

## CLASSICAL DETERMINATE TRUTH I

KENTARO FUJIMOTO AND VOLKER HALBACH

**Abstract.** We introduce and analyze a new axiomatic theory CD of truth. The primitive truth predicate can be applied to sentences containing the truth predicate. The theory is thoroughly classical in the sense that CD is not only formulated in classical logic, but that the axiomatized notion of truth itself is classical: The truth predicate commutes with all quantifiers and connectives, and thus the theory proves that there are no truth value gaps or gluts. To avoid inconsistency, the instances of the T-schema are restricted to *determinate* sentences. Determinateness is introduced as a further primitive predicate and axiomatized. The semantics and proof theory of CD are analyzed.

Ad veritatem autem copulativae requiritur quod utraque pars sit vera, et ideo si quaecumque pars copulativae sit falsa, ipsa copulativa est falsa.

William of Ockham, *Summa Logicae* II.32

Ad veritatem autem [propositionis] disiunctivae requiritur quod aliqua pars sit vera[.]

William of Ockham, *Summa Logicae* II.33

**§1. Classicity and compositional semantics.** Philosophy abounds with general claims expressed with a truth predicate: Philosophers debate whether there are contingent or synthetic a priori truths, and whether there are unprovable or unverifiable truths; they mostly agree that what is known is true, but that some justified true beliefs are not known; they teach their students that the conclusion of a valid argument is true if its premisses are true and try to convince their students that the rules of natural deduction are truth-preserving. Specific instances of these claims can often be stated without a predicate for truth; but the quantified claims make essential use of the truth predicate.

In the arguments for this kind of claim, assumptions about the truth predicate are often used implicitly without further ado. In the present paper we strive to make these assumptions explicit by listing principles about the truth predicate that are jointly consistent with further plausible assumptions about the objects to which truth is ascribed. We presuppose that these objects share the structure of (types of)

---

Received November 5, 2021.

2020 *Mathematics Subject Classification*. Primary 03F03.

*Key words and Phrases*. axiomatic theories of truth, ordinal analysis, liar paradox, determinateness, grounding.

© The Author(s), 2023. Published by Cambridge University Press on behalf of The Association for Symbolic Logic. This is an Open Access article, distributed under the terms of the Creative Commons Attribution licence (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted re-use, distribution, and reproduction in any medium, provided the original work is properly cited.

0022-4812/24/8901-0013

DOI:10.1017/jsl.2023.49



sentences without entering the debate about whether truth applies to propositions, beliefs, sentence tokens, or still other things.<sup>1</sup>

It is widely believed that truth serves its role as a device of generalization in virtue of its disquotational feature; and a large part of the literature on truth is devoted to the development of strong disquotational truth theories. However, in their zeal to devise theories of disquotational truth, some philosophers seem to have lost sight of how we reason with generalizations. When truth is used as a device of generalization, disquotation is not the only principle that is used in reasoning with generalizations. Compositional principles are just as important and used routinely without qualms. We give three examples.

Our first example comes from logic. When motivating axioms or rules for a logical calculus, compositional laws of truth are used. In the case of the elimination rule for conjunction or for the universal quantifier in natural deduction, students usually quickly accept the claim that the rule preserves truth. At this point in the course, students may not have seen the model-theoretic notions of satisfaction and truth or the object/metalanguage distinction at all; nevertheless they have no problems in seeing that a conjunct must be true if the conjunction is true and that a substitution instance is true if the universally quantified sentence is true. The point can also be made about the history of logic: Logicians chose sound logical calculi long before the arrival of the modern model-theoretic notion of truth in Tarski and Vaught [43] in 1956 (as far as we know; see Hodges [26]). Truth preservation of logical rules is often assumed to hold without any restriction. Any restriction to a specific sublanguage is at odds with the universality of logic.

The use of the compositional principles is by no means confined to the realm of logical theorizing. Our second example is from epistemology. Most epistemologists would be happy to endorse the following argument, which could form part of a Gettier example:

Smith believes a disjunction. He is justified in believing one disjunct, which happens to be false, while he does not believe the other disjunct, which happens to be true. Therefore Smith has a justified true belief.

In this argument we reason about the truth (and justification) of beliefs by analyzing and manipulating their logical structures without explicitly specifying exactly what these beliefs are. The validity of the argument depends on the assumption that a disjunction is true if one of its disjuncts is true, as stated by Ockham in the second quote above. Fujimoto [16, 17] called this kind of reasoning *blind deduction*, because we reason about the truth of beliefs or sentences without being able to specify these beliefs or sentences by means of a quotational name or a structural description. Disquotation axioms or T-sentences used as axioms permit only reasoning about the truth of sentences that are explicitly given by such a naming device.

---

<sup>1</sup>There are many alternative approaches. Philosophers have tried to recover the expressive power of the truth predicate by employing propositional quantification. Sophisticated theories of propositions as objects to which truth is ascribed have been developed. Here, we do not have space to discuss these alternatives. See Halbach and Leigh [24] for a discussion.

The use of the compositional principles is so ubiquitous that their use often goes unnoticed as in our third example:

Alfred made a claim that was denied by Kurt (i.e., Kurt asserted the negation of the claim); and everything Kurt claimed is true.  
Therefore not everything Alfred claimed is true.

For the validity of this argument we require the principle that a claim is not true if its negation is true.

In the examples we quantify over *all* sentences, beliefs, or propositions. Blind deductions are not restricted to sentences that are safe, grounded, determinate, healthy, non-circular, or in some other way ‘unproblematic’. Moreover, the full compositional principles are consistent without any restrictions on the underlying logic or syntax theory.

The axioms for truth can and should be added to a base theory that yields at least a comprehensive theory of syntax (in a direct or coded form); the base theory may also go far beyond a theory of syntax. In [22, 23] the second author formulated the truth theory over set theory as the base. Here, however, we start from Peano arithmetic, which is traditionally used as base theory for axiomatic theories of truth; but we consider it only as a simple model case. Expressions are identified with their codes. As usual, we confine ourselves to a truth rather than a satisfaction predicate in the case of arithmetic, as names for all objects are available and satisfaction and variable assignments are not needed. All extensions of Peano arithmetic we consider are formulated in classical logic.

In what follows,  $\mathcal{L}_0$  is the language of arithmetic. The axioms are formulated in an expansion of  $\mathcal{L}_0$  with further predicate symbols T and D. Negation and conjunction are the only connectives, and the universal is the only quantifier in the language; other connectives and the existential quantifier are assumed to be defined. In the axioms below  $\text{Sent}(x)$  expresses that  $x$  is a sentence of  $\mathcal{L}$ . The symbol  $\neg$  expresses the function that, if it is applied to a sentence, returns its negation. If a suitable function symbol is not in  $\mathcal{L}_0$ , this function needs to be expressed by a suitable formula;  $\wedge$  and  $\forall$  are defined analogously. The quantifier  $\forall t$  ranges over (codes of) closed terms and is defined using a suitable formula describing the set of closed terms; and  $x(t/v)$  designates the result of formally substituting  $t$  for the variable  $v$  in  $x$ . The following axioms then express that a negated sentence is true iff the sentence is not, that a conjunct is true iff both conjuncts are, and that a universally quantified sentence is true iff all its substitution instances are:

- T4  $\forall x (\text{Sent}(x) \rightarrow (\text{T}(\neg x) \leftrightarrow \neg \text{T}x)),$   
 T5  $\forall x \forall y (\text{Sent}(x \wedge y) \rightarrow (\text{T}(x \wedge y) \leftrightarrow \text{T}x \wedge \text{T}y)),$   
 T6  $\forall v \forall x (\text{Sent}(\forall v x) \rightarrow (\text{T}(\forall v x) \leftrightarrow \forall t \text{T}x(t/v))).$

The axioms T4–T6 capture a *thoroughly* classical concept of truth: The truth theory is not only formulated in classical logic, but the notion of truth axiomatized by T4–T6 is itself classical. According to the axioms, truth commutes with quantifiers and connectives for *all* sentences, including those with the truth predicate. Theories such as the Kripke–Feferman theory [8, 37] or Cantini’s [6] VF are formulated in classical logic as well, but capture a non-classical notion of truth.

Axiom T4 disproves the existence of truth value gaps and gluts. We define falsity as the truth of the negation, that is, we stipulate:  $Fx :\Leftrightarrow T\neg x$ . If the truth predicate is classical and T4 is assumed, then falsity is equivalent to non-truth. Thus, every sentence is true or false; and no sentence is both true and false.

Disquotation axioms have been thought to be more fundamental than the compositional axioms. Consequently, there have been various attempts to obtain compositional from disquotational principles. Tarski [42, p. 259] tried to derive restricted compositional principles from the T-sentences using what he called the *rule of infinite induction*, which is a generalization of the  $\omega$ -rule. More recently, some authors, including Halbach [19] and Horsten and Leigh [27], have used reflection principles instead, which are of course finitized versions of infinitary rules. We are sceptical about the prospects of justifying compositional from disquotational axioms. There are general objections against the use of infinitary and reflection rules to this end. Another crucial objection is especially relevant when unrestricted compositional principles are employed. Usually, attempts to obtain compositional axioms for *all* sentences of a certain kind from disquotational principles rely on the disquotation principle for all sentences of this kind. For instance, we may try to derive the compositional axioms for all T-free (or T-positive) sentences from all equivalences  $T\phi^\neg \leftrightarrow \phi$ , where  $\phi$  is T-free (or T-positive). However, in the case of the *completely unrestricted* compositional axioms T4–T6, this strategy is not very promising. Of course, the unrestricted compositional axioms T4–T6 can be derived from all instances  $T\phi^\neg \leftrightarrow \phi$ ; but this is because the unrestricted T-schema is inconsistent over arithmetic. Hence, we will not be able to derive all type-free compositional from disquotational principles, even if some kind of infinitary rule is assumed.<sup>2</sup> Therefore we adopt the unproblematic, fully general compositional principles as axioms and do not attempt to derive them from a disquotation schema, which has to be restricted in some way.

**§2. Determinateness and disquotation.** Our policy for restricting disquotation is inspired by disquotationalist and deflationist theories of truth. We expand the base language  $\mathcal{L}_0$  by adding sentences that are needed for semantic ascent and generalization. The sentences of the base language together with those needed for semantic ascent and generalization form the set  $\mathfrak{D}$ . Roughly, if we have a sentence  $\phi$  in  $\mathfrak{D}$ , we also add an equivalent new sentence  $T\phi^\neg$  as a ‘copy’ of  $\phi$ . These copies are required for semantic ascent. We close then under connectives and quantifiers in a way to be explained. This permits generalizations, because we can now express, for instance, that all provable sentences of  $\mathcal{L}_0$  are true. This process is then iterated: We add copies  $T\phi^\neg$  for semantic ascent and then close under connectives and quantifiers in order to be able to generalize. We add only sentences to  $\mathfrak{D}$  that

<sup>2</sup>This claim requires some qualification. It is possible to obtain all compositional axioms from a single consistent instance of the T-schema by a trick due to McGee [32] based on Curry’s paradox. Let  $C$  be the conjunction of all compositional axioms. By the diagonal lemma there is a sentence  $\gamma$  such that  $\gamma \leftrightarrow (T\gamma^\neg \leftrightarrow C)$  is provable. This sentence is logically equivalent to  $(T\gamma^\neg \leftrightarrow \gamma) \leftrightarrow C$ . Therefore the equivalence  $T\gamma^\neg \leftrightarrow \gamma$  is a re-axiomatization of the compositional axioms over arithmetic. We see this as a mere curiosity; but it shows that the claim that one cannot get all type-free compositional axioms from a consistent set of T-sentences needs some qualification.

are required for semantic ascent and generalization, not all sentences of the full language  $\mathcal{L}$ . In particular, not all sentences of the form  $T\ulcorner\phi\urcorner$  are in  $\mathfrak{D}$ . The set  $\mathfrak{D}$  fails to be decidable. In a general setting with ‘contingent’ vocabulary, whether a sentence belongs to  $\mathfrak{D}$  may depend on empirical facts.

Our truth theory CD itself is formulated in the full language  $\mathcal{L}$ ; but the disquotation schema is restricted to sentences of  $\mathfrak{D}$ , that is, to the sentences of the base language and those needed for semantic ascent and generalization. Being a sentence of  $\mathfrak{D}$  is expressed in CD by the primitive predicate symbol  $D$ , which is suitably axiomatized. As disquotation schema we can then use  $D\ulcorner\phi\urcorner \rightarrow (T\ulcorner\phi\urcorner \leftrightarrow \phi)$  for all sentences of  $\mathcal{L}$ , or rather a generalization thereof with parameters. We call the sentences in  $\mathfrak{D}$  *determinate*. If it were not for the additional predicate  $D$ , they would be the sentences that are grounded in Kripke’s [28] sense (with some qualifications).

This approach ensures that truth can be used as a device of generalization over sentences in  $\mathcal{L}_0$  and over generalizations of such generalizations, and so on. Hence this use of truth is fully available in CD. However, disquotation is not available for sentences that cannot be reached by semantic ascent and generalization.

**2.1. Determinateness.** We describe and axiomatize  $\mathfrak{D}$ , the set of determinate sentences in more detail. First, we declare all atomic sentences of the base language  $\mathcal{L}_0$  determinate:

$$D1 \quad \forall s \forall t \ Ds = t.$$

That is, all closed identity statements are determinate. If there were further predicate symbols in the base language, analogous axioms would be added. The axioms below will enable us to prove that all sentences of the base language  $\mathcal{L}_0$  are determinate.

If and only if  $\phi$  is determinate, that is, in  $\mathfrak{D}$ , we add a copy  $T\ulcorner\phi\urcorner$  of  $\phi$  to  $\mathfrak{D}$ . This is stated for arbitrary closed terms  $t$ , not only numerals  $\ulcorner\phi\urcorner$ :

$$D2 \quad \forall t \ (DT\ulcorner t \urcorner \leftrightarrow Dt^\circ).$$

The expression  $t^\circ$  stands for the value of the term  $t$ . We cannot have a function symbol  $^\circ$  in our language, and thus this function needs to be expressed by a suitable formula.

A negated sentence is determinate iff the sentence is:

$$D4 \quad \forall x \ (Sent(x) \rightarrow (D(\ulcorner\neg x\urcorner) \leftrightarrow Dx)).$$

Whether a conjunction is determinate depends on the determinateness of its conjuncts. Clearly, if both conjuncts are determinate, then so is their conjunction; and if both are indeterminate, so is their conjunction. But is a conjunction determinate if only one conjunct is?

Our choice here is motivated by the function of truth as a generalizing device. Typically, a universal generalization would have the form  $\forall x (\phi(x) \rightarrow Tx)$ . The formula  $\phi(x)$  could express that  $x$  is a provable sentence of Peano arithmetic or that  $x$  is sentence of the form  $T\ulcorner T \dots \ulcorner 0=0 \urcorner \dots \urcorner$ . We consider the simple example  $\forall x (x = \ulcorner 0=0 \urcorner \rightarrow Tx)$  whose antecedent is satisfied by a single sentence only; we ‘generalize’ via semantic ascent over all sentences of the form  $0=0$ . Of course, generalization via semantic ascent is not needed for a single sentence; but ‘generalizing’ over a single sentence demonstrates the principle. If  $\mathfrak{D}$  contains the sentences required for generalizing via semantic ascent, then  $\forall x (x = \ulcorner 0=0 \urcorner \rightarrow Tx)$  should be determinate. The determinateness of universal quantified sentences

depends on their instances. Let  $\lambda$  be an indeterminate sentence. Then the instance

$$\ulcorner \lambda \urcorner = \ulcorner 0 = 0 \urcorner \rightarrow T\ulcorner \lambda \urcorner \tag{1}$$

has a determinate antecedent and an indeterminate consequent, by axioms D1 and D2. If the instance (1) were indeterminate, so would be the universal generalization  $\forall x (x = \ulcorner 0 = 0 \urcorner \rightarrow Tx)$  by any reasonably determinateness rules for quantification. Hence, our axioms for D should declare (1) determinate.

The determinateness of (1) should follow from the determinateness of the antecedent  $\ulcorner \lambda \urcorner = \ulcorner 0 = 0 \urcorner$ . The falsity of the determinate antecedent renders the consequent irrelevant; and the truth value of (1) does not depend on the consequent. Thus, we stipulate that a sentence  $\phi \rightarrow \psi$  can inherit its determinateness from a determinate and false antecedent. Equally, we also postulate that it is determinate if it has a true and determinate consequent.

This method of generalization applies equally if we do not only ‘generalize’ over a single sentence such as  $0 = 0$ , but, for instance, over all sentences provable in PA. Moreover, iterated semantic ascent permits generalization over all sentences  $0 = 0$ ,  $T\ulcorner 0 = 0 \urcorner$ ,  $T\ulcorner T\ulcorner 0 = 0 \urcorner \urcorner$ , and so on.

In our official language negation and conjunction are our only connectives.  $\phi \rightarrow \psi$  is conceived as an abbreviation of  $\neg(\phi \wedge \neg\psi)$ . In the presence of D4 we can then state our determinateness axiom for binary connectives in the following way:

- $\phi \wedge \psi$  is determinate iff:
- $\phi$  and  $\psi$  are both determinate, or
- one of the conjuncts is false and determinate.

This is expressed in the following axiom:

$$D5 \quad \forall x \forall y \left( \text{Sent}(x \wedge y) \rightarrow (D(x \wedge y) \leftrightarrow ((Dx \wedge Dy) \vee (Dx \wedge Fx) \vee (Dy \wedge Fy))) \right).$$

By our conventions above,  $Fx$  abbreviates  $T\neg x$ , namely, falsehood of  $x$ ; also note that  $Dx$  and  $D\neg x$  are equivalent by the axiom D4. By axiom T4,  $F$  can be replaced with  $\neg T$ .

Another way to argue for our treatment of binary connectives would be to argue that sentences such as  $\forall x (x = \ulcorner 0 = 0 \urcorner \rightarrow Tx)$  and its instance (1) are only about the determinate sentence  $0 = 0$  and should be treated just like  $T\ulcorner 0 = 0 \urcorner$ . The latter is determinate by axiom D1. Thus, it may be argued, semantic ascent applies to (1) and the universally quantified sentence in the same way it applies to  $T\ulcorner 0 = 0 \urcorner$ . However, making precise the underlying notion of aboutness is notoriously challenging (see Picollo [33, 34]); we do not attempt to pursue this line here.

Universally quantified sentences are treated as conjunctions of their instances:

$$D6 \quad \forall v \forall x \left( \text{Sent}(\forall v x) \rightarrow (D(\forall v x) \leftrightarrow (\forall t Dx(t/v) \vee \exists t (Dx(t/v) \wedge Fx(t/v)))) \right).$$

Injecting the notion of determinateness into the object language makes it possible to apply the truth and determinateness predicates to sentences containing D. We treat both, truth and determinateness, as completely type-free: D and T can be applied meaningfully to any sentence of the language. This is in contrast to traditional conceptions of groundedness, which apply to sentences with truth, while truth cannot be applied to sentences containing groundedness claims. Therefore, in our setting, the question arises whether atomic sentences of the form  $Dt$  are determinate.

Treating  $D$  like an elementary predicate of  $\mathcal{L}_0$ , that is, a predicate from  $\mathcal{L}_0$  would suggest the axiom  $\forall t \, D D t$ . This would mean that we have  $T \ulcorner D t \urcorner \leftrightarrow D t$  for all terms  $t$ . However, when it comes to complex sentences, we axiomatize  $D$  in terms of  $T$  in D5 and D6. Hence, it looks unsafe to declare all sentence  $D t$  determinate;  $D$  should be treated with the same caution as  $T$ : If, and only if  $\phi$  is determinate, we stipulate that  $D \ulcorner \phi \urcorner$  is determinate. Generalizing this to terms other than numerals gives the following axiom analogous to D2:

$$D3 \quad \forall t \, (D D t \leftrightarrow D t^\circ).$$

This concludes the list of determinateness axioms.

**2.2. Disquotation.** We still need to add axioms stipulating that a sentence  $T \ulcorner \phi \urcorner$  in  $\mathcal{D}$  is always obtained by semantic ascent and thus a ‘copy’ of  $\phi$ . Formally, we require that  $D \ulcorner \phi \urcorner \rightarrow (T \ulcorner \phi \urcorner \leftrightarrow \phi)$  is a theorem of our theory CD for every sentence of  $\mathcal{L}$ . We also require that our theory proves a more general, ‘uniform’ version of the disquotation schema; that is, we generalize it by quantifying over the closed terms and postulate for every formula  $\phi(x_1, \dots, x_n)$  with at most  $x_1, \dots, x_n$  the following:

$$DDS \quad \forall t_1 \dots \forall t_n \, (D \ulcorner \phi(t_1, \dots, t_n) \urcorner \rightarrow (T \ulcorner \phi(t_1, \dots, t_n) \urcorner \leftrightarrow \phi(t_1^\circ, \dots, t_n^\circ))).$$

In this uniform determinate disquotation schema the expression  $\ulcorner \phi(t_1, \dots, t_n) \urcorner$  stands for a complex term expressing that  $t_1, \dots, t_n$  are formally substituted for the free variables in  $\phi(x_1, \dots, x_n)$ . This permits us to bind the variables  $t_1, \dots, t_n$  in  $\ulcorner \phi(t_1, \dots, t_n) \urcorner$ . This concludes the list of axioms for our theory CD.

**2.3. Alternative axiomatizations.** In the presence of the compositional axioms T4 –T6, we need to stipulate DDS only for atomic  $\phi$ . Whatever our policy on disquotation is, we can always focus on the atomic instances; T4 –T6 will ensure that also all instances of  $T \ulcorner \phi \urcorner \leftrightarrow \phi$  will be provable for all sentences  $\phi$  built from those atomic formulae. Because we aim at an axiomatization that is as lean as possible for technical reasons, our official axiomatization does not feature schema DDS, but only disquotational axioms concerning atomic sentences. In particular, schema DDS can be replaced with the following three axioms:

- T1  $\forall s \, \forall t \, (T s = t \leftrightarrow s^\circ = t^\circ)$ ,
- T2  $\forall t \, (D t^\circ \rightarrow T D t)$ ,
- T3  $\forall t \, (D t^\circ \rightarrow (T T t \leftrightarrow T t^\circ))$ .

T2 can be derived from the following instance of DDS and axiom D3:

$$\forall t \, (D D t \rightarrow (T D t \leftrightarrow D t^\circ)). \tag{2}$$

Thus the three axioms imply DDS in the presence of the other axioms as is established in Lemma 3.1.

The principles T3, (2), and T2 show that the truth and determinateness predicates interact in a serious manner. In particular,  $D$  is not only a formalization of a metatheoretic notion that has been injected into the object language; predicates of the object language—and especially the truth predicate—can be meaningfully applied to sentences with  $D$ .

From (2) we can prove one direction of the T-sentences  $\phi \rightarrow T \ulcorner \phi \urcorner$  for all sentences without T, including those with  $D$ . Some sentences without T but with  $D$ —such as

$D\ulcorner\lambda\urcorner$  for some indeterminate  $\lambda$ —are not determinate, and thus DDS will not give us the T-schema for some sentences with D. Therefore, we consider the addition of the right-to-left direction to T2:

$$T2^+ \quad \forall t (Dt^\circ \leftrightarrow TDt).$$

This sentence is not provable from the disquotation schema DDS, and is not covered by our policy about  $\mathcal{D}$ . With  $T2^+$  we can derive all T-sentences  $T\ulcorner\phi\urcorner \leftrightarrow \phi$  as long as  $\phi$  is T-free.

