :\leftrightarrow \exists y \text{Proof}_W(x, y)$ .<sup>5</sup> As to notational conventions, we will often abuse of Gödel corners and interchangeably employ them and Feferman's dot convention.

**2.2. Reductions.** The preferred means of reduction throughout the paper will be *many sorted*, —possibly non direct— *relative interpretability*.

Let  $T$  and  $W$  be theories containing  $S_2^1$ .<sup>6</sup> A *relative translation* of  $\mathcal{L}_T$  into  $\mathcal{L}_W$  can be described as a pair  $(\delta, F)$  where  $\delta$  is a  $\mathcal{L}_W$ -formula with one free variable —the domain of the translation— and  $F$  is a (finite) mapping that takes  $n$ -ary relation symbols of  $\mathcal{L}_T$  and gives back formulas of  $\mathcal{L}_W$  with  $n$  free variables.<sup>7</sup> The translation extends to the mapping  $\tau$ :

- $(R)^\tau(x_1, \dots, x_n) :\leftrightarrow F(R)(x_1, \dots, x_n)$ ;
- $\tau$  commutes with the propositional connectives;
- $(\forall x \varphi(x))^\tau :\leftrightarrow \forall x (\delta(x) \rightarrow \varphi^\tau)$  and  $(\exists x \varphi(x))^\tau :\leftrightarrow \exists x (\delta(x) \wedge \varphi^\tau)$ .

An *interpretation*  $K$  is then specified by a triple  $(T, \tau, W)$  such that for all sentences  $\varphi$  of  $\mathcal{L}_T$ ,

$$T \vdash \varphi \Rightarrow W \vdash \varphi^\tau$$

$W$  *locally* interprets  $T$  if and only if every finite subsystem of  $T$  is interpretable in  $W$ . An interpretation is *direct* if and only if it maps identity to identity and it does not relativize quantifiers. We will often not distinguish between an interpretation and the relative translation that supports it.

On occasion, we will employ the notion of *truth definability*, extensively studied by Fujimoto (2010). Truth definability will be mainly applied in the first part of the paper, as it assumes a fixed base theory (at least in its original formulation in Fujimoto (2010)). Let  $U, V$  be theories extending a syntactic base theory  $B$ : it is assumed, moreover, that  $U$ , but not necessarily also  $V$ , results from the extension of  $B$  via truth axioms.  $U$  is *relatively truth definable* in  $V$  if there is a  $\mathcal{L}_B$ -conservative relative interpretation of  $U$  into  $V$ . In other words, a truth definition of  $U$  in  $V$  is a relative interpretation  $\tau_o$  that behaves like the identity mapping when applied to truth-free formulas of  $\mathcal{L}_U$ .

**2.3. Definable Cuts and Incompleteness.** Cuts are initial segments of the natural numbers. We will in particular be interested in initial segments definable in theories  $T$  containing  $S_2^1$ . More precisely, a formula  $\varphi(x)$  is called *inductive* in  $T$  containing  $S_2^1$  if and only if

$$T \vdash \varphi(o) \wedge \forall y (\varphi(y) \rightarrow \varphi(Sy))$$

$\varphi(x)$  is a *T-cut* if and only if, additionally, it defines an initial segment of the  $T$ -numbers. In other words,

$$T \vdash \forall x, y (\varphi(x) \wedge y \leq x \rightarrow \varphi(y))$$

Here it is important to notice that  $o, S, \leq$  are understood in  $T$  via the assumed interpretation of  $S_2^1$  in  $T$ . By a well-know result of Solovay,<sup>8</sup> every inductive formula has a subcut:

<sup>5</sup>The class of  $\exists\Delta_1^b$  formulas is obtained from the class of  $\Delta_1^b$ -formulas by closing it under unbounded existential quantification.

<sup>6</sup>The relation of containment can be understood as subtheory relation or relative interpretability relation.

<sup>7</sup>The definition generalizes naturally to the many-sorted case by considering multiple domains.

<sup>8</sup>Although contained in an unpublished note (cf. Hájek & Pudlák (1993)).

**Lemma 1.** Let  $\varphi(x)$  be inductive in  $T$  containing  $S_2^1$ . Then there exists a  $T$ -cut  $\psi(x)$  such that

$$T \vdash \psi(x) \rightarrow \varphi(x).$$

The  $\psi$  in Lemma 1 is obtained by closing  $\varphi(\bar{n})$  for each  $n$  under transitivity of  $\leq$  so that this holds for all  $m \leq n$  as well. It is often useful, however, to employ a slightly modified notion of definable cut. The modifications are justified by the following Lemma:

**Lemma 2.** Let  $T$  interpret  $S_2^1$  and  $\mathcal{I}$  be a  $T$ -cut. Then we can find a subcut  $\mathcal{J}$  of  $\mathcal{I}$  such that  $T$  proves the following:

- (1)  $\mathcal{J}(x) \wedge \mathcal{J}(y) \rightarrow \mathcal{J}(x + y)$
- (2)  $\mathcal{J}(x) \wedge \mathcal{J}(y) \rightarrow \mathcal{J}(x \times y)$
- (3)  $\mathcal{J}(x) \wedge \mathcal{J}(y) \rightarrow \mathcal{J}(x \# y)$ .

where the smash function  $\#$  is such that  $x \# y = 2^{|x| \times |y|}$ , and  $|x| = \lceil \log_2(x + 1) \rceil$  — that is, the upper integer part of the binary logarithm of  $x + 1$ .

Therefore, in what follows, a definable cut can always be taken to be closed under addition, multiplication and the smash function. Lemma 2 is particularly important as it guarantees that, given a  $T$ -cut  $\mathcal{I}$  — with  $T$  again containing  $S_2^1$ , we can always restrict our attention to a subcut  $\mathcal{J}$  of  $\mathcal{I}$  satisfying the axioms of  $S_2^1$ . This will be extensively employed in §4: having  $S_2^1$  available on a cut enables one to have a meaningful and smooth arithmetization of the syntax in the cut. This will prove to be crucial in many of the interpretations defined later on.

We will extensively make use of the following:

**Lemma 3** (Wilkie, Nelson).  $\mathbb{Q}$  interprets  $S_2^1$  on a definable cut.

$\mathbb{Q}$  does not have any induction: it cannot thus prove the usual properties of  $\leq$ . An important part of the proof of Lemma 3 consists in showing how the interpretation is indeed defined on a initial segment of the  $\mathbb{Q}$ -numbers. This makes possible the downwards preservation of  $\Pi_1$ -sentences (cf. Hájek & Pudlák (1993, V.5(c))) so that, for instance,  $S_2^1 + \text{Con}(U)$  is interpretable in  $\mathbb{Q} + \text{Con}(U)$ .

Another remarkable fact concerning  $S_2^1$  is that Gödel's Second Incompleteness theorems can be meaningfully stated and proved in it. This makes possible a generalization of the Gödel's second incompleteness theorem, due to Pavel Pudlák in its original form. The beautiful version that we present is due to Albert Visser.

**Lemma 4.** Let  $U$  be given by a  $\Sigma_1$ -formula. Then  $U$  does not interpret  $S_2^1 + \text{Con}(U)$ .

*Proof.* Assume it does. There must be, therefore,  $A$  be a finite subsystem of  $U$  that interprets  $S_2^1 + \text{Con}(U)$ . By  $\Sigma_1$ -completeness,  $S_2^1$  proves the formalization of the fact that  $A \subset U$ . Thus  $A$  interprets  $S_2^1 + \text{Con}(A)$ . Since  $A$  and  $S_2^1 + \text{Con}(A)$  are finite, we have

$$S_2^1 \vdash \text{Con}(A) \rightarrow \text{Con}(S_2^1 + \text{Con}(A)).$$

□

## 3. ARITHMETICAL COMPREHENSION

In this section we survey the mutual reductions between truth axioms and predicative comprehension in the standard setting, that is when PA is taken as object theory. After a few useful definitions and lemmata, in §3.1 we consider how Tarskian and positive truth can define the comprehension schema of ACA and variants thereof. We then consider the converse direction in §3.2. §3.3 hints at how to compare iterations of Tarskian truth and predicative comprehension to each other and to systems of type-free truth.

The language  $\mathcal{L}_2$  extends the language  $\mathcal{L} := \{o, S, +, \times, \leq\}$  of arithmetic with an additional sort for ‘sets of natural numbers’, or ‘reals’.

**Definition 1.**

(i) ACA is formulated in  $\mathcal{L}_2$ . It extends PA with

$$(CA) \quad \exists Y \forall u (u \in X \leftrightarrow \varphi(u))$$

where  $\varphi$  does not contain bound set-variables (it may contain set parameters); and

$$(S\text{-Ind}^2) \quad \varphi(o) \wedge \forall y (\varphi(y) \rightarrow \varphi(Sy)) \rightarrow \forall y \varphi(y)$$

for  $\varphi \in \mathcal{L}_2$ .

(ii) To obtain  $ACA_o$  one replaces (S-Ind<sup>2</sup>) with the single  $\mathcal{L}_2$ -sentence

$$(A\text{-Ind}^2) \quad o \in X \wedge \forall y (y \in X \rightarrow Sy \in X) \rightarrow \forall y y \in X$$

**Lemma 5.**  $ACA_o$  is not interpretable in PA.

*Proof.* If  $ACA_o$  were interpretable in PA, then  $ACA_o$  would be interpretable in a suitable finite subsystem  $A$  of PA.

Thus PA would be interpretable in  $A$ . But PA is reflexive — that is, it proves the consistency of all its finite subtheories. Therefore,

$$\begin{aligned} PA &\vdash \text{Con}(A) \rightarrow \text{Con}(ACA_o) \\ PA &\vdash \text{Con}(ACA_o) \end{aligned}$$

But  $PA \subset ACA_o$ , so the last line contradicts Gödel’s second incompleteness theorem.  $\square$

## 3.1. Arithmetical Comprehension from Truth.

**Definition 2** (CT).

(i) CT extends PAT (i.e. PA in  $\mathcal{L} \cup \{\text{Tr}\}$ ) with the universal closure of

$$(CT1) \quad \text{CTerm}_{\mathcal{L}}(x) \wedge \text{CTerm}_{\mathcal{L}}(y) \rightarrow (\text{Tr } x \circ y \leftrightarrow \text{val}(x) \circ \text{val}(y)) \text{ with } \circ \in \{=, \leq\}$$

$$(CT2) \quad \text{Sent}_{\mathcal{L}}(x) \rightarrow (\text{Tr } \neg x \leftrightarrow \neg \text{Tr } x)$$

$$(CT3) \quad \text{Sent}_{\mathcal{L}}(x \wedge y) \rightarrow (\text{Tr } (x \wedge y) \leftrightarrow \text{Tr } x \wedge \text{Tr } y)$$

$$(CT4) \quad \text{Sent}_{\mathcal{L}}(\forall v x) \rightarrow (\text{Tr } \forall v x \leftrightarrow \forall y (\text{CTerm}(y) \rightarrow \text{Tr } \text{sub}(x, v, y)))$$

In (CT4),  $v$  codes a variable.

