



## Reference and Truth

Lavinia Picollo<sup>1</sup> 

Received: 19 July 2019 / Accepted: 30 July 2019 / Published online: 29 August 2019  
© The Author(s) 2019

### Abstract

I apply the notions of alethic reference introduced in previous work in the construction of several classical semantic truth theories. Furthermore, I provide proof-theoretic versions of those notions and use them to formulate axiomatic disquotational truth systems over classical logic. Some of these systems are shown to be sound, proof-theoretically strong, and compare well to the most renowned systems in the literature.

**Keywords** Semantic paradoxes · Disquotation · Reference · Self-reference · Well-foundedness · Formal truth theories

### 1 Introduction

Tarski [32, 33] has bequeathed to us a pessimistic but valuable lesson about truth. Given a language  $L$  with a monadic predicate  $T$  and a name  $\ulcorner\varphi\urcorner$  for each sentence  $\varphi$  of the language, for  $T$  to adequately express the notion of truth of sentences of  $L$ , each sentence  $\varphi$  should be materially equivalent to  $T\ulcorner\varphi\urcorner$ . However, if the underlying logic is classical and the language is ‘expressive enough’, the equivalence must fail in some cases, on pain of triviality. For instance, if  $L$  allows for self-referential expressions and, in particular, contains a liar sentence  $\lambda$ , that says of itself that it is not true,  $\lambda$  and  $\neg T\ulcorner\lambda\urcorner$  will turn out to be equivalent. Thus, the equivalence between  $\lambda$  and  $T\ulcorner\lambda\urcorner$  is untenable. As a consequence, an adequate truth predicate for the whole language is not possible in  $L$ , if working within classical logic.

Tarski proposed to restrict the equivalence between each  $\varphi$  and  $T\ulcorner\varphi\urcorner$  (expressed as a biconditional) – the so-called T-schema – to the  $T$ -free sentences of the language. In this way,  $\lambda$  is easily excluded. Tarskian *typed* truth theories have been pervasive for many years and play a prominent role in diverse areas of logic and philosophy.

---

✉ Lavinia Picollo  
l.picollo@ucl.ac.uk

<sup>1</sup> Department of Philosophy, University College London, 19 Gordon Square, London WC1H 0AW, UK

Nonetheless, the shared view nowadays is that these theories are far too restrictive. After all, many sentences containing the truth predicate – e.g. expressions of the form  $T \ulcorner T \ulcorner \varphi \urcorner \urcorner$  where  $\varphi$  is  $T$ -free – seem to be entirely unproblematic.

Less drastic solutions have been subsequently explored. Theories which weaken the logic but keep some form of equivalence between each sentence and its truth ascription – originating in the work of Kripke [20] and Martin & Woodruff [23] – are very popular. Others have investigated the possibility of remaining classical but imposing less strict restrictions on the admissible instances of the T-schema. The goal is to find a unified criterion that allows us to keep as many unproblematic instances as possible, leaving the paradoxical ones out.

A promising strategy in this direction is to restrict the T-schema to sentences that exhibit ‘safe’ reference patterns. The orthodox view, championed by Russell, Poincaré, and Tarski, has it that the cause of semantic paradox is self-reference. Despite the popularity of this view, the restriction of the T-schema to non-self-referential sentences has not yet been explored. The main reason is that the notions of self-reference and reference have proved to be quite elusive in the past, and are now surrounded by an aura of scepticism.<sup>1</sup> Additionally, the self-reference diagnosis has been (relatively) recently challenged by a new paradox that is *prima facie* free of self-reference: the Visser-Yablo paradox.<sup>2</sup> However, the bearing of this new antinomy is not yet entirely clear: the lack of a precise notion of self-reference has hindered the evaluation of the Visser-Yablo paradox and thus the self-reference diagnosis of paradoxicality.

In [27] I provide a systematic and rigorous account of reference *in the context of truth* – or “alethic reference”, as I call it – and define self-reference and other reference patterns in terms of this notion. I hope the intuitive appeal of the definitions I put forward there helps dissipate the scepticism the corresponding notions are immersed in, at least with regard to the semantic paradoxes. Furthermore, I show that the expressions involved in the Yablo-Visser paradox are not self-referential according to the new notions, refuting the self-reference diagnosis on the roots of paradox. Although this ruins the prospects of restricting the T-schema to non-self-referential expressions, the more general project of restricting it to sentences that exhibit safe reference patterns need not be abandoned, as all paradoxical expressions might share other reference patterns.

The purpose of this paper is to confirm this hypothesis. I provide both semantic and axiomatic theories of truth in which the T-schema is restricted to *well-founded* expressions, and prove they are encompassing and sound. Moreover, I show this condition can be relaxed even further, i.e. that the reference patterns that underlie paradoxical expressions can be given a more fine-grained characterization. In short, I deploy the reference notions introduced in [27] to formulate different restricting criteria for admissible instances of T-schema in terms of their underlying reference patterns and show that this strategy is successful, as it results in classically consistent, sound, and fairly attractive theories of truth.

<sup>1</sup> See Picollo [26] for an overview of the state of the matter.

<sup>2</sup> See Herzberger [15], Visser [34], and Yablo [35, 36].

Section 2 provides a technical introduction to formal theories of truth, followed by an overview of the state of the issue regarding different restrictive criteria for the truth predicate that have been proposed in the literature. I also give a compact presentation of the approach to reference in the context of truth introduced in [27]. Section 3 employs these notions to construct semantic truth theories. Section 4 provides simpler, proof-theoretic versions of the notions of reference (and related concepts) from Section 2 and formulates several axiomatic truth systems based on them, some of which are shown to be sound and proof-theoretically strong. Finally, in Section 5 I conclude by evaluating these systems in light of Leitgeb's [22] criteria for formal theories of truth.

## 2 Preliminaries

### 2.1 Formal Truth Theories

Let  $\mathcal{L}$  be the language of first-order Peano arithmetic (PA) and let  $\mathcal{L}_T$  augment  $\mathcal{L}$  with a monadic predicate T, to express truth.  $\mathcal{L}$  contains  $=, \neg, \wedge, \vee, \rightarrow, \forall, \exists$  as logical constants, an individual constant 0, a monadic function symbol S, two dyadic function symbols  $+$  and  $\times$ , and a finite stock of extra function symbols for primitive recursive (p.r.) functions to be specified. All other logical and non-logical symbols occurring in formulae are to be understood as the usual abbreviations. Let  $\mathbb{N}$  be the standard model for  $\mathcal{L}$ , with  $\omega$  as its domain. Note that, for each  $n \in \omega$ ,  $\mathcal{L}$  has a term  $\bar{n}$  denoting  $n$  (the numeral of  $n$ ), that consists of  $n$  occurrences of S followed by 0.

We work with a fixed (effective and monotonic) coding of expressions of  $\mathcal{L}_T$  by numbers in  $\omega$ .<sup>3</sup> If  $\sigma$  is a string of symbols of  $\mathcal{L}_T$ , we write  $\ulcorner \sigma \urcorner$  for the numeral of its code. We often identify expressions of  $\mathcal{L}_T$  with their codes if there's no room for confusion. Unless otherwise indicated, by "formula" and "sentence" we mean formula of  $\mathcal{L}_T$  and sentence of  $\mathcal{L}_T$ , respectively.

Although  $\mathcal{L}$  speaks primarily about natural numbers, our coding allows it to express many syntactic properties, relations, and functions about the expressions of  $\mathcal{L}_T$ . Thus, our truth theories can be formulated in  $\mathcal{L}_T$ , with background syntactic principles formulated in  $\mathcal{L}$ .

Formal truth theories can be either semantic or axiomatic. Semantic theories consists of a model or family of models  $\langle \mathbb{N}, \Gamma \rangle$  expanding  $\mathbb{N}$  to a model of  $\mathcal{L}_T$ , where  $\Gamma$  is the extension of T in the model. In evaluating them, we look at the truth principles that hold in every model of the family. Epistemic considerations are also relevant, however: it is important that there is a way to know which truth principles belong to the theory, to the extent possible. If the theory is too complex, this would hardly be the case. This is not to say that complex semantic constructions are of no value. On the contrary, they can serve as witnesses of the consistency of collections of truth principles, or play a heuristic role in the formulation of more constructive theories, that is, axiomatic ones.<sup>4</sup>

<sup>3</sup>See my companion paper for more details about the coding.

<sup>4</sup>See Field [5] and Halbach & Horsten [11].

By contrast, axiomatic truth theories result from adding truth-theoretic axioms to a syntax theory, which we assume to be PA.<sup>5</sup> We assume PA contains the defining recursion equations for each extra function symbol in  $\mathcal{L}$ . As is well known, PA is strong enough to represent every recursive relation between numbers and, therefore, expressions of  $\mathcal{L}_T$ , and to weakly represent every recursively enumerable relation. Let PAT consist of the axioms of PA formulated in  $\mathcal{L}_T$  with induction for the whole language. Call an axiomatic truth theory in  $\mathcal{L}_T$  any recursive extension of PAT. Of course, some theories will be highly incomplete and others simply unsound, but the terminology is convenient. As before, the merits of a theory lie in the truth principles the theory entails. To know what principles hold in an axiomatic theory it is enough to find a proof.<sup>6</sup>

## 2.2 Tarski's Theorem (and ways to circumvent it)

Ideally, any truth theory (whether semantic or axiomatic) would satisfy Tarski's condition of material adequacy, according to which all instances of the following schema hold in the theory:

$$(T\text{-schema}) \quad T^{\ulcorner} \varphi^{\urcorner} \leftrightarrow \varphi$$

These are known as "T-biconditionals". Alternatively, we could work with a more general variant of the T-biconditionals, the *uniform* T-biconditionals:

$$(Uniform\ T\text{-schema}) \quad \forall t_1 \dots \forall t_n (T^{\ulcorner} \varphi(t_1, \dots, t_n)^{\urcorner} \leftrightarrow \varphi(t_1^{\circ}, \dots, t_n^{\circ}))$$

Both the T-schema and its uniform version are known as "disquotational" principles. Note that the set of closed terms of  $\mathcal{L}_T$  is p.r., as well as the value relation that holds between each term of the language and the number it denotes (as  $\mathcal{L}$  contains only finitely many function symbols). Let  $CTerm(x)$  and  $x^{\circ} = y$  represent in PA ("represent", from now on) this set and relation, respectively.<sup>7</sup> Moreover, the substitution function that takes a formula  $\varphi$ , a term  $t$ , and a variable  $v$  and returns the formula that results from replacing all free occurrences of  $v$  in  $\varphi$  with  $t$ , is also p.r. and, thus represented by a term  $x(y/z) \in \mathcal{L}$ . We write  $\forall t \varphi$  for  $\forall v (CTerm(v) \rightarrow \varphi)$  and  $\exists t \varphi$  for  $\exists v (CTerm(v) \wedge \varphi)$ , for a suitable variable  $v$ . Finally,  $\ulcorner \varphi(t)^{\urcorner}$  abbreviates  $\ulcorner \varphi^{\urcorner}(t/\ulcorner v^{\urcorner})$ , for some suitable term variable  $t$ , provided that  $v$  is the only free variable in  $\varphi$ . Thus, the instances of the Uniform T-schema quantify over closed terms, entailing all substitution instances of the standard T-schema *uniformly*.

Unfortunately, neither of these principles can be implemented unrestrictedly, as the language is 'expressive enough' to allow for paradoxical expressions such as liar sentences.

Let  $\mathbf{v}$  abbreviate the string of variables  $v_1, \dots, v_n$  different from  $x$  and  $y$ .

<sup>5</sup>Robinson Arithmetic, a weaker subsystem of PA, would also do. I choose PA instead to facilitate the comparison between our truth theories and other systems that are found in the literature.

<sup>6</sup>However, to know whether a schematic principle holds, the schema of a proof is required.

<sup>7</sup>Since  $\mathcal{L}$  contains enough function symbols for the proof of the Strong Diagonal Lemma (cf. Theorem 2) to go through, it cannot have a function symbol for the value relation (which is also a function), on pain of triviality. Nonetheless, I write  $x^{\circ} = y$  instead of, e.g.  $Val(x, y)$ , to preserve readability, as is customary.

**Theorem 1** (Diagonalization) *For every formula  $\varphi(x, \mathbf{v})$  there is a formula  $\psi(\mathbf{v})$  s.t. the (universal closure of the) following equivalence is a theorem of PAT:*<sup>8</sup>

$$\psi(\mathbf{v}) \leftrightarrow \varphi(\ulcorner \psi \urcorner, \mathbf{v}) \tag{1}$$

In equivalences of the form (1),  $\psi(\mathbf{v})$  is said to be a fixed point of  $\varphi(x, \mathbf{v})$ . Let  $\varphi$  in Theorem 1 be  $\neg T x$ . Then there is a sentence  $\lambda$  such that the following is a theorem of PAT:

$$\lambda \leftrightarrow \neg T \ulcorner \lambda \urcorner \tag{2}$$

$\lambda$  is normally understood as saying of itself that it is untrue, as a liar sentence. Given (2), no consistent extension of PAT can contain an instance of the T-schema for  $\lambda$  and, *a fortiori*, full disquotation is untenable. This is what Tarski’s undefinability result consists in. Likewise, if we opt for a semantic account instead, no model  $\langle \mathbb{N}, \Gamma \rangle$  of  $\mathcal{L}_T$  can validate unrestricted versions of disquotation, since all theorems of PAT are true in  $\langle \mathbb{N}, \Gamma \rangle$ , including (2). Thus, we say  $\lambda$  is paradoxical.

To avoid paradox, Tarski opted to restrict disquotation to sentences without T, but more permissive restrictions are possible without stepping into triviality. Thus, we wonder with Leitgeb [21, p. 156], “What kinds of sentences with truth predicate may be inserted plausibly and consistently into the T-schema?” An idea that suggests itself is to restrict our disquotational principles to non-paradoxical expressions, that is, to those that can be consistently inserted in the T-schema. Alas, McGee [24] has shown that maximal consistent sets of T-biconditionals decide each and every sentence of  $\mathcal{L}_T$ , which means they are far too complex for an axiomatization. Moreover, McGee’s result shows there are uncountably many of those sets, so picking one of them amounts to an arbitrary choice. Consider the following 2-liar cycle:

$$\begin{aligned} \lambda_1 &\leftrightarrow \neg T \ulcorner \lambda_2 \urcorner \\ \lambda_2 &\leftrightarrow T \ulcorner \lambda_1 \urcorner \end{aligned} \tag{3}$$

Theorem 1 guarantees that both biconditionals are provable in PAT if we diagonalize the predicate  $\neg T \ulcorner T x \urcorner$ . Given a formula  $\varphi$  with exactly one free variable  $v$ , let  $\ulcorner \varphi(v) \urcorner$  be short for  $\ulcorner \varphi \urcorner(\ulcorner v \urcorner / \ulcorner v \urcorner)$ , where  $\dot{x}$  is a term of  $\mathcal{L}$  for the p.r. function that maps each natural number to the code of its numeral. Since  $v$  is free in  $\ulcorner \varphi(v) \urcorner$ , we can quantify over it. Clearly, the biconditionals in (3) are inconsistent with the T-biconditionals for  $\lambda_1$  and  $\lambda_2$ . However, if one of these T-biconditionals is dropped, consistency is restored. Every maximal consistent set will contain one of them but not the other.

A more promising suggestion has been made by Leitgeb [21]. Roughly, his idea is to restrict the T-schema to grounded sentences, that is, sentences whose truth value ultimately depends on expressions not containing the truth predicate. Since the truth value of the liar and of sentences in liar cycles seems to depend, directly or indirectly, on these very same sentences, the latter are considered ungrounded and

<sup>8</sup>See Carnap [2] and Montague [25]. See Picollo [27] for a proof.

excluded from our disquotational principles. The truth value of other expressions such as  $0 = 0$ ,  $T^{\ulcorner}0 = 0^{\urcorner}$ , and  $T^{\ulcorner}\lambda^{\urcorner} \vee \neg T^{\ulcorner}\lambda^{\urcorner}$ , instead, is ultimately fixed by non-semantic facts, so their corresponding T-biconditionals hold in the theory. Leitgeb’s theory consists of a model  $\langle \mathbb{N}, \Phi_{\text{If}} \rangle$  of  $\mathcal{L}_T$  where only grounded sentences belong to  $\Phi_{\text{If}}$  and all instances of disquotation for grounded sentences are true. The theory is quite natural, elegant, and seemingly free from adhocness. However, it is also fairly complex.<sup>9</sup> This means that an  $\mathbb{N}$ -categorical axiomatization is not possible, i.e. the class of purely truth-theoretic principles the theory entails is not recursively axiomatizable.<sup>10</sup> Schindler [30] provides a nice axiomatic system, but it is naturally far from capturing the original semantic construction.

Others have put forward simpler, syntactic restrictions on disquotation. Halbach [9] explores the restriction of uniform disquotation to formulae in which T occurs only positively – i.e. in the scope of an even number of negations and conditional antecedents. The resulting system is known as PUTB, for “Positive Uniform Tarski Biconditionals”. Schindler [31], in turn, considers restricting the Uniform T-schema to translations of formulae of the language of second-order arithmetic without second-order parameters. Both criteria seem somewhat unnatural and ad hoc.

I would like to propose an alternative path, namely, to restrict disquotation to expressions that exhibit ‘safe’ reference patterns. According to the self-reference diagnosis, all paradoxical expressions share a common reference pattern, i.e. self-reference. This seems to be certainly the case of  $\lambda$  and the liar cycle given by  $\lambda_1$  and  $\lambda_2$ . If this hypothesis were correct, a sensible plan would be to restrict disquotation to non-self-referential expressions. The prior lack of adequate and precise notions of reference and self-reference has prevented us so far from exploring this route. Luckily, this situation has been remedied. In my companion paper [27] I give a systematic and formal account of reference in the context of truth, designed specifically for the study of the reference patterns underlying paradoxical sentences. Moreover, since according to my account reference has a syntactic aspect, restricting disquotation to non-self-referential expressions could turn out to be simple enough for the formulation of axiomatic truth systems.