Determinateness is one of many predicates such as analyticity, knowledge, necessity, and logical validity that are deeply intertwined with truth. Having the T-sentences  $T\ulcorner\phi\urcorner \leftrightarrow \phi$  also for  $\phi$  containing such predicates may be desirable; but we know already that this is not possible for all such  $\phi$  (see Halbach [20]), as axioms like  $T2^+$  may engender inconsistency, depending on the axioms for the other predicates. We cannot think of a better strategy than to consider these axioms on a case-by-case basis and to endorse axioms like  $T2^+$ , as long as they do not yield undesirable consequences. With  $T2^+$  we leave the safety of disquotation for determinate sentences only. In the case of determinateness we show that axiom  $T2^+$  does not affect the  $\omega$ -soundness of our theory and thus  $T2^+$  can be added as a further optional axiom.

The quantifier  $\forall t$  in D2, D3, T2, T3, and  $T2^+$  ranges over all closed terms, including those that do not denote a sentence. It may be desirable to restrict the quantifier in these axioms to terms denoting sentences. This would allow us to make any stipulations about sentences  $Tt$  and  $Dt$  where  $t$  fails to denote a sentence. We have chosen the version above for their simplicity.

**2.4. Extensionality.** In this section we discuss two additional axioms that are irrelevant to most of the metamathematical properties of our system. However, at least the first of these axioms is conceptually important. The issue goes beyond our specific system, and even beyond axiomatic theories, because semantic theories of truth are affected as well.

We expect truth to be an *extensional* notion unlike necessity, apriority, or being known. The compositional axioms prove the extensionality of truth for sentences in the following sense: Substituting a subformula in a sentence  $\phi$  with another subformula with the same truth value does not affect truth or falsity of  $\phi$ .

One would expect extensionality at the level of terms as well. That is, if the terms  $s$  and  $t$  have the same value, a sentence  $\phi(t)$  ought to be true if, and only if  $\phi(s)$  is true. However, this principle fails to be provable from the axioms we have listed so far. Only restricted versions can be proved, for instance, if the values of  $s$  and  $t$  are determinate or  $\chi$  is purely arithmetical. Without axiom R1, we can refute  $T\ulcorner Ts\urcorner \wedge \neg T\ulcorner Tt\urcorner$  under the assumption  $s = t$  only for determinate  $s$  and  $t$ , but not for all terms. Here, we call a term determinate iff its value is. Thus, we state extensionality with respect to terms as an axiom:

$$R1 \quad \forall x \forall v \forall s \forall t ((\text{Sent}(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (Tx(s/v) \leftrightarrow Tx(t/v))).$$

Axiom R1 is closely related to the question whether identity is a logical constant. Using R1, we can prove the truth of  $\forall x \forall y (x = y \rightarrow (\phi(x) \rightarrow \phi(y)))$  for all formulae  $\phi$ , which is a logical truth if identity is a logical constant. If identity is

a logical constant and our theory is to prove all logically valid sentences, then we have to add R1 (or some other axiom).<sup>3</sup>

So far we have focused on truth. For determinateness we also stipulate extensionality:

$$R2 \quad \forall x \forall v \forall s \forall t ((\text{Sent}(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (\text{D}x(s/v) \leftrightarrow \text{D}x(t/v))).$$

Whether determinateness should also be extensional may be less clear. For our purposes the extensionality of determinateness is not decisive. Omitting R2 would not affect the results below.

**§3. Axioms for classical determinate truth.** We now collect the stipulations in the previous sections. We start from the language  $\mathcal{L}_0$  of arithmetic. It contains at least a constant for 0 and a function symbol  $S$  for the successor function, a symbol  $+$  for addition, and a symbol  $\cdot$  for multiplication. For each number  $n$  we use the numeral  $\bar{n}$ , defined as the result of applying  $S$  to 0 for  $n$ -many times, as the canonical name for  $n$ .  $\mathcal{L}_0$  may feature further function and constant symbols, even those for all primitive recursive functions. The identity symbol  $=$  is the only predicate symbol of  $\mathcal{L}_0$ ; but further predicate symbols could easily be added.<sup>4</sup>

The language  $\mathcal{L}_T$  is  $\mathcal{L}_0$  augmented with T;  $\mathcal{L}_D$  is  $\mathcal{L}_0$  augmented with D. Adding both T and D to  $\mathcal{L}_0$  yields the language  $\mathcal{L}$ . Therefore, the following inclusions hold, if the languages are identified with the set of their formulae:

$$\mathcal{L}_0 \subset \mathcal{L}_T \text{ and } \mathcal{L}_D \subset \mathcal{L}.$$

Our system is labelled CD for ‘Classical Determinate Truth’. It is formulated in the language  $\mathcal{L}$  with D and T. The syntax of  $\mathcal{L}$  is appropriately arithmetized within PA. We adopt the notation of Halbach [21] concerning arithmetization of syntax. However, for dealing with new predicate D and other technical reasons, we slightly supplement (and change) his notation as follows:

- $\text{Sent}(x)$ ,  $\text{Sent}_T(x)$ ,  $\text{Sent}_D(x)$ , and  $\text{Sent}_0(x)$  represent the set of codes of  $\mathcal{L}$ -,  $\mathcal{L}_T$ -,  $\mathcal{L}_D$ -, and  $\mathcal{L}_0$ -sentences, respectively.
- $\text{Fml}(x)$ ,  $\text{Fml}_T(x)$ ,  $\text{Fml}_D(x)$ , and  $\text{Fml}_0(x)$ , respectively, represent the set of codes of formula of these languages.
- Similarly, we use  $\text{AtFml}(x)$ ,  $\text{AtFml}_T(x)$ ,  $\text{AtFml}_D(x)$ , and  $\text{AtFml}_0(x)$  for the codes of atomic formulae of these languages, and  $\text{AtSent}(x)$ ,  $\text{AtSent}_T(x)$ ,  $\text{AtSent}_D(x)$ , and  $\text{AtSent}_0(x)$  for the codes of atomic sentences of these languages, respectively.
- $\text{Var}(x)$  and  $\text{Term}(x)$  represent the sets of codes of variables (and codes of  $\mathcal{L}$ -terms (= the set of codes of  $\mathcal{L}_0$ -terms), respectively; recall that  $\text{ClTerm}$  is for the set of codes of *closed*  $\mathcal{L}_0$ -terms, and note that every closed  $\mathcal{L}$ -term is a closed  $\mathcal{L}_0$ -term.
- For each primitive predicate symbol  $R$  with arity  $k$ ,  $\bar{R}$  is a  $k$ -ary function that represents the syntactic operation of applying  $R$  to  $k$ -many terms. For instance,

<sup>3</sup>We thank Anton Broberg for pointing out to us the underderivability of the truth of  $\forall x \forall y (x = y \rightarrow (\phi(x) \rightarrow \phi(y)))$  in CD without axiom R1. His comment motivated the present section.

<sup>4</sup>With modifications the theory can be reformulated in a purely relational language without closed terms such as the language of set theory. For such a setting without function symbols, a satisfaction predicate might be a better fit than a unary truth predicate (see Halbach [22]).

$\mathbb{D}$  is an arithmetical representation of the syntactic operation of applying  $\mathbb{D}$  to a term.

- $x(y/z)$  represents the syntactic substitution of a term  $y$  for a variable  $z$  in an expression  $x$ .

Hence, it can be proved in PA that  $\text{Sent}(\forall v x)$  implies that  $v$  is the code of a variable (i.e.,  $\text{Var}(v)$ ) and that  $x$  codes some formula  $\phi$  with at most one free variable. We use  $\text{Sent}(\forall v x)$  to express that  $x$  is a formula with at most the variable  $v$  free. Often we need to quantify over closed terms and abbreviate  $\forall z (\text{CITerm}(z) \rightarrow \psi(z))$  as  $\forall t \psi(t)$  and also use  $s$  as a further variable ranging over closed terms.

The system CD is given by all axioms of PA with induction in the language  $\mathcal{L}$  with  $\mathbb{T}$  and  $\mathbb{D}$  and the following axioms:

**Truth axioms**

- T1  $\forall s \forall t (\mathbb{T}s = t \leftrightarrow s^\circ = t^\circ)$ ,
- T2  $\forall t (\mathbb{D}t^\circ \rightarrow \mathbb{T}\mathbb{D}t)$ ,
- T3  $\forall t (\mathbb{D}t^\circ \rightarrow (\mathbb{T}\mathbb{T}t \leftrightarrow \mathbb{T}t^\circ))$ ,
- T4  $\forall x (\text{Sent}(x) \rightarrow (\mathbb{T}(\neg x) \leftrightarrow \neg \mathbb{T}x))$ ,
- T5  $\forall x \forall y (\text{Sent}(x \wedge y) \rightarrow (\mathbb{T}(x \wedge y) \leftrightarrow \mathbb{T}x \wedge \mathbb{T}y))$ ,
- T6  $\forall v \forall x (\text{Sent}(\forall v x) \rightarrow (\mathbb{T}(\forall v x) \leftrightarrow \forall t \mathbb{T}x(t/v)))$ .

**Determinateness axioms**

- D1  $\forall s \forall t \mathbb{D}s = t$ ,
- D2  $\forall t (\mathbb{D}\mathbb{T}t \leftrightarrow \mathbb{D}t^\circ)$ ,
- D3  $\forall t (\mathbb{D}\mathbb{D}t \leftrightarrow \mathbb{D}t^\circ)$ ,
- D4  $\forall x (\text{Sent}(x) \rightarrow (\mathbb{D}(\neg x) \leftrightarrow \mathbb{D}x))$ ,
- D5  $\forall x \forall y (\text{Sent}(x \wedge y) \rightarrow (\mathbb{D}(x \wedge y) \leftrightarrow ((\mathbb{D}x \wedge \mathbb{D}y) \vee (\mathbb{D}x \wedge \mathbb{F}x) \vee (\mathbb{D}y \wedge \mathbb{F}y))))$ ,
- D6  $\forall v \forall x (\text{Sent}(\forall v x) \rightarrow ((\mathbb{D}(\forall v x) \leftrightarrow (\forall t \mathbb{D}x(t/v) \vee \exists t (\mathbb{D}x(t/v) \wedge \mathbb{F}x(t/v))))))$ .

As explained above,  $\mathbb{F}x$  abbreviates  $\mathbb{T}\neg x$ , which expresses the falsity of  $x$ .

**Extensionality axioms**

- R1  $\forall x \forall v \forall s \forall t ((\text{Sent}(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (\mathbb{T}x(s/v) \leftrightarrow \mathbb{T}x(t/v)))$ ,
- R2  $\forall x \forall v \forall s \forall t ((\text{Sent}(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (\mathbb{D}x(s/v) \leftrightarrow \mathbb{D}x(t/v)))$ .

We now turn to some subsystems and supersystems of CD in order to assess the significance of axioms T2 and D3. It will be shown that removing one of these two axioms or both does not affect the proof-theoretic strength of CD. That is, the three systems

$$\begin{aligned} \text{CD}_0 &:= \text{CD} - \text{T2} - \text{D3} \\ \text{CD}_1 &:= \text{CD} - \text{T2} \\ \text{CD}_2 &:= \text{CD} - \text{D3} \end{aligned}$$

are proof-theoretically equivalent, although they are different theories.

The addition of axiom  $\text{T2}^+$  boosts the proof-theoretic strength of CD, that is,

$$\begin{aligned} \text{CD}^+ &:= \text{CD} + \text{T2}^+ \\ \text{CD}_2^+ &:= \text{CD}_2 + \text{T2}^+ \end{aligned}$$

is properly stronger than CD. Further variants of CD will be studied in Part II of the paper.

We define disjunction  $\vee$ , conditional  $\rightarrow$ , and the existential quantifier  $\exists$  in the standard manner in terms of negation  $\neg$ , conjunction  $\wedge$ , and universal quantifier  $\forall$ : accordingly, their corresponding operations  $x \vee y$ ,  $x \rightarrow y$ , and  $\exists v x$  are defined as

$$x \vee y := \neg((\neg x) \wedge (\neg y)), \quad x \rightarrow y := (\neg x) \vee y, \quad \text{and} \quad \exists v x := \neg(\forall v(\neg x)).$$

It readily follows from T4 to T6 that T commutes with  $\vee$ ,  $\rightarrow$ , and  $\exists$ . In addition, the following compositional rules of determinateness for these connectives are provable in  $CD_0$ :

$$\forall x \forall y \left( \text{Sent}(x \vee y) \rightarrow (\mathbf{D}(x \vee y) \leftrightarrow ((\mathbf{D}x \wedge \mathbf{D}y) \vee (\mathbf{D}x \wedge \mathbf{T}x) \vee (\mathbf{D}y \wedge \mathbf{T}y))) \right), \quad (3)$$

$$\forall x \forall y \left( \text{Sent}(x \rightarrow y) \rightarrow (\mathbf{D}(x \rightarrow y) \leftrightarrow ((\mathbf{D}x \wedge \mathbf{D}y) \vee (\mathbf{D}x \wedge \mathbf{F}x) \vee (\mathbf{D}y \wedge \mathbf{T}y))) \right), \quad (4)$$

$$\forall v \forall x \left( \text{Sent}(\exists v x) \rightarrow (\mathbf{D}(\exists v x) \leftrightarrow (\forall s \mathbf{D}x(s/v) \vee \exists s(\mathbf{D}x(s/v) \wedge \mathbf{T}x(s/v)))) \right). \quad (5)$$

Concluding this section, we mention some simple observations. First, we show that the uniform determinate disquotation schema is provable in CD, as mentioned above.

**LEMMA 3.1.** *For all  $\mathcal{L}$ -formulae  $\phi(x_1, \dots, x_k)$  with at most  $x_1, \dots, x_k$  free the theory CD proves the following:*

$$\text{DDS } \forall t_1 \dots \forall t_n \left( \mathbf{D}^\Gamma \phi(t_1, \dots, t_n)^\neg \rightarrow (\mathbf{T}^\Gamma \phi(t_1, \dots, t_n)^\neg \leftrightarrow \phi(t_1^\circ, \dots, t_n^\circ)) \right).$$

*In particular,  $\mathbf{D}^\Gamma \phi^\neg \rightarrow (\mathbf{T}^\Gamma \phi^\neg \leftrightarrow \phi)$  is provable for all sentences  $\phi$ .*

The lemma can be proved by a metatheoretic induction on the complexity of  $\phi$ , where the complexity of a formula  $\phi$ ,  $cp(\phi)$  for short, is standardly defined: every atomic formula has the complexity 0;  $cp(\neg\phi) = cp(\forall x\phi) = cp(\phi) + 1$ ;  $cp(\phi \wedge \psi) = \max\{cp(\phi), cp(\psi)\} + 1$ .

By an induction within CD all sentences of the base language can be shown to be determinate.  $\text{Sent}_0(x)$  expresses that  $x$  is a sentence of  $\mathcal{L}_0$ .

**LEMMA 3.2.**  $CD_0 \vdash \forall x (\text{Sent}_0(x) \rightarrow \mathbf{D}x)$ .

Using these two lemmata, we can establish the T-sentences for all sentences of the base language.

**LEMMA 3.3.** *For all  $\mathcal{L}_0$ -formulae  $\phi(x_1, \dots, x_k)$  with only  $x_1, \dots, x_k$  free, we have*

$$CD \vdash \forall t_1 \dots \forall t_k (\mathbf{T}^\Gamma \phi(t_1, \dots, t_k)^\neg \leftrightarrow \phi(t_1^\circ, \dots, t_k^\circ)).$$

Obviously, we cannot have the unrestricted T-schema in any axiomatic theory of truth. Some truth theories feature one direction of the T-schema without any restrictions. However, both directions are incompatible with CD.

**PROPOSITION 3.4.** *The axiom schemata (T-Out) and (T-In) are defined as follows:*

- (T-Out)  $\mathbf{T}^\Gamma \phi^\neg \rightarrow \phi$  for all  $\phi \in \mathcal{L}$ .
- (T-In)  $\phi \rightarrow \mathbf{T}^\Gamma \phi^\neg$  for all  $\phi \in \mathcal{L}$ .

*Then, each of (T-Out) and (T-In) is inconsistent with T4 over arithmetic, and thus with CD.*

PROOF. We first derive a contradiction from CD + (T-Out). Let  $\lambda$  be the liar sentence with  $\lambda \leftrightarrow \neg T \ulcorner \lambda \urcorner$  as above. The liar sentence  $\lambda$  implies  $T \ulcorner \neg \lambda \urcorner$  by T4, from which  $\neg \lambda$  follows by (T-Out). Hence we have  $\neg \lambda$ . However,  $\neg \lambda$  is equivalent to  $T \ulcorner \lambda \urcorner$  and thus implies  $\lambda$  by (T-Out).

The inconsistency of (T-In) is shown in a similar way by using the other direction of T4. ⊥

A consequence of the compositional axioms is the thorough classicality of CD. CD proves that (classical) logic is truth-preserving for the entire language  $\mathcal{L}$ .

LEMMA 3.5. *The following holds for all formulae  $\phi(x)$  and canonical provability predicates constructed from  $\phi(x)$  as defining the axioms:*

$$CD \vdash \forall x (\phi(x) \rightarrow \text{Sent}(x) \wedge Tx) \rightarrow \forall y (\text{Bew}_{\phi(x)}(y) \rightarrow Ty).$$

Here,  $\text{Bew}_{\phi(x)}(y)$  expresses that  $y$  is provable from all sentences  $\psi$  with  $\phi(\ulcorner \psi \urcorner)$ .

PROOF. The claim is shown by induction on the length of a proof. Since T commutes with all logical connectives, quantifiers, and identity in CD, CD proves that every logical axiom is true; note that the extensionality axioms R1 and R2 are needed to show the logical axioms for identity. Every non-logical axiom is true by the assumption. The base step is thereby obtained. The induction step is straightforward, since commutativity ensures that every logical inference rule preserves truth. ⊥

Fischer et al. [12] criticized truth theories such as FS and variants of KF formulated in classical logic, because they are incompatible with reflection principles for logic. In CD the soundness of a calculus for predicate logic is provable. This is a trivial consequence of the lemma above if  $\text{Bew}(x)$  expresses provability in pure logic.

COROLLARY 3.6. *CD proves that every logically valid sentence (in classical logic) is true: namely,  $CD \vdash \forall x (\text{Bew}(x) \wedge \text{Sent}(x) \rightarrow Tx)$ .*

**§4. Semantics.** McGee [31] showed that certain theories of truth that are thoroughly classical are  $\omega$ -inconsistent. In particular, the system FS that also features axioms T1 and T4 –T6 is  $\omega$ -inconsistent. In order to defend CD and its variants we therefore do not only show that they are consistent, but also that they possess  $\omega$ -models.

In this section, we will give an  $\omega$ -model of  $CD^+$ , namely, a model of  $CD^+$  whose  $\mathcal{L}_0$ -reduct is the standard model of arithmetic. The existence of an  $\omega$ -model entails that  $CD^+$  and thus CD are  $\omega$ -consistent.<sup>5</sup>

**4.1. The model.** Let us define  $\mathcal{L}$ -formulae  $\mathcal{D}_i(x)$  ( $1 \leq i \leq 6$ ) as follows:

$$\begin{aligned} \mathcal{D}_1(x) &:\Leftrightarrow \exists s \exists t (x = (s \doteq t)), \\ \mathcal{D}_2(x) &:\Leftrightarrow \exists s (x = \ulcorner \top s \urcorner \wedge \text{Ds}^\circ), \\ \mathcal{D}_3(x) &:\Leftrightarrow \exists s (x = \ulcorner \text{D}s \urcorner \wedge \text{Ds}^\circ), \\ \mathcal{D}_4(x) &:\Leftrightarrow \text{Sent}(x) \wedge \exists y (x = \ulcorner \neg y \urcorner \wedge \text{Dy}), \end{aligned}$$

---

<sup>5</sup>The existence of an  $\omega$ -model of  $CD^+$  in itself also follows from Theorem 7.11.

$$\mathcal{D}_5(x) :\Leftrightarrow \text{Sent}(x) \wedge \exists y \exists z \left( x = y \wedge z \wedge \left( (Dy \wedge Dz) \vee (D\neg y \wedge Fy) \vee (D\neg z \wedge Fz) \right) \right),$$

$$\mathcal{D}_6(x) :\Leftrightarrow \text{Sent}(x) \wedge \exists v \exists y \left( x = \forall v y \wedge \left( \forall s Dy(s/v) \vee \exists s (D\neg(y(s/v)) \wedge F(y(s/v))) \right) \right).$$

We thereby define  $\mathcal{D}(x) :\Leftrightarrow \bigvee_{1 \leq i \leq 6} \mathcal{D}_i(x)$ .  $\mathcal{D}$  describes the closure condition of a determinateness predicate  $D$  (relative to a fixed interpretation of the truth predicate  $T$ ).

By  $(\mathbb{N}, X, Y)$  let us denote the  $\mathcal{L}$ -structure with domain  $\omega$ , the set of natural numbers, in which  $D$  is interpreted by  $X$ ,  $T$  is interpreted by  $Y$ , and all the other symbols receive the standard interpretations. Both  $D$  and  $T$  occur only positively in  $\mathcal{D}$ . Hence, for each  $Y \subset \omega$ ,  $\mathcal{D}$  induces the following monotone operator  $\Gamma_{\mathcal{D}[Y]}: \mathcal{P}(\omega) \rightarrow \mathcal{P}(\omega)$  so that

$$\Gamma_{\mathcal{D}[Y]}(X) = \{n \in \omega \mid (\mathbb{N}, X, Y) \models \mathcal{D}(n)\}.$$

In general, given a monotone operator  $\Gamma: \mathcal{P}(\omega) \rightarrow \mathcal{P}(\omega)$ , we say that  $X \subset \omega$  is  $\Gamma$ -closed if  $\Gamma(X) \subset X$ , and that  $X \subset \omega$  is a  $\Gamma$ -fixed point if  $\Gamma(X) = X$ .

Given an  $\mathcal{L}_0$ -formula or  $\mathcal{L}_0$ -term  $e$ , let us denote its standard interpretation by  $e^{\mathbb{N}}$  and its Gödel number by  $\#e$ . For instance,  $\text{Sent}^{\mathbb{N}}$  denote the set of the Gödel numbers of all  $\mathcal{L}$ -sentences.

LEMMA 4.1. *For every  $X, Y \subset \text{Sent}^{\mathbb{N}}$ , the following are equivalent.<sup>6</sup>*

1.  $X$  is a  $\Gamma_{\mathcal{D}[Y]}$ -fixed point.
2.  $(\mathbb{N}, X, Y)$  is a model of all the determinateness axioms D1–D6.

For each ordinal  $\alpha$ , we define sets  $D_\alpha$  and  $T_\alpha$  as follows:

$$\begin{aligned} D_0 &:= \emptyset, & T_0 &:= \emptyset, \\ D_{\xi+1} &:= \Gamma_{\mathcal{D}[T_\xi]}(D_\xi), & T_{\xi+1} &:= \{\#\phi \in \text{Sent}^{\mathbb{N}} \mid (\mathbb{N}, D_\xi, T_\xi) \models \phi\}, \\ D_\lambda &:= \bigcup_{\xi < \lambda} D_\xi, & T_\lambda &:= \bigcup_{\xi < \lambda} T_\xi \cap D_\xi. \end{aligned}$$

It is obvious from the definition that  $T_\lambda = D_\lambda \cap T_\lambda$  for all limit ordinals  $\lambda$ . We can also show by induction on ordinals that  $D_\xi, T_\xi \subset \text{Sent}^{\mathbb{N}}$  for all ordinals  $\xi$ .

The next two propositions are obvious by definition.