(ii)  $CT \uparrow$  is obtained by allowing only  $\mathcal{L}$ -formulas as instances of the induction scheme of CT;

(iii)  $CT^- := (CT_1) - (CT_4)$  plus:

$$(IInd) \quad \text{tot}(x) \rightarrow \text{Tr}x(\ulcorner o \urcorner) \wedge \forall y (\text{Tr}x(y) \rightarrow \text{Tr}x(Sy)) \rightarrow \forall y \text{Tr}x(y)$$

where

$$(tot) \quad \text{tot}(x) :\leftrightarrow \text{Fml}_{\mathcal{L}}^1(x) \rightarrow \forall y (\text{Tr}x(y) \vee \text{Tr}_{\neg}x(y))$$

and  $\text{Fml}_{\mathcal{L}}^1(x)$  expresses that  $x$  is a formula of  $\mathcal{L}$  with only the first variable free.

$CT^-$  is considered in Fischer (2009). It is finitely axiomatized and thus weaker than  $CT$ : the latter is in fact reflexive.

We summarize some useful facts about  $CT$  and its restricted variants. Most of them, in slightly different form, can also be found for instance in Cantini (1989), Feferman (1991), Halbach (2014).

**Lemma 6.**

- (i) for all  $\varphi(v) \in \mathcal{L}$ ,  $CT \uparrow \vdash \forall x (\text{CTerm}_{\mathcal{L}}(x) \rightarrow (\text{Tr sub}(\ulcorner \varphi \urcorner, v, x) \leftrightarrow \varphi(\text{val}(x))))$ ;
- (ii)  $\forall x (\text{Sent}_{\mathcal{L}}(x) \wedge \text{Pr}_{\text{PA}}(x) \rightarrow \text{Tr}x)$ , where  $\text{Pr}_{\text{PA}}(\cdot)$  expresses canonical provability in PA;
- (iii)  $CT \uparrow \vdash \forall x (\text{Fml}_{\mathcal{L}}^1(x) \rightarrow \text{tot}(x))$
- (iv)  $CT^- \vdash \text{Con}(\text{PA})$
- (v)  $CT \vdash \text{Sent}_{\mathcal{L}}(\forall v z) \wedge \text{CTerm}_{\mathcal{L}}(x) \wedge \text{CTerm}_{\mathcal{L}}(y) \rightarrow (\text{val}(x) = \text{val}(y) \rightarrow (\text{Tr}(\text{sub}(z, v, x) \leftrightarrow \text{Tr}(\text{sub}(z, v, y))))$

*Proof.* (i) is obtained by external induction on  $\varphi$ . For (ii) one employs the induction axioms of  $CT$ . (iii) follows from  $CT_2$ . (iv) is obtained by combining (IInd) and (i) to mimic in  $CT^-$  the  $CT$ -proof of (ii). (v) also follows from a formal induction on the complexity of the formula  $z$ .  $\square$

By extending PAT with the schema of Lemma 6(i) taken as axiom, we obtain the theory UTB. In  $UTB \uparrow$ , like in  $CT \uparrow$ , the truth predicate is not allowed into instances of the induction scheme. The next lemma summarizes some well-known facts concerning the interpretability of typed truth in PA.

**Lemma 7.**

- (i) PA locally interprets  $TB \uparrow$ ,  $UTB \uparrow$ ,  $TB$ ,  $UTB$ ;
- (ii) PA interprets  $TB \uparrow$ ,  $UTB \uparrow$ ,  $TB$ ,  $UTB$ ;
- (iii) PA interprets  $CT \uparrow$ ;
- (iv) PA does not interpret  $CT^-, CT$ .

*Proof.* Ad (i), in a finite subsystem  $S$  of  $UTB$  instances of the uniform disquotation scheme involve sentences  $\varphi \in \mathcal{L}$  such that  $\text{lc}_x(\varphi) \leq n$  for a standard  $n$ , where  $\text{lc}_x(x)$  a primitive recursive function that, applied to a  $\mathcal{L}$ -formula  $\varphi$ , yields the number of its logical symbols. We interpret the truth predicate of  $S$  with the partial truth predicate, definable in PA, for sentences of complexity up to  $n$ . (ii) follows from (i) and Orey's compactness theorem, according to which if  $W$  is locally interpretable in  $T$  and  $T$  is reflexive, then  $T$  interprets  $W$ . The proof of (iii) contained in Fischer (2009) works if instead of the cut-elimination argument of Halbach (1999),<sup>9</sup> one employs the proof of the conservativity of  $CT \uparrow$  over PA

<sup>9</sup>It has been shown to contain a mistake by Kentaro Fujimoto.

contained in Leigh (2014). A new, direct proof is contained in Enayat & Visser (2014, Theorem 5.1). (iv) is straightforward from Lemma 6(iv).  $\square$

**Proposition 1.** CT interprets ACA. Moreover, the interpretation behaves like the identity mapping on the arithmetical vocabulary.

*Proof.* A full proof is contained in Halbach (2014): we present some of its details as they will be useful for later proofs.

We define a substitution function for elements of  $\text{Fml}_{\mathcal{L}}^1$ :

$$\text{sb}(\ulcorner \varphi(v_1) \urcorner, \ulcorner t \urcorner) := \ulcorner \varphi(t) \urcorner$$

We let  $\text{sb}(x, y)$  take a default value, e.g.  $\ulcorner o = o \urcorner$ , when it is not applied to an element of  $\text{Fml}_{\mathcal{L}}^1$  and to a term. The translation  $\iota$  is then crucially defined by the following clauses:

$$\begin{aligned} (x \in y)^\iota &:\leftrightarrow \text{Tr}(\text{sb}(x, y)) && y \text{ is satisfied by } x \\ \iota &\text{ maps the arithmetical nonlogical vocabulary of } \mathcal{L}_2 \text{ into itself;} \\ \iota &\text{ commutes with propositional quantifiers;} \\ (\forall X \varphi)^\iota &:\leftrightarrow \forall X' (\text{Fml}_{\mathcal{L}}^1(X') \rightarrow (\varphi)^\iota) \end{aligned}$$

Set variables and number variables are kept distinct to avoid clashes: this is emphasized by denoting with  $X'$  the  $\mathcal{L}$ -variable resulting from the translation. Strictly speaking,  $X'$  is a variable of the only available sort.

To verify that  $\iota$  behaves as required, we consider the crucial case of comprehension axioms, i.e. we prove

$$(4) \quad (\exists X \forall x (x \in X \leftrightarrow \varphi(x)))^\iota$$

Notice that  $\varphi(x)$  may contain first- and second-order parameters. Then, in CT:

$$(5) \quad \forall x (\varphi'(x, \text{Tr}(\text{sb}(w, X'))) \leftrightarrow \varphi'(x, \text{Tr}(\text{sb}(w, X')))) \quad \text{where } w \text{ is a closed term}$$

$$(6) \quad \forall x (\text{Tr} \ulcorner \varphi'(\dot{x}, \text{sb}(w, X')) \urcorner \leftrightarrow \varphi'(x, \text{Trsb}(w, X'))) \quad \text{Lemma 6(i) and CT}_{1-4}$$

$$(7) \quad \exists Y' \forall x (\text{Tr}(\text{sb}(\dot{x}, Y')) \leftrightarrow \varphi'(x, \text{Tr}(\text{sb}(w, X'))))$$

Essentially, the last line results from the application of the compositional axioms of CT to move the occurrences of the truth predicate in  $\varphi'(x, \text{Trsb}(w, X'))$  in front of the formula. In line (6), we reason for an arbitrary  $\varphi$ : strictly speaking, an external induction on the complexity of the formula is needed in which the different compositional axioms are employed. The last line follows from the properties of substitution.  $\square$

In the proof of Lemma 1, the compositional axioms CT<sub>1-4</sub> are only employed in moving from (5) to (6) and, in particular, to deal with parameters occurring into the comprehension scheme. We then let  $\text{ACA}^{\text{pf}}$  to be exactly as ACA but with no parameters allowed to appear into the comprehension scheme, and  $\text{ACA}\uparrow$  and  $\text{ACA}^{\text{pf}}\uparrow$  their counterparts with induction involving only  $\mathcal{L}$ -formulas. From the proof of Proposition 1 we can thus extract:

**Proposition 2.**

- (i)  $\text{UTB}\uparrow$  interprets  $\text{ACA}^{\text{pf}}\uparrow$
- (ii)  $\text{UTB}$  interprets  $\text{ACA}^{\text{pf}}$ .

In both cases, the translation behaves like id for arithmetical vocabulary.

Proposition 2 tells us that the admission of second-order parameters in the comprehension schema corresponds, on the truth-theoretic side, to a substantial jump from a disquotational theory of truth to a compositional one. The substantiality of this jump can be measured in terms of proof-theoretic strength — and thus in terms of interpretability strength — as Lemma 6 shows. In other words, the presence of parameters in the comprehension schema enables one to go from a theory that is conservative over PA to a theory that proves  $\text{Con}(\text{PA})$ .

On the negative side, we have

**Proposition 3.**

- (i)  $\text{TB}\uparrow$ ,  $\text{UTB}\uparrow$ ,  $\text{CT}\uparrow$  do not interpret  $\text{ACA}_0$ ;
- (ii)  $\text{CT}^-$  interprets  $\text{ACA}_{\Pi_1}$

*Proof.* (i) If one of  $\text{TB}\uparrow$ ,  $\text{UTB}\uparrow$ ,  $\text{CT}\uparrow$  interpreted  $\text{ACA}_0$ , by Lemma 7, PA would interpret  $\text{ACA}_0$ , quod non. (ii) follows from Cantini (1989, Proposition 3.6).  $\square$

Positive inductive definitions can be employed to define the set of arithmetical truths Moschovakis (1974). The idea is to identify this set as the fixed point of a suitable operator on sets of natural numbers extracted by the inductive definition Halbach (2014, §8.7). The theory PT captures the clauses of this inductive definition. PT and its subsystems are studied in Fischer (2009). The peculiarity of positive inductive definitions is that in their clauses the truth predicate can only appear in the scope of an even number of negation symbols.