Unfortunately, the notions I put forward in [27] reveal that the self-reference diagnosis is not correct. As it turns out, there are semantic paradoxes that are free of self-reference. Such is the case of the Visser-Yablo paradox, consisting of an infinite list of sentences, each of which says of the ones coming after that they are untrue. The existence of this list can be proved in PAT by Theorem 1. Diagonalizing the formula  $\forall z (z > w \rightarrow \neg T^{\ulcorner}x(\dot{z}/^{\ulcorner}w^{\urcorner})^{\urcorner})$ , we obtain a predicate  $Y(w)$  such that

$$\forall w (Y(w) \leftrightarrow \forall z (z > w \rightarrow \neg T^{\ulcorner}Y(\dot{z})^{\urcorner})) \tag{4}$$

<sup>9</sup>Leitgeb shows that  $\Phi_{\text{If}}$  is  $\Pi_1^1$ -complete, i.e. the simplest formula expressing this set begins with a string of second-order universal quantifiers followed by a formula of  $\mathcal{L}$ .

<sup>10</sup>The notion of  $\mathbb{N}$ -categoricity has been introduced and put forward as a reasonable criterion for the axiomatizability of semantic theories by Fischer et al. [6].

is provable in PAT. Instantiating  $w$  in each numeral results in the following biconditionals, i.e. the list:

$$\begin{aligned}
 Y(0) &\leftrightarrow \forall z (z > 0 \rightarrow \neg T^{\Gamma}Y(\dot{z})^{\neg}) \\
 Y(\bar{1}) &\leftrightarrow \forall z (z > \bar{1} \rightarrow \neg T^{\Gamma}Y(\dot{z})^{\neg}) \\
 &\dots \\
 Y(\bar{n}) &\leftrightarrow \forall z (z > \bar{n} \rightarrow \neg T^{\Gamma}Y(\dot{z})^{\neg}) \\
 &\dots
 \end{aligned}
 \tag{5}$$

By reductio ad absurdum, the T-biconditionals for  $Y(\bar{n})$  entail  $\neg T^{\Gamma}Y(\bar{n})^{\neg}$  for each  $n \in \omega$ , as well as  $\neg \forall z \neg T^{\Gamma}Y(\dot{z})^{\neg}$ . Nonetheless,  $\forall z \neg T^{\Gamma}Y(\dot{z})^{\neg}$  does not follow in PAT plus the T-biconditionals. This means that the theory is consistent, albeit  $\omega$ -inconsistent.<sup>11</sup> Note, however, that no model  $\langle \mathbb{N}, \Gamma \rangle$  of  $\mathcal{L}_T$  can make all T-biconditionals for each  $Y(\bar{n})$  true at the same time: since each  $\neg T^{\Gamma}Y(\bar{n})^{\neg}$  would have to be true in the model, so would  $\forall z \neg T^{\Gamma}Y(\dot{z})^{\neg}$ . For these reasons, the Visser-Yablo paradox is not considered to be a paradox in the strict sense, but an  $\omega$ -paradox. Despite not directly trivializing our axiomatic theories, it is still problematic, as no semantic truth theory can validate all T-biconditionals for the sentences in the list. Furthermore, although (4) is consistent with the T-biconditionals, an outright inconsistency can be obtained by combining (4) with an instance of *uniform* disquotation for  $Y(w)$ .

According to my account of alethic reference, no sentence in the Visser-Yablo list is self-referential but they are all *unfounded*, as will be seen in Section 2.3. This shows that restricting disquotation to non-self-referential expressions is not a viable project; banning self-reference is not enough. However, as will be shown in Section 3, there are other reference patterns shared by unparadoxical expressions that prove to be sensible restrictions on disquotation. Moreover, their relative simplicity will allow us in Section 4 to formulate proof-theoretic approximations that we then deploy in the formulation of sound and encompassing axiomatic theories of disquotational truth. But before we get to that, I will briefly introduce my general account of alethic reference and related notions.

### 2.3 Alethic reference

Before giving a formal definition of alethic reference, I will briefly discuss its motivation and some of its general features.

Alethic reference (just “reference”, from now on) is introduced in order to study the reference patterns underlying the semantic paradoxes. Accordingly, it is a relation

---

<sup>11</sup>See Hardy [14] and Ketland [18, 19]. A theory formulated in an extension of  $\mathcal{L}$  is said to be  $\omega$ -inconsistent just in case there is a formula  $\varphi(x)$  such that, for each  $n \in \omega$ ,  $\varphi(\bar{n})$  is a theorem and, at the same time, the theory entails  $\neg \forall x \varphi(x)$ . An  $\omega$ -inconsistent theory may be nonetheless consistent, as inferring  $\forall x \varphi(x)$  from the set of all its instances  $\varphi(\bar{n})$  would require an infinitary rule, not admissible in finitary systems such as the ones we are working with.

between sentences of  $\mathcal{L}_T$ . Sentences can refer to one another in two different ways: by mention and by quantification. For a sentence  $\varphi$  to refer by mention – “m-refer” – to a sentence  $\psi$  the former must contain a closed term denoting (the code of)  $\psi$ ; and for it to refer to  $\psi$  by quantification – “q-refer” – (the code of)  $\psi$  must fall under the range of the quantifiers in  $\varphi$ , which might be restricted by predicates in a sense to be specified. Since we are only interested in reference in the context of truth, to determine the m- and q-referents of a sentence, we will focus only on terms occurring in the scope of T. As a consequence, only sentences containing T can (alethically) refer. Finally, it is worth pointing out that since reference is strongly tied to the presence of terms in sentences, the notion cannot be closed under logical equivalence, on pain of triviality (cf. footnote 13). In other words, reference is not extensional, but hyperintensional.

We turn now to the definitions.

**Definition 1** (M-reference) Let  $\varphi$  and  $\psi$  be sentences.  $\varphi$  m-refers to  $\psi$  iff  $\varphi$  contains a subsentence of the form  $Tt$  and  $\mathbb{N} \models t = \ulcorner \psi \urcorner$ .

A sentence m-refers only to those sentences denoted by *closed* terms that occur immediately after T. For instance, let  $\neg$  be a function symbol of  $\mathcal{L}$  representing the p.r. function that maps each formula to its negation – and similarly for other logical connectives – and let  $\text{Bew}_{\text{PA}}(x)$  weakly represent provability in PA in a natural way.<sup>12</sup> According to the definition of m-reference, while  $\neg T\ulcorner 0 = 0 \urcorner = 0^\ulcorner$  m-refers to  $0 = 0$ ,  $T\neg\ulcorner 0 = 0 \urcorner = 0^\ulcorner$  m-refers to *the negation* of  $0 = 0$  – but not to  $0 = 0$  itself – and  $\text{Bew}_{\text{PA}}(\ulcorner 0 = 0 \urcorner)$  doesn’t m-refer to any expression. Thus, proper subterms in the scope of T and T-free sentences don’t play a role in m-reference.

Note also that while  $T\ulcorner 0 = 0 \urcorner = 0^\ulcorner \rightarrow T\ulcorner 0 = 0 \urcorner = 0^\ulcorner$  m-refers to  $0 = 0$ ,  $0 = 0 \rightarrow 0 = 0$  doesn’t, so m-reference is not closed under logical – let alone material – equivalence.<sup>13</sup> For this reason, one should not expect the equivalences delivered by Theorem 1 to bring forth self-m-referential sentences. However, since  $\mathcal{L}$  contains a function symbol for substitution, there is a stronger version of this result that does so, provided that  $Tx$  is a subformula of  $\varphi$ .

**Theorem 2** (Strong Diagonalization) For every formula  $\varphi(x, \mathbf{v})$  there is a term  $t$  s.t.  $t = \ulcorner \varphi(t, \mathbf{v}) \urcorner$  is a theorem of PA.<sup>14</sup>

We say that  $\varphi(t, \mathbf{v})$  is a strong fixed point of  $\varphi(x, \mathbf{v})$ . Strongly diagonalizing, for instance, the predicate  $\neg Tx$ , we obtain in PA the following identity:

$$1 = \ulcorner \neg T1 \urcorner$$

that is, a sentence that m-refers to itself.  $\neg T1$  is a ‘strong’ liar sentence.

<sup>12</sup>See Halbach & Visser [12, 13] for a discussion on natural representations.

<sup>13</sup>If it were, every sentence  $\varphi$  would refer to every other sentence  $\psi$ , as  $\varphi$  and, e.g.  $\varphi \wedge (T\ulcorner \psi \urcorner \rightarrow T\ulcorner \psi \urcorner)$  are logically equivalent, and the latter obviously m-refers to  $\psi$ .

<sup>14</sup>See Jeroslow [17] for a proof.



Q-reference is considerably more complicated. The underlying idea is the following. While, e.g.  $\forall x Tx$  q-refers unrestrictedly to every expression, sentences such as PA’s global reflection principle,

$$(GRfn_{PA}) \quad \forall x (Bew_{PA}(x) \rightarrow Tx)$$

q-refer only to the theorems of PA. In general, universal quantifiers followed by conditionals are to be understood as *restricted*, i.e. ranging over just the sentences satisfying the antecedent, with one proviso: that the antecedents are T-free. In this way, we are able to determine what the sentence q-refers to.<sup>15</sup> Otherwise, the quantifiers are taken to be unrestricted.

Giving a formal definition of q-reference requires some additional setup. We first need to introduce the procedure of *normalizing* formulae of  $\mathcal{L}_T$ ; q-reference for sentences is then defined in terms of their normalizations. Since distinct sentences can have the same normalization, the procedure therefore induces an equivalence relation – *having the same normalization* – under which q-reference is closed. Triviality is avoided, however, since this equivalence relation is more fine grained than logical equivalence.

The normalization of an expression is the result of a series of logically valid transformations that deliver a formula in *alethic disjunctive normal form* (ADNF). Call “prime” any atomic or universal formulae or their negation. A formula is in ADNF just in case it contains no conditionals and no existential or dummy quantifiers, and every subformula of the form  $\forall v \varphi$  is s.t. (i)  $\varphi$  is a disjunction (of length  $\geq 0$ ) of conjunctions (of length  $\geq 0$ ) of primes and, (ii) if  $\varphi$  is of the form  $\psi \vee \chi$ ,  $\psi$  contains all and only the T-free disjuncts of  $\varphi$  (if any) and  $\chi$  all of the T-containing disjuncts (if any). The point of normalization is that one can easily see whether sentences in ADNF take the form of a restricted quantified claim. Consider a sentence  $\forall v \varphi$  in ADNF containing T. If  $\varphi$  is a disjunction  $\psi \vee \chi$  with some T-free disjuncts,  $\forall v \varphi$  can be seen to be equivalent to the restricted quantificational claim  $\forall v (\neg\psi \rightarrow \chi)$ , for the disjuncts of  $\psi$  must be T-free whereas the disjuncts of  $\chi$  must not. In this way, the reference-restricting conditional  $\neg\psi \rightarrow \chi$  is made explicit, and  $\neg\psi$  is guaranteed to encapsulate all truth-free restrictions imposed on  $\forall v$ . If, instead,  $\varphi$  is not a disjunction or contains no T-free disjuncts, the quantifiers in  $\forall v \varphi$  are unrestricted.

Since formulae in ADNF cannot contain conditionals or existential quantifiers, the first step in normalizing an expression is to replace these connectives with negations, conjunctions, disjunctions, and universal quantifiers, making use of the standard definitions. Let  $\tau : \mathcal{L}_T \rightarrow \mathcal{L}_T$  carry out these replacements, that is,  $\tau(\varphi \rightarrow \psi) := \neg\tau(\varphi) \vee \tau(\psi)$  and  $\tau(\exists v \varphi) := \neg\forall v \neg\tau(\varphi)$ . Then, the normalization process is done in stages. It consists of successive transformations of each subfor-

---

<sup>15</sup>Recall there is no standard interpretation of T yet. Our goal is, instead, to define reference to formulate restrictive criteria for disquotation in terms of this notion and *then* learn about truth and its extension.

mula of the form  $\forall v \varphi$  into ADNF, starting from those of lesser *depth*. Let  $\varphi$  contain no conditionals, existential, or dummy quantifiers.

$$dep(\forall v \varphi) = \begin{cases} 1 & \text{if } \varphi \text{ is atomic} \\ dep(\forall v \psi) & \text{if } \varphi := \neg\psi \\ \max\{dep(\forall v \psi), dep(\forall v \chi)\} & \text{if } \varphi := (\psi \wedge \chi) \\ \max\{dep(\forall v \psi), dep(\forall v \chi)\} & \text{if } \varphi := (\psi \vee \chi) \\ dep(\forall u \psi) + 1 & \text{if } \varphi := \forall u \psi \end{cases}$$

For each  $i \in \omega$ , the  $i$ -normalization of  $\varphi$  is the result of successively applying the following transformations to each subformula  $\forall v \psi$  of depth  $i$ :

1. Replace every subformula of the form  $\neg(\psi_1 \vee \psi_2)$  and  $\neg(\psi_1 \wedge \psi_2)$  with  $(\neg\psi_1 \wedge \neg\psi_2)$  and  $(\neg\psi_1 \vee \neg\psi_2)$  resp. until they don't occur any longer, starting with the innermost.
2. Erase all double negations.
3. Replace every subformula of the form  $\psi_1 \wedge (\psi_2 \vee \psi_3)$  and  $(\psi_2 \vee \psi_3) \wedge \psi_1$  with  $(\psi_1 \wedge \psi_2) \vee (\psi_1 \wedge \psi_3)$  and  $(\psi_2 \wedge \psi_1) \vee (\psi_3 \wedge \psi_1)$  resp. until they don't occur any longer, starting with the innermost.
4. In every subformula of the form  $\forall v(\psi_1 \vee \dots \vee \psi_m)$  (where each  $\psi_i$ ,  $1 \leq i \leq m$ , is not itself a disjunction), rearrange the disjuncts into  $\chi_1 \vee \chi_2$  such that the ones not containing T (if any) occur in  $\chi_1$ , whilst the others (if any) occur in  $\chi_2$ .

**Definition 2** (Normalization) The normalization  $\varphi^*$  of a formula  $\varphi$  is the result of erasing all dummy quantifiers in  $\tau(\varphi)$  and then, if there are any quantifiers left, performing successive  $i$ -normalizations starting with  $i = 1$  and stopping after  $i = \max\{dep(\forall v \psi) : \forall v \psi \text{ is a subformula of } \varphi\}$ .

Since every step in the  $i$ -normalization of a formula involves only finitely many transformations and formulae contain finitely many quantifiers, the normalization process always terminates. Moreover, it delivers a logically equivalent expression.

We are finally in a position to provide an adequate and precise definition of q-reference. Let  $\mathbf{n}$  abbreviate  $n_1, \dots, n_m \in \omega$ , and  $\bar{\mathbf{n}}$  abbreviate  $\bar{n}_1, \dots, \bar{n}_m$ .

**Definition 3** (Q-reference) Let  $\varphi, \psi$  be sentences.  $\varphi$  q-refers to  $\psi$  iff  $\varphi^*$  has a subsentence of the form  $\forall \mathbf{v} \chi$  s.t. T occurs in  $\chi$ , and one of the following holds:

1.  $\chi := (\chi_1 \vee \chi_2)$ , T doesn't occur in  $\chi_1$ , and  $\chi_2[\bar{\mathbf{n}}/\mathbf{v}]$  q-refers or newly m-refers to  $\psi$ ,<sup>16</sup> for some  $\mathbf{n} \in \omega$  s.t.  $\mathbb{N} \models \neg\chi_1[\bar{\mathbf{n}}/\mathbf{v}]$ .
2. Either  $\chi := (\chi_1 \vee \chi_2)$  and T occurs in  $\chi_1$  or  $\chi$  is not a disjunction or a universal statement, and  $\chi[\bar{\mathbf{n}}/\mathbf{v}]$  q-refers or newly m-refers to  $\psi$ , for some  $\mathbf{n} \in \omega$ .

In brief, the q-referents of universally quantified claims are the sentences their (possibly restricted) instances m- or q-refer to. More precisely, a sentence whose

<sup>16</sup>To keep m- and q-reference apart, we require here that m-reference is achieved through 'new' closed terms, that is, that the latter are not present in the sentence itself but are a product of the instantiation of the quantifiers at issue.

normalization contains a subsentence of the form  $\forall \mathbf{v} \varphi$  q-refers to what each instance  $\varphi[\mathbf{n}/\mathbf{v}]$  newly m-refers or q-refers to, unless it  $\varphi$  a reference-restricting conditional, in which case only the instances with true antecedents are to be considered.<sup>17</sup> For instance,

$$\forall x \forall y (T \top 0 = 0^\top \wedge Tx \rightarrow y) \tag{6}$$

which is already in ADNF, q-refers to the m-referents of each instance  $T \top 0 = 0^\top \wedge T \bar{m} \rightarrow \bar{n}$ , provided that the terms delivering m-reference are a product of the instantiation of the quantifiers  $\forall x \forall y$ , which must be eliminated together ‘at once’. Thus, (6) q-refers just to every conditional sentence. On the other hand,

$$\forall x (x = \top 0 = 0^\top \rightarrow \forall y (y = \neg x \rightarrow Ty)) \tag{7}$$

which normalizes into  $\forall x (x \neq \top 0 = 0^\top \vee \forall y (y \neq \neg x \vee Ty))$ , q-refers to the q-referents of each instance  $\forall y (y \neq \neg \bar{n} \vee Ty)$  provided that  $\mathbb{N} \models \bar{n} = \top 0 = 0^\top$ : i.e. (7) q-refers to the negation of  $0 = 0$ .

Just as the definition of m-reference accounts for the self-referentiality of certain sentences delivered by the Strong Diagonal Lemma, Definition 3 implies that ‘weakly’ diagonalizing certain predicates – e.g. in which  $Tx$  occurs as a subformula – also delivers self-referential sentences.<sup>18</sup> This implies that  $\lambda$  is self-q-referential. The Visser-Yablo sentences, instead, only q-refer to sentences coming later in the sequence; they are not self-q-referential.

Mention and quantification exhaust the ways in which a sentence can directly refer to other expressions. Thus, we have the following definition of direct reference – “d-reference”.

**Definition 4 (D-reference)** Let  $\varphi, \psi$  be sentences.  $\varphi$  directly refers to  $\psi$  iff it m- or q-refers to  $\psi$ .