LEMMA 4.2. *Let  $\xi$  be any ordinal,  $a$  be a closed  $\mathcal{L}_0$ -term,  $\phi$  and  $\psi$  be  $\mathcal{L}$ -sentences, and  $\theta(x)$  be an  $\mathcal{L}$ -formula with only a single variable  $x$  free.*

1. If  $\xi > 0$  and  $\phi$  is  $\mathcal{L}_0$ -atomic, then  $\#\phi \in D_\xi$ .
2.  $a^{\mathbb{N}} \in D_\xi$  iff  $\#\Gamma a \in D_{\xi+1}$ .
3.  $a^{\mathbb{N}} \in D_\xi$  iff  $\#\mathbf{D}a \in D_{\xi+1}$ .

<sup>6</sup> $\mathcal{D}_5$  and  $\mathcal{D}_6$  are defined so that the condition 1 implies that  $(\mathbb{N}, X, Y)$  satisfies not the original axioms D5 and D6, but the following variants of them:

$$(D5') \quad \forall x \forall y \left( \text{Sent}(x \wedge y) \rightarrow (D(x \wedge y) \leftrightarrow ((Dx \wedge Dy) \vee (D\neg x \wedge Fx) \vee (D\neg y \wedge Fy))) \right),$$

$$(D6') \quad \forall v \forall x \left( \text{Sent}(\forall v x) \rightarrow (D(\forall v x) \leftrightarrow (\forall t Dx(t/v) \vee \exists t (D\neg x(t/v) \wedge Fx(t/v)))) \right);$$

but they are equivalent to the original D5 and D6 anyway due to the axiom D4; we adopt this definition of  $\mathcal{D}_5$  and  $\mathcal{D}_6$  for a purely technical reason.

4.  $\# \phi \in D_\xi$  iff  $\# \neg \phi \in D_{\xi+1}$ .
5.  $\#(\phi \wedge \psi) \in D_{\xi+1}$ , iff either  $\# \phi, \# \psi \in D_\xi$  or  $\# \neg \phi \in D_\xi \cap T_\xi$  or  $\# \neg \psi \in D_\xi \cap T_\xi$ .
6.  $\# \forall x \theta(x) \in D_{\xi+1}$ , iff either of the following holds:
  - $\# \theta(b) \in D_\xi$  for all  $\mathcal{L}_0$ -closed terms  $b$ ;
  - $\# \neg \theta(b) \in D_\xi \cap T_\xi$  for some  $\mathcal{L}_0$ -closed term  $b$ .

PROOF. For example, we have  $\# \mathbf{T}a = \mathbb{T}^\mathbb{N}(\#a)$  and  $((\#a)^\circ)^\mathbb{N} = a^\mathbb{N}$ , from which claim 2 follows, since we have

$$a^\mathbb{N} \in D_\xi \iff (\mathbb{N}, D_\xi, T_\xi) \models \mathcal{D}_2(\# \mathbf{T}a) \iff \# \mathbf{T}a \in D_{\xi+1}. \quad \dashv$$

LEMMA 4.3. Let  $\xi$  be any ordinal,  $a$  and  $b$  be closed  $\mathcal{L}_0$ -terms,  $\phi$  and  $\psi$  be  $\mathcal{L}$ -sentences, and  $\theta(x)$  be an  $\mathcal{L}$ -formula only with some single variable  $x$  free.

1.  $\#(a = b) \in T_{\xi+1}$  iff  $a^\mathbb{N} = b^\mathbb{N}$ .
2.  $a^\mathbb{N} \in T_\xi$  iff  $\# \mathbf{T}a \in T_{\xi+1}$ .
3.  $a^\mathbb{N} \in D_\xi$  iff  $\# \mathbf{D}a \in T_{\xi+1}$ .
4.  $\# \neg \phi \in T_{\xi+1}$  iff  $\phi \notin T_{\xi+1}$ .
5.  $\#(\phi \wedge \psi) \in T_{\xi+1}$  iff  $\# \phi \in T_{\xi+1}$  and  $\# \psi \in T_{\xi+1}$ .
6.  $\# \forall x \theta(x) \in T_{\xi+1}$  iff  $\# \theta(c) \in T_{\xi+1}$  for all closed  $\mathcal{L}_0$ -terms  $c$ .

Let us define a formula  $x \approx y$  as follows:

$$x \approx y \text{ :} \Leftrightarrow x = y \vee \exists z \exists v \exists s \exists t (\text{Sent}(\forall v z) \wedge s^\circ = t^\circ \wedge x = z(s/v) \wedge y = z(t/v)).$$

Namely, for natural numbers  $n, m \in \mathbb{N}$ ,  $n \approx^\mathbb{N} m$  means that either  $n = m$  or  $n$  and  $m$  code numerically equivalent  $\mathcal{L}$ -sentences in the sense that they are substitution instances of the same  $\mathcal{L}$ -formula with closed terms with the same values. Since all the existential quantifiers in the definition can actually be bounded by the values of some primitive recursive functions on  $x$  and  $y$ , the relation  $\approx^\mathbb{N}$  is primitive recursive in the value function (or the index function  $[\cdot]$  of the primitive recursive functions) and thus provably recursive in PA. The extensionality axioms R1 and R2 say that truth and determinateness are invariant across numerically equivalent  $\mathcal{L}$ -sentences, namely, sentences  $\phi$  and  $\psi$  with  $\# \phi \approx^\mathbb{N} \# \psi$ .<sup>7</sup> The following proposition is shown by a straightforward induction on ordinals.

PROPOSITION 4.4. Let  $\xi$  be any ordinal and  $n \approx^\mathbb{N} m$ .

1.  $n \in D_\xi$  iff  $m \in D_\xi$ .
2.  $n \in T_\xi$  iff  $m \in T_\xi$ ; note that if  $\# \phi \approx^\mathbb{N} \# \psi$  then  $(\mathbb{N}, X, Y) \models \phi \leftrightarrow \psi$ .

The next is the main lemma of this subsection.

LEMMA 4.5. For all ordinals  $\xi$  and  $\eta$  with  $\xi \leq \eta$ , the following hold:

1.  $D_\xi \cap T_\xi = D_\xi \cap T_\eta$  and
2.  $D_\xi \subset D_\eta$ .

<sup>7</sup>The first disjunct ‘ $x = y$ ’ is added for a purely technical reason, namely, for incorporating the ordinary logical axioms for atomic sentences of the form  $\mathbf{D}a$  or  $\mathbf{T}a$  to the axioms (Ax2) and (Ax3) of the semi-formal system  $\text{CD}^\infty$ , which will be introduced later in Section 7.

PROOF. By simultaneous induction on  $\xi$ . The base step is trivial. Let  $\xi$  be a limit ordinal and  $\eta \geq \xi$ . For each  $n \in D_\xi$ , there is some  $\zeta < \xi$  such that  $n \in D_\zeta$ , for which we have  $D_\zeta \cap T_\xi = D_\zeta \cap T_\zeta = D_\zeta \cap T_\eta$  and  $D_\zeta \subset D_\eta$  by the induction hypothesis.

Let  $\xi$  be a successor ordinal. Claim 1 is shown by sub-induction on  $\eta$ ; we can assume  $\eta > \xi$ . First, suppose  $\eta$  is a limit. If  $n \in D_\xi \cap T_\xi$ , then  $n \in T_\eta$  by definition. Let  $n \in D_\xi \cap T_\eta$  for the converse. By definition, there is some  $\zeta < \eta$  such that  $n \in D_\zeta \cap T_\zeta$ . If  $\zeta < \xi$ , then  $D_\zeta \cap T_\zeta \subset T_\xi$  by the induction hypothesis for 1; otherwise,  $\xi \leq \zeta < \eta$  and thus  $D_\xi \cap T_\zeta = D_\xi \cap T_\xi$  by the sub-induction hypothesis. Next, suppose  $\eta$  is a successor, and take any  $n = \# \phi \in D_\xi$  ( $\subset \text{Sent}^{\mathbb{N}}$ ) for an  $\mathcal{L}$ -sentence  $\phi$ . If  $\phi$  is an  $\mathcal{L}_0$ -atomic, then the claim follows from Lemmata 4.2.1 and 4.3.1. If  $\phi = \text{Ta}$  for some closed  $\mathcal{L}_0$ -term  $a$ , then we have  $a^{\mathbb{N}} \in D_{\xi-1}$  by Lemma 4.2.2, and thus we obtain

$$n \in T_\xi \stackrel{4.3.2}{\Leftrightarrow} a^{\mathbb{N}} \in T_{\xi-1} \stackrel{\text{IH}}{\Leftrightarrow} a^{\mathbb{N}} \in T_{\eta-1} \stackrel{4.3.2}{\Leftrightarrow} n \in T_\eta,$$

where ‘IH’ denotes the induction hypothesis. If  $\phi = \text{Da}$ , then we have  $a^{\mathbb{N}} \in D_{\xi-1}$  by Lemma 4.2.3, and thus  $n \in T_\xi$  and  $a^{\mathbb{N}} \in D_{\eta-1}$  by Lemma 4.3.3 and the induction hypothesis for 2, respectively, the latter of which implies  $n \in T_\eta$  again by Lemma 4.3.3; hence, both  $n \in T_\xi$  and  $n \in T_\eta$  hold. We move on to the cases for complex  $\phi$ s. Let  $\phi = \psi \wedge \theta$ . By Lemma 4.2.5,  $n \in D_\xi$  implies

$$(\# \psi, \# \theta \in D_{\xi-1}) \vee (\# \neg \psi \in D_{\xi-1} \cap T_{\xi-1}) \vee (\# \neg \theta \in D_{\xi-1} \cap T_{\xi-1}).$$

If the first disjunct holds, then we have

$$n \in T_\xi \stackrel{4.3.5}{\Leftrightarrow} \# \psi, \# \theta \in T_\xi \stackrel{\text{IH}}{\Leftrightarrow} \# \psi, \# \theta \in T_{\xi-1} \stackrel{\text{IH}}{\Leftrightarrow} \# \psi, \# \theta \in T_\eta \stackrel{4.3.5}{\Leftrightarrow} n \in T_\eta.$$

If the second disjunct holds, then we have  $\# \neg \psi \in T_\xi \cap T_\eta$  by the induction hypothesis for 1, from which we can infer

$$\# \neg \psi \in T_\xi \cap T_\eta \stackrel{4.3.4}{\Leftrightarrow} \# \psi \notin T_\xi \text{ and } \# \psi \notin T_\eta \stackrel{4.3.5}{\Rightarrow} \# \phi \notin T_\xi \text{ and } \# \phi \notin T_\eta.$$

The case where the third disjunct holds can be similarly treated. The claim for the remaining cases where  $\phi$  is of the form  $\neg \psi$  or  $\forall x \psi(x)$  can be shown similarly.

Claim 2 for a successor  $\xi$  is also shown by sub-induction on  $\eta$ . As before, we can assume  $\eta > \xi$ , and the limit case is obvious. Let  $\eta$  be a successor ordinal and take any  $n = \# \phi \in D_\xi$  for an  $\mathcal{L}$ -sentence  $\phi$ . If  $\phi$  is  $\mathcal{L}_0$ -atomic, then the claim trivially obtains by Lemma 4.2.1. If either  $\phi = \text{Ta}$  or  $\phi = \text{Da}$ , then we have

$$n \in D_\xi \stackrel{4.2.2 \text{ or } 3}{\Leftrightarrow} a^{\mathbb{N}} \in D_{\xi-1} \stackrel{\text{IH}}{\Rightarrow} a^{\mathbb{N}} \in D_{\eta-1} \stackrel{4.2.2 \text{ or } 3}{\Leftrightarrow} n \in D_\eta.$$

The claim for the other cases for complex  $\phi$ s can be straightforwardly shown. For instance, if  $\phi = \psi \wedge \theta$ , then we have

$$\begin{aligned} n \in D_\xi &\stackrel{4.2.5}{\Leftrightarrow} (\# \psi, \# \theta \in D_{\xi-1}) \vee (\# \neg \psi \in D_{\xi-1} \cap T_{\xi-1}) \vee (\# \neg \theta \in D_{\xi-1} \cap T_{\xi-1}) \\ &\stackrel{\text{IH}}{\Rightarrow} (\# \psi, \# \theta \in D_{\eta-1}) \vee (\# \neg \psi \in D_{\eta-1} \cap T_{\eta-1}) \vee (\# \neg \theta \in D_{\eta-1} \cap T_{\eta-1}) \\ &\stackrel{4.2.5}{\Leftrightarrow} n \in D_\eta. \end{aligned}$$

The other cases can be treated in a similar way. ⊣

LEMMA 4.6. *Let  $\phi$  be an  $\mathcal{L}$ -sentence and  $\xi$  an ordinal.*

1. *If  $\#\phi \in D_\xi$ , then  $\#\phi \in T_\xi$  iff  $(\mathbb{N}, D_\xi, T_\xi) \models \phi$ .*
2. *If  $a^\mathbb{N} = \#\phi \in D_\xi$ , then  $(\mathbb{N}, D_\xi, T_\xi) \models \phi$  iff  $(\mathbb{N}, D_\xi, T_\xi) \models \text{Ta}$ .*

PROOF. 1. Let  $\#\phi \in D_\xi$ . Then, we have

$$\#\phi \in T_\xi \stackrel{4.5.1}{\Leftrightarrow} \#\phi \in T_{\xi+1} \stackrel{\text{def}}{\Leftrightarrow} \#(\mathbb{N}, D_\xi, T_\xi) \models \phi.$$

2. Let  $a^\mathbb{N} = \#\phi \in D_\xi$ . Then, we have

$$(\mathbb{N}, D_\xi, T_\xi) \models \phi \stackrel{\text{by } 1}{\Leftrightarrow} \#\phi \in T_\xi \stackrel{\text{def}}{\Leftrightarrow} (\mathbb{N}, D_\xi, T_\xi) \models \text{Ta}. \quad \dashv$$

- LEMMA 4.7. 1. *For each  $\xi$  and  $\mathcal{L}$ -sentence  $\phi$ , either  $\#\phi \notin T_\xi$  or  $\#\neg\phi \notin T_\xi$ .*  
 2. *For each  $\xi$  and  $\mathcal{L}$ -sentence  $\phi$ , if  $\#\phi \in D_\xi$ , then either  $\#\phi \in T_\xi$  or  $\#\neg\phi \in T_\xi$ .*

PROOF. 1. By induction on  $\xi$ . The claim for a non-limit  $\xi$  immediately follows by definition. Let  $\xi$  be a limit, and suppose for contradiction that  $\#\phi \in T_\xi$  and  $\#\neg\phi \in T_\xi$ . Then, there are some  $\zeta, \eta < \xi$  such that  $\#\phi \in T_\zeta \cap D_\zeta$  and  $\#\neg\phi \in T_\eta \cap D_\eta$ , which is impossible by Lemma 4.5 and the induction hypothesis.

2. By induction on  $\xi$ . If  $\xi = 0$ , then  $D_0 = \emptyset$  and the claim trivially holds. If  $\xi$  is a successor, then the claim immediately follows by Lemma 4.3.4. Finally, if  $\xi$  is a limit and  $\#\phi \in D_\xi$ , then  $\#\phi \in D_\zeta$  for some  $\zeta < \xi$ , and thus either  $\#\phi \in T_\zeta$  or  $\neg\phi \in T_\zeta$  by the induction hypothesis, from which the claim follows by Lemma 4.5.1.  $\dashv$

Let  $\omega_1$  denote the least uncountable ordinal. By Lemma 4.5, we have  $D_\xi \subset D_\eta$  and  $D_\xi \cap T_\xi \subset D_\eta \cap T_\eta$  for all  $\xi \leq \eta$ . Hence, by the standard cardinality consideration, we have the following nice properties of  $D_{\omega_1}$  and  $T_{\omega_1}$ :

$$D_{\omega_1+1} = D_{\omega_1}, \text{ and } D_{\omega_1+1} \cap T_{\omega_1+1} = D_{\omega_1} \cap T_{\omega_1} = T_{\omega_1}; \quad (6)$$

in particular,  $D_{\omega_1}$  is a  $\Gamma_{\mathcal{D}[T_{\omega_1}]}$ -fixed point. Let us denote  $D_{\omega_1}$  and  $T_{\omega_1}$  by  $D_\infty$  and  $T_\infty$ , respectively. We finally give an  $\omega$ -model of CD.

THEOREM 4.8. *Let  $\mathbb{T}_\infty := \{\#\phi \mid (\mathbb{N}, D_\infty, T_\infty) \models \phi\}$ . Then,  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{CD}^+$ . Hence, in particular,  $\text{CD}^+$  and CD have an  $\omega$ -model and thus are  $\omega$ -consistent.*

PROOF. By definition, it immediately follows that  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty)$  is a model of T1, T2<sup>+</sup>, and T4–T6. To see  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{T3}$ , take any  $n \in \text{ClTerm}^\mathbb{N}$  and suppose  $(n^\circ)^\mathbb{N} \in D_\infty$ . Then, there exist some closed  $\mathcal{L}_0$ -term  $a$  and  $\mathcal{L}$ -sentence  $\phi$  such that  $n = \#a$  and  $(n^\circ)^\mathbb{N} = a^\mathbb{N} = \#\phi \in D_\infty (\subset \text{Sent}^\mathbb{N})$ . Hence, we have

$$\begin{aligned} (\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{T}\ddagger n &\stackrel{\text{def}}{\Leftrightarrow} (\mathbb{N}, D_\infty, T_\infty) \models \text{Ta} \stackrel{4.6.2}{\Leftrightarrow} (\mathbb{N}, D_\infty, T_\infty) \models \phi \\ &\stackrel{\text{def}}{\Leftrightarrow} (\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{T}n^\circ. \end{aligned}$$

Since  $D_\infty$  is a  $\Gamma_{\mathcal{D}[T_\infty]}$ -fixed point,  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty)$  satisfies D1–D4 (which does not depend on the interpretation of T). To show  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{D5}$ , take any  $\mathcal{L}$ -sentences  $\phi$  and  $\psi$ . By Lemma 4.6.1, we have  $D_\infty \cap \mathbb{T}_\infty = D_\infty \cap T_\infty$  and thus obtain

$$\begin{aligned} \#\phi \wedge \psi \in D_\infty &\stackrel{4.2.5}{\Leftrightarrow} (\#\phi, \#\psi \in D_\infty) \vee (\#\neg\phi \in D_\infty \cap T_\infty) \vee (\#\neg\psi \in D_\infty \cap T_\infty) \\ &\Leftrightarrow (\#\phi, \#\psi \in D_\infty) \vee (\#\neg\phi \in D_\infty \cap \mathbb{T}_\infty) \vee (\#\neg\psi \in D_\infty \cap \mathbb{T}_\infty). \end{aligned}$$

We can similarly show  $(\mathbb{N}, D_\infty, \mathbb{T}_\infty) \models \text{D6}$  using  $D_\infty \cap T_\infty = D_\infty \cap \mathbb{T}_\infty$ .  $\dashv$

**4.2. The minimality of  $D$ .** In this section, we will show that  $D_\infty$  ( $:= D_{\omega_1}$ ) has a certain ‘minimality’ property.

Let us define  $\mathcal{L}$ -formulae  $\mathcal{T}_j(x)$  ( $1 \leq j \leq 6$ ) as follows:

$$\begin{aligned} \mathcal{T}_1(x) &:\Leftrightarrow \exists s \exists t \left( (x = (s \dot{=} t) \wedge s^\circ = t^\circ) \vee (x = (\neg s \dot{=} t) \wedge s^\circ \neq t^\circ) \right), \\ \mathcal{T}_2(x) &:\Leftrightarrow \exists s \left( (x = \mathbf{T}s \wedge \mathbf{D}s^\circ \wedge \mathbf{T}s^\circ) \vee (x = \neg \mathbf{T}s \wedge \mathbf{D}(\neg s^\circ) \wedge \mathbf{F}s^\circ) \right), \\ \mathcal{T}_3(x) &:\Leftrightarrow \exists s (x = \mathbf{D}s \wedge \mathbf{D}s^\circ), \\ \mathcal{T}_4(x) &:\Leftrightarrow \text{Sent}(x) \wedge \exists y (x = \neg \neg y \wedge \mathbf{T}y), \\ \mathcal{T}_5(x) &:\Leftrightarrow \text{Sent}(x) \wedge \exists y \exists z \left( (x = (y \wedge z) \wedge \mathbf{T}y \wedge \mathbf{T}z) \vee (x = \neg(y \wedge z) \wedge (\mathbf{F}y \vee \mathbf{F}z)) \right), \\ \mathcal{T}_6(x) &:\Leftrightarrow \text{Sent}(x) \wedge \exists v \exists y \left( (x = \forall v y \wedge \forall s \mathbf{T}y(s/v)) \vee (x = \neg \forall v y \wedge \exists s \mathbf{F}y(s/v)) \right). \end{aligned}$$

We define  $\mathcal{T}(x) :\Leftrightarrow \bigvee_{1 \leq j \leq 6} \mathcal{T}_j(x)$ . Note that  $\mathcal{T}$  says nothing about the Gödel numbers of  $\mathcal{L}$ -sentences of the form  $\neg Da$ . Both  $\mathbf{D}$  and  $\mathbf{T}$  occur only positively in  $\mathcal{T}$ . Hence, for each given  $X \subset \omega$ , it induces a monotone operator  $\Gamma_{\mathcal{T}[X]} : \mathcal{P}(\omega) \rightarrow \mathcal{P}(\omega)$  such that

$$\Gamma_{\mathcal{T}[X]}(Y) = \{n \in \omega \mid (\mathbb{N}, X, Y) \models \mathcal{T}(n)\};$$

note that  $\Gamma_{\mathcal{T}[X]}(Y) \subset \text{Sent}^{\mathbb{N}}$  for all  $X, Y \subset \omega$ .

**LEMMA 4.9.** *Let  $X, Y \subset \omega$ . If  $(\mathbb{N}, X, Y) \models \text{CD}$ , then  $Y$  is  $\Gamma_{\mathcal{T}[X]}$ -closed.*

**PROOF.** The proof is routine. Take any  $n \in \Gamma_{\mathcal{T}[X]}(Y) \subset \text{Sent}^{\mathbb{N}}$  and let  $n = \#\phi$  where  $\phi$  is an  $\mathcal{L}$ -sentence. If  $\phi = Da$  for a closed  $\mathcal{L}_0$ -term  $a$ , then we have  $a^{\mathbb{N}} \in X$  and thus  $(\mathbb{N}, X, Y) \models \mathbf{D}a^\circ$ , which implies  $(\mathbb{N}, X, Y) \models \mathbf{T}\mathbf{D}a$ , namely,  $\#\phi \in Y$ , because  $(\mathbb{N}, X, Y) \models \mathbf{T}2$ ; note that it is never the case that  $\#\phi = \neg Da$  for any closed  $\mathcal{L}_0$ -term  $a$ . If  $\phi = \neg Ta$  for a closed  $\mathcal{L}_0$ -term  $a$ , then  $(\neg a)^{\mathbb{N}} \in X \cap Y$  and thus  $a^{\mathbb{N}} \in X \setminus Y$  by  $(\mathbb{N}, X, Y) \models \mathbf{D}4 \wedge \mathbf{T}4$ , which implies  $\#\phi \in Y$  because  $(\mathbb{N}, X, Y) \models \mathbf{T}3 \wedge \mathbf{T}4$ . Next, let  $\phi = \neg(\psi \wedge \theta)$ . Then, either  $\#\neg\psi \in Y$  or  $\#\neg\theta \in Y$ . If the former is the case, then  $\#\psi \notin Y$  by  $(\mathbb{N}, X, Y) \models \mathbf{T}4$  and thus  $\#\phi \notin Y$  by  $(\mathbb{N}, X, Y) \models \mathbf{T}5$ , from which we get  $\#\neg\phi \in Y$  by  $(\mathbb{N}, X, Y) \models \mathbf{T}4$ . We leave the remaining cases to the reader.  $\dashv$

**LEMMA 4.10.** *Let  $X, Y \subset \omega$ . Suppose the following:*

- (a)  $X$  is  $\Gamma_{\mathcal{D}[Y]}$ -closed;
- (b)  $Y$  is  $\Gamma_{\mathcal{T}[X]}$ -closed.