**Definition 3.**

- (i) PT extends PAT with the axioms:

- (PT1)  $\text{CTerm}_{\mathcal{L}}(x) \wedge \text{CTerm}_{\mathcal{L}}(y) \rightarrow (\text{Tr } x \circ y \leftrightarrow \text{val}(x) \circ \text{val}(y))$
- (PT2)  $\text{CTerm}_{\mathcal{L}}(x) \wedge \text{CTerm}_{\mathcal{L}}(y) \rightarrow (\text{Tr } \neg x \circ y \leftrightarrow \neg(\text{val}(x) \circ \text{val}(y)))$
- (PT3)  $\text{Sent}_{\mathcal{L}}(x) \rightarrow (\text{Tr } \neg \neg x \leftrightarrow \text{Tr } x)$
- (PT4)  $\text{Sent}_{\mathcal{L}}(x \wedge y) \rightarrow (\text{Tr } x \wedge y \leftrightarrow (\text{Tr } x \wedge \text{Tr } y))$
- (PT5)  $\text{Sent}_{\mathcal{L}}(x \wedge y) \leftrightarrow (\text{Tr } \neg(x \wedge y) \leftrightarrow (\text{Tr } \neg x \vee \text{Tr } \neg y))$
- (PT6)  $\text{Sent}_{\mathcal{L}}(\forall v x) \leftrightarrow (\text{Tr } \forall v x \leftrightarrow \forall y(\text{CTerm}(y) \rightarrow \text{Tr}(x(y/v))))$
- (PT7)  $\text{Sent}_{\mathcal{L}}(\neg \forall v x) \leftrightarrow (\text{Tr } \neg \forall v x \leftrightarrow \exists y(\text{CTerm}(y) \wedge \text{Tr } \neg x(y/v)))$

- (ii)  $\text{PT}\uparrow$  restricts the induction scheme to  $\mathcal{L}$ -formulas
- (iii) In  $\text{PT}^-$  the induction scheme of PAT is replaced by (IInd).

PT is identical to CT, as PT proves (CT2) by formal induction on the complexity of sentences. The proof of Proposition 1 carries over without modifications: PT thus interprets ACA in the manner described.

$\text{PT}\uparrow$  on the other hand, that is when induction is restricted to  $\mathcal{L}$ -sentences, is a *proper* subtheory of  $\text{CT}\uparrow$ . To see this, one notices that the general theory of inductive definitions entails the existence of fixed points of positive inductive definitions: this fact can be used to define a subset  $S \subseteq \mathcal{M}$  for any  $\mathcal{M} \models \text{PA}$  such that  $(\mathcal{M}, S) \models \text{PT}\uparrow$  (cf. Cantini (1989)). By a theorem of Lachlan (Lachlan (1981)), by contrast, not any model of PA can be expanded



to a model of  $\text{CT}\uparrow$ . This explains the properness of the subtheory relation. Therefore, by Lemma 7, PA interprets  $\text{PT}\uparrow$ .<sup>10</sup>

To get closer to our interests and consider subsystems of second-order arithmetic, we have:

**Proposition 4.**

- (i)  $\text{PT}$  interprets  $\text{ACA}$ ;
- (ii)  $\text{PT}^-$  interprets  $\text{ACA}_0$ .

*Proof.* (i) follows from the identity of  $\text{CT}$  and  $\text{PT}$  and Lemma 1. (ii) is proved by Fischer (2009): the comprehension axioms require external induction on the complexity of the formula involved. Compare (5)-(7) in the proof of Proposition 1: in (6), the hidden induction on the complexity of formulas is replaced here by an induction on its *positive* complexity. Moreover,  $\text{PT}^-$  proves the translation of  $\text{tot}(x)$  where  $x \in \text{Fml}_{\mathcal{L}}^1$ .  $\square$

**3.2. Truth from Comprehension.**

**Proposition 5.**  $\text{ACA}$  defines the truth predicate of  $\text{CT}$ .

*Proof.* The proof is contained in Takeuti (1987). Let  $\text{lcx}(\cdot)$  be as above, and  $\mathcal{T}(X, x)$  mean ‘ $X$  is a truth set for sentences of complexity  $\leq x$ ’, that is

$$\begin{aligned} \mathcal{T}(X, x) :& \leftrightarrow \forall u, v (\text{Cterm}_{\mathcal{L}}(u) \wedge \text{Cterm}_{\mathcal{L}}(v) \rightarrow (u \circ v \in X \leftrightarrow \text{val}(u) \circ \text{val}(v))) \wedge \\ & \forall u (\text{Sent}_{\mathcal{L}}(\neg u) \wedge \text{lcx}(\neg u) \leq x \rightarrow (\neg u \in X \leftrightarrow u \notin X)) \wedge \\ & \forall u, v (\text{Sent}_{\mathcal{L}}(u \wedge v) \wedge \text{lcx}(u \wedge v) \leq x \rightarrow (u \wedge v \in X \leftrightarrow (u \in X \wedge v \in X))) \wedge \\ & \forall u, v (\text{Sent}_{\mathcal{L}}(\forall uv) \wedge \text{lcx}(\forall uv) \leq x \rightarrow (\forall uv \in X \leftrightarrow \forall y (\text{CTerm}_{\mathcal{L}}(y) \rightarrow \text{sb}(y, v) \in X))) \end{aligned}$$

It is first noticed that  $\text{ACA}^{\text{pf}}$  with the induction schema restricted to  $\mathcal{L}_2$ -formulas with no bound set variables proves, by induction on  $x$ , that truth definitions for sentences of restricted complexity are unique:

$$(8) \quad \mathcal{T}(X, x) \wedge \mathcal{T}(Y, x) \wedge \text{Sent}_{\mathcal{L}}(y) \wedge \text{lcx}(y) \leq x \rightarrow (y \in X \leftrightarrow y \in Y)$$

By arithmetical comprehension, one obtains, in  $\text{ACA}^{\text{pf}}\uparrow$ ,

$$(9) \quad \exists X \mathcal{T}(X, 0),$$

because sentences with logical complexity  $\leq 0$  are atomic statements. Moreover, in  $\text{ACA}_0$ , we have, still by arithmetical comprehension, that

$$(10) \quad \forall x (\exists X (\mathcal{T}(X, x)) \rightarrow \exists X (\mathcal{T}(X, x+1)))$$

An instance of the induction schema of  $\text{ACA}$  (indeed  $\text{ACA}_{\Sigma_1^1}$  would suffice, as we only need an instance of  $\Sigma_1^1$ -induction), gives us

$$(11) \quad \forall x \exists X \mathcal{T}(X, x)$$

The required truth definition is given by the formula

$$(12) \quad \tau(x) : \leftrightarrow \exists X (\mathcal{T}(X, \text{lcx}(x)) \wedge x \in X)$$

$\square$

**Corollary 1.**

<sup>10</sup>There is something more that holds: any model of PA can be expanded to a model of  $\text{PT}^-$ .

- (i)  $\text{ACA}^{\text{pf}} \uparrow$  defines the truth predicate of  $\text{UTB} \uparrow$ ;
- (ii)  $\text{ACA}^{\text{pf}}$  defines the truth predicate of  $\text{UTB}$ ;
- (iii)  $\text{ACA}_o$  interprets  $\text{PT}^-$ ;
- (iv)  $\text{ACA}_{\Pi_1^1}$  defines the truth predicate of  $\text{CT}^-$ ;
- (v)  $\text{ACA}$  defines the truth predicate of  $\text{PT}$ .

*Proof.* (i) and (ii) are immediate by inspection of the proof of Proposition 5, by induction on the complexity of the formula involved in the uniform diquotation. A detailed proof of (iii) can be found in Fischer (2009): one replaces  $\text{lcx}(\cdot)$  with a measure of the *positive complexity*  $\text{pcx}(\cdot)$  in an analog of (12). Then it is shown that  $\text{ACA}_o$  proves  $(\text{IInd})$  under the translation: to this end, the induction schema for arithmetical formulas of  $\mathcal{L}_2$ , provable in  $\text{ACA}_o$ , is employed. (iv), in addition, requires  $\Pi_1^1$  induction to prove  $\text{CT}^-$ . (v) is immediate by Proposition (5) and the fact that  $\text{PT}$  is a subtheory of  $\text{CT}$ .  $\square$

**3.3. Extensions and Hierarchies.** The processes of extending  $\text{PA}$  to  $\text{CT}$  and  $\text{ACA}$  can be iterated. To this end, one first assumes a suitable notation for ordinals below  $\Gamma_o$  (cf. Pohlers (2009, Ch. 2)): in other words, any ordinal  $\alpha < \Gamma_o$  can be coded by a natural number and thus represented in  $\text{PA}$  by a unique term  $\bar{\alpha}$ .

To define the theories  $\text{RT}_{<\theta}$  for  $\theta \leq \Gamma_o$  — $\text{RT}$  standing for ‘ramified truth’— one considers languages  $\mathcal{L}_{<\theta} := \mathcal{L}_{\text{PA}} \cup \{T_o, \dots, T_\eta\}$  for  $\eta < \theta$ , where  $T_\alpha := (T, \alpha)$ . The axioms of  $\text{RT}_{<\theta}$  thus stipulate how the truth predicates  $T_\eta$  work for sentences of  $\mathcal{L}_{<\eta}$ , for  $\eta < \theta$  (cf. Halbach (2014, p. 113)).

The definition of iterations of  $\text{ACA}$  to  $\theta \leq \Gamma_o$ , called  $\text{RA}_{<\theta}$  for ‘ramified analysis up to  $\theta$ ’, enjoys different variations. Since we are not so much interested in iterations of predicative comprehension, we refer to Feferman (1964) for the definition of  $\text{RA}_{<\theta}$ .