**Observation 3** For all  $\varphi, \psi, \chi \in \mathcal{L}_T$ :

1. If  $\varphi \in \mathcal{L}$ ,  $\varphi$  doesn’t d-refer to  $\psi$ .
2. If  $\mathbb{N} \models s = t$ ,  $\varphi$  and  $\varphi[s/t]$  d-refer to the same sentences.
3.  $\varphi$  and  $\neg\varphi$  d-refer to the same sentences.
4.  $\varphi \vee \chi$ ,  $\varphi \wedge \chi$ , and  $\varphi \rightarrow \chi$  d-refer to  $\psi$  iff either  $\varphi$  or  $\chi$  do.
5. If  $v$  is not free in  $\varphi$ ,  $\varphi$ ,  $\forall v \varphi$ , and  $\exists v \varphi$  d-refer to the same sentences.
6.  $\forall v \varphi$  and  $\forall u \varphi[u/v]$  d-refer to the same sentences, if  $u$  is free for  $v$  in  $\varphi$ .
7. The following pairs of logical equivalents d-refer to the same sentences:

- $\varphi$  and  $\neg\neg\varphi$ ,
- $\varphi \vee \psi$  and  $\psi \vee \varphi$ ,

<sup>17</sup>Several worked examples, together with additional motivation and discussion, can be found in Picollo [27].

<sup>18</sup>See Picollo [27, Section 3.3].

- $\varphi \wedge \psi$  and  $\psi \wedge \varphi$ ,
- $(\varphi \vee \psi) \vee \chi$  and  $\varphi \vee (\psi \vee \chi)$ ,
- $(\varphi \wedge \psi) \wedge \chi$  and  $\varphi \wedge (\psi \wedge \chi)$ ,
- $\varphi \rightarrow \psi$  and  $\neg\psi \rightarrow \neg\varphi$ ,
- $\varphi \vee (\psi \wedge \chi)$  and  $(\varphi \vee \psi) \wedge (\varphi \vee \chi)$ ,
- $\varphi \wedge (\psi \vee \chi)$  and  $(\varphi \wedge \psi) \vee (\varphi \wedge \chi)$ ,
- $\neg(\varphi \vee \psi)$  and  $\neg\varphi \wedge \neg\psi$ ,
- $\neg(\varphi \wedge \psi)$  and  $\neg\varphi \vee \neg\psi$ ,
- $\varphi \rightarrow \psi$  and  $\neg\varphi \vee \psi$ ,
- $\exists v \varphi$  and  $\neg\forall v \neg\varphi$ .

Definition 4 allows us to characterize many ‘mixed’ reference patterns. Consider, for instance, sentences  $\lambda_1$  and  $\lambda_2$  in the 2-liar cycle in (3). While  $\lambda_1$  m-refers to  $\lambda_2$ , the latter only q-refers to the former. Thus, we can say they directly refer to each other. However, in a very clear sense they also refer to themselves, albeit *indirectly*. Otherwise, we would get a semantic paradox without self-reference on the cheap. The following notions are intended to deal with this and other similar cases.

**Definition 5** (Chain of reference) A (possibly infinite) sequence of sentences s.t. each sentence in the sequence d-refers to the one coming after, if any.

**Definition 6** (Reference) Let  $\varphi, \psi$  be sentences.  $\varphi$  refers to  $\psi$  iff there’s a chain of reference starting with  $\varphi$  and ending with  $\psi$ .

We can employ the notions just introduced to define salient reference patterns, such as the following.

**Definition 7** (Self-reference) A sentence is self-referential iff it refers to itself.

**Definition 8** (Well-foundedness) A sentence  $\varphi$  is well-founded iff all chains of reference starting with  $\varphi$  are finite. Otherwise, we say that  $\varphi$  is unfounded.

Thus, sentences that don’t d-refer to any expressions – e.g. the theorems of PA – are well-founded. And sentences that only refer to well-founded expressions are well-founded too. On the other hand, every self-referential expression is obviously unfounded. But there are also unfounded sentences that don’t refer to themselves, such as the Visser-Yablo sentences in (5). Thus, not all paradoxical expressions are self-referential. However, the notions introduced in this section can still be deployed in the formulation of restrictions to our disquotational truth principles. As will be seen in the next section, all paradoxical expressions are unfounded, for theories in which disquotation is restricted to well-founded sentences will be shown to be ( $\omega$ -) consistent. What is more, we will show that this characterization can be refined even further, prompting more encompassing truth systems.

### 3 Semantic Theories of Truth

In this section we deploy the notions of alethic reference introduced in Section 2.3 to prove the existence of models of  $\mathcal{L}_T$  expanding  $\mathbb{N}$  which verify large and well-motivated sets of instances of disquotation. In other words, we put forward semantic truth theories. In turn, these will serve as witnesses to the consistency and arithmetical soundness of the axiomatic systems we introduce in Section 4.

#### 3.1 Well-Founded Truth

Our first theory consists of a single model  $\langle \mathbb{N}, \Phi_{wf} \rangle$  of  $\mathcal{L}_T$  in which the extension assigned to T,  $\Phi_{wf}$ , contains all and only true well-founded sentences. Thus, all instances of the T-schema for well-founded sentences hold in the model. Moreover, since  $\langle \mathbb{N}, \Phi_{wf} \rangle$  expands the standard model of arithmetic, the Uniform T-schema restricted to formulae with only well-founded numerical instances also holds in  $\langle \mathbb{N}, \Phi_{wf} \rangle$ .

The set  $\Phi_{wf}$  is obtained via a usual Kripke-style construction, by considering a transfinite sequence of sets  $\Phi_\alpha \subseteq \omega$ , with  $\alpha \in \text{On}$  (the class of all ordinals), of which  $\Phi_{wf}$  is shown to be a fixed point. In order to construct this sequence, we first consider another sequence of sets  $\Psi_\alpha$  whose fixed point,  $\Psi_{wf}$ , is the set of well-founded sentences. Thus, this sequence ‘stratifies’ this set. If  $\varphi$  is a sentence, let  $\Theta_\varphi$  be the set of sentences  $\varphi$  d-refers to. For each  $\alpha \in \text{On}$ ,  $\Psi_\alpha$  is defined as follows:

$$\Psi_\alpha := \begin{cases} \emptyset & \text{if } \alpha = 0 \\ \{\varphi : \Theta_\varphi \subseteq \Psi_\beta\} & \text{if } \alpha = \beta + 1 \\ \bigcup_{\beta < \alpha} \Psi_\beta & \text{if } \alpha \text{ is a limit} \end{cases}$$

**Lemma 1** *For all  $\alpha, \beta \in \text{On}$ , if  $\alpha < \beta$ , then  $\Psi_\alpha \subseteq \Psi_\beta$ .*

*Proof* By transfinite induction on  $\beta$ . Let  $\alpha < \beta$  and  $\varphi \in \Psi_\alpha$ . If  $\beta$  is 0 or a limit ordinal, the result follows trivially. Let  $\beta = \xi + 1$ . By inductive hypothesis (i.h.),  $\varphi \in \Psi_\xi$ . Thus,  $\xi \neq 0$ . If  $\xi = \zeta + 1$ ,  $\Theta_\varphi \subseteq \Psi_\zeta$ . Again, by i.h.,  $\Psi_\zeta \subseteq \Psi_\xi$ , so  $\Theta_\varphi \subseteq \Psi_\xi$ . Therefore,  $\varphi \in \Psi_\beta$ . The case in which  $\xi$  is a limit can be proved in a similar way.  $\square$

**Proposition 1** *There is an  $\alpha \in \text{On}$  s.t., for every  $\beta > \alpha$ ,  $\Psi_\alpha = \Psi_\beta$  and  $\Psi_\alpha$  is the set of well-founded sentences.*

*Proof* The first conjunct follows immediately from Lemma 1 and cardinality considerations. Therefore, the sequence reaches a fixed point,  $\Psi_{wf}$ . We show that  $\Psi_{wf}$  is the set of well-founded sentences.

Let  $\varphi \in \Psi_{wf}$ . Thus,  $\varphi \in \Psi_\alpha$ , for some  $\alpha \in \text{On}$ . We show by transfinite induction on  $\alpha$  that all sentences in  $\Psi_\alpha$  are well-founded. Assume that, for every  $\beta < \alpha$ ,  $\Psi_\beta$  contains only well-founded sentences. If  $\alpha = 0$ , the result follows trivially. If  $\alpha = \xi + 1$ , then  $\xi$  is a set of well-founded sentences, by i.h.. Then, all members of  $\Psi_\alpha$  d-refer just to well-founded sentences and are also well-founded. If  $\alpha$  is a limit ordinal, the result follows trivially from the i.h. as well.

Now let  $\varphi$  be a well-founded sentence. By Definition 8, either  $\varphi$  doesn't refer or every chain of reference starting with  $\varphi$  is finite. Thus, the following function from the set of well-founded sentences to On is well defined:

$$f(\varphi) := \begin{cases} 0 & \text{if } \Theta_\varphi = \emptyset \\ \sup\{f(\psi) + 1 : \varphi \text{ d-refers to } \psi\} & \text{otherwise} \end{cases}$$

I show that  $\varphi \in \Psi_{f(\varphi)+1}$  by transfinite induction on  $f(\varphi)$ . If  $f(\varphi) = 0$ ,  $\Theta_\varphi = \emptyset$ , so  $\Theta_\varphi \subseteq \Psi_0$ , so  $\varphi \in \Psi_1$ . Assume the result holds for every  $\alpha < f(\varphi)$  and let  $f(\varphi) = \beta + 1$ . The case in which  $f(\varphi)$  is a limit is similar. Then,  $\beta$  is the supremum of all  $f(\psi)$  s.t.  $\varphi$  d-refers to  $\psi$ . By i.h. and Lemma 1,  $\Theta_\varphi \subseteq \Psi_{\beta+1}$ , which means  $\varphi \in \Psi_{f(\varphi)+1}$ . By Lemma 1 and the first conjunct of this proposition, we can conclude that  $\varphi \in \Psi_{\text{wf}}$ . □

We have now paved the way for the construction of the sequence of  $\Phi_\alpha$ , with  $\alpha \in \text{On}$ , that will give us the desired extension of the truth predicate. At each ordinal  $\alpha$ ,  $\Phi_\alpha$  contains only sentences that are well-founded at this stage, i.e. that belong to  $\Psi_\alpha$ .

$$\Phi_\alpha := \begin{cases} \emptyset & \text{if } \alpha = 0 \\ \{\varphi : \langle \mathbb{N}, \Phi_\beta \rangle \models \varphi\} \cap \Psi_\alpha & \text{if } \alpha = \beta + 1 \\ \bigcup_{\beta < \alpha} \Phi_\beta & \text{if } \alpha \text{ is a limit ordinal} \end{cases}$$

To show that the sequence reaches a fixed point we first need to prove two lemmata. According to the first, the truth value of a statement in a model is not affected by the mere addition or removal of sentences it doesn't d-refer to from the extension of the truth predicate. This establishes a link between direct reference and Leitgeb's [21] notion of dependence.<sup>19</sup> More specifically, it follows that  $\varphi$  depends on  $\Theta_\varphi$ . The converse doesn't hold, that is,  $\Theta_\varphi$  is not a subset of every set  $\varphi$  depends on as, e.g.  $\text{T}\ulcorner\lambda\urcorner \rightarrow \text{T}\ulcorner\lambda\urcorner$  d-refers to  $\lambda$  but depends on  $\emptyset$ . Just like aboutness, dependence is not as tied to the syntactic structure of sentences as reference is.

**Lemma 2** *Let  $\Gamma$  be a set of sentences. If  $\Theta_\varphi \subseteq \Gamma$ , then for every set of sentences  $\Delta$ ,*

$$\langle \mathbb{N}, \Delta \rangle \models \varphi \text{ iff } \langle \mathbb{N}, \Gamma \cap \Delta \rangle \models \varphi.$$

*Proof* Since every sentence is logically equivalent to its normalization (cf. Definition 2), we can prove the result for the normalization of  $\varphi$ ,  $\varphi^*$ , by induction on its complexity (number of logical operators). If  $\varphi^* \in \mathcal{L}$ ,  $\varphi^*$  is true in every expansion of  $\mathbb{N}$  or in none, so the result follows trivially. Thus, we assume  $\varphi^*$  contains T.

If  $\varphi^*$  is an atomic sentence, then it's of the form  $\text{T}t$ , where either  $t$  denotes a sentence in  $\mathbb{N}$  or it doesn't. If it doesn't, then  $\varphi^*$  is false both in  $\langle \mathbb{N}, \Delta \rangle$  and in  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle$ . If  $t$  denotes a sentence  $\psi$ , then, by Definition 1,  $\psi \in \Theta_\varphi \subseteq \Gamma$ . Thus,  $\psi \in \Delta$  iff  $\psi \in \Gamma \cap \Delta$ , so we have the desired result.

---

<sup>19</sup>A sentence  $\varphi$  is said to depend on a set of sentences  $\Gamma$  iff, for every set of sentences  $\Delta$ ,  $\langle \mathbb{N}, \Delta \rangle \models \varphi$  iff  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle \models \varphi$ .

Assume the claim holds for every sentence of lower complexity than  $\varphi^*$ . Since d-reference is closed under logical connectives (cf. Observation 3), if  $\varphi$  is a negation, conjunction, or disjunction, the result follows trivially from the i.h.

Let  $\varphi^* := \forall \mathbf{v} \psi$ , where  $\psi$  is not a universal statement.  $\varphi^*$  contains no dummy quantifiers. Let  $\psi := (\psi_1 \vee \psi_2)$  s.t. T occurs in every disjunct of  $\psi_2$  but in none of  $\psi_1$ . Otherwise, the proof is similar but with less complications. Since  $\psi_1 \in \mathcal{L}$ , for each  $\mathbf{n} \in \omega$  either  $\mathbb{N} \models \psi_1[\bar{\mathbf{n}}/\mathbf{v}]$  or  $\mathbb{N} \models \neg\psi_1[\bar{\mathbf{n}}/\mathbf{v}]$ . Note that if  $\mathbb{N} \models \psi_1[\bar{\mathbf{n}}/\mathbf{v}]$ , then  $\langle \mathbb{N}, \Delta \rangle \models \psi[\bar{\mathbf{n}}/\mathbf{v}]$  for every set of sentences  $\Delta$ .

Assume then that  $\mathbb{N} \models \neg\psi_1[\bar{\mathbf{n}}/\mathbf{v}]$ . Thus, for every  $\Delta$ ,  $\langle \mathbb{N}, \Delta \rangle \models \psi[\bar{\mathbf{n}}/\mathbf{v}]$  iff  $\langle \mathbb{N}, \Delta \rangle \models \psi_2[\bar{\mathbf{n}}/\mathbf{v}]$ . Since  $\Theta_{\forall \mathbf{v} \psi} \subseteq \Gamma$  and, by Definition 3,  $\Theta_{\forall \mathbf{v} \psi}$  is the union of all  $\Theta_{\psi_2[\bar{\mathbf{n}}/\mathbf{v}]}$  s.t.  $\mathbb{N} \models \neg\psi_1[\bar{\mathbf{n}}/\mathbf{v}]$ , we have that  $\Theta_{\psi_2[\bar{\mathbf{n}}/\mathbf{v}]} \subseteq \Gamma$ . By i.h.,  $\langle \mathbb{N}, \Delta \rangle \models \psi_2[\bar{\mathbf{n}}/\mathbf{v}]$  iff  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle \models \psi_2[\bar{\mathbf{n}}/\mathbf{v}]$ . By our last assumption, this entails that  $\langle \mathbb{N}, \Delta \rangle \models (\psi_1 \vee \psi_2)[\bar{\mathbf{n}}/\mathbf{v}]$  iff  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle \models (\psi_1 \vee \psi_2)[\bar{\mathbf{n}}/\mathbf{v}]$  or, what is the same,  $\langle \mathbb{N}, \Delta \rangle \models \psi[\bar{\mathbf{n}}/\mathbf{v}]$  iff  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle \models \psi[\bar{\mathbf{n}}/\mathbf{v}]$ .

Thus, in any case we also have that  $\langle \mathbb{N}, \Delta \rangle \models \forall \mathbf{v} \psi$  iff  $\langle \mathbb{N}, \Gamma \cap \Delta \rangle \models \forall \mathbf{v} \psi$ . □

**Lemma 3** *For all  $\alpha, \beta \in \text{On}$ , if  $\alpha < \beta$ , then  $\Phi_\alpha \subseteq \Phi_\beta$ .*

*Proof* By transfinite induction on  $\beta$ . If  $\beta$  is not a successor ordinal, the result follows trivially. Let  $\beta = \xi + 1$ , and assume that, for all  $\zeta < \xi$ ,  $\Phi_\zeta \subseteq \Phi_\xi$  and  $\alpha < \beta$ . Thus,  $\alpha \leq \xi$ . If  $\alpha = \xi$ , trivially  $\Phi_\alpha \subseteq \Phi_\xi$ , and if  $\alpha < \xi$ ,  $\Phi_\alpha \subseteq \Phi_\xi$  by i.h.. Thus, it remains to be shown that  $\Phi_\xi \subseteq \Phi_\beta$ .

Let  $\varphi \in \Phi_\xi$ . Therefore,  $\Phi_\xi \neq \emptyset$ , so  $\xi$  is either a successor or a limit ordinal. The proof for both cases is similar. I only show it for  $\xi = \zeta + 1$ . Since  $\varphi \in \Psi_\xi$ ,  $\varphi$  d-refers only to sentences in  $\Psi_\zeta$  and  $\langle \mathbb{N}, \Phi_\zeta \rangle \models \varphi$ . Moreover, note that  $\Phi_\zeta = \Psi_\zeta \cap \Phi_\xi$ : by construction of  $\Phi_\zeta$ ,  $\Phi_\zeta \subseteq \Psi_\zeta$  and, by i.h.,  $\Phi_\zeta \subseteq \Phi_\xi$ . For the other direction, assume for contradiction that  $\psi \in \Psi_\zeta \cap \Phi_\xi$  and  $\psi \notin \Phi_\zeta$ . Thus,  $\zeta \neq \emptyset$ . Let  $\zeta = \theta + 1$ ; the case in which  $\zeta$  is a limit ordinal is analogous. Since  $\psi \notin \Phi_\zeta$  but  $\psi \in \Psi_\zeta$ , we have that  $\langle \mathbb{N}, \Phi_\theta \rangle \not\models \psi$ , that is,  $\langle \mathbb{N}, \Phi_\theta \rangle \models \neg\psi$ . Since d-reference is closed under negation (cf. Observation 3), we also have that  $\neg\psi \in \Psi_\zeta$ . Thus,  $\neg\psi \in \Phi_\zeta$ . By i.h.,  $\neg\psi \in \Phi_\xi$ . As  $\psi \in \Phi_\xi$  too, both  $\langle \mathbb{N}, \Phi_\zeta \rangle \models \psi$  and  $\langle \mathbb{N}, \Phi_\zeta \rangle \models \neg\psi$ , which is impossible.