*Then we have  $D_\infty \subset X$  and  $T_\infty \subset Y$  and thus, by (6) above,  $T_\infty \subset D_\infty \cap Y$ .*

**PROOF.** It suffices to show the following under the supposition of (a) and (b): for all  $\mathcal{L}$ -sentences  $\phi$  and ordinals  $\alpha$ ,

$$\#\phi \in D_\alpha \rightarrow \#\phi \in X, \tag{7}$$

$$\#\phi \in D_\alpha \cap T_\alpha \rightarrow \#\phi \in Y, \tag{8}$$

$$\#\phi \in D_\alpha \setminus T_\alpha \rightarrow \#\neg\phi \in Y. \tag{9}$$

They are shown by simultaneous induction on  $\alpha$ . The base step is trivial. If  $\alpha$  is a limit, the claim follows from the induction hypothesis using Lemma 4.5. Let  $\alpha$  be

a successor. The claim for the case where  $\phi$  is an  $\mathcal{L}_0$ -atomic is obvious. We will go through the remaining cases.

Let  $\phi = Ta$  and suppose  $\#\phi \in D_\alpha$ . We have  $a^{\mathbb{N}} \in D_{\alpha-1} (\subset \text{Sent}^{\mathbb{N}})$  by Lemma 4.2.2. Then, we have  $a^{\mathbb{N}} \in X$  by the induction hypothesis and thus  $\#\phi \in X$  by (a); hence (7) holds for  $\phi$ . The other claims (8) and (9) are obtained as follows:

$$\begin{aligned} \#\phi \in T_\alpha &\stackrel{4.3.2}{\Leftrightarrow} a^{\mathbb{N}} \in T_{\alpha-1} \stackrel{\text{IH}}{\Rightarrow} a^{\mathbb{N}} \in X \cap Y \stackrel{(b)}{\Rightarrow} \#\phi \in Y; \\ \#\phi \notin T_\alpha &\stackrel{4.3.2}{\Leftrightarrow} a^{\mathbb{N}} \notin T_{\alpha-1} \stackrel{\text{IH}}{\Rightarrow} (\neg a)^{\mathbb{N}} \in X \cap Y \stackrel{(b)}{\Rightarrow} \#\neg\phi \in Y. \end{aligned}$$

Let  $\phi = Da$  and suppose  $\#\phi \in D_\alpha$ . We have  $a^{\mathbb{N}} \in D_{\alpha-1}$  by Lemma 4.2.3. Then, we get  $a^{\mathbb{N}} \in X$  by the induction hypothesis, which implies  $\#\phi \in Y$  by (b); hence, (7) and (8) hold for  $\phi$ . The other claim (9) trivially holds, since  $\#\phi \notin T_\alpha$  would imply  $a^{\mathbb{N}} \notin D_{\alpha-1}$  by Lemma 4.3.3, which is absurd.

Let  $\phi = \neg\psi$  and suppose  $\#\phi \in D_\alpha$ . We have  $\#\psi \in D_{\alpha-1}$  by Lemma 4.2.4. Then, it follows that  $\#\psi \in X$  by the induction hypothesis and thus  $\#\phi \in X$  by (a). The other claims (8) and (9) are obtained as follows:

$$\begin{aligned} \#\phi \in T_\alpha &\stackrel{4.3.4}{\Leftrightarrow} \#\psi \notin T_\alpha \stackrel{4.5.1}{\Leftrightarrow} \#\psi \notin T_{\alpha-1} \stackrel{\text{IH}}{\Rightarrow} \#\phi \in Y; \\ \#\phi \notin T_\alpha &\stackrel{4.3.4}{\Leftrightarrow} \#\psi \in T_\alpha \stackrel{4.5.1}{\Leftrightarrow} \#\psi \in T_{\alpha-1} \stackrel{\text{IH}}{\Rightarrow} \#\psi \in Y \stackrel{(b)}{\Rightarrow} \#\neg\phi \in Y. \end{aligned}$$

Let  $\phi = \psi \wedge \theta$  and suppose  $\#\phi \in D_\alpha$ . By Lemma 4.2.5, there are three cases to be separately considered. First assume  $\#\psi, \#\theta \in D_{\alpha-1}$ . Then, we get  $\#\psi \in X$  and  $\#\theta \in X$  by the induction hypothesis, and thus  $\#\phi \in X$  by (a); hence, (7) holds for  $\phi$ . We next obtain (8) as follows:

$$\#\phi \in T_\alpha \stackrel{4.3.5}{\Leftrightarrow} \#\psi, \#\theta \in T_\alpha \stackrel{4.5.1}{\Leftrightarrow} \#\psi, \#\theta \in T_{\alpha-1} \stackrel{\text{IH}}{\Rightarrow} \#\psi, \#\theta \in Y \stackrel{(b)}{\Rightarrow} \#\phi \in Y.$$

For the remaining claim (9), we infer

$$\begin{aligned} \#\phi \notin T_\alpha &\stackrel{4.3.5}{\Leftrightarrow} \#\psi \notin T_\alpha \text{ or } \#\theta \notin T_\alpha \stackrel{4.5.1}{\Leftrightarrow} \#\psi \notin T_{\alpha-1} \text{ or } \#\theta \notin T_{\alpha-1} \\ &\stackrel{\text{IH}}{\Rightarrow} \#\neg\psi \in Y \text{ or } \#\neg\theta \in Y \stackrel{(b)}{\Rightarrow} \#\neg\phi \in Y. \end{aligned}$$

Second assume  $\#\neg\psi \in D_{\alpha-1} \cap T_{\alpha-1}$ . We have  $\#\neg\psi \in X \cap Y$  by the induction hypothesis, which implies  $\#\phi \in X$  and  $\#\neg\phi \in Y$  by (a) and (b); hence, (7) and (9) hold for  $\phi$ . Note that  $\#\phi \in T_\alpha$  can never be the case under the current assumption, since it would imply  $\#\neg\psi \notin T_\alpha$  by Lemmata 4.3.4 and 4.3.5 and thus  $\#\neg\psi \notin T_{\alpha-1}$  by Lemma 4.5.1; hence, (8) trivially holds for  $\phi$ . The case where  $\#\neg\theta \in D_{\alpha-1} \cap T_{\alpha-1}$  can be similarly treated.

The remaining case where  $\phi = \forall x\psi(x)$  can be similarly treated to the last case, and we omit the details. ⊖

REMARK 4.11. As we have remarked,  $D_\infty$  is  $\Gamma_{\mathcal{D}[T_\infty]}$ -closed. We can also show that  $T_\infty$  is  $\Gamma_{\mathcal{T}[D_\infty]}$ -closed; the proof is routine. Hence, it follows from Lemma 4.10 that  $D_\infty$  and  $T_\infty$  are simultaneously inductively defined by the following monotone operators from  $\mathcal{P}(\mathbb{N}) \times \mathcal{P}(\mathbb{N})$  to  $\mathcal{P}(\mathbb{N})$ :

$$\Gamma_{\mathcal{D}}(X, Y) = \Gamma_{\mathcal{D}[Y]}(X) \quad \text{and} \quad \Gamma_{\mathcal{T}}(X, Y) = \Gamma_{\mathcal{T}[X]}(Y).$$

LEMMA 4.12. *Let  $X, Y \subset \omega$ . Suppose the same conditions (a) and (b) as in Lemma 4.10, as well as the following additional condition:*

- (c) for all  $\mathcal{L}$ -sentences  $\phi$ , if  $\#\phi \in X$ , then either  $\#\phi \notin Y$  or  $\#\neg\phi \notin Y$ .

*Then the following holds.*

1.  $D_\infty \cap Y = T_\infty$ .
2.  $D_\infty$  is the least  $\Gamma_{\mathcal{D}[Y]}$ -closed set.

PROOF. We have  $T_\infty \subset D_\infty$  by (6) and  $T_\infty \subset Y$  by Lemma 4.10; hence,  $T_\infty \subset D_\infty \cap Y$ . For the converse, take any  $\#\phi \notin T_\infty$ . If  $\#\phi \in D_\infty$ , then we have  $\#\neg\phi \in T_\infty \subset Y$  by Lemma 4.7.2 and thus  $\#\phi \notin Y$  by (c).

For claim 2, we first show that  $D_\infty$  is  $\Gamma_{\mathcal{D}[Y]}$ -closed. We observe that  $\mathsf{T}$  only appears in the clauses  $\mathcal{D}_5$  and  $\mathcal{D}_6$  of  $\mathcal{D}$  in the form  $(\mathsf{D}\neg x \wedge \mathsf{T}\neg x)$ . Hence, since  $D_\infty \cap Y = D_\infty \cap T_\infty$  by claim 1, we have  $\Gamma_{\mathcal{D}[Y]}(D_\infty) = \Gamma_{\mathcal{D}[T_\infty]}(D_\infty) \subset D_\infty$ . Now, take any  $\Gamma_{\mathcal{D}[Y]}$ -closed set  $Z$ . Since the intersection of  $\Gamma_{\mathcal{D}[Y]}$ -closed sets is also  $\Gamma_{\mathcal{D}[Y]}$ -closed (by the monotonicity of  $\Gamma_{\mathcal{D}[Y]}$ ),  $D_\infty \cap Z$  is  $\Gamma_{\mathcal{D}[Y]}$ -closed. Also, since  $\mathsf{D}$  occurs in  $\mathcal{T}$  only positively, it also follows that  $Y$  is  $\Gamma_{\mathcal{T}[D_\infty \cap Z]}$ -closed; for,  $D_\infty \subset X$  by Lemma 4.10 and thus  $\Gamma_{\mathcal{T}[D_\infty \cap Z]}(Y) \subset \Gamma_{\mathcal{T}[D_\infty]}(Y) \subset \Gamma_{\mathcal{T}[X]}(Y) \subset Y$ . Hence,  $D_\infty \cap Z$  and  $Y$  satisfy the conditions (a) and (b) of Lemma 4.10, and thus we obtain  $D_\infty \subset Z$ . ⊥

The next theorem immediately follows from Lemmata 4.1, 4.9, and 4.12.

THEOREM 4.13. *Let  $X, Y \subset \mathbb{N}$ . If  $(\mathbb{N}, X, Y) \models \text{CD}$ , then the following hold:*

1.  $D_\infty \subset X$ .
2.  $D_\infty \cap Y = T_\infty$ .
3.  $D_\infty$  is the least  $\Gamma_{\mathcal{D}[Y]}$ -fixed point.

This theorem says that every sentence in  $D_\infty$  is determinate in any  $\omega$ -model of CD, the truth and the falsity are invariable on  $D_\infty$  across all  $\omega$ -models of CD, and  $D_\infty$  is the least fixed point of  $\Gamma_{\mathcal{D}[Y]}$  (and thus is the least set  $Z$  with  $(\mathbb{N}, Z, Y) \models \text{D1} - \text{D6}$ ) whenever  $(\mathbb{N}, X, Y) \models \text{CD}$  for some  $X$ . This mathematical fact, as well as the philosophical story that motivated CD, suggests an axiom expressing such a minimality property of  $\mathsf{D}$ , which yields a new theory  $\text{CD}_\mu$ . This axiom and the resulting system will be briefly explained in Section 9, but their full analysis is left for the Part II.

**§5. A lower bound of the strength of CD.** In this section, we will show that the system  $\text{RA}_{<\varepsilon_0}$  of ramified analysis up to  $\varepsilon_0$  and the system  $\text{RT}_{<\varepsilon_0}$  of  $\varepsilon_0$ -iterated truth (see [21] for its definition) are arithmetically conservative over CD, that is, every arithmetical theorem of the former system is provable in the latter system, and also that  $\text{RA}_{<\varepsilon_0}$  and  $\text{RT}_{<\varepsilon_0}$  are arithmetically conservative over  $\text{CD}^+$ . The reductions will be achieved by constructing a relative interpretation of intermediate systems KF and  $\text{CT}[\text{KF}]$  in CD and  $\text{CD}^+$ , respectively.

**5.1. The systems KF and CT.** For the sake of the reader’s convenience, we repeat the definition of the system KF of Feferman [8] (in our notation).

DEFINITION 5.1. The  $\mathcal{L}_T$ -system KF is defined as PA with full induction in  $\mathcal{L}_T$  augmented with the following axioms:

- K1  $\forall s \forall t (T(s \doteq t) \leftrightarrow s^\circ = t^\circ) \wedge \forall s \forall t (F(s \doteq t) \leftrightarrow s^\circ \neq t^\circ)$ ,
- K2  $\forall s (TTs \leftrightarrow Ts^\circ) \wedge \forall s (FTs \leftrightarrow Fs^\circ)$ ,
- K3  $\forall x (\text{Sent}_T(x) \rightarrow (T(\neg x) \leftrightarrow Tx))$ ,
- K4  $\forall x \forall y (\text{Sent}_T(x \wedge y) \rightarrow (T(x \wedge y) \leftrightarrow (Tx \wedge Ty)))$ ,
- K5  $\forall x \forall y (\text{Sent}_T(x \wedge y) \rightarrow (F(x \wedge y) \leftrightarrow (Fx \vee Fy)))$ ,
- K6  $\forall v \forall x (\text{Sent}_T(\forall v x) \rightarrow (T(\forall v x) \leftrightarrow \forall s Tx(s/v)))$ ,
- K7  $\forall v \forall x (\text{Sent}_T(\forall v x) \rightarrow (F(\forall v x) \leftrightarrow \exists s Fx(s/v)))$ ,
- R3  $\forall x \forall v \forall s \forall t ((\text{Sent}_T(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (Tx(s/v) \leftrightarrow Tx(t/v)))$ .

Note that, while  $Fx$  and  $\neg Tx$  are equivalent in CD, they are not equivalent in KF, and we have to distinguish them when working in KF. It is easily shown that R3 is redundant and provable from the other axioms (see [5, Lemma 3.1]). We will occasionally consider two extra axioms Cons and Comp, which are defined as follows:

- Cons:  $\forall x (\text{Sent}_T(x) \rightarrow \neg(Tx \wedge Fx))$ ,
- Comp:  $\forall x (\text{Sent}_T(x) \rightarrow (Tx \vee Fx))$ .

For each  $\varphi \in \mathcal{L}_T$ , let  $\varphi^c$  denote Cantini’s ‘dual’ translation of  $\varphi$ : namely,  $T^c a := \neg Fa$  for each  $\mathcal{L}_0$ -term  $a$ , and  $c$  preserves the  $\mathcal{L}_0$ -vocabulary (as well as the logical connectives and quantifiers); see [5, Section 4].  $c$  is a relative interpretation of KF + Cons in KF + Comp and *vice versa*. Cantini [5] showed that KF + Cons and KF + Comp are both proof-theoretically equivalent to KF.

LEMMA 5.2. KF proves the following.

1.  $\forall x \forall y (\text{Sent}_T(x \vee y) \rightarrow ((T(x \vee y) \leftrightarrow (Tx \vee Ty)) \wedge (F(x \vee y) \leftrightarrow (Fx \wedge Fy))))$ .
2.  $\forall x \forall y (\text{Sent}_T(x \rightarrow y) \rightarrow ((T(x \rightarrow y) \leftrightarrow (\neg Fx \rightarrow Ty)) \wedge (F(x \rightarrow y) \leftrightarrow (Tx \wedge Fy))))$ .
3.  $\forall v \forall x (\forall v x \in \text{Sent}_T \rightarrow ((T(\exists v x) \leftrightarrow \exists s Tx(s/v)) \wedge F(\exists v x) \leftrightarrow \forall s Fx(s/v)))$ .

DEFINITION 5.3. Let  $\mathcal{L}_1$  be a first-order language that extends  $\mathcal{L}_0$  only with new predicate symbols. We set a new language  $\mathcal{L}_1^+$  to be  $\mathcal{L}_1 \cup \{\mathbb{T}\}$  where  $\mathbb{T}$  is a fresh unary predicate symbol: the new predicate  $\mathbb{T}$  will be interpreted as the Tarskian typed compositional truth for the language  $\mathcal{L}_1$ . Let  $\text{Sent}_1$  be a representation of the set of (codes of)  $\mathcal{L}_1$ -sentences. Given any  $\mathcal{L}_1$ -system S, the  $\mathcal{L}_1^+$ -system CT[S] is defined as S with all axiom schemata of S (possibly including other schemata than induction) extended for  $\mathcal{L}_1^+$  together with the following axioms expressing Tarski’s ‘inductive clauses’ of truth.

- T1  $\forall \vec{s} (\mathbb{T}R\vec{s} \leftrightarrow R\vec{s}^\circ)$ , for all  $\mathcal{L}_1$ -atomic formulae  $R\vec{x}$ .
- T2  $\forall x (\text{Sent}_1(x) \rightarrow (\mathbb{T}(\neg x) \leftrightarrow \neg \mathbb{T}x))$ .
- T3  $\forall x \forall y (\text{Sent}_1(x \wedge y) \rightarrow (\mathbb{T}(x \wedge y) \leftrightarrow \mathbb{T}x \wedge \mathbb{T}y))$ .
- T4  $\forall v \forall x (\forall v x \in \text{Sent}_1 \rightarrow (\mathbb{T}(\forall v x) \leftrightarrow \forall s \mathbb{T}x(s/v)))$ .
- R4  $\forall x \forall v \forall s \forall t ((\text{Sent}_1(\forall v x) \wedge s^\circ = t^\circ) \rightarrow (\mathbb{T}x(s/v) \leftrightarrow \mathbb{T}x(t/v)))$ .

Here we abbreviate a sequence of variables (or terms),  $x_0, \dots, x_k$  (or  $a_1, \dots, a_k$ ), by  $\vec{x}$  (or  $\vec{a}$ ) for saving space. If  $\mathcal{L}_1 \supset \mathcal{L}_T$ , then T1 contains  $\forall s (\mathbb{T}Ts \leftrightarrow Ts^\circ)$ . It can be

shown in a parallel manner to the case of R3 in KF that R4 is derivable from the other axioms.

LEMMA 5.4. *Let S be as above. CT[S] proves the following theorems:*

1.  $\forall x \forall y (\text{Sent}_1(x \vee y) \rightarrow (\mathbb{T}(x \vee y) \leftrightarrow (\mathbb{T}x \vee \mathbb{T}y)))$ .
2.  $\forall x \forall y (\text{Sent}_1(x \rightarrow y) \rightarrow (\mathbb{T}(x \rightarrow y) \leftrightarrow (\mathbb{T}x \rightarrow \mathbb{T}y)))$ .
3.  $\forall v \forall x (\text{Sent}_1(\forall vx) \rightarrow (\mathbb{T}(\exists vx) \leftrightarrow \exists s \mathbb{T}x(s/v)))$ .

**5.2. Reduction of KF to CD<sub>0</sub>.** The translation  $\natural$  of  $\mathcal{L}_T$  into  $\mathcal{L}$  replaces every atomic formula  $\mathbb{T}^{\natural}a$  (for a term  $a$ ) with  $Ta \wedge Da$ ;  $\natural$  does not change anything else.

In what follows, we will occasionally write  $s \in \text{CTerm}$  instead of  $\text{CTerm}(s)$ ,  $x \in \text{Sent}$  instead of  $\text{Sent}(x)$ , and similarly for other formulae.

LEMMA 5.5.  *$\natural$  is a relative interpretation of KF in CD<sub>0</sub>.*

PROOF. We will only exhibit the proofs that  $\natural$  preserves the axioms K2, K5, and K7. For K2, take any  $s \in \text{CTerm}$ . Then we have

$$\begin{aligned} \mathbb{T}_{\neg} \mathbb{T}s \wedge \mathbb{D}_{\neg} \mathbb{T}s &\stackrel{T4 \ \& \ D4}{\Leftrightarrow} \neg \mathbb{T}\mathbb{T}s \wedge \mathbb{D}\mathbb{T}s &\stackrel{D2}{\Leftrightarrow} &\neg \mathbb{T}\mathbb{T}s \wedge \mathbb{D}s^{\circ} &\stackrel{T3}{\Leftrightarrow} &\neg \mathbb{T}s^{\circ} \wedge \mathbb{D}s^{\circ} \\ &&&&&&&&&&&&&\stackrel{T4 \ \& \ D4}{\Leftrightarrow} &\mathbb{T}_{\neg} s^{\circ} \wedge \mathbb{D}_{\neg} s^{\circ}. \end{aligned}$$

The other conjunct can be shown similarly. For K5, take any  $x, y \in \text{Sent}_T$ ; then we have

$$\begin{aligned} \mathbb{T}_{\neg}(x \wedge y) \wedge \mathbb{D}_{\neg}(x \wedge y) &\stackrel{T4 \ \& \ T5}{\Leftrightarrow} (Fx \vee Fy) \wedge \mathbb{D}_{\neg}(x \wedge y) \\ &\stackrel{D4 \ \& \ D5}{\Leftrightarrow} (Fx \vee Fy) \wedge ((\mathbb{D}x \wedge \mathbb{D}y) \vee (\mathbb{D}x \wedge Fx) \vee (\mathbb{D}y \wedge Fy)) \\ &\stackrel{\text{by logic}}{\Leftrightarrow} (\mathbb{D}x \wedge Fx) \vee (\mathbb{D}y \wedge Fy) \\ &\stackrel{D4}{\Leftrightarrow} (\mathbb{D}_{\neg}x \wedge \mathbb{T}_{\neg}x) \vee (\mathbb{D}_{\neg}y \wedge \mathbb{T}_{\neg}y). \end{aligned}$$

For K7, take any  $v$  and  $x$  such that  $\forall vx \in \text{Sent}_T$ ; then we have

$$\begin{aligned} \mathbb{T}_{\neg}(\forall vx) \wedge \mathbb{D}_{\neg}(\forall vx) &\stackrel{T4 \ \& \ T6}{\Leftrightarrow} \exists s Fx(s/v) \wedge \mathbb{D}_{\neg}(\forall vx) \\ &\stackrel{D4 \ \& \ D6}{\Leftrightarrow} \exists s Fx(s/v) \wedge (\forall s \mathbb{D}x(s/v) \vee \exists s (\mathbb{D}x(s/v) \wedge Fx(s/v))) \\ &\stackrel{\text{by logic}}{\Leftrightarrow} \exists s (\mathbb{D}x(s/v) \wedge Fx(s/v)) \\ &\stackrel{D4}{\Leftrightarrow} \exists s (\mathbb{D}_{\neg}x(s/v) \wedge \mathbb{T}_{\neg}x(s/v)). \end{aligned}$$

The remaining cases are similarly and even more easily shown. ⊢

Because  $\natural$  does not affect arithmetical sentences, the next corollary follows.

COROLLARY 5.6. *KF is arithmetically conservative over CD<sub>0</sub> and thus CD.*

Finally, since KF and  $\text{RA}_{<\varepsilon_0}$  are known to have the same arithmetical theorems by results due to Cantini [5] and Feferman [8],  $\text{RA}_{<\varepsilon_0}$  is arithmetically conservative over CD.