The fundamental fact that links ramified truth and analysis is the following, whose proof is sketched in Feferman (1991):

**Proposition 6** (Feferman). For  $\alpha \leq \Gamma_o$ ,  $\text{RT}_{<\alpha}$  and  $\text{RA}_{<\alpha}$  are mutually interpretable.

In the light of Proposition 6, one may for instance obtain a reduction of the classical axiomatization of Kripke’s theory of truth  $\text{KF}$  (cf. Feferman (1991), Halbach (2014)) to  $\text{RT}_{<\epsilon_o}$ : in particular,  $\text{KF}$  defines all truth predicates of  $\text{RT}_{<\epsilon_o}$ . Also, from Cantini (1989, §9) one can extract an interpretation of  $\text{KF}$  in  $\text{RT}_{<\epsilon_o}$ . Given the mutual truth definability between the remarkable disquotational system  $\text{PUTB}$  of positive uniform (type free) disquotation and  $\text{KF}$  (cf. Halbach (2009)), Proposition 6 closely relates  $\text{PUTB}$  and iterations of predicative comprehension. It also shows the close relationships between  $\text{KF}$ ,  $\text{PUTB}$ , and  $\widehat{\text{ID}}_1$ .<sup>11</sup>

Many other results and reductions between type-free truth and second order-arithmetic have been investigated in recent years. As we anticipated in the introduction, however, the strategy put forward in the next section does not enjoy a plausible extension to the type-free case yet: we lack natural ways of extending the strategy proposed in the next section allow self-applications of the operations on theories to be defined. Therefore we decided to omit a systematic survey of the relationships between type-free systems and subsystems of second-order arithmetic.

<sup>11</sup> $\widehat{\text{ID}}_1$  is the theory postulating the existence of fixed-points for arbitrary positive arithmetical operators. For a detailed definition of the full  $\text{ID}_1$ , cf. Pohlers (2009, Ch. 9, Def. 9.1).

## 4. TRUTH AND PREDICATIVE COMPREHENSION AS FUNCTORS

As anticipated in the introduction, we will now consider the reductions surveyed in §3 from a different angle. This approach finds its roots in a philosophical attempt to distinguish patterns of reasoning belonging to the syntactico/truth-theoretic component of a theory of truth on one side and its mathematical/object-theoretic component on the other. This separation has been used, also by the author, to analyze the connections between axiomatic truth and truth-theoretic deflationism (cf. Halbach (2014), Heck (2015), Nicolai (2015a)). The fundamental insight offered by this approach is thus that, unlike what happens in the standard construction, the assumption of a full, Tarskian theory of truth does not immediately lead to new object-theoretic, or mathematical insights, although it may do so via a suitably defined interpretation.

In §4.1 we define and introduce the basic properties of  $T[\cdot]$  and  $\text{Tp}[\cdot]$  that operate on theories  $U$  by applying Tarskian truth and positive, typed truth to them. Although inductive reasoning involving truth-theoretic and syntactic notions is not directly available in  $T[\cdot]$  (and  $\text{Tp}[\cdot]$ ), it can be mimicked in it by exploiting the properties of definable cuts introduced in §2: this will be the core of §4.2. The arithmetized model corresponding to the formalization—in weak arithmetical context—of the construction of the term model used in Henkin’s proof of the completeness theorem, comes with a truth predicate. This truth predicate is used in §4.3 to interpret the truth predicates of  $T[\cdot]$  and  $\text{Tp}[\cdot]$ . In §4.4 we turn to predicative comprehension and investigate the functors  $\text{PC}[\cdot]$  and  $\text{PCS}[\cdot]$ ; in §4.5 we finally relate them to  $T[\cdot]$  and  $\text{Tp}[\cdot]$ .

In the following sections the object theory  $U$  will be assumed to be formulated in a relational language. On occasion we will require sequentiality or the capability of interpreting  $S_2^1$ .

**4.1. Typed-Truth as an Operation on Theories.** The functor  $T[\cdot]$  applies to an arbitrary RE theory  $U$  and yields the three-sorted theory  $T[U]$  in a language  $\mathcal{L}_T$  with sorts in  $\{\mathfrak{s}, \mathfrak{o}, \mathfrak{sq}\}$  (for ‘syntax’, ‘object-theory’, ‘sequences/variable assignments’), constants proper of the language  $\mathcal{L}_\#$  of  $S_2^1,^{12}$  the function symbol  $\cdot(\cdot)$  of type  $(\mathfrak{sq}, \mathfrak{s}) \rightarrow \mathfrak{o}$  and the predicate symbol  $\text{Sat}$  of sort  $(\mathfrak{s}, \mathfrak{sq})$ . The former will give rise to expressions of the form  $v_i^{\mathfrak{sq}q}(v_i^{\mathfrak{s}}) = v_i^{\mathfrak{o}}$ , stating that the  $v_i^{\mathfrak{s}}$ -th element of a variable assignment  $v_i^{\mathfrak{sq}q}$  is the  $U$ -object  $v_i^{\mathfrak{o}}$ , whereas the latter will be characterized as a satisfaction predicate. Greek letters  $\varphi, \psi, \dots$  are taken to range over formulas of  $\mathcal{L}_U$ .

The axioms of  $T[U]$ , besides the axioms of  $U$ , are

( $S_2^1$ ) axioms of  $S_2^1$

$$(\text{sq}) \quad \exists v_j^{\mathfrak{sq}q} (\forall v_k^{\mathfrak{s}} (v_k^{\mathfrak{s}} \neq v_i^{\mathfrak{s}} \rightarrow v_i^{\mathfrak{sq}q}(v_k^{\mathfrak{s}}) = v_j^{\mathfrak{sq}q}(v_k^{\mathfrak{s}})) \wedge v_j^{\mathfrak{sq}q}(v_i^{\mathfrak{s}}) = v_i^{\mathfrak{o}})$$

$$(\text{tat}) \quad \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner R(v_1^{\mathfrak{o}}, \dots, v_n^{\mathfrak{o}}) \urcorner) \leftrightarrow R(v_i^{\mathfrak{sq}q}(\ulcorner v_1^{\mathfrak{o}} \urcorner), \dots, v_i^{\mathfrak{sq}q}(\ulcorner v_n^{\mathfrak{o}} \urcorner))$$

for every relation symbol  $R$  in  $\mathcal{L}_U$

$$(\text{t}\neg) \quad \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner \neg \varphi \urcorner) \leftrightarrow \neg \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner \varphi \urcorner)$$

$$(\text{tv}) \quad \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner \varphi \vee \psi \urcorner) \leftrightarrow \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner \varphi \urcorner) \vee \text{Sat}(v_i^{\mathfrak{sq}q}, \ulcorner \psi \urcorner)$$

<sup>12</sup>For simplicity, we may require them to be unofficial abbreviations of their relational counterparts.

$$(tq) \text{ Sat}(v_i^{sq}, \ulcorner \exists v_i^o \varphi \urcorner) \leftrightarrow (\exists v_j^{sq} \overset{v_i^o}{\sim} v_i^{sq})(\text{Sat}(v_j^{sq}, \ulcorner \varphi \urcorner))$$

where

$$v_j^{sq} \overset{v_i^o}{\sim} v_i^{sq} \leftrightarrow \forall v_j^s (v_j^s \neq v_i^s \rightarrow v_j^{sq}(v_j^s) = v_i^{sq}(v_j^s))$$

(sq) states minimal conditions on the existence and on the behaviour of sequences. (tq) simply formalizes Tarski's clause on quantified formulas:  $\exists v_i \varphi$  is satisfied by a variable assignment if and only if there is a second variable assignment, differing from the former only in what they assign to  $v_i$  (this is captured by the abbreviation  $v_j^{sq} \overset{v_i^o}{\sim} v_i^{sq}$ ), satisfying  $\varphi$ .

Intuitively, by applying the operator  $T[\cdot]$  to  $U$ , one expands a model of  $U$  with a universe of syntactic objects —here numbers in the sense of  $S_2^1$ — and a universe of 'mixed' objects, sequences of  $\mathcal{L}_U$ -variable assignments living in a disjoint domain.

We define an alternative functor:  $\text{Tp}[U]$  adds *positive* compositional axioms similar to the axioms of  $\text{PT}\uparrow$  in Definition 3 but in a many-sorted way compatible with the idea behind the construction of  $T[U]$ . Again  $\text{Tp}[U]$  is formulated in  $\mathcal{L}_T$ ; its axioms, besides the axioms of  $U$ ,  $S_2^1$ , (sq), (tat), (tv), are

$$(fat) \text{ Sat}(v_i^{sq}, \ulcorner \neg R(v_1^o, \dots, v_n^o) \urcorner) \leftrightarrow \neg (R(v_i^{sq}(\ulcorner v_1^o \urcorner), \dots, v_i^{sq}(\ulcorner v_n^o \urcorner)))$$

for every relation symbol  $R$  in  $\mathcal{L}_U$

$$(tdn) \text{ Sat}(v_i^{sq}, \ulcorner \neg \neg \ulcorner \varphi \urcorner \urcorner) \leftrightarrow \text{Sat}(v_i^{sq}, \ulcorner \varphi \urcorner)$$

$$(fv) \text{ Sat}(v_i^{sq}, \ulcorner \neg \ulcorner \varphi \vee \psi \urcorner \urcorner) \leftrightarrow (\text{Sat}(v_i^{sq}, \ulcorner \neg \ulcorner \varphi \urcorner \urcorner) \wedge \text{Sat}(v_i^{sq}, \ulcorner \neg \ulcorner \psi \urcorner \urcorner))$$

$$(fq) \text{ Sat}(v_i^{sq}, \ulcorner \neg \ulcorner \exists v_i^o \varphi \urcorner \urcorner) \leftrightarrow (\forall v_j^{sq} \overset{v_i^o}{\sim} v_i^{sq})(\text{Sat}(v_j^{sq}, \ulcorner \neg \ulcorner \varphi \urcorner \urcorner))$$

By external induction on the positive complexity of  $\varphi(\vec{v}_i^o)$ , we have:

**Lemma 8.**  $\text{Tp}[U]$  (and thus  $T[U]$ ) proves, for all  $\varphi(\vec{v}_i^o)$  in  $\mathcal{L}_U$ ,

$$(i) \text{ Sat}(v_i^{sq}, \ulcorner \varphi(v_i^o, \dots, v_n^o) \urcorner) \leftrightarrow \varphi(v_i^{sq}(\ulcorner v_i^o \urcorner), \dots, v_i^{sq}(\ulcorner v_n^o \urcorner))$$

We will often refer to variables of sort  $o, s, sq$  as  $u, v, w, x, y, z, \dots, i, j, k, l, m, n, \dots$ , and  $a, b, c, \dots$  respectively.