Therefore,  $\langle \mathbb{N}, \Psi_\zeta \cap \Phi_\xi \rangle \models \varphi$ . By Lemma 2,  $\langle \mathbb{N}, \Phi_\xi \rangle \models \varphi$ . Also, since  $\varphi$  d-refers only to sentences in  $\Psi_\zeta$ ,  $\varphi \in \Psi_\xi$ , and by Lemma 1,  $\varphi \in \Psi_\beta$ . Thus,  $\varphi \in \Phi_\beta$ . □

**Proposition 2** *There is an ordinal  $\alpha \in \text{On}$  s.t.  $\Phi_\alpha = \Phi_{\alpha+1}$ .*

*Proof* By Lemma 3 and cardinality considerations. □

The sequence of sets  $\Phi_\alpha$  stabilizes at some ordinal, it has a fixed point. Let  $\Phi_{\text{wf}}$  be this fixed point. The following result establishes that all instances of the T-schema for well-founded sentences hold in  $\langle \mathbb{N}, \Phi_{\text{wf}} \rangle$ .

**Proposition 3** *For every well-founded sentence  $\varphi$ ,  $\langle \mathbb{N}, \Phi_{\text{wf}} \rangle \models T^\top \varphi^\top \leftrightarrow \varphi$ .*

*Proof* Let  $\varphi$  be a well-founded expression. By Proposition 1, there is an ordinal  $\alpha$  s.t.  $\varphi \in \Psi_\alpha$ . Thus,  $\alpha \neq \emptyset$ , so  $\alpha$  is either a successor or a limit ordinal. I only prove the result for the latter case, as the proof for the former is analogous. If  $\alpha$  is a limit, by Lemma 3 there is a successor ordinal  $\beta < \alpha$  s.t.  $\varphi \in \Psi_\beta$ . As d-reference is closed under negation, we also have that  $\neg\varphi \in \Psi_\beta$ . Let  $\beta = \xi + 1$ . Either  $\langle \mathbb{N}, \Phi_\xi \rangle \models \varphi$  or  $\langle \mathbb{N}, \Phi_\xi \rangle \models \neg\varphi$ , so either  $\varphi \in \Phi_\beta$  or  $\neg\varphi \in \Phi_\beta$ . Therefore, by Lemma 3,  $\varphi \in \Phi_{\text{wf}}$  or  $\neg\varphi \in \Phi_{\text{wf}}$ . Let  $\zeta$  be the smallest ordinal s.t.  $\Phi_\zeta = \Phi_{\text{wf}}$ , i.e. the fixed-point ordinal.

If  $\varphi \in \Phi_\zeta$ , then  $\langle \mathbb{N}, \Phi_\zeta \rangle \models T^\Gamma\varphi^\neg$ , but also  $\varphi \in \Phi_{\zeta+1}$ , since  $\Phi_\zeta = \Phi_{\zeta+1}$ . Thus,  $\varphi \in \{\varphi : \langle \mathbb{N}, \Phi_\zeta \rangle \models \varphi\} \cap \Psi_\zeta$ , which means that  $\langle \mathbb{N}, \Phi_\zeta \rangle \models \varphi$ . So, we get  $\langle \mathbb{N}, \Phi_\zeta \rangle \models T^\Gamma\varphi^\neg \leftrightarrow \varphi$ .

If, instead,  $\neg\varphi \in \Phi_\zeta$ , then  $\neg\varphi \in \Phi_{\zeta+1}$ , so  $\neg\varphi \in \{\varphi : \langle \mathbb{N}, \Phi_\zeta \rangle \models \varphi\} \cap \Psi_\zeta$ . This implies that  $\langle \mathbb{N}, \Phi_\zeta \rangle \models \neg\varphi$  and, thus,  $\langle \mathbb{N}, \Phi_\zeta \rangle \not\models \varphi$ . Therefore,  $\varphi \notin \Phi_{\zeta+1}$ , which means that  $\varphi \notin \Phi_\zeta$ , i.e.  $\langle \mathbb{N}, \Phi_\zeta \rangle \not\models T^\Gamma\varphi^\neg$ . Again, we have that  $\langle \mathbb{N}, \Phi_\zeta \rangle \models T^\Gamma\varphi^\neg \leftrightarrow \varphi$ . □

As a consequence, for every well-founded sentence, either it belongs to  $\Phi_{\text{wf}}$  or its negation does. Note also that, by Lemma 1 and Proposition 1,  $\Phi_{\text{wf}}$  contains only well-founded sentences.

### 3.2 More Permissive Criteria

Well-foundedness is not the only restriction we can impose on instances of the T-schema in an expansion of  $\mathbb{N}$  to  $\mathcal{L}_T$  using the notions of alethic reference: more permissive criteria may also be adopted. For instance, note that in the case of the Visser-Yablo paradox, every expression on the list is unfounded, but also there are infinitely many of them and, additionally, each of them d-refers to infinitely many others. It can be shown that, if only finitely many sentences on the list are considered, no  $\omega$ -inconsistency arises. More generally, let  $\Delta$  by a *finite* set of non-self-referential but unfounded sentences. One can easily find a model  $\langle \mathbb{N}, \Gamma \rangle$  in which all instances of the T-schema for sentences in  $\Psi_{\text{wf}} \cup \Delta$  hold.

**Proposition 4** *If  $\Delta$  is a finite set of non-self-referential unfounded sentences, there's a  $\Gamma \subseteq \Psi_{\text{wf}} \cup \Delta$  s.t. all instances of the T-schema for sentences in  $\Psi_{\text{wf}} \cup \Delta$  are true in  $\langle \mathbb{N}, \Gamma \rangle$ .*

*Proof* Let  $f : \Psi_{\text{wf}} \cup \Delta \rightarrow \omega$  be s.t.

$$f(\varphi) := \begin{cases} 0 & \text{if } \Delta \cap \Theta_\varphi = \emptyset \\ \max\{f(\psi) : \psi \in \Delta \cap \Theta_\varphi\} + 1 & \text{otherwise} \end{cases}$$

$f$  is well defined, as  $\Delta$  is finite and contains no self-referential sentences. Thus, every complete chain of reference restricted to members of  $\Delta$  ends in a sentence that is not in  $\Delta$ . For each  $n \in \omega$ , let  $\Delta_n := \{\varphi \in \Delta : f(\varphi) = n\}$ . Since  $\Delta$  is finite, only finitely many of these sets are non-empty. Let  $k \in \omega$  be the greatest number s.t.  $\Delta_k \neq \emptyset$ . Thus,  $\Delta = \bigcup_{i \leq k} \Delta_i$ .



For each  $n \in \omega$ , let

$$\Xi_n := \begin{cases} \Phi_{wf} \cup \{\varphi \in \Delta_0 : \langle \mathbb{N}, \Phi_{wf} \rangle \models \varphi\} & \text{if } n = 0 \\ \Xi_m \cup \{\varphi \in \Delta_m : \langle \mathbb{N}, \Xi_m \rangle \models \varphi\} & \text{if } n = m + 1 \end{cases}$$

I show that all instances of the T-schema for sentences in  $\Psi_{wf} \cup \Delta$  are true in  $\langle \mathbb{N}, \Xi_k \rangle$ . By Lemma 2, all instances of the T-schema for well-founded sentences are still true in  $\langle \mathbb{N}, \Xi_k \rangle$ , for  $\Phi_{wf} \subseteq \Xi_k$  and no well-founded sentence may refer to members of  $\Delta$ , as they are unfounded. Now let  $\varphi \in \Delta$ , which means there's an  $n \leq k$  s.t.  $\varphi \in \Delta_n$ . Let  $\Xi_{-1} := \Phi_{wf}$ . Therefore, we have that  $\langle \mathbb{N}, \Xi_k \rangle \models \varphi$  iff  $\langle \mathbb{N}, \Xi_{n-1} \rangle \models \varphi$  iff  $\langle \mathbb{N}, \Xi_n \rangle \models T^\Gamma \varphi^\neg$  iff, by Lemma 2,  $\langle \mathbb{N}, \Xi_k \rangle \models T^\Gamma \varphi^\neg$ , as  $\varphi$  cannot d-refer to any member of  $\Xi_k$  that isn't in  $\Xi_n$ . Consequently, the T-schema also holds for all sentences in  $\Delta$ .  $\square$

More interestingly, it can be shown that for any *infinite* set  $\Delta$  of non-self-referential unfounded sentences, each of which d-refers only to a *finite* number of expressions, there is a model  $\langle \mathbb{N}, \Gamma \rangle$  in which all instances of the T-schema for sentences in  $\Psi_{wf} \cup \Delta$  are true.

**Proposition 5** *If  $\Delta$  is an infinite set of non-self-referential unfounded sentences, each of which d-refers only to a finite number of expressions, there is a  $\Gamma \subseteq \Psi_{wf} \cup \Delta$  s.t., for every  $\varphi \in \Psi_{wf} \cup \Delta$ ,  $\langle \mathbb{N}, \Gamma \rangle \models T^\Gamma \varphi^\neg \leftrightarrow \varphi$ .<sup>20</sup>*

*Proof* Let  $\Delta := \{\varphi_0, \varphi_1, \dots\}$ ,  $\Delta_n := \{\varphi_0, \dots, \varphi_n\}$ , and

$$\mathcal{G}_n := \{\Xi : \Phi_{wf} \subseteq \Xi \subseteq \Phi_{wf} \cup \bigcup_{i \leq n} (\Theta_{\varphi_i} \cup \{\varphi_i\}) \ \& \ \langle \mathbb{N}, \Xi \rangle \models T^\Gamma \varphi^\neg \leftrightarrow \varphi, \text{ for } \varphi \in \Delta_n\}$$

Each  $\mathcal{G}_n$  is finite, for each  $\bigcup_{i \leq n} (\Theta_{\varphi_i} \cup \{\varphi_i\})$  is. I show that, for every  $n \in \omega$ , we can find sets  $\Xi_1 \subseteq \dots \subseteq \Xi_n$  s.t.  $\Xi_i \in \mathcal{G}_i$  for each  $i \leq n$ .

By Proposition 4, for each  $n \in \omega$  there's a model  $\langle \mathbb{N}, \Gamma_n \rangle$  with  $\Phi_{wf} \subseteq \Gamma_n \subseteq \Phi_{wf} \cup \bigcup_{i \leq n} (\Theta_{\varphi_i} \cup \{\varphi_i\})$  in which all instances of the T-schema for sentences in  $\Psi_{wf} \cup \Delta_n$  are true. For  $i \leq n$  let

$$\Gamma_n^i := \Gamma_n \cap (\Phi_{wf} \cup \bigcup_{j \leq i} (\Theta_{\varphi_j} \cup \{\varphi_j\}))$$

Thus, for each  $i \leq n$ ,  $\Gamma_n^i \subseteq \Gamma_n^{i+1}$ . By Lemma 2, we have that all instances of the T-schema for sentences in  $\Delta_i$  hold in  $\langle \mathbb{N}, \Gamma_n^i \rangle$ , so  $\Gamma_n^i \in \mathcal{G}_i$ .

Let  $\mathcal{G}$  be the smallest tree consisting of sequences of members of  $\{\emptyset\} \cup \bigcup_{n \in \omega} \mathcal{G}_n$  s.t.:

- $\langle \emptyset \rangle \in \mathcal{G}$ ,
- if  $\langle \Xi_0, \dots, \Xi_n \rangle \in \mathcal{G}$  and  $\Xi_n \subseteq \Xi_{n+1} \in \mathcal{G}_{n+1}$ , then  $\langle \Xi_0, \dots, \Xi_n, \Xi_{n+1} \rangle \in \mathcal{G}$ .

<sup>20</sup>This result has been proved in collaboration with Thomas Schindler.

Notice that  $\mathcal{G}$  is locally finite, that is, each node has finitely many children, for each  $\mathcal{G}_n$  is finite. Moreover, for every  $n \in \omega$ ,  $\mathcal{G}$  contains a sequence of length  $n$ , e.g.  $\langle \Gamma_n^0, \dots, \Gamma_n^n \rangle$ , so  $\mathcal{G}$  is an infinite graph. Finally, all sequences start with  $\emptyset$ , so all nodes are connected by  $\emptyset$ . By König’s Lemma,<sup>21</sup>  $\mathcal{G}$  contains an infinite sequence  $\langle \Xi_0, \dots, \Xi_n, \dots \rangle$  s.t., for each  $i \in \omega$ ,  $\Xi_i \in \mathcal{G}_i$ .

Let  $\Gamma := \bigcup_{n \in \omega} \Xi_n$ . By Lemma 2, all instances of the T-schema for well-founded sentences are still true in  $\langle \mathbb{N}, \Gamma \rangle$  as, for each  $n \in \omega$ ,  $\Xi_n \cap \Psi_{\text{wf}} = \Phi_{\text{wf}}$ . Let  $\varphi_i \in \Delta$  and let  $m > i$  be the least natural number s.t., for every  $k \geq m$ ,  $\Xi_k$  doesn’t contain elements of  $\Theta_{\varphi_i}$  that are not already in  $\Xi_m$  – i.e.  $\Xi_k \cap \Theta_{\varphi_i} = \Xi_m$ . Such an  $m$  must exist, for  $\Theta_{\varphi_i}$  is finite and sequences in  $\mathcal{G}$  are monotonic. Since  $i < m$  and  $\Xi_m \in \mathcal{G}_m$ ,  $\langle \mathbb{N}, \Xi_m \rangle \models \text{T}^\Gamma \varphi_i^{-1} \leftrightarrow \varphi_i$ . Given the choice of  $m$ , Lemma 2 entails that  $\langle \mathbb{N}, \Gamma \rangle \models \text{T}^\Gamma \varphi_i^{-1} \leftrightarrow \varphi_i$ . □

Propositions 4 and 5 equip us with a more precise criterion for paradox. Whilst the former shows that infinitely many sentences are required for non-self-referential ( $\omega$ -)paradox, by Proposition 5 we also know that at least some of the infinitely many paradoxical sentences must directly refer to infinitely many expressions.

Based on the notions of alethic reference and the results provided in this section, different reference patters may be employed to obtain a plethora of semantic truth theories, some more permissive, philosophically motivated, or elegant than others. Our restriction to well-founded sentences seems to fare quite well in these respects, although we have seen more permissive ones may also be adopted.

Unfortunately, though, reference and, *a fortiori*, well-foundedness and other reference patterns are far too complex to guide the development of axiomatic truth theories in a straightforward manner. For they are not definable in  $\mathcal{L}$ , as they make essential use of the notion of satisfaction in the standard model (cf. Definition 3). In the next section I provide mirroring notions of reference tailored specifically for that purpose; these are simplified counterparts of those introduced in the previous section.

### 4 Axiomatic Truth

The first part of this section is concerned with proof-theoretic versions of the notions presented in Section 2.3, that is, reference concepts that are relative to a particular proof-system. Section 4.2 explores a naïve truth theory formulated purely in terms of these proof-theoretic notions, shows it is unsound, and thus motivates the restriction of disquotation to reference stable expressions to avoid paradox. Finally, in Section 4.3 I formulate several axiomatic theories of disquotational truth deploying the concepts introduced in Sections 4.1 and 4.2. The resulting systems are shown to be sound and proof-theoretically strong in comparison with other systems that exists in the literature.

---

<sup>21</sup>König’s Lemma establishes that every connected, locally finite, infinite graph contains an infinite path.

### 4.1 Proof-Theoretic Reference

Note that the set of sentences of  $\mathcal{L}_T$  is p.r., as is as the occurrence relation that holds between a string of symbols of the language and each of its substrings. Let  $\text{Sent}(x)$  and  $\text{Occ}(x, y) \in \mathcal{L}$  express this set and relation, respectively, in a natural way.<sup>22</sup> Thus, m-reference is also a p.r. relation (cf. Definition 1), as it can be expressed in  $\mathcal{L}$  by the  $\Delta_0^0$  formula

$$\text{MRef}(x, y) := \text{Sent}(x \wedge y) \wedge \exists t (\text{Occ}(Tt, x) \wedge t^\circ = y)$$

where  $T$  is a function symbol of  $\mathcal{L}$  for the p.r. function that maps each term  $t$  to the formula  $Tt$  – and likewise for other predicates that will come up later.

Q-reference, instead, is much more complex, as it’s not even arithmetical (cf. Definition 3), i.e. expressible in  $\mathcal{L}$ . As is well known, truth-in- $\mathbb{N}$  can be expressed by a  $\Delta_1^1$  formula, and this is best possible. Note that q-reference and truth-in- $\mathbb{N}$  are inter-definable over a small class of arithmetical notions. On the one hand, truth-in- $\mathbb{N}$  is the only notion occurring in the recursive definition of q-reference that involves second-order quantification, as the value and the occurrence relations, the normalization function, and syntactic properties of expressions such as being a disjunction or a universal statement are all p.r. and, thus, arithmetical. Therefore, q-reference is expressible by a  $\Delta_1^1$  formula. On the other hand, truth-in- $\mathbb{N}$  can be defined for each sentence  $\varphi \in \mathcal{L}$  in terms of q-reference as follows:  $\mathbb{N} \models \varphi$  iff  $\forall x (x = \ulcorner \varphi \urcorner \wedge \varphi \rightarrow Tx)$  q-refers to  $\varphi$ . Thus, q-reference is  $\Delta_1^1$ , namely, *hyperarithmetical*.