**5.3. Reduction of CT[[KF]] to CD<sub>2</sub><sup>+</sup>.** Let  $k$  be an  $\mathcal{L}_0$ -definable function that represents the arithmetization of the translation  $\mathfrak{h}$  of  $\mathcal{L}_T$  in  $\mathcal{L}$ : that is, for all  $x \in \omega$ ,

$$kx := \begin{cases} x, & \text{if } x \in \text{AtFml}_0, \\ (\mathbb{T}y) \wedge (\mathbb{D}y), & \text{if } x = \mathbb{T}y \text{ and } y \in \text{Term}, \\ \neg ky, & \text{if } x = \neg y \in \text{Fml}_T, \\ (ky) \wedge (kz), & \text{if } x = y \wedge z \in \text{Fml}_T, \\ \forall v(ky), & \text{if } x = \forall v y \in \text{Fml}_T, \\ 0, & \text{otherwise.} \end{cases} \tag{10}$$

With the standard coding,  $k$  is primitive recursive and thus definable in PA, and PA proves  $\forall x(\text{Fml}_T(x) \rightarrow \text{Fml}(kx))$ . Furthermore, with the notation of [5], PA also proves  $\forall z(\text{Free}(z, x) \leftrightarrow \text{Free}(z, kx))$ , where  $\text{Free}(a, b)$  says that  $a$  is a code of a variable that is free in the formula coded by  $b$ ; in particular, we have  $\forall x(\text{Sent}_T(x) \leftrightarrow \text{Sent}(kx))$ .

We thereby extend  $\mathfrak{h}$  to a translation of  $\mathcal{L}_T^+$  to  $\mathcal{L}$  by stipulating the following:

$$\mathbb{T}^{\mathfrak{h}}x := \text{Tk}x.$$

We will use the same symbol  $\mathfrak{h}$  for the two translations for simplicity.

**LEMMA 5.7.**  $\text{PA} \vdash (\forall x \in \text{Fml}_T) \forall s (\forall v \in \text{Var})(k(x(s/v)) = (kx)(s/v))$ .

**PROOF.** By a routine induction on the complexity of  $x$ . ⊢

**LEMMA 5.8.** *The following claims are provable in CD<sub>2</sub>.*

1.  $\forall x (\text{Sent}_T(x) \rightarrow (\text{Tk}(\neg x) \leftrightarrow \neg \text{Tk}x))$ ; by T4 and (10).
2.  $\forall x \forall y (\text{Sent}_T(x \wedge y) \rightarrow (\text{Tk}(x \wedge y) \leftrightarrow (\text{Tk}x \wedge \text{Tk}y)))$ ; by T5 and (10).
3.  $\forall x \forall v (\text{Sent}_T(\forall v x) \rightarrow (\text{Tk}(\forall v x) \leftrightarrow \forall s \text{Tk}(x(s/v))))$ ; by T6, (10), and Lemma 5.7.

**THEOREM 5.9.**  $\mathfrak{h}$  is a relative interpretation of CT[[KF]] in CD<sub>2</sub><sup>+</sup>.

**PROOF.** It immediately follows from the definition of  $\mathfrak{h}$  and the last lemma that the  $\mathfrak{h}$ -translations of the axioms  $\mathbb{T}1$  for all  $\mathcal{L}_0$ -atomics and  $\mathbb{T}2$ – $\mathbb{T}4$  are provable in CD<sub>2</sub>.  $\text{R4}^{\mathfrak{h}}$  readily follows from R1 and Lemma 5.7 (though this step is actually redundant because R4 is derivable from the other axioms of CT[[KF]]). For the remaining case, i.e., the axiom  $\mathbb{T}1$  for an atomic of the form  $\text{T}x$ , take any  $s \in \text{CTerm}$ . Then we have

$$\mathbb{T}^{\mathfrak{h}}(\mathbb{T}s) \Leftrightarrow \text{T}(\mathbb{T}s \wedge \mathbb{D}s) \stackrel{\text{T5}}{\Leftrightarrow} \mathbb{T}\mathbb{T}s \wedge \mathbb{T}\mathbb{D}s \stackrel{\text{T2}^+}{\Leftrightarrow} \mathbb{T}\mathbb{T}s \wedge \mathbb{D}s^\circ \stackrel{\text{T3}}{\Leftrightarrow} \mathbb{T}s^\circ \wedge \mathbb{D}s^\circ \Leftrightarrow \mathbb{T}^{\mathfrak{h}}s^\circ;$$

note that the use of  $\text{T2}^+$  in the third equivalence is crucial here. ⊢

**COROLLARY 5.10.** CT[[KF]] is arithmetically conservative over CD<sub>2</sub><sup>+</sup> and thus CD<sub>2</sub><sup>+</sup>

Now, the following fact is essentially due to Cantini [5].

**FACT 5.11.** CT[[KF + Cons]] (= CT[[KF]] + Cons) is arithmetically equivalent to  $\text{RA}_{<\varepsilon_0}$ , that is, the two systems have the same arithmetical theorems.

Hence,  $\text{RT}_{<\varepsilon_0}$  is arithmetically conservative over CD<sup>+</sup>.

**§6. An upper bound of the strength of CD.** In this section, we will show that CD is proof-theoretically reducible to KF. This will be achieved via partial cut-elimination for a semi-formal system for CD.

**6.1. Semi-formal system  $CD^\infty$ .** In this subsection, we will introduce a Tait-style semi-formal system  $CD^\infty$ .

Throughout this Section 6, capital Greek letters  $\Gamma$  and  $\Delta$  will be used as variables ranging over finite sets of  $\mathcal{L}$ -sentences, and we will only consider  $\mathcal{L}$ -sentences in negation-normal form (or in what is sometimes called ‘Tait form’), which are constructed from closed  $\mathcal{L}_0$ -literals by conjunction, disjunction, and universal and existential quantifiers; we will abuse the notation and denote the negation operation on negation-normal forms, which is often written as ‘ $\sim$ ’ in the literature, also by  $\neg$ , e.g.,  $\neg(Ta \wedge Da) = \neg Ta \vee \neg Da$ . Following the convention, for a finite set  $\Gamma$  of  $\mathcal{L}$ -sentences and an  $\mathcal{L}$ -sentence  $A$ , we will write  $\Gamma, A$  for  $\Gamma \cup \{A\}$ .

The derivation relation  $CD^\infty \frac{\delta}{p} \Gamma$  in the semi-formal system  $CD^\infty$  for ordinals  $\delta$  and  $p < \omega$  and a set  $\Gamma$  of  $\mathcal{L}$ -sentences, which intuitively means that at least one  $\phi \in \Gamma$  is verified by a derivation of length  $\delta$  with cut-rank  $p$ , is defined as follows.

**Axioms.** Let  $a, b$ , and  $c$  be any closed  $\mathcal{L}_0$ -terms,  $\Gamma$  any set of  $\mathcal{L}$ -sentences,  $\delta$  any ordinal, and  $p < \omega$  any finite ordinal.

- (Ax1) If  $A$  is a true closed  $\mathcal{L}_0$ -literal, then  $CD^\infty \frac{\delta}{p} \Gamma, A$ .
- (Ax2) If  $a^{\mathbb{N}} \approx^{\mathbb{N}} b^{\mathbb{N}}$  (see p.13 for the definition of  $\approx$ ), then  $CD^\infty \frac{\delta}{p} \Gamma, Da, \neg Db$ .
- (Ax3) If  $a^{\mathbb{N}} \approx^{\mathbb{N}} b^{\mathbb{N}}$ , then  $CD^\infty \frac{\delta}{p} \Gamma, Ta, \neg Tb$ .
- (D1) If  $b^{\mathbb{N}}, c^{\mathbb{N}} \in \text{CTerm}^{\mathbb{N}}$  and  $a^{\mathbb{N}} = (b \dot{=} c)^{\mathbb{N}}$ , then  $CD^\infty \frac{\delta}{p} \Gamma, Da$ .

**Logical Rules.** Let  $A$  and  $B$  be  $\mathcal{L}$ -sentences,  $C(x)$  an  $\mathcal{L}$ -formula with at most one free variable,  $\delta$  an ordinal, and  $p < \omega$ .

- ( $\vee 1$ ) If  $CD^\infty \frac{\delta'}{p} \Gamma, A \vee B, A$  for some ordinal  $\delta' < \delta$ , then  $CD^\infty \frac{\delta}{p} \Gamma, A \vee B$ .
- ( $\vee 2$ ) If  $CD^\infty \frac{\delta'}{p} \Gamma, A \vee B, B$  for some ordinal  $\delta' < \delta$ , then  $CD^\infty \frac{\delta}{p} \Gamma, A \vee B$ .
- ( $\wedge$ ) If  $CD^\infty \frac{\delta_0}{p} \Gamma, A \wedge B, A$  and  $CD^\infty \frac{\delta_1}{p} \Gamma, A \wedge B, B$  for some ordinals  $\delta_0, \delta_1 < \delta$ , then  $CD^\infty \frac{\delta}{p} \Gamma, A \wedge B$ .
- ( $\exists$ ) If  $CD^\infty \frac{\delta'}{p} \Gamma, \exists x C(x), C(a)$  for some closed  $\mathcal{L}_0$ -term  $a$  and ordinal  $\delta' < \delta$ , then  $CD^\infty \frac{\delta}{p} \Gamma, \exists x C(x)$ .
- ( $\forall$ ) If there exists some ordinal  $\delta_a < \delta$  for each closed  $\mathcal{L}_0$ -term  $a$  such that  $CD^\infty \frac{\delta_a}{p} \Gamma, \forall x C(x), C(a)$ , then  $CD^\infty \frac{\delta}{p} \Gamma, \forall x C(x)$ .
- (cut) If  $CD^\infty \frac{\delta_0}{p} \Gamma, A$  and  $CD^\infty \frac{\delta_1}{p} \Gamma, \neg A$  for some ordinals  $\delta_0, \delta_1 < \delta$  and  $\mathcal{L}$ -sentence  $A$  with complexity  $< p$ , then  $CD^\infty \frac{\delta}{p} \Gamma$ .

**Rules for positive occurrences of D.** For ordinals  $\delta$  and  $p < \omega$  and a closed term  $a$ ,  $CD^\infty \frac{\delta}{p} \Gamma, Da$  holds, if one of the following holds for some ordinal  $\delta' < \delta^8$ :

<sup>8</sup>Note that  $D5^\pm_\infty$  and  $D6^\pm_\infty$  are the rule versions of the variants  $D5'$  and  $D6'$  of  $D5$  and  $D6$  (see fn. 4.1 for their definitions). Recall that they are equivalent to the original  $D5$  and  $D6$  due to the axiom  $D4$ .

- (D2<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\top b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Da}, \text{Db}^{\circ}$ .
- (D3<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\text{D}b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Da}, \text{Db}^{\circ}$ .
- (D4<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\neg b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Da}, \text{Db}$ .
- (D5<sup>pos</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \wedge c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Da}, (\text{Db} \wedge \text{Dc}) \vee (\text{D}\neg b \wedge \text{Fb}) \vee (\text{D}\neg c \wedge \text{Fc})$ .
- (D6<sup>pos</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (\forall bc)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Da}, \forall s \text{Dc}(s/b) \vee \exists s(\text{D}\neg c(s/b) \wedge \text{Fc}(s/b))$ .

**Rules for negative occurrences of D.** For ordinals  $\delta$  and  $p < \omega$  and a closed term  $a$ ,  $\text{CD}^{\infty} \mid_{\frac{\delta}{p}} \Gamma, \neg \text{Da}$  holds, if one of the following holds for some ordinal  $\delta' < \delta$ :

- (D2<sup>neg</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\top b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \neg \text{Da}, \neg \text{Db}^{\circ}$ .
- (D3<sup>neg</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\text{D}b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \neg \text{Da}, \neg \text{Db}^{\circ}$ .
- (D4<sup>neg</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\neg b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \neg \text{Da}, \neg \text{Db}$ .
- (D5<sup>neg</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \wedge c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \neg \text{Da}, (\neg \text{Db} \vee \neg \text{Dc}) \wedge (\neg \text{D}\neg b \vee \neg \text{Fb}) \wedge (\neg \text{D}\neg c \vee \neg \text{Fc})$ .
- (D6<sup>neg</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (\forall bc)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \neg \text{Da}, \exists s \neg \text{Dc}(s/b) \wedge \forall s(\neg \text{D}\neg c(s/b) \vee \neg \text{Fc}(s/b))$ .

**Rules for positive occurrences of T.** For ordinals  $\delta$  and  $p < \omega$  and a closed term  $a$ ,  $\text{CD}^{\infty} \mid_{\frac{\delta}{p}} \Gamma, \text{Ta}$  holds, if one of the following holds for some ordinal  $\delta' < \delta$ :

- (T1<sup>pos</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \equiv c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, b^{\circ} = c^{\circ}$ .
- (T2<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\text{D}b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, \text{Db}^{\circ}$ .
- (T3<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\top b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, \text{Tb}^{\circ} \wedge \text{Db}^{\circ}$ .
- (T4<sup>pos</sup><sub>∞</sub>) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\neg b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, \neg \text{Tb}$ .
- (T5<sup>pos</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \wedge c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, \text{Tb} \wedge \text{Tc}$ .
- (T6<sup>pos</sup><sub>∞</sub>) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (\forall bc)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $\text{CD}^{\infty} \mid_{\frac{\delta'}{p}} \Gamma, \text{Ta}, \forall s \text{Tc}(s/b)$ .

**Rules for a negative occurrence of T.** For ordinals  $\delta$  and  $p < \omega$  and a closed term  $a$ ,  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, \neg Ta$  holds, if one of the following holds for some ordinal  $\delta' < \delta$ :

- (T1<sup>neg</sup> <sub>$\infty$</sub> ) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \equiv c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta'}{p}} \Gamma, \neg Ta, b^\circ \neq c^\circ$ .
- (T3<sup>neg</sup> <sub>$\infty$</sub> ) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\top b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta'}{p}} \Gamma, \neg Ta, \neg Tb^\circ \wedge Db^\circ$ .
- (T4<sup>neg</sup> <sub>$\infty$</sub> ) There is a closed term  $b$  such that  $a^{\mathbb{N}} = (\neg b)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta'}{p}} \Gamma, \neg Ta, Tb$ .
- (T5<sup>neg</sup> <sub>$\infty$</sub> ) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (b \wedge c)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta'}{p}} \Gamma, \neg Ta, \neg Tb \vee \neg Tc$ .
- (T6<sup>neg</sup> <sub>$\infty$</sub> ) There are some closed terms  $b$  and  $c$  such that  $a^{\mathbb{N}} = (\forall bc)^{\mathbb{N}} \in \text{Sent}^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta'}{p}} \Gamma, Ta, \exists s \neg Tc(s/b)$ .

Note that there is no rule for a negative occurrence of T corresponding to (T2<sup>pos</sup> <sub>$\infty$</sub> ).

LEMMA 6.1. *The following basic properties of  $CD^\infty$  can be shown in the standard manner, and we omit their proofs.*

1. (Structural Lemma)  
For all  $\gamma \leq \delta, q \leq p$ , and  $\Delta \subset \Gamma$ , if  $CD^\infty \mid_{\frac{\gamma}{q}} \Delta$ , then  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma$ .
2. (Numerical Equivalence Lemma)  
For all closed terms  $a$  and  $b$ , if  $a^{\mathbb{N}} = b^{\mathbb{N}}$  and  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A(a)$ , then  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A(b)$ .
3. (Tautology Lemma)  
 $CD^\infty \mid_{\frac{2 \cdot cp(A)}{0}} A, \neg A$ ; recall that  $cp(A)$  denotes the complexity of  $A$  (see page 12).
4. ( $\wedge$ -Inversion)  
If  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A \wedge B$ , then  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A$ , and  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, B$ .
5. ( $\vee$ -Exportation)  
If  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A \vee B$ , then  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A, B$ .
6. ( $\forall$ -Inversion)  
If  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, \forall x A(x)$ , then  $CD^\infty \mid_{\frac{\delta}{p}} \Gamma, A(a)$  for all closed terms  $a$ .

The next lemma is nearly obvious from the design of  $CD^\infty$ , and the proof is routine.

LEMMA 6.2. *For every  $\mathcal{L}$ -sentence  $\psi$ , if  $CD \vdash \psi$ , then  $CD^\infty \mid_{\frac{\omega+k}{n}} \psi$  for some  $k, n < \omega$ ; actually, we can primitive recursively calculate  $k$  and  $n$  from a given derivation of  $\psi$ .*

The following three lemmata can also be shown standardly.

LEMMA 6.3 (Reduction Lemma). *If  $CD^\infty \mid_{\frac{\alpha}{1+n+1}} \Gamma$ , then  $CD^\infty \mid_{\frac{2\alpha}{1+n}} \Gamma$ .*

LEMMA 6.4 (Partial Cut-elimination Lemma). *If  $CD^\infty \mid_{\frac{\alpha}{1+n}} \Gamma$ , then  $CD^\infty \mid_{\frac{\omega_n(\alpha)}{1}} \Gamma$ , where  $\omega_n(\alpha)$  is recursively defined as follows:  $\omega_0(\alpha) = \alpha$  and  $\omega_{k+1}(\alpha) = \omega^{\omega_k(\alpha)}$ .*

LEMMA 6.5. *For every  $\mathcal{L}$ -sentence  $\psi$ , if  $CD \vdash \psi$ , then  $CD^\infty \mid_{\top} \psi$  for some  $\delta < \varepsilon_0$ .*

Note that  $CD^\infty \uparrow_1^\delta \psi$  means that there is a derivation of  $\psi$  in  $CD^\infty$  in which (cut) is only applied to closed  $\mathcal{L}$ -literals; such a derivation corresponds to what Cantini [5] calls a *quasi-normal* derivation in his semi-formal system  $TKF^\infty$ .

**6.2. Reduction of CD to KF.** Let  $\sigma$ ,  $\tau$ , and  $\xi$  be ordinals. Recall that we only consider  $\mathcal{L}$ -sentences in negation-normal form throughout the Section 6. Given an  $\mathcal{L}$ -sentence  $\phi$ , we write  $\models \phi[\sigma, \tau, \xi]$  when  $\phi$  is true under the following interpretation:

- the quantifiers range over  $\omega$ ;
- the  $\mathcal{L}_0$ -vocabulary receives the standard interpretations;
- each negative occurrence of D in  $\phi$  is interpreted by  $D_{\sigma+1}$ ;
- each positive occurrence of D in  $\phi$  is interpreted by  $D_{\tau+1}$ ;
- each (either negative or positive) occurrence of T in  $\phi$  is interpreted by  $T_{\xi+1}$ .

We extend this notation to finite sets  $\Gamma$  of  $\mathcal{L}$ -sentences and write  $\models \Gamma[\sigma, \tau, \xi]$  when  $\models \phi[\sigma, \tau, \xi]$  for some  $\phi \in \Gamma$ . By Lemma 4.5.2, the sequence  $\langle D_\xi \mid \xi \in On \rangle$  is a monotonically increasing, where  $On$  denotes the class of ordinals, and thus the next lemma can be shown in the standard manner.

LEMMA 6.6 (Persistence Lemma). *Let  $\phi$  be an  $\mathcal{L}$ -sentences and  $\Gamma$  a finite set of  $\mathcal{L}$ -sentences. For every ordinals  $\sigma_1 \leq \sigma_0$ ,  $\tau_0 \leq \tau_1$ , and  $\xi$ , the following hold.*

1. *If  $\models \phi[\sigma_0, \tau_0, \xi]$ , then  $\models \phi[\sigma_1, \tau_1, \xi]$ .*
2. *If  $\models \Gamma[\sigma_0, \tau_0, \xi]$ , then  $\models \Gamma[\sigma_1, \tau_1, \xi]$ .*

The next is the main lemma of the present section.

LEMMA 6.7. *For all ordinals  $\delta$ , if  $CD^\infty \uparrow_1^\delta \Gamma$ , then for all ordinals  $\sigma$  and  $\xi$ , if  $\sigma + 2^\delta \leq \xi$ , then  $\models \Gamma[\sigma, \sigma + 2^\delta, \xi]$ .*

PROOF. The claim is shown by induction on derivation. Throughout this proof,  $a$ ,  $b$ ,  $c$ , and  $d$  always denote closed  $\mathcal{L}_0$ -terms.

Suppose that the last inference is made by either of the following axioms:

$$\frac{}{CD^\infty \uparrow_1^\delta \Gamma, A} \text{ (Ax1)} \qquad \frac{}{CD^\infty \uparrow_1^\delta \Gamma, Ta, \neg Tb} \text{ (Ax3)} \qquad \frac{}{CD^\infty \uparrow_1^\delta \Gamma, Dc} \text{ (D1)},$$

where  $A$  is a true closed  $\mathcal{L}_0$ -literal,  $a^{\mathbb{N}} \approx^{\mathbb{N}} b^{\mathbb{N}}$ , and  $c^{\mathbb{N}}$  is the Gödel number of a closed  $\mathcal{L}_0$ -atomic. The claim for the first case is obvious; the claim for the second case follows from Proposition 4.4.2, and the claim for the third case follows from Proposition 4.2.1. Next, suppose the last inference is made by the axiom (Ax2):

$$\frac{}{CD^\infty \uparrow_1^\delta \Gamma, Da, \neg Db} \text{ (Ax2)},$$

where  $a^{\mathbb{N}} \approx^{\mathbb{N}} b^{\mathbb{N}}$ . Take any ordinal  $\sigma$  and  $\xi$ . Since  $D_{\sigma+1} \subset D_{\sigma+2^\delta+1}$  by Lemma 4.5.1,  $a^{\mathbb{N}} \notin D_{\sigma+2^\delta+1}$  implies  $a^{\mathbb{N}} \notin D_{\sigma+1}$  and thus  $b^{\mathbb{N}} \notin D_{\sigma+1}$  by Proposition 4.4.1. Hence, either  $a^{\mathbb{N}} \in D_{\sigma+2^\delta+1}$  or  $b^{\mathbb{N}} \notin D_{\sigma+1}$  holds, and thus  $\models \{Da, \neg Db\}[\sigma, \sigma + 2^\delta, \xi]$ .

We move on to the induction step. The case where we use a logical rule in the last inference, except for cut, can be treated in the usual manner. For instance, suppose the last inference is made by  $(\forall)$ , namely, (our variant of) the  $\omega$ -rule:

$$\frac{CD^\infty \uparrow_1^\delta \Gamma, \forall x A(x), A(a), \text{ for each closed } \mathcal{L}_0\text{-term } a,}{CD^\infty \uparrow_1^\delta \Gamma, \forall x A(x)}$$

where  $\delta_a < \delta$  for each  $a$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . It follows from the induction hypothesis that for each closed term  $a$ , we have the following:

$$\models \Gamma \cup \{\forall x A(x)\} \cup \{A(a)\} [\sigma, \sigma + 2^{\delta_a}, \xi].$$

If  $\models \Gamma \cup \{\forall x A(x)\} [\sigma, \sigma + 2^{\delta_a}, \xi]$  for any  $a$ , then we have  $\models \Gamma \cup \{\forall x A(x)\} [\sigma, \sigma + 2^\delta, \xi]$  by Lemma 6.6. Otherwise, for all closed terms  $a$ , we have  $\models A(a) [\sigma, \sigma + 2^{\delta_a}, \xi]$  and thus  $\models A(a) [\sigma, \sigma + 2^\delta, \xi]$  by Lemma 6.6, which implies  $\models \forall x A(x) [\sigma, \sigma + 2^\delta, \xi]$ .