**4.2. Consistency from Truth.** A remarkable fact concerning the method of shortening of cuts, introduced in §2, is that it enables one to prove the consistency of the object theory  $U$  relativized to a definable cut, once we are able to prove, in the theory of truth, that all axioms of  $U$  are true.<sup>13</sup>

Therefore we define

$$\text{AxT}_U := \forall a \forall k (\text{AxL}_U(k) \vee \text{AxU}(k) \rightarrow \text{Sat}(a, k))$$

where  $\text{AxL}_U(k)$  and  $\text{AxU}(k)$  are  $\Delta_1^b$ -representations of the logical and the nonlogical axioms of  $U$  in  $S_2^1$ .  $\text{AxT}_U$  thus states that all logical and nonlogical axioms of  $U$  are true, i.e. satisfied by all sequences. This definition of  $\text{AxT}_U$  enables us to simplify the proof of a similar result given in Nicolai (2015); to this extent, we recall that  $U$  is taken to be formulated in a Hilbert-style calculus in which Modus Ponens is the only rule of inference. It should be noticed, however, that the result is not dependent on the choice of a specific logical calculus: this

<sup>13</sup>This fact has been brought to my attention by Richard Heck (cf. Heck (2015)), although already well-known in other contexts.

only simplifies the description. Thus we omit details on the specific formulation of  $U$  in the statements of the results.

To simplify the notation, from now on we set:

$$\begin{aligned} \mathsf{T}^+[U] &:= \mathsf{T}[U] + \mathsf{AxT}_U \\ \mathsf{Tp}^+[U] &:= \mathsf{Tp}[U] + \mathsf{AxT}_U \end{aligned}$$

**Lemma 9.**  $\mathsf{T}[U]^+$  proves the consistency of  $U$  on a cut.

*Proof.* We recall that  $\mathsf{Prf}_U(k)$  is taken to be a  $\Delta_1^b$ -formula in  $\mathsf{S}_2^1$  expressing that  $k$  is the code of a  $U$ -proof, and  $\mathsf{lst}(k)$  a  $\Sigma_1^b$  function yielding, when applied to a sequence  $k$ , the last element of  $k$ .

We break down the proof in several steps.

We first consider the following formula with only  $n$  free.

$$(13) \quad \mathcal{K}(n) : \leftrightarrow (\forall m \leq n) (\mathsf{Prf}_U(m) \rightarrow \forall a \mathsf{Sat}(a, \mathsf{lst}(m)))$$

It expresses that any proof smaller than  $n$  has a true conclusion. We show that it is inductive. The claim clearly holds for  $n = 0$ , by assumption. Assuming it holds for arbitrary proofs  $k, l \leq n$  and sequence  $a$ , we want to prove it for proofs  $m \leq n + 1$ . In the interesting case, we have

$$\mathsf{lst}(m) = \mathsf{mp}(\mathsf{lst}(k), \mathsf{lst}(l)),$$

where  $\mathsf{mp}(m, n)$  is a  $\Sigma_1^b$ -function yielding the result of applying Modus Ponens to  $m$  and  $n$ . In particular, here we require  $\mathsf{lst}(k)$  of the form  $\mathsf{lst}(l) \rightarrow \mathsf{lst}(m)$ . The claim is thus obtained by applying (t $\rightarrow$ ) and (tv) to  $\mathsf{Sat}(a, \mathsf{lst}(k))$  and combining it with  $\mathsf{Sat}(a, \mathsf{lst}(l))$ .

We shorten  $\mathcal{K}$  to a cut  $\mathcal{J}$  such that, for formulas  $j$ , we have

$$\exists m (\mathcal{J}(m) \wedge \mathsf{Proof}_U(m, j)) \rightarrow \forall a \mathsf{Sat}(a, j)$$

We can safely assume that  $\mathsf{S}_2^1$  is available in  $\mathcal{J}$ .

Assuming the monotonicity of the coding — that is codes of elements of sequences are smaller of the code of the whole sequence — if  $m \in \mathcal{J}$ , also  $j$  will be. Finally we reason in the standard way: let  $\perp$  code, in the  $\mathsf{S}_2^1$ -numbers, an absurdity in the sense of  $U$ . If  $\mathsf{Pr}_U^{\mathcal{J}}(\perp)$ , then  $\mathsf{Sat}(a, \perp)$ . By Lemma 8, we obtain  $\neg \mathsf{Pr}_U^{\mathcal{J}}(\perp)$ , that is  $\neg \exists m (\mathcal{J}(m) \wedge \mathsf{Proof}_U(m, \perp))$ , that is  $\mathsf{Con}^{\mathcal{J}}(U)$ . □

By Lemma 8, the parametrized Tarski-biconditionals are provable in  $\mathsf{Tp}[U]$ . To prove that all logical axioms of  $A$  are true on a cut, we need extra-care. In a word, we have to make sure that syntax and truth interact well on this cut. The problem is essentially that, in order to prove the truth of a logical axiom schema, say  $\forall v_o \varphi \rightarrow \varphi(y/v_o)$ , one has to employ an instance of the very same principle, e.g

$$(14) \quad \forall m \mathsf{Sat}(a, k(m/\ulcorner v_o \urcorner)) \rightarrow \mathsf{Sat}(a, k(t/\ulcorner v_o \urcorner))$$

to be able to apply (t $\rightarrow$ ), (tv), (tq) and obtain

$$(15) \quad \mathsf{Sat}(a, \ulcorner \forall v_o \varphi \rightarrow \varphi(y/v_o) \urcorner)$$

The step from (14) to (15) requires, for instance, that variable assignments and formal substitution behave as required under the scope of  $\mathsf{Sat}$ . Since these lemmata are usually proved by induction, we need to resort to a definable cut in which this behaviour is preserved. For

details, we refer to Heck (2015) and Nicolai (2015). Once this is done, we have (recall that  $A, B, C, \dots$  range over finitely axiomatized theories)

**Corollary 2.**  $\mathsf{T}[A]$  proves the consistency of  $A$  on a cut.

Surprisingly enough, a similar results can be transferred to  $\mathsf{Tp}^+[U]$ . We have already seen that CT and PT, unlike their variants featuring arithmetical induction only, are extensionally identical. Induction extended to the truth predicate is in fact needed to prove the full axiom for negation. The technology of shortening of cuts enables us to mimic the role played by the extended induction schema of PT in deriving (CT2).

**Lemma 10.** There is a  $\mathsf{Tp}[U]$ -definable cut in which  $(\mathsf{t}\neg)$  holds.

*Proof.* We reason in  $\mathsf{Tp}[U]$ . Let  $\mathsf{lc}(x)$  a  $\Sigma_1^b$ -function that, applied to a  $\mathcal{L}_U$ -formula  $\varphi$ , yields the number of its logical symbols. By Lemma 1 the claim is obtained once we show that the following is inductive:

$$\mathcal{K}_o(n) := \forall k \forall a (\mathsf{Fml}_{\mathcal{L}_U}(k) \wedge \mathsf{lc}(k) \leq n \rightarrow (\mathsf{Sat}(a, \neg k) \leftrightarrow \neg \mathsf{Sat}(a, k)))$$

If  $n = o$ , we use (tat) and (fat). Assuming the claim is true for formulas  $k$  with  $\mathsf{lc}(k) \leq n$ , we obtain the claim by the truth/falsity clauses for connectives and quantifiers and the ‘inductive’ assumption. For instance, if  $k = \neg k_o$ , we have:

$$\begin{aligned} \neg \mathsf{Sat}(a, \neg k_o) &\leftrightarrow \neg \neg \mathsf{Sat}(a, k_o) && \text{by assumption} \\ &\leftrightarrow \mathsf{Sat}(a, k_o) \\ &\leftrightarrow \mathsf{Sat}(a, \neg \neg k_o) && (\mathsf{tdn}) \end{aligned}$$

□

**Corollary 3.**  $\mathsf{Tp}[U]$  interprets  $\mathsf{T}[U]$  on a definable cut.

More generally, Lemma 10 gives a strategy to interpret  $\mathsf{CT}\uparrow$  into  $\mathsf{PT}\uparrow$ , the standard systems of Tarskian and positive truth considered in §3.

Once we have  $(\mathsf{t}\neg)$  available on suitable initial segment  $\mathcal{K}_o$ , we can run the proof of Lemma 9 paying attention to work in  $\mathsf{Tp}^+[U]$ -definable cuts that shorten  $\mathcal{K}_o$  and in which the axioms of  $\mathsf{S}_2^1$  are satisfied. We thus have:

**Proposition 7.**  $\mathsf{Tp}[U]^+$  proves the consistency of  $U$  on a cut.

*Proof.* Instead of  $\mathcal{K}(n)$  in (13), we can start with the formula

$$(16) \quad \mathcal{K}_1(n) := (\forall m \leq n) (\mathsf{Prf}_U(m) \wedge \mathcal{K}_o(\mathsf{lst}(m)) \rightarrow \forall a \mathsf{Sat}(a, \mathsf{lst}(m)))$$

By the same argument used in the case of  $\mathcal{K}(n)$  in the proof of Lemma 9, it can be shown that  $\mathcal{K}_1(n)$  is inductive (only here the properties of  $\mathcal{K}_o$  are used in the inductive step). The proof then proceeds unchanged. □

By Lemma 8, and taking care of truth of all logical axioms of  $A$  on a suitable definable cut as suggested above, we have

**Corollary 4.**  $\mathsf{Tp}[A]$  proves the consistency of  $A$  on a cut.

Finally, we state the desired property of the result of applying  $\mathsf{T}^+[\cdot]$  and  $\mathsf{Tp}^+[\cdot]$  to arbitrary base theories:

**Corollary 5.**

- (i)  $T^+[U]$  and  $\text{Tp}^+[U]$  interpret  $S_2^1 + \text{Con}(U)$ , and thus  $Q + \text{Con}(U)$ ;
- (ii)  $T[A]$  and  $\text{Tp}[A]$  interpret  $S_2^1 + \text{Con}(A)$ , and thus  $Q + \text{Con}(A)$ .