By the usual complexity considerations, this result extends to d-reference, reference, chains of reference, self-reference, and also well-foundedness (cf. Definitions 4 and 5–8), which means there is no formula in  $\mathcal{L}$  that could be employed to restrict the instances of the T-schema or the Uniform T-schema in a first-order axiomatic system of disquotational truth. In other words, there is no straightforward way of axiomatizing the semantic theories introduced in Section 3.

Although there might be interesting, rather indirect, ways of axiomatizing this model, I opt instead to consider simpler, proof-theoretic versions of the concepts introduced in Section 2.3; in particular of well-foundedness. The strategy consists in replacing in Definition 3 the notion of truth-in- $\mathbb{N}$  with that of provability-in-PA, our base theory. Thus, all  $\Delta_1^1$  notions occurring in Definition 3 are replaced with  $\Sigma_1^0$  ones. All further reference notions defined in terms of q-reference are modified accordingly. The new concepts can be seen as approximations to the original ones. They turn out to be expressible in  $\mathcal{L}$  by formulae of fairly low quantificational complexity. This allows for the formulation of generous and well-motivated restrictions on disquotation, which result in powerful axiomatic truth theories, as will be seen soon.

Note that replacing in the definition of m-reference truth-in- $\mathbb{N}$  with provability-in-PA would not alter the extension of the relation, for all identity statements that

---

<sup>22</sup>That is, by means of  $\Delta_0^0$  formulae. Recall that every p.r. property (or relation) can be expressed in  $\mathcal{L}$  by a  $\Delta_0^0$  expression and every recursively enumerable property by a  $\Sigma_1^0$  expression. As PA decides all  $\Delta_0^0$  and proves all true  $\Sigma_1^0$  sentences – i.e. PA is  $\Sigma_1^0$ -complete – those expressions represent and weakly represent, respectively, the properties they express.

are true in  $\mathbb{N}$  are provable in PA and vice versa. As mentioned before, m-reference is already a p.r. notion. On the other hand, substituting provability-in-PA for truth-in- $\mathbb{N}$  in the definition of q-reference makes a significant difference: it results in a much simpler, recursively enumerable notion. As is well known, provability-in-PA is itself recursively enumerable. Thus, if we replace truth-in- $\mathbb{N}$  by provability-in-PA in the recursive definition of q-reference, it can be shown by induction on the depth of  $\chi$  (cf. Section 2.3) that the altered notion is recursively enumerable as well, since both the definiendum and provability-in-PA occur only positively in the definiens, provability-in-PA is only preceded by a string of existential quantifiers, and all other notions occurring in the definition are p.r., as it has already been argued. Therefore, the resulting proof-theoretic notion of q-reference *relative to PA* can be expressed in  $\mathcal{L}$  by a  $\Sigma_1^0$  formula,  $\text{QRef}_{\text{PA}}(x, y)$ .

We can find arithmetical formulae expressing proof-theoretic versions of direct reference, reference *simpliciter*, self-reference, and well-foundedness relative to PA in terms of  $\text{MRef}(x, y)$  and  $\text{QRef}_{\text{PA}}(x, y)$ . For d-reference (cf. Definition 4), let  $\text{DRef}_{\text{PA}}(x, y) := \text{MRef}(x, y) \vee \text{QRef}_{\text{PA}}(x, y)$ , which is trivially  $\Sigma_1^0$ .

For reference *simpliciter*, we first need to take a detour to define chain of reference. This requires some subtlety, for infinite sequences cannot be coded by natural numbers. However, finite sequences can. Moreover, note that, according to Definition 6, for  $\varphi$  to refer to  $\psi$  there must be a chain of reference starting with the former and ending with the latter, i.e. a *finite* chain. Thus, we can define reference in  $\mathcal{L}$  exclusively in terms of finite chains. Let  $\text{Seq}(x) \in \Delta_0^0$  express in  $\mathcal{L}$  the p.r. property of being a (finite) sequence, and let  $\text{len}(x)$  and  $(x)_i \in \Delta_0^0$  express the p.r. function that maps each sequence to its length and the one that maps each sequence  $x$  and number  $i$  to the  $i$ -th entry of  $x$ , if  $i$  is smaller than the length of  $x$ , respectively. The  $\Sigma_1^0$  formula

$$\text{CRef}_{\text{PA}}(x) := \text{Seq}(x) \wedge \forall y < \text{len}(x) \forall z < y \text{DRef}_{\text{PA}}((x)_z, (x)_{S_z})$$

expresses in  $\mathcal{L}$  the property of being a finite chain of reference relative to PA (cf. Definition 5). Thus, reference relative to PA can be expressed in  $\mathcal{L}$  by the  $\Sigma_1^0$  formula

$$\text{Ref}_{\text{PA}}(x, y) := \exists z (\text{CRef}_{\text{PA}}(z) \wedge (z)_0 = x \wedge (z)_{\text{len}(z)-1} = y)$$

Accordingly, self-reference relative to PA (cf. Definition 7) is expressed in  $\mathcal{L}$  by  $\text{SRef}_{\text{PA}}(x) := \text{Ref}_{\text{PA}}(x, x) \in \Sigma_1^0$ . Well-foundedness, on the other hand, requires a bit more work, as it's defined in terms of possibly infinite chains of reference (cf. Definition 8). However, this is not a hindrance, for there is an equivalent formulation of the definition in terms of finite chains only. We say a finite chain of reference  $\varphi_1, \dots, \varphi_n$  can be *extended* just in case there is a sentence  $\varphi_{n+1}$  such that  $\varphi_1, \dots, \varphi_n, \varphi_{n+1}$  is also a chain of reference. Clearly, a sentence  $\varphi$  is well-founded just in case every finite chain of reference starting with  $\varphi$  can be extended only a finite number of times. Thus, well-foundedness relative to PA is expressed by the  $\Pi_3^0$  formula

$$\begin{aligned} \text{Wf}_{\text{PA}}(x) &:= \forall y (\text{CRef}_{\text{PA}}(y) \wedge (y)_0 = x \rightarrow \\ &\exists z \forall w (\text{CRef}_{\text{PA}}(w) \wedge \text{len}(y) < \text{len}(w) \wedge \forall k < \text{len}(y) (y)_k = (w)_k \rightarrow \text{len}(w) < z)) \end{aligned}$$

This cumbersome formula states that for each finite chain of reference  $y$  starting with  $x$  there is a limit  $z$  to the length of every chain that extends  $y$ .

### 4.2 Unstable Reference

One would think our job is almost done now. It just remains to relativize disquotation to  $Wf_{PA}(x)$  and then extend PAT with the resulting instances (see Section 2.1), that is,

$$Wf_{PA}(\ulcorner\varphi\urcorner) \rightarrow (T\ulcorner\varphi\urcorner \leftrightarrow \varphi) \tag{8}$$

for each sentence  $\varphi \in \mathcal{L}_T$ , to obtain the desired truth system. Unfortunately, such a theory would be too weak and, what is worse, unsound.

Although many sentences turn out to be well-founded relative to PA, only very few are *provably* so in PA itself. Consider, for instance, the following:

$$\forall x (x = \ulcorner 0 = 0 \urcorner \rightarrow Tx) \tag{9}$$

Despite being obviously harmless, as it only refers to  $0 = 0$ , PA cannot prove that (9) isn't self-referential, for that would mean that PA proves

$$\neg Bew_{PA}(\ulcorner \forall x (x = \ulcorner 0 = 0 \urcorner \rightarrow Tx) \urcorner) = \ulcorner 0 = 0 \urcorner$$

In other words, PA would prove its own consistency, which is not possible by Gödel's second incompleteness theorem. In general, PA is not able to prove that a given sentence doesn't q-refer to an arbitrary expression, for it cannot prove that a given number or sequence of numbers provably fails to satisfy a given formula. Therefore, PA can only establish positive cases of q-reference plus those negative ones in which q-reference is nowhere restricted by a conditional – e.g. in  $\forall x T\neg x$  – or no (non-dummy) quantifiers occur in the formula.

An idea that suggests itself is to inform PA that it doesn't prove false statements, i.e. that it is sound. This can be done by extending PA with all instances of PA's *uniform reflection principle*, that is,

$$(URfn_{PA}) \quad \forall t_1 \dots t_n (Bew_{PA}(\ulcorner \varphi(t_1, \dots, t_n) \urcorner) \rightarrow \varphi(t_1^o, \dots, t_n^o))$$

where  $\varphi \in \mathcal{L}$ . Let URfn(PA) be the resulting theory. URfn(PA) is trivially sound, as all instances of URfn<sub>PA</sub> are true in  $\mathbb{N}$ . Nonetheless, if we extend URfn(PA) with all instances of (8), we obtain a trivial system.

**Observation 4** URfn(PA)  $\cup$   $\{Wf_{PA}(\ulcorner\varphi\urcorner) \rightarrow (T\ulcorner\varphi\urcorner \leftrightarrow \varphi) : \varphi \in \mathcal{L}_T\}$  is inconsistent.

*Proof* ‘Weakly’ diagonalizing (cf. Theorem 1) the predicate  $\neg Bew_{PA}(x)$  we obtain a sentence  $\gamma$  that is provably equivalent in PA to  $\neg Bew_{PA}(\ulcorner\gamma\urcorner)$  – i.e. PA's Gödel sentence. In turn, Strong Diagonalization (cf. Theorem 2) applied to  $\forall x (x = y \wedge \gamma \rightarrow \neg Tx)$  delivers a term  $l^*$  s.t.

$$PA \vdash l^* = \ulcorner \forall x (x = l^* \wedge \gamma \rightarrow \neg Tx) \urcorner$$

Note that  $\forall x (x = I^* \wedge \gamma \rightarrow \neg Tx)$  doesn't refer to anything in PA, for PA doesn't prove  $\text{Bew}_{\text{PA}}(\ulcorner \gamma \urcorner)$ . Moreover,  $\text{URfn}(\text{PA})$  knows this, as it famously implies  $\gamma$  and, therefore, that PA doesn't prove  $\gamma$ . As a consequence,  $\forall x \neg \text{Bew}_{\text{PA}}(\ulcorner x = I^* \wedge \gamma \urcorner)$  is a theorem of  $\text{URfn}(\text{PA})$ , and so is  $\text{Wf}_{\text{PA}}(I^*)$ . By an argument similar to that employed in the liar paradox, one can show that the corresponding instance of the T-schema for  $\forall x (x = I^* \wedge \gamma \rightarrow \neg Tx)$  entails a contradiction in  $\text{URfn}(\text{PA})$ .  $\square$

An immediate consequence of Observation 4 is that the theory that extends PAT with all instances of (8) is unsound. The reason behind this perplexing result is that being well-founded relative to PA does not mean that a sentence is actually well-founded, but rather that PA has no 'evidence' for its unfoundedness, i.e. it doesn't know that  $I^* = I^* \wedge \gamma$ . The proof of Observation 4 shows that, if more evidence is available, such as  $\gamma$ , a sentence might turn out to be unfounded relative to the extended system,  $\text{URfn}(\text{PA})$ . This is a consequence of the incompleteness of PA.

Relativizing q-reference and all the other reference concepts that depend on it to  $\text{URfn}(\text{PA})$  or a stronger system instead of PA will obviously not solve the problem, as any recursively axiomatizable theory will also be incomplete. The issue would emerge all the same, only at a higher level. A better way of bypassing the obstacle is to focus on sentences for which new evidence does not make a difference in the expressions they refer to. I call these sentences "reference-stable" or "r-stable", for short.<sup>23</sup> R-unstable expressions bear a certain analogy with blind truth ascriptions and contingent liars: in both cases we don't know what they express and, a fortiori, if they are paradoxical or not. Only for r-stable sentences we can be sure that their reference patterns are safe.

To properly define this class of sentences, we first need to introduce the idea of a directly-reference-stable sentence, or "dr-stable" for short. These are sentences that only contain reference-restricting conditionals with 'simple' antecedents, that is, antecedents that, if satisfied by a sequence of numbers, PA knows so. Thus, these antecedents must be provably equivalent in PA to  $\Sigma_1^0$  formulae.

**Definition 9 (Dr-stability)** Let  $\varphi$  be a sentence.  $\varphi$  is dr-stable iff in  $\varphi^*$  all subformulae of the form  $\forall v (\chi_1 \vee \chi_2)$  where T occurs in  $\chi_2$  but not in  $\chi_1$  are s.t.  $\chi_1$  is provably equivalent in PA to a  $\Pi_1^0$  formula.

For instance, whereas  $\forall x (x = I^* \rightarrow \neg Tx)$  and  $\text{GRfn}_{\text{PA}}$  – i.e.  $\forall x (\text{Bew}_{\text{PA}}(x) \rightarrow Tx)$  – are both dr-stable,  $\forall x (\neg \text{Bew}_{\text{PA}}(x) \rightarrow Tx)$  and  $\forall x (x = I^* \wedge \gamma \rightarrow \neg Tx)$  aren't. Since truth-in- $\mathbb{N}$  and provability-in-PA coincide for sentences provably equivalent in

---

<sup>23</sup> A complementary strategy would be to replace PA as our base theory with an extension of it by means of a hierarchy of iterated reflection principles structurally identical to  $\text{URfn}_{\text{PA}}$  into the transfinite. This would increase the number of expressions that are rightly declared well-founded and, consequently, the number of instances of disquotation that hold in the ultimate truth system. However, this would introduce unnecessary complications, as we will see that  $\text{URfn}(\text{PA})$  can already serve as base theory of considerably powerful truth systems.

PA to a  $\Sigma_1^0$  statement, q-reference in PA and q-reference *simpliciter* coincide in the case of dr-stable sentences. Thus, so do d-reference in PA and d-reference *simpliciter*, so we can drop the qualifier in these cases. The same cannot be said of reference itself. For a dr-stable sentence  $\varphi$  might refer to an expression  $\psi$  that is not itself dr-stable, so that the latter d-refers to a sentence  $\chi$  but doesn't d-refer to it in PA. Our original sentence  $\varphi$ , then, also refers to  $\chi$  but not in PA. Thus, we need the following definition.

**Definition 10** (R-stability) Let  $\varphi$  be a sentence.  $\varphi$  is r-stable iff it's dr-stable and only refers to dr-stable sentences.

Since  $\text{GRfn}_{\text{PA}}$  only refers to T-free sentences, it must be r-stable. On the other hand,  $\forall x (x = 1^* \rightarrow \neg \text{T}x)$  refers to  $\forall x (x = 1^* \wedge \gamma \rightarrow \neg \text{T}x)$  and, therefore, isn't r-stable. Since r-stable sentences are dr-stable and only refer to dr-stable sentences, it follows from our previous considerations that Definition 10 can be equivalently restated by replacing “refers” with “refers in PA”. Thus, reference in PA and reference *simpliciter* concur in the case of r-stable sentences. Moreover, note that dr-stability is expressible in  $\mathcal{L}$  by a  $\Sigma_1^0$  formula,  $\text{DRSt}(x)$ , as it's defined in terms of p.r. syntactic properties plus the property of being provably equivalent in PA to a  $\Pi_1^0$  formula. Therefore, by the usual complexity considerations, r-stability is expressible in  $\mathcal{L}$  by the following  $\Pi_2$  formula:

$$\text{RSt}(x) := \text{DRSt}(x) \wedge \forall y (\text{Ref}_{\text{PA}}(x, y) \rightarrow \text{DRSt}(y))$$

### 4.3 Restricting Disquotation

We now have all the resources needed to formulate sound axiomatic theories of disquotational truth. Let  $\text{URfn}(\text{PAT})$  be  $\text{PAT} + (\text{URfn}_{\text{PA}})$  and let  $\text{WFUTB}$  be the theory extending  $\text{URfn}(\text{PAT})$  with the following restricted version of the Uniform T-schema, where  $\mathbf{t}$  abbreviates  $t_1, \dots, t_n$ :

$$\forall \mathbf{t} (\text{RSt}(\ulcorner \varphi(\mathbf{t}) \urcorner) \wedge \text{Wf}_{\text{PA}}(\ulcorner \varphi(\mathbf{t}) \urcorner) \rightarrow (\text{T}\ulcorner \varphi(\mathbf{t}) \urcorner \leftrightarrow \varphi(\mathbf{t}^\circ)))$$

$\text{WFUTB}$  – for “Well-founded Uniform Tarski Biconditionals” – uniformly entails all instances of the T-schema for sentences that are provably r-stable and well-founded in PA.

**Proposition 6** *WFUTB is  $\omega$ -consistent.*

*Proof* Since every r-stable sentence that is well-founded in PA is also well-founded *simpliciter*, by Proposition 3 the axioms of  $\text{WFUTB}$  are all true in  $\langle \mathbb{N}, \Phi_{\text{wf}} \rangle$ .  $\square$

Despite being a disquotational theory,  $\text{WFUTB}$  is proof-theoretically strong, as it can relatively interpret the theory of ramified truth up to  $\Gamma_0$ ,  $\text{RT}_{<\Gamma_0}$ . This consists of iterations of the theory of typed compositional truth over  $\text{PAT}$  up to the ordinal  $\Gamma_0$ , the

Feferman-Schütte ordinal.<sup>24</sup> Each of these iterations requires its own truth predicate, bringing about a hierarchy of corresponding languages.

As is well known, natural numbers can effectively codify ordinals up to  $\Gamma_0$ .<sup>25</sup> Assuming a fixed (effective) coding, if  $\alpha < \Gamma_0$ , we write  $\bar{\alpha}$  for the numeral of the code of  $\alpha$ . For clarity, I often identify ordinals with their codes if there is no room for confusion. Based on this coding, PA can talk about ordinals and represent some of their properties. For instance, the set of ordinals below  $\Gamma_0$  can be represented by a formula of  $\mathcal{L}$ ,  $\text{Ord}(x)$ , as well as their ordering, say, by the formula  $x < y$ . Let  $\forall\alpha\varphi$  and  $\exists\alpha\varphi$  be short for  $\forall v(\text{Ord}(v) \rightarrow \varphi)$  and  $\exists v(\text{Ord}(v) \wedge \varphi)$ , respectively, where  $v$  is a suitable variable. As is also well known, PA can prove all instances of transfinite induction up to  $\epsilon_0 (< \Gamma_0)$ , i.e.

$$(TI_{\xi}) \quad \forall\alpha (\forall\beta < \alpha \varphi(\beta) \rightarrow \varphi(\alpha)) \rightarrow \forall\alpha < \bar{\xi} \varphi(\alpha)$$

for  $\varphi \in \mathcal{L}$  and  $\xi < \epsilon_0$ .<sup>26</sup> This means PA can prove that all ordinals below  $\epsilon_0$  are well ordered.