Next, the claim for the cases where the last inference is made by one of the rules for determinateness can be straightforwardly shown using Lemmata 4.2 and 4.5. For example, suppose the last inference is of the following form:

$$\frac{\text{CD}^\infty \frac{\gamma}{1} \Gamma, \neg Da, (\neg Db \vee \neg Dc) \wedge (\neg D\neg b \vee \neg Fb) \wedge (\neg D\neg c \vee \neg Fc)}{\text{CD}^\infty \frac{\delta}{1} \Gamma, \neg Da} \text{ (D5}^\infty\text{neg)},$$

where  $\gamma < \delta$  and  $a^\mathbb{N} = (b \wedge c)^\mathbb{N} \in \text{Sent}^\mathbb{N}$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . If we have  $\models \Gamma \cup \{\neg Da\} [\sigma, \sigma + 2^\gamma, \xi]$ , then we get  $\models \Gamma \cup \{\neg Da\} [\sigma, \sigma + 2^\delta, \xi]$  by Lemma 6.6. Otherwise, we have the following by the induction hypothesis:

$$\models (\neg Db \vee \neg Dc) \wedge (\neg D\neg b \vee \neg Fb) \wedge (\neg D\neg c \vee \neg Fc) [\sigma, \sigma + 2^\gamma, \xi].$$

Since  $D_{\sigma+1} \supset D_\sigma$  and  $D_\sigma \cap T_{\xi+1} = D_\sigma \cap T_\sigma$  by Lemma 4.5, this implies

$$b^\mathbb{N} \notin D_\sigma \text{ or } c^\mathbb{N} \notin D_\sigma, \text{ and } (\neg b)^\mathbb{N} \notin D_\sigma \cap T_\sigma, \text{ and } (\neg c)^\mathbb{N} \notin D_\sigma \cap T_\sigma.$$

Hence, by Lemma 4.2.5, we obtain  $a^\mathbb{N} \notin D_{\sigma+1}$  and thus  $\models \neg Da [\sigma, \sigma + 2^\delta, \xi]$ . Let us see one more example. Suppose the last inference is of the following form:

$$\frac{\text{CD}^\infty \frac{\gamma}{1} \Gamma, Da, \forall s Dc(s/b) \vee \exists s (D\neg c(s/b) \wedge Fc(s/b))}{\text{CD}^\infty \frac{\delta}{1} \Gamma, Da} \text{ (D6}^\infty\text{pos)},$$

where  $\gamma < \delta$  and  $a^\mathbb{N} = (\forall bc)^\mathbb{N} \in \text{Sent}^\mathbb{N}$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . As before, if  $\models \Gamma \cup \{Da\} [\sigma, \sigma + 2^\gamma, \xi]$ , then the claim follows by Lemma 6.6. Otherwise, by the induction hypothesis, we have

$$\models \forall s Dc(s/b) \vee \exists s (D\neg c(s/b) \wedge Fc(s/b)) [\sigma, \sigma + 2^\gamma, \xi].$$

Since  $\sigma + 2^\gamma + 1 \leq \sigma + 2^\delta \leq \xi$ , it follows from Lemma 4.5 that

$$(c(d/b))^\mathbb{N} \in D_{\sigma+2^\delta} \text{ for all closed terms } d \text{ with } d^\mathbb{N} \in \text{CI}^\mathbb{N}\text{Term}^\mathbb{N},$$

$$\text{or } (\neg c(d/b))^\mathbb{N} \in D_{\sigma+2^\delta} \cap T_{\sigma+2^\delta} \text{ for some } d \text{ with } d^\mathbb{N} \in \text{CI}^\mathbb{N}\text{Term}^\mathbb{N},$$

which implies  $a^\mathbb{N} \in D_{\sigma+2^\delta+1}$  by Lemma 4.2.6 and thus  $\models Da [\sigma, \sigma + 2^\delta, \xi]$ . We leave the other cases to the reader.

Now, let us assume that the last inference is made by a truth rule. Suppose the last inference is made by either (T4<sup>pos</sup><sub>∞</sub>) or (T3<sup>neg</sup><sub>∞</sub>):

$$\frac{\text{CD}^\infty \frac{\gamma}{1} \Gamma, Ta, \neg Tb}{\text{CD}^\infty \frac{\delta}{1} \Gamma, Ta} \text{ (T4}^\infty\text{pos)} \qquad \frac{\text{CD}^\infty \frac{\gamma}{1} \Gamma, \neg Ta, Tb}{\text{CD}^\infty \frac{\delta}{1} \Gamma, \neg Ta} \text{ (T4}^\infty\text{neg)},$$

where  $\gamma < \delta$  and  $a^\mathbb{N} = (\neg b)^\mathbb{N} \in \text{Sent}^\mathbb{N}$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . We will prove the claim only for the former case; the latter case can be similarly shown.

If  $\models \neg Tb[\sigma, \sigma + 2^\gamma, \xi]$ , then  $b^{\mathbb{N}} \notin T_{\xi+1}$  and thus  $a^{\mathbb{N}} \in T_{\xi+1}$  by Lemma 4.3.4, which implies  $\models Ta[\sigma, \sigma + 2^\delta, \xi]$ . Otherwise, we have  $\models \Gamma \cup \{Ta\}[\sigma, \sigma + 2^\gamma, \xi]$  by the induction hypothesis, and the claim follows by Lemma 6.6 as before. Next, suppose that we make either of the following inferences at the last step:

$$\frac{CD^\infty \frac{\gamma}{1} \Gamma, Ta, Tb^\circ \wedge Db^\circ}{CD^\infty \frac{\delta}{1} \Gamma, Ta} (T3_\infty^{pos}) \quad \frac{CD^\infty \frac{\gamma}{1} \Gamma, \neg Ta, \neg Tb^\circ \wedge Db^\circ}{CD^\infty \frac{\delta}{1} \Gamma, \neg Ta} (T3_\infty^{neg}),$$

where  $\gamma < \delta$  and  $a^{\mathbb{N}} = (Tb)^\mathbb{N} \in \text{Sent}^\mathbb{N}$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . For the former case, if  $\models Tb^\circ \wedge Db^\circ[\sigma, \sigma + 2^\gamma, \xi]$ , then, by Lemma 4.5, we have

$$(b^\circ)^\mathbb{N} \in D_{\sigma+2^\gamma+1} \cap T_{\xi+1} \subset D_{\sigma+2^\delta} \cap T_{\xi+1} = D_{\sigma+2^\delta} \cap T_{\sigma+2^\delta},$$

and thus, by Lemmata 4.2.2, 4.3.2, and 4.5, we obtain

$$a^{\mathbb{N}} \in D_{\sigma+2^\delta+1} \cap T_{\sigma+2^\delta+1} = D_{\sigma+2^\delta+1} \cap T_{\xi+1},$$

which implies  $\models Ta[\sigma, \sigma + 2^\delta, \xi]$ ; otherwise, we have  $\models \Gamma \cup \{Ta\}[\sigma, \sigma + 2^\gamma, \xi]$  by the induction hypothesis, and the claim follows by Lemma 6.6 as before. Similarly, for the latter case, if  $\models \neg Tb^\circ \wedge Db^\circ[\sigma, \sigma + 2^\gamma, \xi]$ , then, by Lemma 4.5, we have

$$(b^\circ)^\mathbb{N} \in D_{\sigma+2^\gamma+1} \setminus T_{\xi+1} \subset D_{\sigma+2^\delta} \setminus T_{\xi+1} = D_{\sigma+2^\delta} \setminus T_{\sigma+2^\delta},$$

and thus, again by Lemmata 4.2.2, 4.3.2, and 4.5, we obtain

$$a^{\mathbb{N}} \in D_{\sigma+2^\delta+1} \setminus T_{\sigma+2^\delta+1} = D_{\sigma+2^\delta+1} \setminus T_{\xi+1},$$

which implies  $\models \neg Ta[\sigma, \sigma + 2^\delta, \xi]$ ; otherwise, we have  $\models \Gamma \cup \{\neg Ta\}[\sigma, \sigma + 2^\gamma, \xi]$  by the induction hypothesis, and the claim follows by Lemma 6.6 as before. Thirdly, suppose that the last inference is made by the rule  $(T2_\infty^{pos})$ :

$$\frac{CD^\infty \frac{\gamma}{1} \Gamma, Ta, Db^\circ}{CD^\infty \frac{\delta}{1} \Gamma, Ta} (T2_\infty^{pos}),$$

where  $\gamma < \delta$  and  $a^{\mathbb{N}} = (Db)^\mathbb{N} \in \text{Sent}^\mathbb{N}$ . Take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . Suppose  $\models Db^\circ[\sigma, \sigma + 2^\gamma, \xi]$ . We have  $(b^\circ)^\mathbb{N} \in D_{\sigma+2^\gamma+1}$  and thus  $(b^\circ)^\mathbb{N} \in D_\xi$  by Lemma 4.5.2. Hence, we get  $a^{\mathbb{N}} \in T_{\xi+1}$  by Lemma 4.3.3 and thus  $\models Ta[\sigma, \sigma + 2^\delta, \xi]$ . Otherwise, the claim follows by Lemma 6.6 as before. The remaining cases for the other truth rules can be dealt with similarly.

Finally, let us assume that the last inference is made by (cut). We only consider the crucial case where the cut formulae are of the forms  $Da$  and  $\neg Da$ :

$$\frac{CD^\infty \frac{\gamma_0}{1} \Gamma, Da \quad CD^\infty \frac{\gamma_1}{1} \Gamma, \neg Da}{CD^\infty \frac{\delta}{1} \Gamma},$$

where  $\gamma_0, \gamma_1 < \delta$ . Let  $\gamma = \max\{\gamma_0, \gamma_1\}$  and take any  $\sigma$  and  $\xi$  with  $\sigma + 2^\delta \leq \xi$ . Suppose  $\not\models \Gamma[\sigma, \sigma + 2^\delta, \xi]$  for contradiction. Since  $2^\gamma + 2^\gamma \leq 2^\delta$ , it would follow that

$$\not\models \Gamma[\sigma, \sigma + 2^\gamma, \xi] \quad \text{and} \quad \not\models \Gamma[\sigma + 2^\gamma, \sigma + 2^\gamma + 2^\gamma, \xi]$$

by Lemma 6.6. Hence, by the induction hypothesis, we would have

$$\models Da[\sigma, \sigma + 2^\gamma, \xi] \quad \text{and} \quad \models \neg Da[\sigma + 2^\gamma, \sigma + 2^\gamma + 2^\gamma, \xi],$$

which means  $a^{\mathbb{N}} \in D_{\sigma+2\gamma+1}$  and  $a^{\mathbb{N}} \notin D_{\sigma+2\gamma+1}$ ; a contradiction. The remaining case in which the cut formulae are closed  $\mathcal{L}_T$ -literals can be straightforwardly treated.  $\dashv$

Combining this with Lemma 6.5, we obtain the next corollary.

**COROLLARY 6.8.** *For each  $\mathcal{L}$ -sentence  $\phi$ , if  $CD \vdash \phi$ , then there is some  $\delta < \varepsilon_0$  such that  $\models \psi[0, 2^\delta, 2^\delta]$ .*

Now, the argument of this section can be adequately formalized within KF (or  $RA_{<\varepsilon_0}$ ), in which we can formalize the construction of  $D_\xi$ s and  $T_\xi$ s for sufficiently large ordinals ( $< \varepsilon_0$ ) and prove (the arithmetization of) Corollary 6.8 in terms of them. Hence, in particular, we obtain the following.

**COROLLARY 6.9.** *CD is arithmetically conservative over KF.*

Combining this corollary with Corollary 5.6, we obtain our main theorem.

**THEOREM 6.10.**  *$CD_0, CD_1, CD_2$ , and CD are arithmetically equivalent to KF,  $RA_{<\varepsilon_0}$ , and  $RT_{<\varepsilon_0}$ , that is, they have exactly the same arithmetical theorems.*

In fact, our proofs establish stronger results: All the systems are proof-theoretically reducible to each other in the sense of Feferman [7, 9].

**§7. Upper bound of the strength of  $CD^+$ .** We will show that  $CD^+$  is arithmetically conservative over  $RT_{<\varepsilon_0}$ . This will be shown by constructing a relative interpretation of  $CD^+$  in  $CT[[KF + Cons]] (= CT[[KF]] + Cons)$  that preserves the arithmetical part.

**7.1. Total and consistent predicates.** We begin with the following definition:

$$D^+(x) :\Leftrightarrow (Fx \leftrightarrow \neg Tx).$$

Thus,  $D^+(x)$  is equivalent to  $(Tx \vee Fx) \wedge (\neg Tx \vee \neg Fx)$ . We obviously have

$$\begin{aligned} KF + Cons &\vdash (\forall x \in \text{Sent}_T)(D^+(x) \leftrightarrow (Tx \vee Fx)), \\ KF + Comp &\vdash (\forall x \in \text{Sent}_T)(D^+(x) \leftrightarrow (\neg Tx \vee \neg Fx)). \end{aligned} \tag{11}$$

**LEMMA 7.1.** *The following are provable in KF.*

1.  $\forall s (D^+(Ts) \leftrightarrow D^+(s^\circ))$ ; use K2.
2.  $(\forall x \in \text{Sent}_T)(D^+(\neg x) \leftrightarrow D^+x)$ ; use K3.

**LEMMA 7.2.** *The following are provable in  $KF + Cons$ .*

1.  $\forall x \forall y ( \text{Sent}_T(x \wedge y) \rightarrow (D^+(x \wedge y) \leftrightarrow ((D^+x \wedge D^+y) \vee (D^+x \wedge Fx) \vee (D^+y \wedge Fy))) )$ .
2.  $\forall x \forall y ( \text{Sent}_T(x \vee y) \rightarrow (D^+(x \vee y) \leftrightarrow ((D^+x \wedge D^+y) \vee (D^+x \wedge Tx) \vee (D^+y \wedge Ty))) )$ .
3.  $\forall v \forall x ( \text{Sent}_T(\forall vx) \rightarrow (D^+(\forall vx) \leftrightarrow (\forall s D^+x(s/v) \vee \exists s (D^+x(s/v) \wedge Fx(s/v))) ) )$ .
4.  $\forall v \forall x ( \text{Sent}_T(\exists vx) \rightarrow (D^+(\exists vx) \leftrightarrow (\forall s D^+x(s/v) \vee \exists s (D^+x(s/v) \wedge Tx(s/v))) ) )$ .

**PROOF.** We work within  $KF + Cons$ ; recall that we have  $D^+(x) \leftrightarrow (Tx \vee Fx)$  for all  $x \in \text{Sent}_T$  in  $KF + Cons$ . For claim 1, take any  $x, y \in \text{Sent}_T$ ; then we have

$$\begin{aligned} T(x \wedge y) \vee F(x \wedge y) &\stackrel{K4 \ \& \ K5}{\Leftrightarrow} (Tx \wedge Ty) \vee (Fx \vee Fy) \\ &\stackrel{\text{by logic}}{\Leftrightarrow} (D^+x \wedge D^+y) \vee (D^+x \wedge Fx) \vee (D^+y \wedge Fy). \end{aligned}$$

Claim 2 can be shown similarly to 1 using Lemma 7.1.2 and K3. For claim 3, take any  $v$  and  $x$  with  $\forall vx \in \text{Sent}_T$ ; then we have

$$\begin{aligned} T(\forall vx) \vee F(\forall vx) &\stackrel{K6 \ \& \ K7}{\Leftrightarrow} \forall sTx(s/v) \vee \exists sFx(a/v) \\ &\stackrel{\text{by logic}}{\Leftrightarrow} \forall sD^+x(s/v) \vee \exists s(D^+x(s/v) \wedge Fx(s/v)). \end{aligned}$$

Claim 4 can be shown similarly to 3 using Lemma 7.1.2 and K3. ⊢

LEMMA 7.3.  $\text{KF} \vdash \forall s(D^+(Ts \vee Fs) \leftrightarrow D^+s^\circ) \wedge \forall s(D^+(\neg Ts \vee \neg Fs) \leftrightarrow D^+s^\circ)$ .

PROOF. Take any  $s \in \text{CTerm}$ . The following four equivalences are provable in KF:

$$\begin{array}{llll} T(Ts \vee Fs) & \stackrel{5.2.1 \ \& \ K2}{\Leftrightarrow} & Ts^\circ \vee Fs^\circ & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & T(\neg Fs \vee \neg Ts); \\ \neg F(Ts \vee Fs) & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & \neg Fs^\circ \vee \neg Ts^\circ & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & \neg F(\neg Fs \vee \neg Ts); \\ F(Ts \vee Fs) & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & Fs^\circ \wedge Ts^\circ & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & F(\neg Fs \vee \neg Ts); \\ \neg T(Ts \vee Fs) & \stackrel{5.2.1 \ \& \ K2}{\Leftrightarrow} & \neg Ts^\circ \wedge \neg Fs^\circ & \stackrel{5.2.1 \ \& \ K2 \ \& \ K3}{\Leftrightarrow} & \neg T(\neg Fs \vee \neg Ts). \end{array}$$

The claim follows from these by logic. ⊢

**7.2. Reduction of  $\text{CD}^+$  to  $\text{CT}[\text{KF} + \text{Cons}]$ .** In the present subsection, we give an interpretation of  $\text{CD}^+$  in  $\text{CT}[\text{KF} + \text{Cons}]$  that preserves  $\mathcal{L}_0$ . For this purpose, we need a few preliminary definitions.

By the recursion theorem we obtain a recursive function  $f$  that satisfies the following condition:

$$f x = \begin{cases} x, & \text{if } x \in \text{AtFml}_0, \\ T f y, & \text{if } x = T y \text{ for } y \in \text{Term}, \\ T f y \vee F f y, & \text{if } x = D y \text{ for } y \in \text{Term}, \\ \neg f y, & \text{if } x = \neg y \in \text{Fml}, \\ (f y) \wedge (f z), & \text{if } x = y \wedge z \in \text{Fml}, \\ \forall v (f v), & \text{if } x = \forall v y \in \text{Fml}, \\ \ulcorner 0 \neq 0 \urcorner, & \text{otherwise.} \end{cases} \tag{12}$$

Here  $f$  is an arithmetical representation of the syntactic operation of applying the function  $f$  to terms so that  $(f s)^\circ = f(s^\circ)$  is provable in PA for each  $s \in \text{CTerm}$ .<sup>9</sup> We can show that  $f$  is provably total in PA and thus  $f$  is a PA-definable function. Obviously we have  $f x \in \text{Fml}_T$  for every  $x \in \text{Fml}$  (provably in PA), and, as in the case for the function  $k$ , we have  $\forall z(\text{Free}(z, x) \leftrightarrow \text{Free}(z, f x))$ ; hence, we have  $f x \in \text{Sent}_T$  for every  $x \in \text{Sent}$  (provably in PA).

The next is shown by routine induction on the complexity of  $x$  (cf. Lemma 5.7).

LEMMA 7.4.  $\text{PA} \vdash (\forall x \in \text{Fml}) \forall s (\forall v \in \text{Var})(f(x(s/v)) = (f x)(s/v))$ .

<sup>9</sup>The language  $\mathcal{L}_0$  may not possess a function symbol for  $f$  depending on how one defines  $\mathcal{L}_0$ . In such a case,  $f$  is defined by some  $\mathcal{L}_0$ -formula  $\phi(x, y)$ , and  $T f z$  means  $\ulcorner \forall w(\phi(x, w) \rightarrow T w) \urcorner(z/\ulcorner x \urcorner)$ . Then, we can easily verify that KF proves that  $T \ulcorner \forall w(\phi(x, w) \rightarrow T w) \urcorner(z/\ulcorner x \urcorner)$  is equivalent to  $T f(z^\circ)$ .

Now, we define a translation  $\sharp$  of  $\mathcal{L}$  into  $\mathcal{L}_T^+$  as follows: for each term  $a$ ,

$$D^\sharp a := Tfa \vee Ffa \quad \text{and} \quad T^\sharp a := Tfa;$$

$\sharp$  preserves all the  $\mathcal{L}_0$ -vocabulary. By (11), we have

$$KF + \text{Cons} \vdash (\forall x \in \text{Sent})(D^\sharp x \leftrightarrow D^+(fx)). \tag{13}$$

We will show that  $\sharp$  is a relative interpretation of CD in  $\text{CT}[[KF + \text{Cons}]]$ .

LEMMA 7.5. *The following are provable in  $KF + \text{Cons}$ .*

1.  $\forall s \forall t D^\sharp(s=t)$ ; by (12) and K1.
2.  $\forall s (D^\sharp(Ts) \leftrightarrow D^\sharp(s^\circ))$ ; by (12), (13), and Lemma 7.1.1.
3.  $(\forall x \in \text{Sent})(D^\sharp(\neg x) \leftrightarrow D^\sharp x)$ ; by (12), (13), and Lemma 7.1.2.

Hence, the  $\sharp$ -translations of D1, D2, and D4 are provable in  $KF + \text{Cons}$ .

LEMMA 7.6.  $KF + \text{Cons} \vdash \forall s (D^\sharp(Ds) \leftrightarrow D^\sharp(s^\circ))$ . Hence, the  $\sharp$ -translations of D3 is provable in  $KF + \text{Cons}$ .

PROOF. The claim is readily observed from the following equivalences:

$$D^\sharp(Ds) \stackrel{(13)}{\Leftrightarrow} D^+f(Ds) \stackrel{(12)}{\Leftrightarrow} D^+(Tfs \vee Ffs) \stackrel{7,3}{\Leftrightarrow} D^+(fs)^\circ \Leftrightarrow D^+f(s^\circ) \stackrel{(13)}{\Leftrightarrow} D^\sharp(s^\circ). \quad \dashv$$

LEMMA 7.7. *The following are provable in  $\text{CT}[[KF + \text{Cons}]]$ .*

1.  $\forall s \forall t (T^\sharp(s=t) \leftrightarrow s^\circ = m^\circ)$ .
2.  $\forall s (T^\sharp(Ds) \leftrightarrow D^\sharp(s^\circ))$ .
3.  $(\forall x \in \text{Sent})(T^\sharp(\neg x) \leftrightarrow \neg T^\sharp x)$ .
4.  $(\forall x, y \in \text{Sent})(T^\sharp(x \wedge y) \leftrightarrow (T^\sharp x \wedge T^\sharp y))$ .
5.  $\forall v \forall x (\forall v x \in \text{Sent} \rightarrow (T^\sharp(\forall v x) \leftrightarrow \forall s T^\sharp x(s/v)))$ .

Hence, the  $\sharp$ -translation of T1 and T2 and T4–T6 are provable in  $\text{CT}[[KF + \text{Cons}]]$ .

PROOF. We work within  $\text{CT}[[KF + \text{Cons}]]$ . Claim 1 is obvious by T1. For claim 2, take any  $s \in \text{CIterm}$ , and then we infer

$$Tf(Ds) \stackrel{(12)}{\Leftrightarrow} T(Tfs \vee Ffs) \stackrel{5,4,1}{\Leftrightarrow} T(Tfs) \vee T(Ffs) \stackrel{T1}{\Leftrightarrow} Tf(s^\circ) \vee Ff(s^\circ) \Leftrightarrow D^\sharp(s^\circ).$$

Claims 3 and 4 readily follow from T2 and T3, respectively, together with the definition (12) of  $f$ . For claim 5, let  $v$  and  $x$  be such that  $\forall v x \in \text{Sent}$ . Then, we have

$$Tf(\forall v x) \stackrel{(12)}{\Leftrightarrow} T\forall v f(x) \stackrel{T4}{\Leftrightarrow} \forall s T(fx)(s/v) \stackrel{7,4}{\Leftrightarrow} \forall s Tf(x(s/v)). \quad \dashv$$

LEMMA 7.8.  $\text{CT}[[KF + \text{Cons}]] \vdash \forall x (D^\sharp(x) \rightarrow (Tf(x) \leftrightarrow Tf(x)))$ .