*Proof.* Ad (i), we simply take the cut  $\mathcal{J}$  and the cut resulting from  $\mathcal{K}_1$  from Proposition 7 as domains of our interpretations. We notice that, by Lemma 2,  $S_2^1$  holds in them. (ii) is a weaker claim than (i).  $\square$

**4.3. Truth from Consistency.** We now turn to the question of how to recover the semantic machinery of  $T[\cdot]$  and  $\text{Tp}[\cdot]$  via the assumption of the consistency of the base theory. Feferman offered in Feferman (1960) a full formalization of Henkin's term model construction for a first-order, recursive, set of sentences  $S$  in extensions of PA plus a  $\Pi_1$  assertion of the consistency of  $S$ . Following Visser (1991) and Visser (2009b), we refer to the resulting arithmetized model as the 'Henkin-Feferman construction'. The technology of definable cuts enables us to formalize this term model in very weak subsystems of arithmetic complemented with the consistency statement for  $S$ .

**Proposition 8.**  $S_2^1 + \text{Con}(U)$  interprets  $U$ .

*Proof Sketch.* A full proof can be found in Visser (1991). As usual, we start with  $\mathcal{L}_U$  and extend it to  $\mathcal{L}_U^h$  by adding a constant  $c_{\exists v_i \varphi}$  for every  $\exists v_i \varphi$  in the extended language: notice that we have enough induction in  $S_2^1 + \text{Con}(U)$  to formalize this inductive argument. The Henkin theory, from which the arithmetized model for  $U$  is read off, is also defined in stages:

$$U_0^h := U$$

$$U_{n+1}^h := \begin{cases} U_n^h \cup \{\varphi\}, & \text{if } n = \ulcorner \varphi \urcorner \text{ (with } \varphi \in \mathcal{L}_U^h \text{) and } \text{Con}(U_n^h \cup \{\varphi\}), \\ U_n^h \cup \{\varphi\} \cup \{\psi(c_\varphi)\}, & \text{if } \varphi \text{ is } \exists v \psi \text{ and } \text{Con}(U_n^h \cup \{\varphi\}) \\ U_n^h & \text{otherwise} \end{cases}$$

Unlike the definition of  $\mathcal{L}_U^h, \Sigma_1^b$ -induction is not sufficient to prove the existence of the union of all  $U_n^h$ 's. However, it turns out that the formula defining the  $U_n^h$ 's is inductive, so it can be shortened to a cut  $\mathcal{I}$ . The required, complete theory in  $\mathcal{L}_U^h$ , for sentences in  $\text{Sent}_{\mathcal{L}_U^h}^{\mathcal{I}}$  can thus be taken to be  $\mathcal{F} = \bigcup_{n \in \mathcal{I}} U_n^h$ .

From  $\mathcal{F}$ , we define the required interpretation  $\mathfrak{F}$ , whose domain  $\delta_{\mathfrak{F}}$  is the set of codes of Henkin constants of  $\mathcal{L}_U^h$  in  $\mathcal{I}$ . To the theory  $\mathcal{F}$  there corresponds a truth predicate  $S$ . We set

$$(17) \quad R^{\mathfrak{F}}(v^0) := \leftrightarrow \delta_{\mathfrak{F}}(x) \wedge R(x) \in S \quad \text{for any } R \in \mathcal{L}_U$$

We notice that in (17) and in what follows, we write  $\varphi(x) \in S$  to express that the result of formally substituting  $v_1$  for  $x$  in  $\varphi$  falls into the extension of  $S$ . Crucially,  $S_2^1 + \text{Con}(U)$  proves:

- (a) for  $\varphi \in \text{Sent}_{\mathcal{L}_U^h}^{\mathcal{I}}$ ,  $(S^{\ulcorner \neg \varphi \urcorner}) \leftrightarrow \neg(S^{\ulcorner \varphi \urcorner})$ ;
- (b) for  $\varphi, \psi \in \text{Sent}_{\mathcal{L}_U^h}^{\mathcal{I}}$ ,  $S^{\ulcorner \varphi \vee \psi \urcorner} \leftrightarrow (S^{\ulcorner \varphi \urcorner} \vee S^{\ulcorner \psi \urcorner})$ ;
- (c) for  $\exists v \varphi(v) \in \text{Sent}_{\mathcal{L}_U^h}^{\mathcal{I}}$  and  $v^0 \in \text{Var}^{\mathcal{I}}$ ,  $S^{\ulcorner \exists v^0 \varphi \urcorner} \leftrightarrow \exists x (\delta_{\mathfrak{F}}(x) \wedge (\varphi(x) \in S))$ ;
- (d) for all  $\varphi \in \text{Sent}_{\mathcal{L}_U^h}^{\mathcal{I}}$ ,  $(\text{Pr}_U(\ulcorner \varphi \urcorner) \rightarrow S^{\ulcorner \varphi \urcorner})$

Notice, in (d), that reflection holds for sentences of the *non extended language*  $\mathcal{L}_U$  in the cut  $\mathcal{I}$ .  $\square$

By the interpretability of  $S_2^1 + \text{Con}(U)$  in  $\text{Q} + \text{Con}(U)$  on a definable cut, we have:

**Corollary 6.**  $\text{Q} + \text{Con}(U)$  interprets  $U$ .

Now we show that the truth predicate  $S(\cdot)$  defined in the proof of Proposition 8 enables us to interpret the satisfaction predicate of  $\text{T}^+[U]$ , and thus of  $\text{Tp}^+[U]$ . Therefore we obtain Tarskian and positive, typed truth for the base theory from the assertion of the consistency of  $U$ .

**Proposition 9.**  $S_2^1 + \text{Con}(U)$  interprets  $\text{T}^+[U]$ .

*Proof.* We define the translation  $\mathfrak{H}$ . It is assume a suitable renaming of bound variables omitted for readability;  $S$  is the truth predicate defined in the proof of Proposition 8.

- (18)  $(R)^{\mathfrak{H}}(u) :\leftrightarrow R(u) \in S$  for predicates  $R$  of sort  $\mathfrak{o}$
- (19)  $(P)^{\mathfrak{H}}(u) :\leftrightarrow P(u)$  for predicates  $P$  of sort  $\mathfrak{s}$ ;
- (20)  $(\text{Sat})^{\mathfrak{H}}(u, v) :\leftrightarrow S(\text{sb}(v, u))$ ; where  $\text{sb}$  substitutes elements of  $u$  in  $v$
- (21)  $((\cdot)^{\mathfrak{H}})^{\mathfrak{H}}(u, v, w) :\leftrightarrow ((u)_v = w \wedge v < \text{lh}(u)) \vee (v \geq \text{lh}(u) \wedge w = \mathfrak{o})$
- (22)  $(\exists x A)^{\mathfrak{H}} :\leftrightarrow \exists x (\delta_{\mathfrak{F}}(x) \wedge A)$  where  $\delta_{\mathfrak{F}}$  is as in Prop. 8;
- (23)  $(\exists k A)^{\mathfrak{H}} :\leftrightarrow \exists k (\mathcal{I}(k) \wedge A)$  with  $\mathcal{I}$  again as in Prop. 8
- (24)  $(\exists a A)^{\mathfrak{H}} :\leftrightarrow \exists a, a$  is a (finite) sequence and  $(\forall y \in a)(\delta_{\mathfrak{H}}(y))$

In (21),  $(x)_y$  is an efficient version of the  $\beta$  function that outputs the  $y^{\text{th}}$  element of the finite sequence  $x$  (cf. ?). In (22)-(24) we have kept the suggestive notation concerning variables even though, strictly speaking, we are not considering the language  $\mathcal{L}_T$  anymore.

A full verification that  $\mathfrak{H}$  is indeed a relative interpretation can be found in Nicolai (2015). Crucially, for  $\text{AxT}_U$ , one notices that, by Prop. 8(d),

$$(25) \quad \forall u (\text{Sent}_{\mathcal{L}_U}(u) \wedge \mathcal{I}(u) \wedge \text{Ax}_U(u) \rightarrow S(u))$$

Therefore, it suffices to notice that in  $S_2^1 + \text{Con}(U)$

$$(26) \quad \text{for all } u \in \text{Sent}_{\mathcal{L}_U}^{\mathcal{I}}, \text{Ax}_U^{\mathfrak{H}}(u) \rightarrow \text{Ax}_U(u).$$

□

**Corollary 7.**

- (i)  $S_2^1 + \text{Con}(U)$  interprets  $\text{Tp}^+[U]$ ;
- (ii)  $\text{Q} + \text{Con}(U)$  interprets  $\text{T}^+[U]$ , and thus  $\text{Tp}^+[U]$ .

*Proof.* (i) is immediate from the previous Proposition. (ii) follows from Lemma 3. □

Corollary 7 clearly holds for  $\text{T}[A]$  and  $\text{Tp}[A]$ .

**4.4. The functors  $\text{PC}[\cdot]$  and  $\text{PCS}[\cdot]$ .** The facts reported in this subsection — except the connections between these facts and the results above — are essentially due to Visser (2009b). We show that, given an arbitrary theory  $U$ ,

The functor  $\text{PC}[\cdot]$  adds predicative comprehension to an arbitrary  $U$ . In order to be completely general, as our overall program requires, we need to be able to act on arbitrary  $U$ . To mention a famous example, we should be able to generalize the transformation of PA into  $\text{ACA}_0$  to arbitrary theories. This can be done in the following way.



We first render the language two-sorted: objects over which quantifiers of  $\mathcal{L}_U$  range are taken to be of sort  $\sigma$ . We add to the language variables  $X$  of sort  $\mathfrak{c}$ , standing for sets (or concepts, or unary predicates) of objects of sort  $\sigma$  and a binary predicate  $A$  of type  $(\mathfrak{c}, \sigma)$  such that  $Xx := \leftrightarrow A(X, x)$ . To obtain  $\text{PC}[U]$ , the axioms of  $U$  formulated in the new language are extended with the scheme of predicative comprehension that, in a pedantic version, reads

$$(ca) \quad \exists u^{\mathfrak{c}} \forall v^{\sigma} (A(u^{\mathfrak{c}}, v^{\sigma}) \leftrightarrow \varphi(v^{\sigma}, \vec{w}^{\sigma}, \vec{y}^{\mathfrak{c}}))$$

If  $U$  is sequential, then  $\text{PC}[U]$  can be finitely axiomatized (Visser (2009b, Thm 3.4)).