Let  $\mathcal{L}_{<\alpha}$  be  $\mathcal{L}$  and, for every  $\alpha$  such that  $0 < \alpha \leq \Gamma_0$ , let  $\mathcal{L}_{<\alpha}$  be the result of extending  $\mathcal{L}$  with monadic predicates  $T_{\beta}$ , for each  $\beta < \alpha$ .  $\text{PAT}_{<\alpha}$  consists of the axioms of PA formulated in  $\mathcal{L}_{<\alpha}$ , with the induction schema extended to the whole language. For each  $\alpha < \Gamma_0$ , let  $\text{Sent}(x, \alpha) \in \mathcal{L}$  – or  $\text{Sent}_{\alpha}(x)$ , for short – represent the set of sentences of  $\mathcal{L}_{<\alpha}$ . We define a cumulative hierarchy of compositional truth theories formulated in these languages, as follows.

Let  $\text{RT}_{<0}$  be PA and, for every  $\alpha$  such that  $0 < \alpha \leq \Gamma_0$ , let  $\text{RT}_{<\alpha}$  be the theory formulated in  $\mathcal{L}_{<\alpha}$  that extends  $\text{PAT}_{<\alpha}$  with every instance of following axiom-schemata, for each  $\xi < \beta < \alpha$ :

- (RT $_{\beta}$ 1)  $\forall s \forall t (T_{\beta}(s \doteq t) \leftrightarrow s^{\circ} = t^{\circ})$
- (RT $_{\beta}$ 2)  $\forall x ( \text{Sent}_{\bar{\beta}}(x) \rightarrow (T_{\beta} \neg x \leftrightarrow \neg T_{\beta} x) )$
- (RT $_{\beta}$ 3)  $\forall x \forall y ( \text{Sent}_{\bar{\beta}}(x \wedge y) \rightarrow (T_{\beta}(x \wedge y) \leftrightarrow (T_{\beta} x \wedge T_{\beta} y)) )$
- (RT $_{\beta}$ 4)  $\forall x \forall y ( \text{Sent}_{\bar{\beta}}(x \vee y) \rightarrow (T_{\beta}(x \vee y) \leftrightarrow (T_{\beta} x \vee T_{\beta} y)) )$
- (RT $_{\beta}$ 5)  $\forall x \forall y ( \text{Sent}_{\bar{\beta}}(x \rightarrow y) \rightarrow (T_{\beta}(x \rightarrow y) \leftrightarrow (T_{\beta} x \rightarrow T_{\beta} y)) )$
- (RT $_{\beta}$ 6)  $\forall x \forall y ( \text{Sent}_{\bar{\beta}}(\forall y x) \rightarrow (T_{\beta}(\forall y x) \leftrightarrow \forall t T_{\beta} x(t/y)) )$
- (RT $_{\beta}$ 7)  $\forall x \forall y ( \text{Sent}_{\bar{\beta}}(\exists y x) \rightarrow (T_{\beta}(\exists y x) \leftrightarrow \exists t T_{\beta} x(t/y)) )$
- (RT $_{\beta}$ 8)  $\forall t ( \text{Sent}_{\bar{\xi}}(t^{\circ}) \rightarrow (T_{\beta}(T_{\xi} t) \leftrightarrow T_{\xi} t^{\circ}) )$
- (RT $_{\beta}$ 9)  $\forall t \forall \zeta < \bar{\beta} ( \text{Sent}_{\zeta}(t^{\circ}) \rightarrow (T_{\beta}(T_{\zeta} t) \leftrightarrow T_{\beta} t^{\circ}) )$

where  $\forall$  and  $\exists$  are function symbols for the p.r. functions mapping a variable  $v$  and a formula  $\varphi$  to  $\forall v \varphi$  and  $\exists v \varphi$ , respectively. Note that the axioms never quantify over the subindex  $\beta$  of  $T_{\beta}$ , on pain of triviality. But they do quantify over the predicates themselves up to a given ordinal  $\beta$  in  $\text{RT}_{\beta}$ 9. This set is p.r. for each  $\beta < \epsilon_0$  and

<sup>24</sup>See Halbach [10, chap. 8].

<sup>25</sup>More comprehensive ordinal notation systems are also possible, but this is enough for our purposes.

<sup>26</sup>See Pohlers [28, chap. 3].



representable also for  $\beta \geq \epsilon_0$  insofar as all arithmetical instances of  $TI_\beta$  hold in the theory.

$RT_{<1}$  is just the theory of typed compositional truth CT where all occurrences of T have been replaced with  $T_0$  (also within corner quotes), for  $RT_08$  and  $RT_09$  are vacuously true. Similarly, for  $1 < \alpha < \Gamma_0$ , axioms  $RT_\beta 1$ , with  $\beta < \alpha$ , establish all instances of disquotation for identity statements for  $T_\beta$ , and  $RT_{\beta 2}$ - $RT_{\beta 7}$  lay down the compositional character of these truth predicates.  $RT_{\beta 8}$ , in turn, provides instances of disquotation for each truth ascription that belongs to lower levels of the hierarchy, supplementing  $RT_{\beta 1}$ . Finally,  $RT_{\beta 9}$  establishes that the hierarchy is cumulative.

**Proposition 7** *The theory of ramified truth up to  $\Gamma_0$ ,  $RT_{<\Gamma_0}$ , is relatively interpretable in WFUTB.*

*Proof* I first show that  $\mathcal{L}_T$  contains a formula  $\vartheta_{\overline{\beta}}(x)$  that behaves in WFUTB like  $T_\beta(x)$  in  $RT_{<\epsilon_0}$ , for each  $\beta < \epsilon_0$ , following Halbach’s [10] demonstration of his Lemma 15.24. I then extend this result to every  $\beta < \Gamma_0$ .

If  $\varphi(x)$  is a formula of  $\mathcal{L}_T$  containing T, let  $\mathcal{L}_\varphi \subseteq \mathcal{L}_T$  be the language that extends  $\mathcal{L}$  with  $\varphi(x)$  as if it were just a predicate symbol, that is, T occurs in formulae of  $\mathcal{L}_\varphi$  only within  $\varphi(t)$ , where  $t$  is a term.<sup>27</sup> The relation that holds between a sentence  $\psi$  and a formula  $\varphi$  just in case  $\psi$  is a sentence of  $\mathcal{L}_\varphi$  is p.r. Therefore, it can be represented by a  $\Delta_0^0$  expression,  $Sent\downarrow(x, y)$ . Strongly diagonalizing the formula

$$\begin{aligned} & \exists s \exists t (x = (s \equiv t) \wedge s^\circ = t^\circ) \vee \\ & \exists z ((Sent_{\overline{0}}(z) \vee \exists \zeta \prec y \text{ Sent}\downarrow(z, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = \neg z \wedge \neg Tk(\dot{y}/\Gamma y^\neg)(\dot{z}/\Gamma x^\neg)) \vee \\ & \exists z \exists w ((Sent_{\overline{0}}(z \wedge w) \vee \exists \zeta \prec y \text{ Sent}\downarrow(z \wedge w, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = (z \wedge w) \wedge \\ & \quad (Tk(\dot{y}/\Gamma y^\neg)(\dot{z}/\Gamma x^\neg) \wedge Tk(\dot{y}/\Gamma y^\neg)(\dot{w}/\Gamma x^\neg))) \vee \\ & \exists z \exists w ((Sent_{\overline{0}}(z \wedge w) \vee \exists \zeta \prec y \text{ Sent}\downarrow(z \wedge w, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = (z \vee w) \wedge \\ & \quad (Tk(\dot{y}/\Gamma y^\neg)(\dot{z}/\Gamma x^\neg) \vee Tk(\dot{y}/\Gamma y^\neg)(\dot{w}/\Gamma x^\neg))) \vee \\ & \exists z \exists w ((Sent_{\overline{0}}(z \rightarrow w) \vee \exists \zeta \prec y \text{ Sent}\downarrow(z \rightarrow w, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = (z \rightarrow w) \wedge \\ & \quad (Tk(\dot{y}/\Gamma y^\neg)(\dot{z}/\Gamma x^\neg) \rightarrow Tk(\dot{y}/\Gamma y^\neg)(\dot{w}/\Gamma x^\neg))) \vee \\ & \exists z \exists w ((Sent_{\overline{0}}(\forall z w) \vee \exists \zeta \prec y \text{ Sent}\downarrow(\forall z w, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = \forall z w \wedge \\ & \quad \forall t Tk(\dot{y}/\Gamma y^\neg)(\Gamma \dot{w}(\dot{t}/\dot{z})^\neg/\Gamma x^\neg)) \vee \\ & \exists z \exists w ((Sent_{\overline{0}}(\exists z w) \vee \exists \zeta \prec y \text{ Sent}\downarrow(\exists z w, k(\dot{\zeta}/\Gamma y^\neg))) \wedge x = \exists z w \wedge \\ & \quad \exists t Tk(\dot{y}/\Gamma y^\neg)(\Gamma \dot{w}(\dot{t}/\dot{z})^\neg/\Gamma x^\neg)) \vee \\ & \exists t \exists \xi \prec y ((Sent_{\overline{0}}(t^\circ) \vee \exists \xi \prec \zeta \text{ Sent}\downarrow(t^\circ, k(\dot{\xi}/\Gamma y^\neg))) \wedge x = k(\dot{\zeta}/\Gamma y^\neg)(t/\Gamma x^\neg) \wedge \\ & \quad Tk(\dot{\zeta}/\Gamma y^\neg)(t/\Gamma x^\neg)) \end{aligned}$$

<sup>27</sup>More precisely, we add the following recursion clause to the definition of well-formed formula: if  $t$  is a term, then  $\varphi(t)$  is a well-formed formula.

over the free variable  $k$ , we obtain a predicate  $\vartheta(x, y) - \vartheta_y(x) - \text{s.t.}$ <sup>28</sup>

$$\begin{aligned} \text{PAT} \vdash \vartheta_y(x) \leftrightarrow & \exists s \exists t (x = (s \dot{=} t) \wedge s^\circ = t^\circ) \vee & (10) \\ & \exists z (\text{Sent}_{<y}(z) \wedge x = \neg z \wedge \neg \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{z})^\Gamma) \vee \\ & \exists z \exists w (\text{Sent}_{<y}(z \wedge w) \wedge x = (z \wedge w) \wedge (\text{T}^\Gamma \vartheta_{\dot{y}}(\dot{z})^\Gamma \wedge \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{w})^\Gamma)) \vee \\ & \exists z \exists w (\text{Sent}_{<y}(z \vee w) \wedge x = (z \vee w) \wedge (\text{T}^\Gamma \vartheta_{\dot{y}}(\dot{z})^\Gamma \vee \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{w})^\Gamma)) \vee \\ & \exists z \exists w (\text{Sent}_{<y}(z \rightarrow w) \wedge x = (z \rightarrow w) \wedge (\text{T}^\Gamma \vartheta_{\dot{y}}(\dot{z})^\Gamma \rightarrow \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{w})^\Gamma)) \vee \\ & \exists z \exists w (\text{Sent}_{<y}(\forall z w) \wedge x = \forall z w \wedge \forall t \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{w}(\dot{i}/\dot{z}))^\Gamma) \vee \\ & \exists z \exists w (\text{Sent}_{<y}(\exists z w) \wedge x = \exists z w \wedge \exists t \text{T}^\Gamma \vartheta_{\dot{y}}(\dot{w}(\dot{i}/\dot{z}))^\Gamma) \vee \\ & \exists t \exists \zeta < y (\text{Sent}_{<\zeta}(t^\circ) \wedge x = \ulcorner \vartheta_{\dot{\zeta}}(t)^\Gamma \wedge \text{T}^\Gamma \vartheta_{\dot{\zeta}}(t)^\Gamma \end{aligned}$$

where  $\text{Sent}_{<s}(t) := \text{Sent}_0(t) \vee \exists \zeta < s \text{ Sent} \upharpoonright (t, \ulcorner \vartheta_{\dot{\zeta}}(x)^\Gamma) \in \Delta_0^0$ ,  $s$  and  $t$  are terms and  $\zeta$  is a suitable ordinal variable.<sup>29</sup> Let  $\mathcal{L}_{<\vartheta_\alpha}$  be the language that results from merging together every  $\mathcal{L}_{\vartheta_\beta}$ , with  $\beta < \alpha$ . Thus, for every  $\alpha < \epsilon_0$ ,  $\text{Sent}_{<\alpha}(x)$  represents the set of sentences of  $\mathcal{L}_{<\vartheta_\alpha}$ .

Given that (10) was obtained by Strong Diagonalization, we can identify the predicate  $\vartheta_y(x)$  with the right-hand side of the biconditional, except all occurrences of  $\ulcorner \vartheta_y(x)^\Gamma$  are actually occurrences of a more complex term that denotes  $\vartheta_y(x)$ . Since alethic reference is closed under Leibniz’s Law, this difference won’t actually make a difference.

To show that, for each  $\beta < \epsilon_0$ ,  $\vartheta_{\bar{\beta}}(x)$  behaves like  $\text{T}_\beta(x)$ , we first need to remove the occurrences of T in (10), that is, we need to show that WFUTB contains an instance of the Uniform T-schema for each  $\vartheta_\beta(x)$ , with  $\beta < \epsilon_0$ . In other words, I give a uniform proof that every instance of each of these formulae is r-stable and well-founded in PA.

Note first that the normalization of  $\vartheta_y(x)$  is a disjunction of negated universally quantified statements followed by just one reference-restricting conditional, except for the first disjunct that doesn’t contain T. Moreover, the antecedents of these conditionals are all  $\Delta_0^0$ . This means, on the one hand, that each  $\vartheta_{\bar{\beta}}(t)$  is dr-stable, so

$$\text{PA} \vdash \forall \beta \forall t \text{DRSt}(\ulcorner \vartheta_{\bar{\beta}}(t)^\Gamma) \tag{11}$$

that is, if  $\vartheta_{\bar{\beta}}(t)$  q-refers to a sentence in PA, then PA knows about it. But also, the fact that all the antecedents are  $\Delta_0^0$  implies that URfn(PA) knows about all *negative* cases of q-reference for each  $\vartheta_{\bar{\beta}}(t)$  as well. For instance, if  $\text{Sent}_{<\bar{\beta}}(z \wedge w) \wedge t = (z \wedge w)$  isn’t true of  $m, n \in \omega$ , then PA proves the negation of  $\text{Sent}_{<\bar{\beta}}(\bar{m} \wedge \bar{n}) \wedge t = (\bar{m} \wedge \bar{n})$ , which means that URfn(PA) proves that PA doesn’t prove  $\text{Sent}_{<\bar{\beta}}(\bar{m} \wedge \bar{n}) \wedge t = (\bar{m} \wedge \bar{n})$ .

<sup>28</sup>Note that the fourth disjunct on the right-hand side of the following biconditional requires that the conjunction of  $z$  and  $w$  is a sentence of  $\mathcal{L}_{<y}$ , instead of their disjunction. This guarantees that  $z$  and  $w$  are sentences of the language themselves, and not two disjuncts conforming an expression of the form  $\vartheta_k(t)$ , with  $k < y$ , which is always a disjunction. This way, only one disjunct in  $\vartheta_y(x)$  can be satisfied at a time.

<sup>29</sup>Although strictly speaking this formula is not  $\Delta_0^0$ , as it quantifies over the possibly infinitely many ordinals smaller than  $\zeta$ , it’s easy to see that only finitely many cases need to be checked for each  $t$  and each  $s$ , for every sentence of  $\mathcal{L}_\uparrow$  is of finite length. Thus,  $\text{Sent}_{<s}(t)$  is provably equivalent in PA to a  $\Delta_0^0$  expression.

Furthermore, since each of the antecedents holds exactly of one natural number or ordered pair of natural numbers, URfn(PA) can also prove general facts about the sentences each  $\vartheta_{\bar{\beta}}(t)$  q-refers to. For example, it follows in URfn(PA) that, if  $t$  denotes an identity statement or doesn't denote a sentence of  $\mathcal{L}_{<\vartheta_{\beta}}$ , then  $\vartheta_{\bar{\beta}}(t)$  doesn't refer to any expression.

By (10) it also follows easily in URfn(PA) that each  $\vartheta_{\bar{\beta}}(t)$  d-refers just to sentences of the form  $\vartheta_{\bar{\zeta}}(s)$ , where  $\zeta < \beta$  or  $s$  denotes either a sentence in  $\mathcal{L}_{<\vartheta_{\beta}}$  of lower complexity than  $t$ 's or a sentence in  $\mathcal{L}_{<\vartheta_{\zeta}}$ , for some  $\zeta < \beta$ . Thus, if  $\alpha < \epsilon_0$ , we can prove that

$$\text{URfn(PA)} \vdash \forall \beta < \bar{\alpha} \forall t (\text{RSt}(\ulcorner \vartheta_{\bar{\beta}} : (t) \urcorner) \wedge \text{Wf}_{\text{PA}}(\ulcorner \vartheta_{\bar{\beta}} : (t) \urcorner)) \tag{12}$$

by internal transfinite induction up to  $\alpha$  on  $\beta$ . We reason informally in URfn(PA). Assume  $\forall \zeta < \beta \forall t (\text{RSt}(\ulcorner \vartheta_{\bar{\zeta}} : (t) \urcorner) \wedge \text{Wf}_{\text{PA}}(\ulcorner \vartheta_{\bar{\zeta}} : (t) \urcorner))$ . We derive  $\forall t (\text{RSt}(\ulcorner \vartheta_{\bar{\beta}} : (t) \urcorner) \wedge \text{Wf}_{\text{PA}}(\ulcorner \vartheta_{\bar{\beta}} : (t) \urcorner))$  by an internal complete induction on the complexity of the sentence denoted by  $t^\circ$ .