PROOF. If  $x \notin \text{Sent}$  then  $f(x) = \top \neq 0$  and thus  $Tf(x) \wedge Tf(x)$ . Let  $x \in \text{Sent}$ . The claim is shown by induction on the complexity of  $x$ . For the base step, assume  $x \in \text{AtSent}$ . The case where  $x \in \text{AtSent}_0$  is obvious. Next, if  $x = Ts$  for some  $s \in \text{CIterm}$ , then we have

$$Tf(x) \stackrel{(12)}{\Leftrightarrow} T(Tfs) \stackrel{T1}{\Leftrightarrow} Tf(s^\circ) \stackrel{K2}{\Leftrightarrow} TTfs \stackrel{(12)}{\Leftrightarrow} Tf(x).$$

Lastly, let  $x = Ds$  for some  $s \in \text{CTerm}$ . Then, we have

$$\mathbb{T}f(x) \stackrel{(12) \& 5.4.1}{\Leftrightarrow} \mathbb{T}(\mathbb{T}f.s) \vee \mathbb{T}(Ff.s) \stackrel{\mathbb{T}1}{\Leftrightarrow} \mathbb{T}fs^\circ \vee Ffs^\circ \stackrel{\mathbb{K}2}{\Leftrightarrow} \mathbb{T}(\mathbb{T}f.s) \vee \mathbb{T}(Ff.s) \stackrel{5.2.1 \& (12)}{\Leftrightarrow} \mathbb{T}f(x).$$

For the induction step, we first let  $x = \neg y$  and assume  $D^\sharp x$ . We have  $D^\sharp y$  by Lemma 7.5.3 and thus  $D^+f(y)$  by (13). Hence, we have

$$\mathbb{T}f(x) \stackrel{(12)}{\Leftrightarrow} \mathbb{T}\neg f(y) \stackrel{\mathbb{T}2}{\Leftrightarrow} \neg \mathbb{T}f(y) \stackrel{\mathbb{H}}{\Leftrightarrow} \neg \mathbb{T}f(y) \stackrel{D^+f(y)}{\Leftrightarrow} Ff(y) \stackrel{(12)}{\Leftrightarrow} \mathbb{T}f(x);$$

note that the assumption  $D^\sharp x$  is crucial here. The other cases can be shown similarly, but without the need to use the assumption  $D^\sharp x$ .  $\dashv$

LEMMA 7.9.  $\text{CT}[\mathbb{K}F + \text{Cons}] \vdash \forall s (D^\sharp(s^\circ) \rightarrow (\mathbb{T}^\sharp(Ts) \leftrightarrow \mathbb{T}^\sharp(s^\circ)))$ . Hence, the  $\sharp$ -translation of T3 is provable in  $\text{CT}[\mathbb{K}F + \text{Cons}]$ .

PROOF. Take any  $s \in \text{CTerm}$  and suppose  $D^\sharp(s^\circ)$ . Then, we have

$$\mathbb{T}f(Ts) \stackrel{(12)}{\Leftrightarrow} \mathbb{T}\mathbb{T}f.s \stackrel{\mathbb{T}1}{\Leftrightarrow} \mathbb{T}f(s^\circ) \stackrel{7.8}{\Leftrightarrow} \mathbb{T}f(s^\circ). \quad \dashv$$

LEMMA 7.10.  $\text{CT}[\mathbb{K}F + \text{Cons}]$  proves the following.

1.  $(\forall x, y \in \text{Sent})(D^\sharp(x \wedge y) \leftrightarrow ((D^\sharp x \wedge D^\sharp y) \vee (D^\sharp x \wedge F^\sharp x) \vee (D^\sharp y \wedge F^\sharp y)))$ .
2.  $\forall v \forall x (\forall vx \in \text{Sent} \rightarrow (D^\sharp(\forall vx) \leftrightarrow (\forall s D^\sharp x(s/v) \vee \exists s (D^\sharp x(s/v) \wedge F^\sharp x(s/v)))))$ .

Hence, the  $\sharp$ -translations of D5 and D6 are provable in  $\mathbb{K}F + \text{Cons}$ .

PROOF. 1. Let  $x, y \in \text{Sent}$ . Then, we have

$$D^\sharp(x \wedge y) \stackrel{(12) \& (13)}{\Leftrightarrow} D^+(fx \wedge fy) \stackrel{7.2.1}{\Leftrightarrow} (D^+fx \wedge D^+fy) \vee (D^+fx \wedge Ffx) \vee (D^+fy \wedge Ffy) \stackrel{(12) \& (13) \& 7.8}{\Leftrightarrow} (D^\sharp x \wedge D^\sharp y) \vee (D^\sharp x \wedge F^\sharp x) \vee (D^\sharp y \wedge F^\sharp y).$$

2. Let  $\forall vx \in \text{Sent}$ . Then, we similarly obtain

$$D^\sharp(\forall vx) \stackrel{(12) \& (13) \& 7.2.3}{\Leftrightarrow} (\forall s D^+(fx)(s/v) \vee \exists s (D^+(fx)(s/n) \wedge Ff(x)(s/v))) \stackrel{(12) \& (13) \& 7.4 \& 7.8}{\Leftrightarrow} (\forall s D^\sharp x(s/v) \vee \exists s (D^\sharp x(s/v) \wedge F^\sharp x(s/v))). \quad \dashv$$

By Lemmata 7.5–7.7, 7.9, and 7.10, we obtain the next theorem.

THEOREM 7.11. The translation  $\sharp$  is a relative interpretation of  $\text{CD}^+$  in  $\text{CT}[\mathbb{K}F + \text{Cons}]$ . Hence, in particular,  $\text{CD}^+$  is arithmetically conservative over  $\text{CT}[\mathbb{K}F + \text{Cons}]$ .

Combining this theorem with Corollary 5.10 and Fact 5.11, we obtain the next.

THEOREM 7.12.  $\text{CD}_2^+$  and  $\text{CD}^+$  are arithmetically equivalent to  $\text{RA}_{<\varepsilon_{E_0}}$  and  $\text{RT}_{<\varepsilon_{E_0}}$ .

As before, our proofs actually establish that all these systems are proof-theoretically reducible to each other in the sense of Feferman [7, 9].

**§8. Some supplementary results.**

**8.1. An interpretation of  $CD^+$  in  $CT[[KF + Comp]]$ .** Preliminarily, we give an interpretation of  $CD^+$  in  $CT[[KF + Comp]]$ .<sup>10</sup>

Similarly to the function  $f$ , introduced in (12) above, we define a recursive function  $h$  by the recursion theorem so that the following condition obtains:

$$h(x) := \begin{cases} x, & \text{if } x \in \text{AtFml}_0, \\ \top h y, & \text{if } x = \top y \text{ for } y \in \text{Term}, \\ \neg \top h y \vee \neg F f y, & \text{if } x = D y \text{ for } y \in \text{Term}, \\ \neg h(y), & \text{if } x = \neg y \in \text{Fml}, \\ h(y) \wedge h(z), & \text{if } x = y \wedge z \in \text{Fml}, \\ \forall v. h(y), & \text{if } x = \forall v y \in \text{Fml}, \\ 0, & \text{otherwise.} \end{cases} \tag{14}$$

Using  $h$ , the translation  $b$  of  $\mathcal{L}$  into  $\mathcal{L}_T$  is defined as follows: for each term  $a$ ,

$$D^b a = \neg T h a \vee \neg F h a \quad \text{and} \quad T^b a = T h(a);$$

all the  $\mathcal{L}_0$ -vocabulary, logical connectives, and quantifiers are preserved by  $b$ . By (11), we have  $KF + Comp \vdash \forall x (D^b x \leftrightarrow D^+(hx))$ .

The proof of the next theorem is completely dual to that of Theorem 7.11, and we omit the details.

**THEOREM 8.1.**  $b$  is a relative interpretation of  $CD$  in  $CT[[KF + Comp]]$ .

**8.2. Liars and truth tellers.** In previous section we compared  $CD$  and its variants with other truth theories with respect to their strength. Truth theories are often compared by the way liar and truth teller sentences behave. In the present section we establish a few results about their behaviour in  $CD$  and its variants.

Let  $l$  and  $t$  be closed  $\mathcal{L}_0$ -terms with the following properties:

$$PA \vdash l = \neg T l \quad \text{and} \quad PA \vdash t = T t. \tag{15}$$

Hence, we have

$$CD \vdash T l \leftrightarrow \neg T T l \quad \text{and} \quad CD \vdash T t \leftrightarrow T T t. \tag{16}$$

Therefore,  $\neg T l$  and  $T t$  are a liar and a truth-teller sentence, respectively.<sup>11</sup>

<sup>10</sup>Cantini's [5] 'dual' interpretation of  $KF + Cons$  in  $KF + Comp$  induces an interpretation of  $CD^+$  in  $CT[[KF + Comp]]$ , but we will need a different interpretation for our purposes, namely,  $b$  defined here.

<sup>11</sup>If the  $\mathcal{L}_0$ -vocabulary does not have  $l$  and  $t$  as genuine closed terms, they can still be taken as definable closed terms: in such a case, we take  $\mathcal{L}_0$ -formulae  $\phi(x)$  and  $\psi(x)$  with exactly one free variable such that

$$PA \vdash \forall x (\phi(x) \rightarrow x = \neg \forall x (\phi(x) \rightarrow \neg T x)) \quad \text{and} \quad PA \vdash \forall x (\psi(x) \rightarrow x = \neg \forall x (\psi(x) \rightarrow T x)),$$

and (16) means that

$$CD \vdash \forall x (\phi(x) \rightarrow (T x \leftrightarrow \neg T \neg \forall x (\phi(x) \rightarrow \neg T x))) \quad \text{and} \quad CD \vdash \forall x (\psi(x) \rightarrow (T x \leftrightarrow \neg T \neg \forall x (\psi(x) \rightarrow T x))),$$

both of which are readily verified. All the results in this section are true in terms of these paraphrases.

PROPOSITION 8.2.

1.  $CD \vdash \neg DI$ .
2.  $CD^+ \not\vdash Dt$ .
3.  $CD^+ \not\vdash Tt$  and  $CD^+ \not\vdash \neg Tt$ .
4.  $CD^+ \not\vdash \neg Tl$  and  $CD^+ \not\vdash Tl$ .

PROOF. 1. Provably in  $CD$ , if  $DI$  were the case, then we would have  $Tl \leftrightarrow T \ulcorner Tl \urcorner$  by  $T3$ , which contradicts (16).

2. Using (15), we can show that  $t \notin D_\xi$  for all  $\xi$  by straightforward induction on  $\xi$ .

3. For the function  $f$  defined in Section 7.2, we observe that  $PA$  proves the following:

$$f(t) = f(\ulcorner Tt \urcorner) = \ulcorner Tf(t) \urcorner \quad \text{and} \quad h(t) = h(\ulcorner Tt \urcorner) = \ulcorner Th(t) \urcorner; \quad (17)$$

hence, both  $f(t)$  and  $h(t)$  are also truth-tellers. Let  $I$  be the least Kripkean fixed point, namely, the least set  $I \subset \omega$  such that  $(\mathbb{N}, I) \models KF$ , where  $I$  interprets the truth predicate  $T$ . It is known that  $(\mathbb{N}, I) \models Cons$ . By the standard argument, we can show that  $f(t) \notin I$  and  $\neg h(t) \notin I$ . Let  $J := \{\#\neg\phi \mid \#\phi \notin I\}$ ; hence,  $(\mathbb{N}, J)$  is the  $\mathcal{L}_T$ -structure induced from  $(\mathbb{N}, I)$  by Cantini's dual interpretation, and thus  $(\mathbb{N}, J) \models KF + Comp$ . Then, we have

$$(\mathbb{N}, I) \not\models Tf(t) \quad \text{and} \quad (\mathbb{N}, J) \models Th(t). \quad (18)$$

Now, by (17), we have the following:

$$CT \ulcorner [KF] \urcorner \vdash T f(t) \leftrightarrow Tf(t) \quad \text{and} \quad CT \ulcorner [KF] \urcorner \vdash T h(t) \leftrightarrow Th(t). \quad (19)$$

Let  $I^+ := \{\#\phi \mid (\mathbb{N}, I) \models \phi\}$  and  $J^+ := \{\#\phi \mid (\mathbb{N}, J) \models \phi\}$ ; then, the  $\mathcal{L}_T^+$ -structures  $(\mathbb{N}, I, I^+)$  and  $(\mathbb{N}, J, J^+)$ , in which  $T$  is interpreted by  $I^+$  and  $J^+$  respectively, are models of  $CT \ulcorner [KF + Cons] \urcorner$  and  $CT \ulcorner [KF + Comp] \urcorner$ , respectively. It follows from (18) and (19) that

$$(\mathbb{N}, I, I^+) \not\models Tf(t) \quad \text{and} \quad (\mathbb{N}, J, J^+) \models \neg Th(t).$$

Hence,  $CD^+ \not\vdash Tt$  and  $CD^+ \not\vdash \neg Tt$  by Theorems 7.11 and 8.1.

4. Observe that  $PA$  proves the following:

$$f(l) = f(\ulcorner \neg Tl \urcorner) = \ulcorner \neg Tf(l) \urcorner \quad \text{and} \quad h(l) = h(\ulcorner \neg Tl \urcorner) = \ulcorner \neg Th(l) \urcorner; \quad (20)$$

namely,  $f(l)$  and  $h(l)$  are also liar sentences. Hence, we obtain

$$KF \vdash Tf(l) \leftrightarrow Ff(l) \quad \text{and} \quad KF \vdash Th(l) \leftrightarrow Fh(l), \quad (21)$$

which implies  $KF + Cons \vdash \neg Tf(l)$  and  $KF + Comp \vdash Th(l)$ . Finally, it follows from (20) that  $CT \ulcorner [KF] \urcorner$  proves

$$Tf(l) \leftrightarrow T \ulcorner \neg Tf(l) \urcorner \leftrightarrow \neg Tf(l) \quad \text{and} \quad Th(l) \leftrightarrow T \ulcorner \neg Th(l) \urcorner \leftrightarrow \neg Th(l).$$

Hence, we have  $CT \ulcorner [KF + Cons] \urcorner \vdash Tf(l)$  and  $CT \ulcorner [KF + Comp] \urcorner \vdash \neg Th(l)$ , which implies  $CD^+ \not\vdash \neg Tl$  and  $CD^+ \not\vdash Tl$  by Theorems 7.11 and 8.1.  $\dashv$

PROPOSITION 8.3. *Let us write  $\lambda$  for  $\neg\text{TI}$ . We have the following:*

1.  $\text{CD}^+ \not\vdash \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$  and  $\text{CD}^+ \not\vdash \neg \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$ ; hence, the truth of some logical truth is neither provable nor refutable in  $\text{CD}$ .
2.  $\text{CD}^+ \not\vdash \text{T}^\# (\text{T}^\# \neg \lambda^{\neg\neg} \leftrightarrow \neg \text{T}^\# \lambda^{\neg\neg})^{\neg}$ ; hence, not all  $\text{CD}$ -axioms are provably true in  $\text{CD}$ .
3.  $\text{CD}^+ \vdash \neg \text{D}^\# (\text{T}^\# \neg \lambda^{\neg\neg} \leftrightarrow \neg \text{T}^\# \lambda^{\neg\neg})^{\neg}$ ; hence, not all  $\text{CD}$ -axioms are provably determinate in  $\text{CD}$ .

PROOF. 1.  $\text{CT}[\text{KF}]$  proves the following equivalences:

$$\begin{aligned} \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg} &\Leftrightarrow \text{T}^\# \text{T}^\# f(\lambda \vee \neg \lambda^{\neg\neg})^{\neg} \Leftrightarrow \text{T}^\# f(\lambda \vee \neg \lambda^{\neg\neg}) \Leftrightarrow \text{T}^\# f(\lambda) \vee \text{F}^\# f(\lambda), \\ \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg} &\Leftrightarrow \text{T}^\# \text{Th}(\lambda \vee \neg \lambda^{\neg\neg})^{\neg} \Leftrightarrow \text{Th}(\lambda \vee \neg \lambda^{\neg\neg}) \Leftrightarrow \text{Th}(\lambda) \vee \text{Fh}(\lambda). \end{aligned}$$

By (15) and (21),  $\text{KF} + \text{Cons} \vdash \neg(\text{T}^\# f(\lambda) \vee \text{F}^\# f(\lambda))$  and  $\text{KF} + \text{Comp} \vdash \text{Th}(\lambda) \vee \text{Fh}(\lambda)$ . Hence,  $\text{CT}[\text{KF} + \text{Cons}] \vdash \neg \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$  and  $\text{CT}[\text{KF} + \text{Comp}] \vdash \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$ ; thus,  $\text{CD}^+ \not\vdash \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$  and  $\text{CD}^+ \not\vdash \neg \text{T}^\# \text{T}^\# \lambda \vee \neg \lambda^{\neg\neg}$  by Theorems 7.11 and 8.1.

2.  $\text{CT}[\text{KF}]$  derives the following:

$$\begin{aligned} \text{T}^\# (\text{T}^\# \neg \lambda^{\neg\neg} \leftrightarrow \neg \text{T}^\# \lambda^{\neg\neg})^{\neg} &\Leftrightarrow \text{T}^\# f(\text{T}^\# \text{T}^\# \neg \lambda^{\neg\neg}) \Leftrightarrow \neg \text{T}^\# f(\text{T}^\# \lambda^{\neg\neg}) \\ &\Leftrightarrow \text{F}^\# f(\lambda) \Leftrightarrow \neg \text{T}^\# f(\lambda) \\ &\stackrel{(15) \& (21)}{\Leftrightarrow} \text{T}^\# f(\lambda) \Leftrightarrow \neg \text{T}^\# f(\lambda). \end{aligned}$$

Hence, we have  $\text{CT}[\text{KF} + \text{Cons}] \vdash \neg \text{T}^\# (\text{T}^\# \neg \lambda^{\neg\neg} \leftrightarrow \neg \text{T}^\# \lambda^{\neg\neg})^{\neg}$ , from which the claim follows by Theorems 7.11.

3. Provably in  $\text{CD}$ , if  $\text{D}^\# (\text{T}^\# \neg \lambda^{\neg\neg} \leftrightarrow \neg \text{T}^\# \lambda^{\neg\neg})^{\neg}$ , then  $\text{D}^\# \neg \text{T}^\# \lambda^{\neg\neg}$  and thus  $\text{D}^\#$  by (15), which contradicts Proposition 8.2.1. ⊥

**8.3. Additional axioms.** In this subsection, we will consider a few conservative extensions of  $\text{CD}^+$  with additional axioms about iterations of truth.

LEMMA 8.4.

1.  $\text{CT}[\text{KF} + \text{Cons}] \vdash \forall s (\text{T}^\# \text{T}^\# s \rightarrow \text{T}^\# s^\circ)$ .
2.  $\text{CT}[\text{KF} + \text{Cons}] \vdash \forall s (\text{F}^\# s^\circ \rightarrow \text{F}^\# \text{T}^\# s)$ .

PROOF. We work within  $\text{CT}[\text{KF} + \text{Cons}]$ . Take any  $s \in \text{CTerm}$ . We have

$$\text{T}^\# f(\text{T}^\# s) \stackrel{(12) \& \text{T1}}{\Leftrightarrow} \text{T}^\# f(s^\circ) \stackrel{\text{by logic}}{\Leftrightarrow} \text{D}^+ f(s^\circ) \wedge \text{T}^\# f(s^\circ) \stackrel{7,8}{\Rightarrow} \text{T}^\# f(s^\circ);$$

recall that  $\text{T}^\# a := \text{T}^\# f a$  and  $\text{D}^\# a := \text{D}^+ f a$ . Similarly, for each  $s \in \text{CTerm}$ , we have

$$\begin{aligned} \neg \text{T}^\# f(\neg \text{T}^\# s) &\stackrel{(12) \& \text{T2}}{\Leftrightarrow} \text{T}^\# \neg f_s \stackrel{\text{T1}}{\Leftrightarrow} \text{T}^\# f(s^\circ) \stackrel{\text{by logic}}{\Leftrightarrow} \text{D}^+ f(s^\circ) \wedge \text{T}^\# f(s^\circ) \stackrel{7,8}{\Rightarrow} \text{T}^\# f(s^\circ) \\ &\stackrel{(12) \& \text{T2}}{\Leftrightarrow} \neg \text{T}^\# f(\neg s^\circ). \end{aligned}$$

This completes the proof. ⊥

By the dual argument, we can also show the following.

LEMMA 8.5.

1.  $\text{CT}[\text{KF} + \text{Comp}] \vdash \forall s (\text{F}^\# \text{T}^\# s \rightarrow \text{F}^\# s^\circ)$ .
2.  $\text{CT}[\text{KF} + \text{Comp}] \vdash \forall s (\text{T}^\# s^\circ \rightarrow \text{T}^\# \text{T}^\# s)$ .

Let  $\text{num}(n)$  represent the function  $n \mapsto \#\bar{n}$  and thus it is the code of the numeral for  $n$ . Following the notation of [21], given a formula  $\phi(x)$  with a distinguished free variable  $x$ ,  $\ulcorner \phi(\dot{z}) \urcorner$  denotes  $\ulcorner \phi(x) \urcorner(\text{num}(z)/\ulcorner x \urcorner)$  and also  $\ulcorner \phi(\dot{s}) \urcorner$  denotes  $\ulcorner \phi(x) \urcorner(\ulcorner s \urcorner/\ulcorner x \urcorner)$  for  $s \in \text{CTerm}$ . We will occasionally write  $\ulcorner \phi(\dot{z}) \urcorner$  or  $\ulcorner \phi(\dot{s}) \urcorner$  without specifying the distinguished free variable of  $\phi$ , but in doing so we always assume that  $\phi$  contains some free variable to be substituted for.

LEMMA 8.6.  $\text{CT}\llbracket\text{KF}\rrbracket \vdash \forall s(\text{T}^\# \ulcorner \text{T}\dot{s} \urcorner \leftrightarrow \text{T}^\# \ulcorner \text{T}s \urcorner) \wedge \forall s(\text{T}^b \ulcorner \text{T}\dot{s} \urcorner \leftrightarrow \text{T}^b \ulcorner \text{T}s \urcorner)$ .

PROOF. Take any  $s \in \text{CTerm}$ . The first conjunct is shown as follows:

$$\text{T}f(\ulcorner \text{T}\dot{s} \urcorner) \Leftrightarrow \text{T}\ulcorner \text{T}f(\dot{s}) \urcorner \Leftrightarrow \text{T}f(\ulcorner \text{T}s \urcorner) \Leftrightarrow \text{T}\dot{s}f \Leftrightarrow \text{T}f s^\circ \Leftrightarrow \text{T}f(\ulcorner \text{T}s \urcorner).$$

The second conjunct can be shown similarly by just replacing  $f$  above with  $h$ .  $\dashv$

We conclude with the next theorem, which immediately follows from the last three lemmata and Theorems 7.11 and 8.1.

THEOREM 8.7. *The following two systems are both conservative over  $\text{CD}^+$ :*

1.  $\text{CD}^+ + \forall s(\text{T}\dot{s} \rightarrow \text{T}s^\circ) + \forall s(\text{F}s^\circ \rightarrow \text{F}\dot{s}) + \forall s(\text{T}\ulcorner \text{T}\dot{s} \urcorner \leftrightarrow \text{T}\dot{s})$ .
2.  $\text{CD}^+ + \forall s(\text{F}\dot{s} \rightarrow \text{F}s^\circ) + \forall s(\text{T}s^\circ \rightarrow \text{T}\dot{s}) + \forall s(\text{T}\ulcorner \text{T}\dot{s} \urcorner \leftrightarrow \text{T}\dot{s})$ .

Hence, either set of the new axioms can be consistently added to  $\text{CD}$ .

**§9. The minimality axiom.** The axioms for determinateness postulate that certain sentences are determinate; but  $\text{CD}$  lacks axioms that tell us that no other sentences are determinate. For instance, our axioms for determinateness do not rule out that truth teller sentences or similar sentences are determinate. Only in the case of sentences, such as liar sentences, where the assumption of determinateness leads to inconsistency, can we prove in  $\text{CD}$  that they are not determinate.

Theorem 4.13 says that no matter what model  $(\mathbb{N}, X, Y)$  of  $\text{CD}$  we choose,  $D_\infty$  is the least fixed point of  $\Gamma_{\mathcal{D}[Y]}$  and can be inductively defined ‘from the bottom up’ in  $Y$ . Hence,  $D_\infty$  seems to be a natural candidate for the ‘intended’ interpretation of the determinateness predicate  $D$ . In this section, we present an axiomatization of this conception of determinateness, which allows us to refute the determinateness of truth tellers, and state the proof-theoretic strength of the resulting system; however, the full analysis of the system goes beyond limited space of the present paper and is left for the Part II.