In order to describe the next results in full generality, a slight modification of  $\text{PC}[\cdot]$  has to be considered. The required parallelism between consistency and predicative comprehension will break down if we considered arbitrary representations of the axiom set of  $U$ . A powerful theorem of Vaught states that any RE, sequential theory  $U$  can be axiomatized by a scheme. This scheme can be taken to be of the form  $\Psi(\vec{B})$ , where the  $B$ 's are *schematic variables* for formulas of  $\mathcal{L}_U$ . Vaught's result states that  $U$  and the theory  $U_{\forall}$  resulting from the process just described are extensionally identical (cf. Vaught (1967)).

To define the theory  $\text{PCS}[U]$ , we formulate  $U_{\forall}$  in a language with sorts in  $\{\sigma, \mathfrak{c}\}$  and translate  $\Psi(\vec{B})$  into  $\Psi(\vec{X})$ : we call this 'translation' of  $U_{\forall}$  as given by  $\Psi(\vec{X})$  the *v-class form* of  $U$  and denote it with  $U^{\mathfrak{c}}$ . To obtain  $\text{PCS}[U]$ , one then adds to  $U^{\mathfrak{c}}$  both (ca) and the universal closure of  $\Psi(\vec{X})$ . A modification of the strategy required in Lemma 1 shows that we can find a definable  $\text{PCS}[U]$ -cut on which there are no proofs of  $U$ -contradictions. In other words:

**Lemma 11.** Let  $U_{\forall}$  be sequential. Then  $\text{PCS}[U]$  proves the consistency of  $U_{\forall}$  on a cut.

*Proof Sketch.* We recall the main steps of Theorem 7.1 in Visser (2009b), and reason in  $\text{PCS}[U]$ . Let  $\text{lc}(\cdot)$  be as above. By the sequentiality of  $U_{\forall}$ , syntax is formalized in a standard way (e.g. by the natural numbers). We consider a reformulation modification of the definition of  $\mathcal{T}(X, x)$  in Proposition 5, which we call  $\mathcal{T}^*(X, x)$ .

Now in Lemma 5 we were able to use (ca) and the induction schema of ACA that the uniqueness condition for truth sets. In  $\text{PCS}[U]$  we have no such induction. We have to use again the shortening techniques. In particular, we can show that the (pseudo-) class

$$(27) \quad \mathcal{M} := \{x \mid \exists! X \mathcal{T}^*(X, x)\},$$

that is the class of all numbers which bound the complexity of formulas occurring in our truth sets, is indeed inductive. For any  $x$ , this unique  $X$  is denoted with  $X^x$ . We shorten  $\mathcal{M}$  to a cut  $\mathcal{N}$ .

A satisfaction predicate is now extracted from  $\mathcal{N}$ .  $\text{Fml}_{\mathcal{L}_U}^{1, \mathcal{N}}$  is the class of formulas of  $\mathcal{L}_U$ , lying in  $\mathcal{N}$ , with one free variable. It is

$$(28) \quad \text{St}(s, \ulcorner \varphi \urcorner) := \leftrightarrow (\exists x \in \mathcal{N})((s, \ulcorner \varphi \urcorner) \in X^x)$$

$\text{St}$  enjoys the nice property that, for any  $\varphi$  in  $\text{Fml}_{\mathcal{L}_U}^{1, \mathcal{J}}$  and assignment  $s$ , one can find —again by (ca)— a  $Y_{\varphi}$  such that, for all  $y$  and  $v_i$ ,

$$(29) \quad Y_{\varphi} y \text{ if and only if } \text{St}(s(y/v_i), \ulcorner \varphi \urcorner)$$

In fact this  $Y_{\varphi}$  can simply be taken to be the set of  $y$  that satisfy  $\varphi$  in the sense of  $X^{\text{lc}(\varphi)}$ . Since  $\varphi \text{Fml}_{\mathcal{L}_U}^{1, \mathcal{J}}$  we also know that  $Y_{\varphi}$  exists.

By employing an argument similar to the one employed in Lemma 9, we consider the set of  $U_V$ -proofs  $\mathcal{P}$  whose members  $\pi \in \mathcal{P}$  are such that  $\text{lh}(\pi) \leq x$ ,  $\forall i \leq \text{lh}(\pi)$ ,  $(\pi)_i \in \text{Fml}_{\mathcal{L}_U}^{1, \mathcal{N}}$  and such that  $\forall s \text{ St}(s, \text{lst}(\pi))$ . The set  $\mathcal{A}$  of such  $x$  is inductive: crucially, for the axiom  $\Psi(\vec{\psi})$  of  $U_V$  with  $\vec{\psi} \in \text{Fml}_{\mathcal{L}_U}^{1, \mathcal{N}}$ , we require

$$(30) \quad \text{St}(s(y/v_i), \ulcorner \Psi(\vec{\psi}) \urcorner)$$

By definition of  $\text{PCS}[U]$ , we have  $\forall \vec{Y} \Psi(\vec{Y})$  and thus  $\Psi(\vec{\psi})$  for an arbitrary  $\psi \in \text{Fml}_{\mathcal{L}_U}^{1, \mathcal{N}}$ . By (29), we associate to each formula in  $\vec{\psi}$  the appropriate set  $Y_{\psi_i}$ . Thus  $\Psi(\vec{Y}_{\psi})$ . By the properties of  $\text{St}$ , we obtain (30). A similar argument holds for the logical axiom schemata. The inductive step is unproblematic and it is similar to the inductive step in the proof of Lemma 9.

We thus shorten  $\mathcal{A}$  to a cut  $\mathcal{A}_o \subset \mathcal{A}$ .

It follows that, as in Lemma 9, no proof of contradiction can be reached in  $\mathcal{A}_o$ : i.e. we have  $\text{Con}^{\mathcal{A}_o}(U_V)$ .  $\square$

**Corollary 8.** Let  $U_V$  be sequential. Then  $\text{PCS}[U]$  interprets  $S_2^1 + \text{Con}(U_V)$ .

We now show, with the help of Proposition 8 that also the converse direction holds. In other words, that we can obtain predicative comprehension from consistency.

**Lemma 12.**  $S_2^1 + \text{Con}(U_V)$  interprets  $\text{PCS}[U]$ .

*Proof.* We first define the domain(s) of our interpretation  $\mathcal{I}$ . The domain of quantifiers ranging over variables of sort  $\mathfrak{o}$ ,  $\delta_{\mathcal{I}}$ , is just  $\delta_{\mathfrak{S}}$  (cf. Proposition 8), that is the domain of Henkin constant in a  $S_2^1 + \text{Con}(U)$ -definable initial segment of the natural numbers. The domain of class variables, with  $\mathcal{L}_U^c$  denoting again the extension of  $\mathcal{L}_U$  with a countable set of Henkin constants, is

$$\delta_{\mathcal{I}}^c := \{x \mid x \in \text{Fml}_{\mathcal{L}_U^c}^1 \cap \mathcal{I} \text{ and the free variable of } x \text{ is of sort } \mathfrak{o}\}$$

We recall that  $\mathcal{I}$  is the cut defined in the proof of Proposition 8. In the official definition of the language of  $\text{PCS}[U]$ ,  $Xx$  reads  $A(X, x)$ . Therefore we translate:

$$(A(X, x))^{\mathcal{I}} := \text{S}(\ulcorner \varphi \urcorner(x/v)), \quad \text{if } X \text{ is mapped into } \varphi(v)$$

Unsurprisingly, the clause just states that  $Xx$  is translated as ‘ $\varphi$  is satisfied by  $x$ ’. We call  $\varphi_X$  the element of  $\delta_{\mathcal{I}}^c$  associated to  $X$  by the translation.

Crucially,  $S_2^1 + \text{Con}(U)$  proves, by meta-induction on the complexity of  $\Phi(\vec{Y})$ ,<sup>14</sup> and by the crucial contribution of Proposition 8(d):

$$(31) \quad (\forall \vec{Y} \in \delta_{\mathcal{I}}^c) ((\Phi(\vec{Y}))^{\mathcal{I}} \leftrightarrow \text{S}(\ulcorner \Phi(\vec{\varphi}_Y) \urcorner))$$

The base case is basically given by the definition of the translation. The other cases are easy by employing the properties of  $\text{S}$ .

To verify that  $\mathcal{I}$  is indeed a relative interpretation, we prove in  $S_2^1 + \text{Con}(U)$  the translation of the axioms of  $\text{PCS}[U]$ : in particular (i) axioms of the form  $\forall \vec{Y} \Phi(\vec{Y})$  and (ii) all instances of comprehension.

(i) for  $\Phi(\vec{X})$  with  $\vec{X}$  arbitrary, by Lemma 8 we know that  $\text{Pr}_{U_V}(\ulcorner \Phi(\vec{\varphi}_X) \urcorner)$ . Thus by reflection (Lemma 8(d)),  $\text{S}(\ulcorner \Phi(\vec{\varphi}_X) \urcorner)$ . By (31),  $(\Phi(\vec{X}))^{\mathcal{I}}$ . By the properties of interpretations, since  $\vec{X}$  is arbitrary,  $(\forall \vec{X} \Phi(\vec{X}))^{\mathcal{I}}$ .

<sup>14</sup> $\Phi(\vec{Y})$  does not contain any  $\mathfrak{c}$ -quantifiers.

Ad (ii), to any  $\varphi(u, \vec{Z})$  with no bound set variables, we associate to it  $\varphi_X(v_i, \vec{\varphi}_Z)$  by the definition of  $\mathcal{J}$ . Since  $\varphi$  is a standard formula, we know  $\varphi_X(v_i, \vec{\varphi}_Z)$  will be in  $\delta_5^x$ . But  $\varphi_X(u, \vec{\varphi}_Z)$  is just  $(Xu)^{\vec{\mathcal{J}}}$ , and it is equivalent to  $S(\varphi_X(u, \vec{\varphi}_Z))$  by the properties of  $\mathcal{J}$ . By (31), we conclude  $(\varphi(u, \vec{Z}))^{\vec{\mathcal{J}}}$ . This shows that, given any  $\varphi(u, \vec{Z})$ , the translation of the comprehension axiom of  $\text{PCS}[U]$  is satisfied.  $\square$

Thus, by Lemma 3,

**Corollary 9.** Let  $U$  be as above.  $\text{Q} + \text{Con}(U_V)$  interprets  $\text{PCS}[U]$ .

**4.5. Truth and Comprehension, via Consistency.** Finally we can combine the claims introduced in the last two sections.