If  $t^\circ$  doesn't denote a sentence of  $\mathcal{L}_{<\vartheta_{\beta}}$  or denotes an identity statement, the result follows trivially, given (11). Assume  $\text{RSt}(\ulcorner \vartheta_{\bar{\beta}} : (s) \urcorner) \wedge \text{Wf}_{\text{PA}}(\ulcorner \vartheta_{\bar{\beta}} : (s) \urcorner)$  for every  $s$  s.t.  $s^\circ$  denotes a sentence of lower complexity than  $t^\circ$ 's. If  $t^\circ$  denotes a negation  $\neg\varphi \in \mathcal{L}_{<\vartheta_{\beta}}$ , by (10)  $\vartheta_{\bar{\beta}}(t^\circ)$  d-refers to  $\vartheta_{\bar{\beta}}(\ulcorner \varphi \urcorner)$  that, by i.h., is both r-stable and well-founded in PA. By (11), so is  $\vartheta_{\bar{\beta}}(t^\circ)$ . The cases for the other connectives and the quantifiers are similar.<sup>30</sup> Let  $t^\circ$  denote  $\vartheta_{\bar{\zeta}}(s)$ , for some  $\zeta < \beta$ . By the last line of (10),  $\vartheta_{\bar{\beta}}(t^\circ)$  d-refers to  $\vartheta_{\bar{\zeta}}(s)$  which, by our main i.h. and (11), is r-stable and well-founded in PA. This completes the proof of (12).

As a consequence, for every  $\alpha < \epsilon_0$  the following principle of disquotation holds in WFUTB:

$$\forall \beta < \bar{\alpha} \forall t (\text{T} \ulcorner \vartheta_{\bar{\beta}} : (t) \urcorner \leftrightarrow \vartheta_{\bar{\beta}}(t^\circ)) \tag{13}$$

We are now in a position to show that  $\text{RT}_{<\epsilon_0}$  is relatively interpretable in WFUTB.

Consider the following translation function  $\eta : \mathcal{L}_{<\Gamma_0} \rightarrow \mathcal{L}_{\text{T}}$ :

$$\eta(\varphi) := \begin{cases} \varphi & \text{if } \varphi := s = t \\ \neg\eta(\psi) & \text{if } \varphi := \neg\psi \\ \eta(\psi) \wedge \eta(\chi) & \text{if } \varphi := (\psi \wedge \chi) \\ \eta(\psi) \vee \eta(\chi) & \text{if } \varphi := (\psi \vee \chi) \\ \eta(\psi) \rightarrow \eta(\chi) & \text{if } \varphi := (\psi \rightarrow \chi) \\ \forall v \eta(\psi) & \text{if } \varphi := \forall v \psi \\ \exists v \eta(\psi) & \text{if } \varphi := \exists v \psi \\ \vartheta_{\bar{\alpha}}(\eta(t)) & \text{if } \varphi := \text{T}_{\alpha} t \end{cases}$$

where  $\eta$  is a term of  $\mathcal{L}$  representing  $\eta$ . We know such a function exists and is p.r. by Kleene's Recursion Theorem. This way,  $\eta$  not only translates truth predicates occurring in a formula but also those that occur inside corner quotes, as well as those inside corner quotes inside corner quotes, etc. For instance, the translation of  $\text{T}_{\bar{\omega}}(\ulcorner \text{T}_0(\ulcorner 0 = 0 \urcorner) \urcorner)$  is not  $\vartheta_{\bar{\omega}}(\ulcorner \text{T}_0(\ulcorner 0 = 0 \urcorner) \urcorner)$  but  $\vartheta_{\bar{\omega}}(\ulcorner \vartheta_{\bar{0}}(\ulcorner 0 = 0 \urcorner) \urcorner)$ .

<sup>30</sup>A word of caution about the disjunction case is in order, however. Note that, if  $t^\circ$  denotes a disjunction, not only the restricting antecedent on the fourth line of (10) might be true, but also the one on the last line, since each formula of the form  $\vartheta_{\bar{\zeta}}(s)$  is itself a disjunction.

In what follows, I show that WFUTB proves the translation of every axiom of  $RT_{<\epsilon_0}$ . Actually, I prove *general, quantified* versions of the translations to facilitate the derivation of  $RT_{\beta 9}$ . I only deal with  $RT_{\beta 1}$ ,  $RT_{\beta 6}$ ,  $RT_{\beta 8}$ , and  $RT_{\beta 9}$ . The other cases can be proved in a similar fashion. Let  $\alpha < \epsilon_0$ . We reason informally in WFUTB under the assumption that  $\beta < \bar{\alpha}$ . By (10) we have that

$$\forall \beta < \bar{\alpha} (\vartheta_{\beta}(s \doteq t) \leftrightarrow s^{\circ} = t^{\circ})$$

which entails the translation of  $RT_{\beta 1}$ . For the other axioms, note that

$$PA \vdash \forall \beta < \bar{\alpha} \forall x (\text{Sent}_{\beta}(x) \leftrightarrow \text{Sent}_{<\beta}(\eta(x))) \tag{14}$$

by  $TI_{\alpha}$ . Thus, we have that

$$\text{Sent}_{\beta}(\forall y x) \rightarrow \text{Sent}_{<\beta}(\eta(\forall y x)) \tag{14}$$

$$\begin{aligned} &\rightarrow \text{Sent}_{<\beta}(\forall y \eta(x)) \text{Def. } \eta \\ &\rightarrow (\vartheta_{\beta}(\forall y \eta(x)) \leftrightarrow \forall t T^{\Gamma} \vartheta_{\beta}(\eta(\dot{x})(\dot{t}/\dot{y}))^{\neg}) \end{aligned} \tag{10}$$

$$\rightarrow (\vartheta_{\beta}(\forall y \eta(x)) \leftrightarrow \forall t \vartheta_{\beta}(\eta(x)(t/y))) \tag{13}$$

which implies the translation of  $RT_{\beta 6}$ . For  $RT_{\beta 8}$ 's, assume  $\xi < \beta$ :

$$\text{Sent}_{\xi}(t^{\circ}) \rightarrow \text{Sent}_{<\xi}(\eta(t^{\circ})) \tag{14}$$

$$\rightarrow (\vartheta_{\beta}(\ulcorner \vartheta_{\xi}(\eta(t))^{\neg} \urcorner) \leftrightarrow T^{\Gamma} \vartheta_{\xi}(\eta(t))^{\neg}) \tag{10}$$

$$\rightarrow (\vartheta_{\beta}(\ulcorner \vartheta_{\xi}(\eta(t))^{\neg} \urcorner) \leftrightarrow \vartheta_{\xi}(\eta(t^{\circ}))) \tag{13}$$

Finally, we prove the translation of  $RT_{\beta 9}$ , that is,<sup>31</sup>

$$\forall t \forall \zeta < \bar{\beta} (\text{Sent}_{\zeta}(t^{\circ}) \rightarrow (\vartheta_{\bar{\beta}}(\ulcorner \vartheta_{\zeta}(\eta(t))^{\neg} \urcorner) \leftrightarrow \vartheta_{\bar{\beta}}(\eta(t^{\circ}))))$$

Assume  $\zeta < \beta$  and  $\text{Sent}_{\zeta}(t^{\circ})$ . By (14), we have that  $\text{Sent}_{<\zeta}(\eta(t^{\circ}))$ . By our proof of the translation of  $RT_{\beta 8}$ , in which  $\xi$  is a *variable*, we know that  $\vartheta_{\beta}(\ulcorner \vartheta_{\zeta}(\eta(t))^{\neg} \urcorner) \leftrightarrow \vartheta_{\zeta}(\eta(t^{\circ}))$ . Thus, we just need to show that  $\vartheta_{\zeta}(\eta(t^{\circ})) \leftrightarrow \vartheta_{\beta}(\eta(t^{\circ}))$ , for every term  $t$ . I prove it by an internal complete induction on the complexity of  $t^{\circ}$ .

If  $t^{\circ} = (t_1 \doteq t_2)$ , we have that

$$\begin{aligned} \vartheta_{\zeta}(\eta(t^{\circ})) &\leftrightarrow \vartheta_{\zeta}(t_1 \doteq t_2) && \text{Def. } \eta \\ &\leftrightarrow t_1^{\circ} = t_2^{\circ} && \text{Proof of } RT_{\beta 1} \\ &\leftrightarrow \vartheta_{\beta}(t_1 \doteq t_2) && \text{Proof of } RT_{\beta 1} \\ &\leftrightarrow \vartheta_{\beta}(\eta(t^{\circ})) && \text{Def. } \eta \end{aligned}$$

<sup>31</sup>Note that  $RT_{\beta 9}$  does not have a ‘corresponding’ disjunct in (10), as it would clash with that for  $RT_{\beta 8}$  in the sense that there would be two simultaneously satisfiable disjuncts, blocking the proof of  $RT_{\beta 8}$ .

which completes the base case. Assume the result holds for every term denoting a sentence of lower complexity than  $t^\circ$ 's. If  $t^\circ = \neg s$ , we have the following:

$$\begin{aligned}
 \vartheta_\zeta(\eta(t^\circ)) &\leftrightarrow \vartheta_\zeta(\neg\eta(s)) && \text{Def. } \eta \\
 &\leftrightarrow \neg\vartheta_\zeta(\eta(s)) && \text{Proof of RT}_{\beta}2 \\
 &\leftrightarrow \neg\vartheta_\beta(\eta(s)) && \text{Ind. hyp.} \\
 &\leftrightarrow \vartheta_\beta(\neg\eta(s)) && \text{Proof of RT}_{\beta}2 \\
 &\leftrightarrow \vartheta_\beta(\eta(t^\circ)) && \text{Def. } \eta
 \end{aligned}$$

The cases for the other logical operators follow from the i.h. in a similar way. For the last case, assume  $t^\circ = \ulcorner T_\xi(s) \urcorner$ , for some  $\xi < \zeta$ . We then have that

$$\begin{aligned}
 \vartheta_\zeta(\eta(t^\circ)) &\leftrightarrow \vartheta_\zeta(\ulcorner \vartheta_\xi(\eta(s)) \urcorner) && \text{Def. } \eta \\
 &\leftrightarrow \vartheta_\xi(\eta(s)) && \text{Proof of RT}_{\beta}8 \\
 &\leftrightarrow \vartheta_\beta(\ulcorner \vartheta_\xi(\eta(s)) \urcorner) && \text{Proof of RT}_{\beta}8 \\
 &\leftrightarrow \vartheta_\beta(\eta(t^\circ)) && \text{Def. } \eta
 \end{aligned}$$

This completes the proof of the relative interpretability of  $\text{RT}_{<\epsilon_0}$  in WFUTB. Therefore, WFUTB entails every arithmetical theorem of  $\text{RT}_{<\epsilon_0}$ . This includes all instances of transfinite induction for formulae of  $\mathcal{L}$  up to  $\alpha$ , for some  $\alpha > \epsilon_0$ , i.e. WFUTB can prove that larger segments of the ordinal notation system are well-ordered. This implies that our proofs can be extended to show that  $\text{RT}_{<\alpha}$  is relatively interpreted in WFUTB, as only arithmetical transfinite induction has been employed so far. In turn,  $\text{RT}_{<\alpha}$  allows us to prove further instances of arithmetical transfinite induction, which means WFUTB can actually relatively interpret more iterations of compositionality, and so on. The progression reaches a fixed point at  $\Gamma_0$ , completing the proof.<sup>32</sup> □

Whilst Proposition 6 establishes the soundness of WFUTB, Proposition 7 shows that the theory is well positioned with respect to some of the better-known truth systems in terms of proof-theoretic strength. For instance, the Kripke-Feferman theory KF, that classically axiomatizes the family of Kripke’s fixed-point models over the Strong Kleene evaluation scheme, is relatively interpretable already in  $\text{RT}_{<\epsilon_0}$ . So is Halbach’s PUTB. In turn, FS, for “Friedman-Sheard”, is even weaker than the latter.<sup>33</sup> Proposition 7 shows that WFUTB is stronger than these four renowned systems.<sup>34</sup>

<sup>32</sup>See Feferman [4], Halbach [10, sec. 22.2].

<sup>33</sup>See Feferman [4] for KF, Halbach [9] for PUTB, and Friedman & Sheard [7] and Halbach [8] for FS.

<sup>34</sup>These results make WFUTB and other (classical disquotational) systems that will be introduced next attractive candidates for minimalist theories of truth. As it’s often quoted, according to Horwich [16, p. 42] “the principles governing our selection of excluded instances are, in order of priority: (a) that the minimal theory not engender ‘liar-type’ contradictions; (b) that the set of excluded instances be as small as possible; and – perhaps just as important as (b) – (c) that there be a constructive specification of the excluded instances that is as simple as possible.”

Yet, WFUTB can be soundly strengthened in several ways. For instance, we can close WFUTB’s truth predicate under provable equivalence in URfn(PAT), that is, we can replace WFUTB’s truth-theoretic axiom schema with the following:

$$\forall \mathbf{t} \forall x (\text{RSt}(x(\mathbf{t})) \wedge \text{Wf}_{\text{PA}}(x(\mathbf{t})) \wedge \text{Bew}_{\text{URfn(PAT)}}(\ulcorner \varphi(\mathbf{t}) \urcorner \leftrightarrow x(\mathbf{t})) \rightarrow (\text{T}\ulcorner \varphi(\mathbf{t}) \urcorner \leftrightarrow \varphi(\mathbf{t}^\circ)))$$

where  $\varphi$  is a formula with exactly  $n$  free variables  $v_1, \dots, v_n$ ,  $x(\mathbf{t})$  is short for  $x(t_1/\ulcorner v_1 \urcorner) \dots (t_n/\ulcorner v_n \urcorner)$ , and  $\text{Bew}_{\text{URfn(PAT)}}(x)$  weakly represents provability in URfn(PAT). Call the strengthened system WFUTB<sup>+</sup>. In WFUTB<sup>+</sup>, not only r-stable well-founded sentences have corresponding instances of disquotation, but also those that are provably equivalent in URfn(PAT) to r-stable well-founded sentences. This includes, among other expressions, unfounded logical truths and falsities such as  $\forall x (Tx \rightarrow Tx)$  and  $\exists x (0 \neq 0 \wedge Tx)$ , which are obviously harmless but are excluded from WFUTB’s truth principles.

The  $\omega$ -consistency of WFUTB<sup>+</sup> follows from Lemma 2, which establishes that every sentence  $\varphi$  depends on  $\Theta_\varphi$ , the set of sentences  $\varphi$  d-refers to. This implies that well-founded sentences can be shown to be grounded in the sense of Leitgeb.<sup>35</sup> Since dependence is closed under equivalence in every expansion of  $\mathbb{N}$  to  $\mathcal{L}_T$  and, a fortiori, in URfn(PAT), so is Leitgeb’s truth predicate, which means that his semantic theory  $\langle \mathbb{N}, \Phi_{\text{If}} \rangle$  is a model of WFUTB<sup>+</sup>.<sup>36</sup>

Alternatively, we could extend WFUTB along the lines suggested by the results established in Section 3.2. By Proposition 4, any theory extending WFUTB with finitely many instances of the T-schema for r-stable non-self-referential non-well-founded sentences is  $\omega$ -consistent. For example, we could add instances of disquotation for finite subsets of the Visser-Yablo sentences in (5) to WFUTB without stepping into  $\omega$ -inconsistencies.

In turn, Proposition 5 allows us to  $\omega$ -consistently extend WFUTB with all instances of the Uniform T-schema for arbitrary r-stable non-self-referential non-well-founded sentences, provided that they d-refer to finitely many expressions. Call the resulting theory FUUTB, for “Finitely Unfounded Uniform Tarski Biconditionals”. This theory leaves the Visser-Yablo sentences out but allows for instances of disquotation for, e.g. sentences in  $\omega$ -chains such as the ‘truth-teller’ sequence, given by an infinite list of sentences, each of which says only of the one coming after that is true.<sup>37</sup>

Unlike WFUTB’s, the truth predicate of the extensions considered in the previous paragraph cannot always be closed under equivalence, on pain of unsoundness. To see it, notice that  $\forall x (x = 1 \rightarrow \neg Tx)$  is logically equivalent to the strong liar  $\neg \text{TI}$  ( $= 1$ ) but also r-stable, non-self-referential – it only refers to  $\neg \text{TI}$ , which just refers to itself – and it d-refers to finitely many expressions.

Allow me to consider one last disquotational theory, NSRTB, for “Non-self-referential Tarski Biconditionals”. NSRTB extends URfn(PAT) with all instances of the following schema:

$$\text{RSt}(\ulcorner \varphi \urcorner) \wedge \neg \text{SRef}_{\text{PA}}(\ulcorner \varphi \urcorner) \rightarrow (\text{T}\ulcorner \varphi \urcorner \leftrightarrow \varphi) \tag{15}$$

<sup>35</sup>See Leitgeb [21, Lemma 11].

<sup>36</sup>See Leitgeb [21, Lemma 5].

<sup>37</sup>See Picollo [26, Prop. 3.3].

where  $\varphi$  is a sentence. Every r-stable non-self-referential sentence has a corresponding instance of disquotation in NSRTB, including all sentences on the Visser-Yablo list. As a consequence, the theory is  $\omega$ -inconsistent (cf. Section 2.2) and, therefore, unsound. However, it can be easily seen that NSRTB is not inconsistent, by Proposition 4 and an application of compactness. Thus, we may extract the general conclusion that sets of non-self-referential sentences of  $\mathcal{L}_T$  cannot be paradoxical simpliciter but only  $\omega$ -paradoxical.

But note also that, if we adopt the uniform version of (15) instead, that is,

$$\forall \mathbf{t} (\text{RSt}(\ulcorner \varphi(\mathbf{t}) \urcorner) \wedge \neg \text{SRef}_{\text{PA}}(\ulcorner \varphi(\mathbf{t}) \urcorner) \rightarrow (\text{T}\ulcorner \varphi(\mathbf{t}) \urcorner \leftrightarrow \varphi(\mathbf{t}^\circ)))$$

the resulting system is downright inconsistent, for we can show in  $\text{URfn}(\text{PA})$  that all instances of the Visser-Yablo predicate  $Y(v)$  satisfy the antecedent of this principle.<sup>38</sup> Since no new instances of disquotation for sentences are allowed by the uniform version, one might wonder what has gone wrong. According to Priest [29], Beall [1], Cook [3], and others, the inconsistency is a consequence of admitting instances of (uniform) disquotation, not for self-referential sentences, but for self-referential or circular predicates, e.g.  $Y(v)$ . Although the Visser-Yablo sentences are not self-referential, the argument goes, they are formulated in terms of a circular predicate and, therefore, are circular themselves. To evaluate this claim, adequate notions of reference, self-reference, etc. that apply to formulae, and not just sentences, are required. I leave the task of extending the notions introduced in Section 2.3 to formulae for another time. Let me just say that, provided a natural extension is possible, the claim seems quite plausible.