We start with predicative comprehension. By Lemmata 11 and 12, we have

**Corollary 10.** Let  $U_V$  be sequential. Then  $\text{PCS}[U]$  is mutually interpretable with  $\text{Q} + \text{Con}(U_V)$ .

**Corollary 11.** If  $U$  is finitely axiomatized and sequential, then  $\text{PC}[U]$  is mutually interpretable with  $\text{Q} + \text{Con}(U)$ .

*Proof.* If  $U$  is finitely axiomatized, we don't need to resort to the schematic version of  $U$  to prove Lemmata 11 and 12.  $\square$

Now we move to Tarskian and positive typed truth. By Proposition 9, Corollary 7, and Corollary 5,

**Proposition 10.**

- (i) Let  $U$  be finitely axiomatized, then  $\text{T}[U]$  and  $\text{Tp}[U]$  are mutually interpretable with  $\text{Q} + \text{Con}(U)$ ;
- (ii) Let  $U$  be arbitrary, then  $\text{T}^+[U]$  and  $\text{Tp}^+[U]$  are mutually interpretable with  $\text{Q} + \text{Con}(U)$ .

Predicative comprehension and Tarskian truth, via consistency, enjoy a mutual reduction for a wide range of choices of the object theory:

**Proposition 11.**

- (i) Let  $U$  be finitely axiomatized and sequential. Then  $\text{PC}[U]$  is mutually interpretable with  $\text{T}[U]$  and  $\text{Tp}[U]$ ;
- (ii) Let  $U_V$  be sequential. Then  $\text{PCS}[U]$  is mutually interpretable with  $\text{T}^+[U]$  and  $\text{Tp}^+[U]$ .

## 5. CONCLUSION

Proof-theoretic reductions, such as the relative interpretation of PA in ZF, have often been considered as examples of ontological reductions — of natural numbers to sets, in the case at issue. Ideally, we would require a reduction of ACA to CT to replace our commitment to arithmetically definable sets with a semantic commitment to a Tarskian, compositional truth predicate. It might be thought, for instance, that any plausible criterion of ontological reduction of what is implicit in the acceptance of  $T$  to what is implicit in the acceptance of  $W$  should at least require  $T$  to be relatively interpretable in  $W$ . In such scenario, the reductions investigated in this work — and summarized in §4.5 — seem to offer the possibility of freely moving from one's ontological commitment to subsets of the domain of the object theory to

ideological commitment to concepts giving more structure to but not enlarging the original domain, such as truth and satisfaction.

These claims may be contrasted by reflecting on the combinatorial nature of a relative interpretation, which appears to be primarily a syntactical reduction: there seems to be no need to resort to a world-language relation to make sense of the interpretation of ACA in CT, for instance. It might simply be taken to be no more than a complex relation between syntactically individuated entities, that is languages and theories. The significance of the results presented in this work, however, is largely independent from one's stance on the relationships between proof-theoretic and ontological reductions. Even from a more neutral stance our study seems to reinforce the belief that set existence principles and principles governing primitive predicates for truth and satisfaction are deeply intertwined and that a general criterion of theory choice — once we accept the base theory  $U$ , how much comprehension should we add to it? Which truth axioms for it are preferable? — should consider them as interdependent. The analysis of the operations on theories  $T^+[\cdot]$ ,  $TP^+[\cdot]$  and  $PCS[\cdot]$  in §4 generalizes in fact the standard approach surveyed in §3 in several respects and strengthens our conviction that predicative comprehension and suitable truth axioms are indeed parallel logico-mathematical devices.

One might in fact suspect that the reductions between ACA and CT rely on principles that have not much to do with truth and predicative comprehension: by considering for instance  $ACA_0$  and  $CT\uparrow$ , resulting from ACA and CT by restricting induction, the symmetry between truth and predicative comprehension suggested by the case with full induction is lost.  $ACA_0$  is not interpretable in  $CT\uparrow$ . Again part of the reason for this lies with the reflexivity of PA, which is a peculiarity of the base theory and has nothing to do with truth or with predicative comprehension.  $T^+[\cdot]$  and  $PCS[\cdot]$ , in this respect, fare reasonably better. Their construction is rooted in the idea that  $T^+[U]$  genuinely applies a syntactico/truth-theoretic package to the base theory, and that  $PCS[U]$  generalizes to arbitrary theories the method for obtaining  $ACA_0$  from PA: this approach tries to minimize and control peculiar features of the object theory, such as the reflexivity of PA, for instance. Furthermore,  $T^+[PA]$  and  $PCS[PA]$  are comparable, as Proposition 11 shows, whereas a similar characterization for  $CT\uparrow$ +'all axioms of PA are true' still seems to be missing.<sup>15</sup> This suggests that the operations on the theories considered here offer a more uniform, smoother comparison of semantic and set existence assumptions.<sup>16</sup>

Another attractive feature of the proposed results is that they represent different but equivalent ways to uncover the 'lower bound' of our implicit commitment to the acceptance of a base theory  $U$ . Corollaries 10-11 and Propositions 10-11 tell us that the application of  $T^+[\cdot]$ ,  $TP^+[\cdot]$  and  $PCS[\cdot]$  correspond to an intensional consistency statement for  $U$ , which is not reducible — neither provable nor interpretable — to  $U$  by the incompleteness phenomena. Considered in this instrumental fashion, positive, typed truth and Tarskian truth on one side and predicative comprehension on the other represent equivalent routes for making explicit our commitment to  $U$  via the assertion of its consistency.

<sup>15</sup>In particular, it is known that  $CT\uparrow$ +'all axioms of PA are true' interprets  $ACA_0$ , but the converse claim seems to be an open problem.

<sup>16</sup>As we have seen, more generality requires different means of reduction: conservativeness for suitable class of formulas (e.g.  $\Pi_2^0$  or formulas of the base language), for instance, becomes a trivial requirement as all the theories resulting from the application of our functors will be conservative over the base theory. Relative interpretability becomes the preferred means of reduction.

Further technical and philosophical developments are of course possible. On the philosophical side, a unified account of theory choice for systems of truth is much needed, given the explosion the research on axiomatic truth has had in recent years. From the technical standpoint, an obvious extension of the work would have to take into account hierarchies of applications of the operators: the difficulty of the hierarchical approach, in this context, seems to be represented by the fact that methods for a hierarchical analysis of  $T[\cdot]$ , for instance, are foreseeable only when a single object theory is fixed. This would lead to a weakening of the approach that draws much of its strength from its general applicability. At any rate this is only tentative and we defer a treatment of these extensions to forthcoming works.

## REFERENCES

- Buss, S. (1986), *Bounded Arithmetic*, Bibliopolis, Naples.
- Buss, S. (ed.) (1998). *Handbook of Proof Theory*. Elsevier.
- Cantini, A. (1989). Notes on Formal Theories of Truth. *Archives for Mathematical Logic* 35: 97–130.
- Enayat, A. and A. Visser. (2014). New Constructions of Satisfaction Classes. Logic Preprint Groups Series.
- Enderton, H. (2001). *A Mathematical Introduction to Logic*. Harcourt Academic Press.
- Feferman, S. (1960), Arithmetization of metamathematics in a general setting, *Fundamenta Mathematicae*, 49, 35–91.
- Feferman, S. (1964). Systems of Predicative Analysis. *Journal of Symbolic Logic* 27: 1–30.
- Feferman, S. (1991). Reflecting on Incompleteness. *The Journal of Symbolic Logic* 56: 1–49.
- Fischer, M. (2009). Minimal Truth and Interpretability. *The Review of Symbolic Logic* 4(2), 2009: 799–815.
- Fujimoto, K. (2010). Relative Truth Definability of Axiomatic Truth Theories. *The Bulletin of Symbolic Logic* 16(3): 305–344.
- Hájek, P. and P. Pudlák (1993). *Metamathematics of First-Order Arithmetic*. Springer.
- Halbach, V. (1999). Conservative Theories of Classical Truth. *Studia Logica* 62: 353–370.
- Halbach, V. (2009). Reducing Compositional to Disquotational Truth. *Review of Symbolic Logic* 2: 786–798.
- Halbach, V. (2014). *Axiomatic Theories of Truth (Revised Edition)*. Cambridge University Press 2014.
- Heck, R. (2015). Consistency and the Theory of Truth. to appear in *The Review of Symbolic Logic*.
- Lachlan, A. (1981). ‘Full Satisfaction Classes and Recursive Saturation’. *Canadian Mathematical Bulletin* 24: 295–297.
- Leigh, G. (2014). Conservativity for Theories of Compositional Truth via Cut-Elimination. To appear in *The Journal of Symbolic Logic*.
- Mostowski, A. (1950), Some impredicative definitions in the axiomatic set-theory. *Fundamentae Mathematicae* 37: 111–124.
- Moschovakis, Y. (1974). *Elementary Induction in Abstract Structures*. American Elsevier.
- Nicolai, C. (2015). A Note on Typed Truth and Consistency Assertions. *The Journal of Philosophical Logic* (online first version). DOI: 10.1007/s10992-015-9366-6.
- Nicolai, C. (2015a). Deflationary Truth and the Ontology of Expressions. *Synthese* (online first version), DOI: 10.1007/s11229-015-0729-x.
- Pohlers, W. (2009). *Proof Theory. The First Step into Impredicativity*. Springer.
- Takeuti, G. (1987). *Proof Theory*. North-Holland, Amsterdam.
- Vaught, R. A. (1967). Axiomatizability by a Schema. *The Journal of Symbolic Logic* 32: 473–479.
- Visser, A. (1991), ‘The Formalization of Interpretability’, *Studia Logica* 50(1): 81–106, 1991.
- Visser, A. (2009), ‘Can we make the second Incompleteness Theorem Coordinate Free?’, *Journal of Logic and Computation* 21(4), pp.: 543–560.
- Visser, A. (2009). The predicative Frege hierarchy. *Annals of Pure and Applied Logic* 160: 129–153.

Visser, A. (2012). Vaught's Theorem on Axiomatizability by a Scheme. *Bulletin of Symbolic Logic* 18: 382–402.