## 5 Conclusions

In this paper I have explored a number of semantic and axiomatic truth theories motivated by the notions of reference I introduced in [27]. I have shown the latter to be proof-theoretically strong compared to most well-known systems in the literature, despite being purely disquotational. However, proof-theoretic strength is not necessarily a sign (or the only sign) of theoretical value. To assess the worth of the systems in themselves and compared to others, in this concluding section I test them against the criteria discussed in Leitgeb [22].

The criteria are the following:

- (a) Truth should be expressed by a predicate, and a theory of syntax should be available.
- (b) If a theory of truth is added to mathematical or empirical theories, it should be possible to prove the latter true.
- (c) The truth predicate should not be subject to any type restrictions.
- (d) T-biconditionals should be derivable unrestrictedly.
- (e) Truth should be compositional.
- (f) The theory should allow for standard interpretations.

<sup>38</sup>See Ketland [18, 19] and (4).

- (g) The outer logic and the inner logic should coincide.
- (h) The outer logic should be classical.

As Leitgeb shows, it is not possible to meet all these criteria at once. Every formal account of truth will necessarily have to give up some of these requirements. It will be instructive to see how the systems we have considered here fare.

Clearly, all of our systems – WFUTB, WFUTB<sup>+</sup>, FUUTB, and NSRTB – meet the first criterion, for T is a predicate and arithmetic plays the role of the syntax theory in the background. The same holds of other popular theories considered in the previous section, i.e. PUTB, KF, FS, and systems of ramified truth.

(b), on the other hand, is met only partially. For every theorem  $\varphi$  of our base theory – URfn(PA) – each of our systems can be shown to entail  $T^\top\varphi^\top$ . This is because  $\varphi$  is T-free and, therefore, well-founded. However, none of the systems can prove what Leitgeb demands, i.e. a general statement asserting that every theorem of URfn(PA) is true. This would amount to showing that the systems entail URfn(PA)’s global reflection principle:

$$(GRfn_{URfn(PA)}) \quad \forall x (Bew_{URfn(PA)}(x) \rightarrow Tx)$$

KF, FS, and all systems of ramified truth are known to entail similar principles for their respective base theories. PUTB, on the other hand, doesn’t, as shown by Halbach [9, §6]. We can adapt his proof to show that  $GRfn_{URfn(PA)}$  is not provable in any of our systems either. Using Kleene’s Recursion Theorem let us define a p.r. function  $s$ , represented in PA by the function symbol  $\mathfrak{s}$ , as follows:

$$s(\varphi, n) := \begin{cases} \varphi & \text{if } \varphi := s = t \\ T\mathfrak{s}(t, \bar{n}) \wedge lh(t) \leq \bar{n} & \text{if } \varphi := Tt \\ \neg s(\psi, n) & \text{if } \varphi := \neg\psi \\ s(\psi, n) \wedge s(\chi, n) & \text{if } \varphi := (\psi \wedge \chi) \\ s(\psi, n) \vee s(\chi, n) & \text{if } \varphi := (\psi \vee \chi) \\ s(\psi, n) \rightarrow s(\chi, n) & \text{if } \varphi := (\psi \rightarrow \chi) \\ \forall v s(\psi, n) & \text{if } \varphi := \forall v \psi \\ \exists v s(\psi, n) & \text{if } \varphi := \exists v \psi \end{cases}$$

The function symbol  $lh(x)$  represents the function that maps every formula to the number of logical operators occurring in it. Then, we can show the following:

**Lemma 4** *If  $WFUTB \vdash \varphi$ , then there is an  $n \in \omega$  s.t.  $WFUTB \vdash s(\varphi, n)$ , and similarly for  $WFUTB^+$ , FUUTB, and NSRTB.*

I omit the proof, as it’s roughly the same as the proof of Halbach’s Lemma 6.1. An important lesson to draw from this lemma is that none of the systems entails URfn(PA)’s global reflection principle. If they did, then they would also entail

$$\forall x (Bew_{URfn(PA)}(x) \rightarrow T\mathfrak{s}(x, \bar{n}) \wedge lh(x) \leq \bar{n})$$

for some natural number  $n$ , which would fix an upper bound to the number of logical operators occurring in every theorem of URfn(PA).

Moving on to (c), it can be easily seen that all the systems meet this requirement, for in all of them T applies to sentences containing T itself – e.g.  $T^\top T^\top 0 = 0^\top$  is



provable in the four systems. So do KF, FS, and PUTB. Systems of ramified truth obviously fail this criterion.

By contrast, (d) obviously doesn't hold of any of the systems we are considering, as the liar sentence  $\lambda$ , for example, doesn't have a corresponding instance of disquotation, on pain of triviality. Nonetheless, some of our systems fare better than others. While WFUTB contains 'the least number' of T-biconditionals, NSRTB is the most encompassing one, and the other two systems lie somewhere in between.

Criterion (e) is not met by our systems either. Suppose one of the systems entailed the following compositional principle:

$$(T\wedge) \quad \forall x\forall y (\text{Sent}_{\mathcal{L}}(x\wedge y) \rightarrow (T(x\wedge y) \leftrightarrow (Tx \wedge Ty)))$$

where  $\text{Sent}_{\mathcal{L}}(x)$  is a formula of  $\mathcal{L}$  representing the property of being a sentence of this language. By Lemma 4, there should be an  $n \in \omega$  s.t. the following is also derivable:

$$\forall x\forall y (\text{Sent}_{\mathcal{L}}(x\wedge y) \rightarrow (T\mathfrak{s}(x\wedge y, \bar{n}) \wedge \text{lh}(x\wedge y) \leq \bar{n} \leftrightarrow T\mathfrak{s}(x, \bar{n}) \wedge \text{lh}(x) \leq \bar{n} \wedge T\mathfrak{s}(y, \bar{n}) \wedge \text{lh}(y) \leq \bar{n}))$$

Since  $\forall x (\text{Sent}_{\mathcal{L}}(x) \rightarrow x = \mathfrak{s}(x, \bar{n}))$  is provable in PA for every  $n \in \omega$ , we also have that

$$\forall x\forall y (\text{Sent}_{\mathcal{L}}(x\wedge y) \rightarrow (T(x\wedge y) \wedge \text{lh}(x\wedge y) \leq \bar{n} \leftrightarrow Tx \wedge \text{lh}(x) \leq \bar{n} \wedge Ty \wedge \text{lh}(y) \leq \bar{n}))$$

By  $T\wedge$ , this formula entails the following:

$$\forall x\forall y (\text{Sent}_{\mathcal{L}}(x\wedge y) \wedge T(x\wedge y) \rightarrow (\text{lh}(x) \leq \bar{n} \wedge \text{lh}(y) \leq \bar{n} \rightarrow \text{lh}(x\wedge y) \leq \bar{n}))$$

Since there are theorems of PA with more than  $n$  logical operators that take the form of a conjunction, each of whose conjuncts has at most  $n$  logical operators, and all our truth systems prove that each single theorem of arithmetic is true, we get a contradiction. Similar arguments can be given for other compositional principles.

Thus, none of the systems puts forward a compositional notion of truth, not even regarding sentences of  $\mathcal{L}$ . For the same reasons, PUTB's truth predicate is also non-compositional. By contrast, KF, FS, and ramified truth theories are axiomatized by compositional principles. Note, however, that our systems can be consistently extended with compositional principles for  $\mathcal{L}$ : since every sentence of this language is well-founded, our systems prove all instances of each compositional principle already, which means that these principles are all true in the models we used as witnesses to the consistency of the systems. This also shows that  $\text{GRfn}_{\text{URfn}(\text{PA})}$  can be consistently added to our systems, for it's entailed by the compositional principles for  $\mathcal{L}$ .<sup>39</sup> Moreover,  $\omega$ -consistency is preserved in every case.

By contrast, not all of our systems can be consistently extended with compositional principles for expressions containing T. Although unrestricted compositionality might be an unreasonable requirement, one would expect that each truth system is compatible at least with compositional principles for the class of sentences that have a corresponding instance of disquotation in the system. This is certainly the case of

<sup>39</sup>See, for instance, Halbach [10, Theorem 8.32].

WFUTB and WFUTB<sup>+</sup>. This result follows from the fact that the instances of each principle of compositionality for grounded sentences in the sense of Leitgeb are all true in Leitgeb's semantic theory of grounded truth,  $\langle \mathbb{N}, \Phi_{\text{If}} \rangle$ , as dependence is closed under logical connectives and equivalence in every expansion of  $\mathbb{N}$  to  $\mathcal{L}_{\text{T}}$ .<sup>40</sup> Thus, the compositional principles for truth relativized to the predicate WFUTB<sup>+</sup> deployed in the restriction of disquotation are also true in  $\langle \mathbb{N}, \Phi_{\text{If}} \rangle$  (cf. Lemma 2). Since the latter is a model of WFUTB<sup>+</sup>, extending this theory with compositionality relativized to this predicate preserves  $\omega$ -consistency. *A fortiori*, the same can be said of WFUTB and its respective restricting predicate.

On the other hand, neither FUUTB nor NSRTB can be extended with appropriately relativized compositional axioms. If we strongly diagonalize the predicate  $T\neg x$  we obtain a term  $l'$  s.t. the identity statement  $l' = \ulcorner T\neg l' \urcorner$  is provable in PA. Note that  $T\neg l'$  d-refers just to the negation of the sentence denoted by  $l'$ , i.e. to  $\neg T\neg l'$ .  $T\neg l'$  is therefore non-self-referential and d-refers to finitely many expressions. Moreover, it is r-stable. Thus, both FUUTB and NSRTB contain an instance of disquotation for this sentence, namely,  $\ulcorner T\neg l' \urcorner \leftrightarrow T\neg l'$ . By Leibniz's Law, this entails  $Tl' \leftrightarrow T\neg l'$ , which is incompatible with the compositional principle that says that truth commutes with negation.

Regarding Leitgeb's criterion (f), all our systems except NSRTB fare well. While the latter is  $\omega$ -inconsistent, as we saw in the last section, the other three systems all have models extending  $\mathbb{N}$ , which is what Leitgeb had in mind when he demanded that a formal truth theory based on arithmetic allowed for standard interpretations. While KF, PUTB, and ramified truth theories also meet (f), FS famously doesn't.

Finally, we consider criteria (g) and (h) together. According to Leitgeb, the outer logic of a truth theory is given by the logical laws the theory can prove, whereas its inner logic consists of the logical laws the theory proves to be true. Since the systems we are considering are classical, (h) is met by all of them. (g), instead, is only satisfied by WFUTB<sup>+</sup>. To see this, note that all logical truths are provably equivalent to each other in first-order logic and, thus, also in URfn(PAT). Since, say,  $0 = 0$  is well-founded and r-stable, it has an instance of disquotation in WFUTB<sup>+</sup>. This implies that all logical truths have instances of disquotation as well in the system. Given that they are all provable, we can conclude that they are also provably true. In contrast, none of the other systems has an instance of disquotation for, e.g.  $\forall x (Tx \rightarrow Tx)$ . This is a logical truth, yet it q-refers to every sentence. Thus, it is self-referential, so it's not declared true by WFUTB, FUUTB, or NSRTB. While FS is known for its matching inner and outer logics, KF, PUTB, and systems of ramified truth don't meet this criterion.

It's time to take stock. As we have seen, all the axiomatic systems we have considered in this paper fare equally well regarding criteria (a), (d), and (h). Due to their pure disquotational character, like PUTB the new systems fail requirements (b) and (e). However, unlike it, both WFUTB and WFUTB<sup>+</sup> can be soundly extended with compositional principles for the language with the truth predicate. Moreover, unlike the systems of ramified truth, ours are untyped; unlike FS, WFUTB, WFUTB<sup>+</sup>, and

<sup>40</sup>See Leitgeb [21, Lemma 5].

FUUTB, ours allow for standard interpretations; and, unlike KF and PUTB, the inner and the outer logics of WFUTB<sup>+</sup> do coincide.

This shows that the axiomatic systems introduced in the previous section, especially WFUTB<sup>+</sup>, compare favourably to the best known axiomatic truth theories in the literature against Leitgeb's criteria. Additionally, their proof-theoretic power places them above many popular systems.

**Acknowledgements** I am deeply indebted to Volker Halbach, with whom I had countless fruitful discussions on reference and self-reference over the last six years. I would also like to particularly thank Dan Waxman for extremely helpful comments on the final drafts, Thomas Schindler, for great suggestions and encouragement, and two anonymous referees for serious improvements in clarity and exposition. I should mention as well Eduardo Barrio, Catrin Campbell-Moore, Luca Castaldo, Roy T. Cook, Benedict Eastaugh, Martin Fischer, Hannes Leitgeb, Øystein Linnebo, Carlo Nicolai, Graham Priest, Johannes Stern, Albert Visser, the Buenos Aires Logic Group, the MCMP logic community, and the Oxford logic group. Finally, I would like to thank the Alexander von Humboldt Foundation and, especially, the Deutsche Forschungsgemeinschaft (DFG) for generously funding the research projects "Reference patterns of paradox" (PI 1294/1-1) and "The Logics of Truth: Operational and Substructural Approaches" (GZ HJ 5/1-1, AOBJ 617612).

**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. Beall, J.C. (2001). Is Yablo's Paradox non-circular? *Analysis*, 61, 176–187.
2. Carnap, R. (1937). *Logische Syntax der Sprache*. London: Routledge.
3. Cook, R.T. (2006). There are non-circular paradoxes (but Yablo's isn't one of them!). *The Monist*, 89(1), 118–149.
4. Feferman, S. (1991). Reflecting on incompleteness. *Journal of Symbolic Logic*, 56, 1–49.
5. Field, H. (2008). *Saving truth from paradox*. New York: Oxford University Press.
6. Fischer, M., Halbach, V., Stern, J., Kreiner, J. (2015). Axiomatizing semantic theories of truth? *Review of Symbolic Logic*, 8, 257–278.
7. Friedman, H., & Sheard, M. (1987). An axiomatic approach to self-referential truth. *Annals of Pure and Applied Logic*, 33, 1–21.
8. Halbach, V. (1994). A system of complete and consistent truth. *Notre Dame Journal of Formal Logic*, 35, 311–327.
9. Halbach, V. (2009). Reducing compositional to disquotational truth. *Review of Symbolic Logic*, 2, 786–798.
10. Halbach, V. (2014). *Axiomatic theories of truth*, 2nd edn. Cambridge: Cambridge University Press.
11. Halbach, V., & Horsten, L. (2006). Axiomatizing kripke's theory of truth. *Journal of Symbolic Logic*, 71, 677–712.
12. Halbach, V., & Visser, A. (2014). Self-reference in arithmetic I. *Review of Symbolic Logic*, 7, 671–691.
13. Halbach, V., & Visser, A. (2014). Self-reference in arithmetic II. *Review of Symbolic Logic*, 7, 692–712.
14. Hardy, J. (1995). Is Yablo's paradox liar-like? *Analysis*, 55(3), 197–198.
15. Herzberger, H. (1970). Paradoxes of grounding in semantics. *Journal of Philosophical Logic*, 67, 145–167.
16. Horwich, P. (1998). *Truth*, 2nd edn. Oxford: Blackwell.
17. Jeroslow, R.G. (1973). Redundancies in the Hilbert-Bernays derivability conditions for gödel's second incompleteness theorem. *Journal of Symbolic Logic*, 38, 359–367.
18. Ketland, J. (2004). Bueno and colyvan on yablo's paradox. *Analysis*, 64, 165–172.

19. Ketland, J. (2005). Yablo's paradox and  $\omega$ -inconsistency. *Synthese*, 145, 295–307.
20. Kripke, S. (1975). Outline of a theory of truth. *Journal of Philosophy*, 72, 690–716.
21. Leitgeb, H. (2005). What truth depends on. *Journal of Philosophical Logic*, 34, 155–192.
22. Leitgeb, H. (2007). What theories of truth should be like (but cannot be). *Philosophy Compass*, 2(2), 276–290.
23. Martin, R.L., & Woodruff, P.W. (1975). On Representing True-in-L in L. *Philosophia*, 5, 213–217.
24. McGee, V. (1992). Maximal consistent sets of instances of Tarski's schema. *Journal of Philosophical Logic*, 21, 235–241.
25. Montague, R. (1962). Theories incomparable with respect to relative interpretability. *Journal of Symbolic Logic*, 27, 195–211.
26. Picollo, L. (2018). Reference in arithmetic. Review of Symbolic Logic, pp. 1–31. <https://doi.org/10.1017/S1755020317000351>.
27. Picollo, L. Alethic Reference. *Journal of Philosophical Logic* (to appear).
28. Pohlers, W. (2009). *Proof theory: the first step into impredicativity*. Berlin-Heidelberg: Springer.
29. Priest, G. (1997). Yablo's paradox. *Analysis*, 57, 236–242.
30. Schindler, T. (2014). Axioms for grounded truth. *Review of Symbolic Logic*, 7, 73–83.
31. Schindler, T. (2015). A disquotational theory of truth as strong as  $Z_2^-$ . *Journal of Philosophical Logic*, 44, 395–410.
32. Tarski, A. (1935). The concept of truth in formalized languages. In *Logic, semantics, metamathematics* (pp. 152–278). Oxford: Clarendon Press.
33. Tarski, A. (1944). The semantic conception of truth: and the foundations of semantics. *Philosophy and Phenomenological Research*, 4, 341–376.
34. Visser, A. (1989). Semantics and the liar paradox. In Gabbay, D.M., & Günthner, F. (Eds.) *Handbook of philosophical logic*, (Vol. 4 pp. 617–706). Dordrecht: Reidel.
35. Yablo, S. (1985). Truth and reflexion. *Journal of Philosophical Logic*, 14, 297–349.
36. Yablo, S. (1993). Paradox without self-reference. *Analysis*, 53, 251–252.

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