The system  $\text{ID}_1$  of (positive arithmetical) inductive definitions contains an axiom schema expressing that each fixed point expressed by a predicate of  $\text{ID}_1$  is the least, inductively defined, one. We use the same strategy for expressing the leastness of  $D$ . We remind the reader of the definition of  $\text{ID}_1$ . The language  $\mathcal{L}_{\text{fix}}$  of  $\text{ID}_1$  has a predicate symbol  $J_{\mathcal{A}}$  for each second-order arithmetical formula  $\mathcal{A}(x, X)$  with only the displayed variables free in which  $X$  occurs only positively.  $\text{ID}_1$  comprises the following axioms for each predicate  $J_{\mathcal{A}}$  together with the axioms of  $\text{PA}$  and full induction for  $\mathcal{L}_{\text{fix}}$ :

- ID1  $\forall x(\mathcal{A}(x, J_{\mathcal{A}}) \rightarrow J_{\mathcal{A}}x)$ ,
- ID2  $\forall x(\mathcal{A}(x, \Phi) \rightarrow \Phi(x)) \rightarrow \forall x(J_{\mathcal{A}}x \rightarrow \Phi(x))$ , for all  $\Phi \in \mathcal{L}_{\text{fix}}$ ,

where  $\mathcal{A}(x, \Phi)$  for an  $\mathcal{L}_{\text{fix}}$ -formula  $\Phi(u)$  with a designated variable  $u$  is the result of substituting  $\Phi(a)$  for each occurrence of  $a \in X$  in  $\mathcal{A}(x, X)$  for each term  $a$

(with renaming of bound variables as necessary to avoid collision). Inspired by ID2, Burgess [4] proposes the following axiom schema as an axiomatic characterization of the *groundedness* of the truth predicate T in Kripke’s [28] sense:

$$\forall x(\mathcal{B}(x, \Phi) \rightarrow \Phi(x)) \rightarrow \forall x(\text{T}x \rightarrow \Phi(x)), \text{ for all } \Phi \in \mathcal{L}_T,$$

where  $\mathcal{B}(x, X)$  is taken so that  $\forall x(\mathcal{B}(x, X) \rightarrow X)$  expresses that  $X$  is closed under the strong Kleene evaluation schema of truth values. Burgess’s schema expresses that the extension of T is the least Kripkean fixed point under the strong Kleene evaluation schema; the system obtained by augmenting KF with this schema is called KFB. We follow the lead of Burgess’s idea toward an axiomatic characterization of D as the least set that satisfies the closure conditions of determinateness relative to the truth predicate T (but recall that the least such set is invariant across the choice of T).

For each  $\mathcal{L}$ -formula  $\Phi(u)$  with a designated variable  $u$ , we write  $\mathcal{D}(x, \Phi)$  for the result of substituting  $\Phi(a)$  for each occurrence of  $Da$  in  $\mathcal{D}(x)$  for each term  $a$  (with renaming of bound variables as necessary to avoid collision). For the uniformity of notation, we will write  $\mathcal{D}(x, D)$  for  $\mathcal{D}(x)$  in (and only in) the present section. Thereby we introduce two new axioms.

$$D_\mu 1 \quad \forall x(\mathcal{D}(x, D) \rightarrow Dx).$$

$$D_\mu 2 \quad \forall x(\mathcal{D}(x, \Phi) \rightarrow \Phi(x)) \rightarrow \forall x(Dx \rightarrow \Phi(x)), \text{ for all } \Phi \in \mathcal{L}.$$

Note that  $D_\mu 2$  is not a single axiom but an axioms schema.

DEFINITION 9.1. We define the system  $CD_\mu$  as  $CD + D_\mu 1 + D_\mu 2$ . We can easily show that  $D_\mu 1$  is derivable from CD and thus  $CD_\mu$  is identical with  $CD + D_\mu 2$ . The system  $CD_\mu^+$  is defined as  $CD^+ + D_\mu 2$ .

The next follows from Theorems 4.8 and 4.13.

COROLLARY 9.2.  $(\mathbb{N}, D_\infty, \mathbb{T}) \models CD_\mu^+$ .

The system  $CD_\mu$  has some desirable properties. For example, we have the following.

PROPOSITION 9.3.  $CD_\mu \vdash \neg Dt$ ; see Section 8.2 for the definition of  $t$ .

PROOF. Suppose  $Dt$  for contradiction. Let  $\Phi(x) := Dx \wedge x \neq t$ . Since  $\mathcal{D}(t, \Phi)$  implies  $\Phi(t)$  by (15), we have  $\neg \mathcal{D}(t, \Phi)$  and thus  $\forall x(\mathcal{D}(x, \Phi) \rightarrow \Phi(x))$  by  $D_\mu 1$ , which would imply  $\forall x(Dx \rightarrow \Phi(x))$ ; a contradiction.  $\dashv$

In fact, all the underivability results proved in Section 8.2 still hold in  $CD_\mu$  and  $CD_\mu^+$ , and the additional axioms considered in Section 8.3 are also consistent with  $CD_\mu$  and  $CD_\mu^+$ .

The full proof-theoretic analyses of these new systems are left for the Part II. We state the main results below without giving their proofs.

THEOREM 9.4.

1.  $CD_\mu$  is proof-theoretically equivalent to  $ID_1$ .
2.  $CD_\mu^+$  is proof-theoretically equivalent to  $BID_1^2$  (see [36] for its definition and Fujimoto [15] for its proof-theoretic analysis), and thus its proof-theoretic ordinal (if appropriately defined) is  $\psi_\Omega(\varepsilon_{\varepsilon_0})$ .

**§10. Comparison with other axiomatizations of truth.** We think that CD is more promising than many other systems found in the literature. This section contains brief comparisons with such systems. We leave a more thorough discussion of the philosophical aspects for another paper.

**10.1. Comparison with the Friedman–Sheard system.** The Friedman–Sheard system FS features all the truth axioms of CD, except for the truth iteration axiom T3 and all axioms involving D, which is not in the language of FS. Axiom T3 is replaced with two rules called Necessitation NEC and Conecessitation CONEC:

$$\frac{\phi}{T^{\ulcorner}\phi^{\urcorner}} \text{ NEC} \qquad \frac{T^{\ulcorner}\phi^{\urcorner}}{\phi} \text{ CONEC.}$$

Adding these two rules for all sentences to T1 and T4 – T6 yields the Friedman–Sheard system FS, which was analyzed by Friedman and Sheard [13] with a different axiomatization. McGee [31] showed that a subsystem of FS is  $\omega$ -inconsistent and thus does not have a standard model, that is, the standard model of arithmetic cannot be expanded to a model of FS by specifying a suitable extension for the truth predicate. Halbach [18] determined the proof-theoretic strength of FS as ramified analysis  $RA_{<\omega}$  up to  $\omega$  or  $\omega$ -times iterated typed truth (see Halbach [21, Section 14.2]).

A neat feature of FS is its symmetry: What is provable in FS and what is provably true in FS coincide. Moreover, if  $\phi$  is a classical tautology such as  $\lambda \vee \neg\lambda$  (for a liar sentence  $\lambda$ ), the sentence

$$\underbrace{T^{\ulcorner}T^{\ulcorner}\dots T^{\ulcorner}}_{n \text{ iterations of } T} \ulcorner\phi^{\urcorner} \dots \urcorner^{\urcorner}$$

is provable in FS. In CD  $T^{\ulcorner}\phi^{\urcorner}$  is provable for every classical tautology, but  $T^{\ulcorner}T^{\ulcorner}\lambda \vee \neg\lambda^{\urcorner}$  and further iterations are not, as was shown in Proposition 8.3. In FS no transfinite iterations of truth can be proved. Hence neither is FS a subsystem of CD nor is CD a subsystem of FS. Moreover, neither can the truth predicate of FS be defined in CD nor can the truth predicate of CD be defined in FS. This follows from general arguments by Fujimoto [14] (see also Halbach [21]). We will show in Part II that CD can be consistently closed under NEC and CONEC. The result is a system that properly contains FS and is consequently  $\omega$ -inconsistent.

The system FS has rightly been criticized by various authors. Barrio and Picollo [2] list some of the problematic features of FS caused by the  $\omega$ -inconsistency. In particular, the rule NEC can naturally be replaced with an  $\omega$ -times iterated reflection axiom. If the reflection axiom is iterated into the transfinite, an inconsistency ensues (Halbach and Horsten [3] and Halbach [21, Corollary 14.39]). These shortcomings of FS prompted our search for an alternative to FS without giving up the thorough classicality of FS and endorsing the internal non-classicality of systems such as KF.

Some alternative systems have been developed that prove the truth of classical tautologies, but are  $\omega$ -consistent such as Cantini’s [6] VF and Stern’s [41] IT and their variants. For these systems  $\omega$ -models can be obtained via constructions involving supervaluations.

**10.2. Comparison with Kripke–Feferman.** Feferman [8] defined a determinateness predicate  $D$  in terms of a truth and a falsity predicate within the Kripke–Feferman

system KF.<sup>12</sup> Understanding the falsity of a sentence as the truth of its negation, Feferman defined  $D$  as the formula we called  $D^+$  above:

$$Dx \leftrightarrow \text{Sent}(x) \wedge (Tx \vee Fx) \wedge \neg(Tx \wedge Fx).$$

The schema  $D^\top \phi^\top \rightarrow (T^\top \phi^\top \leftrightarrow \phi)$  becomes provable in KF with this definition. The schema had been mentioned by Kripke [28] already.

Reinhardt [37, 38] takes the determinateness predicate to express *meaningful applicability* or *significance*; and truth is only meaningfully applicable to determinate sentences.<sup>13</sup> In particular, the liar sentence is not in the range of significance of the truth predicate. Reinhardt concludes that the system KF, which is formulated in classical logic, is not sound in the sense that it proves only true theorems, where truth is understood in the sense of the truth predicate  $T$  of KF. According to Reinhardt, one might hope that the sentences  $\phi$  with  $\text{KF} \vdash T^\top \phi^\top$  are true and significant. The set of these sentences is not closed under classical logic. The theorems of KF itself need not be significant or trustworthy; only those whose truth can be proved are. Generally, it is very hard to avoid the provability of sentences that are not ‘significant’ or not ‘healthy’, as Reinhardt [37] and Bacon [1] have argued.

In contrast to Reinhardt, we do not think that the range of significance of the truth predicate is restricted in any way; the determinateness predicate cannot be seen as indicating the range of significance of  $T$ . For instance, truth should commute with connectives for *all* sentences, not just with those in some range of significance. However, some sentences, including liar sentences, are sensitive to the addition of another layer of truth. Stacking an additional layer of truth onto the liar sentence will change its semantic status; but that does not mean that the truth predicate cannot meaningfully applied to it. We would only run into the Reinhardt–Bacon problem by postulating that only insensitive sentences ought to be provable. We reject this kind of soundness condition. We endorse sensitive sentences such as the compositional axioms or classical tautologies as theorems. However, we may not be able to ascend semantically from such sentences and add a further level of truth. It is not possible to ascend semantically from the compositional axiom; they are already at the highest level of generality and further semantic ascend and generalization is not possible.

**10.3. Comparison with Feferman’s DT.** Feferman [10, p. 205] originally introduced KF as an instrument to explain ‘what notions and principles one ought to accept if one accepts the basic notions and principles of the theory’, the foundational question which had long been one of the central themes of Feferman’s work.<sup>14</sup>

As a theory of truth *per se*, Feferman later proposed another system DT. Again, each predicate is assumed to have a *domain of significance* and to be meaningfully applicable only to objects in that domain. In the case of the truth predicate  $T$ , its domain  $\mathcal{D}$  of significance is taken to consist of the sentences that are *meaningful*

<sup>12</sup>The system is not called KF in [8]. Moreover, different systems have been called KF. See Halbach [21] for more on the history of KF and its variants.

<sup>13</sup>That predicates have a domain of significance has been part of the philosophy of type theory since [39]. Feferman [10] explicitly refers to Russell.

<sup>14</sup>However, Feferman [10, p. 205] expressed ‘I always thought that the KF axioms were a bit artificial for that purpose’ and abandoned KF as such an instrument in the end. In place of KF, he proposed a new notion of *unfolding* of a schematic system in place of KF for the purpose in question.

and determinate, and all the principles of truth are only applicable to the sentences in the domain  $\mathcal{D}$ . Accordingly, the system DT comprises two groups of axioms: the axioms characterizing the domain  $\mathcal{D}$  and the axioms expressing the T-schema and the compositional axioms for all sentences in  $\mathcal{D}$ .

Feferman’s axioms for determinateness differ from ours in particular with respect to the determinateness of sentences formed with binary connectives. In his system a disjunction is determinate if, and only if both disjuncts are; and a conjunction is determinate if, and only if both conjuncts are. We discussed our choice of axiom D5 on page 6 and justified its endorsement with the generalizing function of truth. Many harmless generalizations would become indeterminate (and consequently the system proof-theoretically weak) if generalization were expressed with conjunction and disjunction and Feferman’s concept of determinateness. To overcome the problem, he introduced a special conditional  $\rightarrow$  as a primitive logical connective so that the truth and determinateness of a conditional sentence  $\phi \rightarrow \psi$  is characterized independently and differently from its usual definition  $\neg(\phi \wedge \neg\psi)$  in terms of negation and conjunction (or  $\neg\phi \vee \psi$  in terms of negation and disjunction) by the following axioms:

$$\forall x \forall y (\text{Sent}(x) \wedge \text{Sent}(y) \rightarrow (\mathcal{D}(x \rightarrow y) \leftrightarrow (\mathcal{D}x \wedge (\text{T}x \rightarrow \mathcal{D}y))))),$$

$$\forall x \forall y (\text{Sent}(x) \wedge \text{Sent}(y) \wedge \mathcal{D}(x \rightarrow y) \rightarrow (\text{T}(x \rightarrow y) \leftrightarrow (\text{T}x \rightarrow \text{T}y))).$$

Feferman’s approach diverges from ours in two crucial respects: First, he took  $\mathcal{D}$  as the range of all the principles of truth comprehensively and restricted every principle of truth to the class of determinate sentences  $\mathcal{D}$ , whereas we take the class of determinate sentences only as the range of the disquotation schema and postulate the compositional axioms unrestrictedly for every sentence. Secondly, he took  $\mathcal{D}$  as definable in terms of truth by stipulating  $\mathcal{D}x \leftrightarrow \text{T}x \vee \text{F}x$ ; but this definition does not yield the desired properties of  $\mathcal{D}$  in our theory and we consequently introduce determinateness as a primitive notion.

Fujimoto [14] observed that DT is identical with the system FKF + Cons of non-classical partial truth whose evaluation rule is given by the Aczel–Feferman evaluation schema. The Aczel–Feferman evaluation schema is the same as the weak Kleene evaluation schema, with the exception of the evaluation rule of the conditional. FKF is a variant of KF with the Aczel–Feferman evaluation schema: FKF is to the Aczel–Feferman evaluation schema what KF is to the strong Kleene evaluation schema.

**10.4. Comparison with the Leitgeb–Schindler system.** Schindler [40] defined a group of systems that resemble ours. They differ from CD in restricting the compositional principles T4–T6 to *grounded* sentences. The predicate symbol  $G$  for groundedness plays a role comparable to that of  $\mathcal{D}$  in CD. The full system is called CG. All the unrestricted compositional axioms T4–T6 are provable except for the right-to-left direction of the negation axiom T4. That is,  $\forall x (\text{Sent}(x) \rightarrow (\neg\text{T}x \rightarrow \text{T}(\neg x)))$  is not provable in CG. In this sense Schindler’s system fails to be thoroughly classical.

One of the axioms of CG is a ‘definitional’ axiom of groundedness [40, p. 77]:

$$\forall x (Gx \leftrightarrow (\text{T}x \vee \text{T}\neg x)).$$

This axiom of CG is not a theorem (with D for  $G$ ) of our theory CD, because the full negation axiom T4 and thus  $\forall x (\text{Sent}(x) \rightarrow (Tx \vee T\neg x))$  are theorems of CD and therefore Schindler's definitional axiom would imply  $\forall x (\text{Sent}(x) \rightarrow Gx)$ . Consequently, Schindler's system does not really require a primitive predicate for groundedness, while the definition of determinateness as truth or falsity fails in our system.

On Schindler's approach, T applies provably only to sentences without  $G$ , because it lacks an axiom analogous to our T2. This is in line with the intended interpretation of  $G$  as Leitgeb's [30] notion of groundedness. Also  $G$  cannot be iterated, and axiom D3 is absent from Schindler's list of axioms. However, contrary to Schindler's official formulation of the system, we could view  $Gx$  as a metalinguistic abbreviation for  $Tx \vee T\neg x$ . Then T provably applies to sentences containing  $G$ .

**10.5. Comparison with Halbach's system PUTB.** The truth-theoretic axioms of the system PUTB are given by all instances  $T^\Gamma \phi^\neg \leftrightarrow \phi$  of the T-schema where the truth predicate occurs only positively in  $\phi$  and parameters are allowed in  $\phi$ .

Positiveness has the advantage over determinateness that it is a very simple syntactic decidable property that is definable in the language of arithmetic. If in the schema

$$\text{DDS} \quad \forall t_1 \dots \forall t_n \left( D^\Gamma \phi(t_1, \dots, t_n)^\neg \rightarrow (T^\Gamma \phi(t_1, \dots, t_n)^\neg \leftrightarrow \phi(t_1^\circ, \dots, t_n^\circ)) \right),$$

the symbol D expresses that T occurs only positively, no primitive predicate D is needed, as positiveness can be expressed by an arithmetical formula. Consequently, because D is arithmetical, our axiom T2<sup>+</sup>, that is,  $\forall s (TD_s \leftrightarrow Ds^\circ)$  becomes provable.

The theory PUTB does not prove the compositional axiom T4–T6. To obtain a system closer to CD, one could add the compositional axioms to PUTB or derive them from reflection principles, as Horsten and Leigh [27] suggest.

However, T-positiveness and determinateness differ significantly: In particular, truth-teller sentences are provably T-positive, while they are not provably determinate in our sense; in fact we can refute their determinateness in  $CD_\mu$  as was shown in proposition 9.3. The proof-theoretic strength of PUTB relies crucially on indeterminate instances that allow one to mimic positive inductive definitions. Overall the restriction to determinate instead of positive sentences is better motivated.

**10.6. Comparison with Picollo's system WFUTB.** Picollo [34, 35] defines a notion of referential well-foundedness, which is related to determinateness in our sense. Here we do not go into the details of its definition and only highlight some differences to our approach. First, the notion of well-foundedness is defined in such a way that it becomes arithmetically definable. Thus only a truth predicate needs to be added and no separate predicate for referential well-foundedness as our primitive predicate D. Roughly, the main axiom schema of her system WFUTB then states the T-sentences for sentences that are referentially well-founded in her sense (or PA- provably equivalent to such a sentence).<sup>15</sup> WFUTB does not feature

<sup>15</sup>Strictly, speaking there are two conditions on the permissible instances of the T-schema: The uniform T-sentences are postulated for instances that are provably r-stable and well-founded in PA. By *referential*

any compositional axioms, and it can be shown that they are not provable in it, although the system was shown by Picollo to be proof-theoretically at least as strong as ramified analysis  $RA_{<\Gamma_0}$  up to  $\Gamma_0$  and thus much stronger than PUTB and compositional theories such as KF.

**10.7. Comparison with Field's systems Int.** The most fundamental decision for the truth theorist is whether to sacrifice classical logic for transparency or transparency for classical logic. By transparency we mean here some equivalence of  $\phi$  with  $T^\Gamma \phi^\neg$  for all  $\phi$ . Field [11] sacrifices classical logic; we sacrifice transparency. Field saves truth from paradox; we save logic from paradox. We regain transparency for determinate sentences; Field regains classical logic for what he calls *strongly classical* sentences.

Field employs a new primitive predicate *ScI* for strongly classical truth. Although Field's and our approach are pulling in exactly opposite directions, Field's axioms for *ScI* and ours for *D* have a striking resemblance. As Field in footnote 6 mentions, he and we arrived at our axiomatizations independently. Of course, a serious comparison of Int and CD leads back to the most fundamental decision that truth theorists face, and we do not enter the discussion here.

**§11. Further perspectives.** In the part II of this paper, we will give proof-theoretic analysis of variants of CD. In particular, among many others, we will give a proof of Theorem 9.4. We conclude the paper by listing two open problems.

- (I) We may consider CD and its variants with the schema of induction restricted to the arithmetical sentences. In many cases, the restriction of induction to the arithmetical sentences results in a proof-theoretically conservative theory (over PA). We conjecture that the same holds for CD. Conservativeness proofs in the analogous case of KF can be given in a model-theoretic way by showing how to extend a given model of PA to a model of KF with arithmetical induction only. This is not possible in the case of CD because Lachlan's theorem [29] applies and nonstandard models that can be expanded to CD with restricted induction have to be recursively saturated.
- (II) It may be of interest to replace the axiom D3 of CD with an alternative axiom  $\forall s DD_s$ . Together with DDS this would yield our additional axiom T2<sup>+</sup>. Thus, adding  $\forall s DD_s$  may be more natural than adding T2<sup>+</sup>. We do not know how strong CD becomes if  $\forall s DD_s$  is added.

**Acknowledgments.** We are grateful to Dora Achourioti, Anton Broberg, Hannes Leitgeb, Beau Mount, Carlo Nicolai, Lavinia Picollo, Lorenzo Rossi, Thomas Schindler, Johannes Stern, Philip Welch, and the referee for comments and suggestions.

**Funding.** V.H. would like to thank the Leverhulme Trust for supporting his work with a Research Fellowship.

---

*well-foundedness* we mean the conjunction of the two conditions. Unlike most notions of groundedness and determinateness, referential well-foundedness is sensitive to the base theory.

## REFERENCES

- [1] A. BACON, *Can the classical logician avoid the revenge paradoxes?* *Philosophical Review*, vol. 124 (2015), pp. 299–352.
- [2] E. BARRIO and L. PICOLLO, *Notes on  $\omega$ -inconsistent theories of truth in second-order languages.* *Review of Symbolic Logic*, vol. 6 (2013), pp. 733–741.
- [3] J. C. BEALL and B. ARMOUR-GARB, *Deflationism and Paradox*, Clarendon Press, Oxford, 2005.
- [4] J. P. BURGESS, *Friedman and the axiomatization of Kripke's theory of truth*, *Foundational Adventures: Essays in Honor of Harvey M. Friedman*, 2009, paper delivered at the Ohio State University conference in honor of the 60th birthday of Harvey Friedman, <http://foundationaladventures.com/>.
- [5] A. CANTINI, *Notes on formal theories of truth.* *Zeitschrift für mathematische Logik und Grundlagen der Mathematik*, vol. 35 (1989), pp. 97–130.
- [6] ———, *A theory of formal truth arithmetically equivalent to  $ID_1$* , this JOURNAL, vol. 55 (1990), pp. 244–259.
- [7] S. FEFERMAN, *Hilbert's program relativized: Proof-theoretical and foundational reductions*, this JOURNAL, vol. 53 (1988), pp. 364–384.
- [8] ———, *Reflecting on incompleteness*, this JOURNAL, vol. 56 (1991), pp. 1–49.
- [9] S. FEFERMAN, *What rests on what? The proof-theoretic analysis of mathematics*, *Proceedings of the 15th International Wittgenstein Symposium*, Hölder-Pichler-Tempsky, Wien, 1992, pp. 47–171.
- [10] S. FEFERMAN, *Axioms for determinateness and truth.* *Review of Symbolic Logic*, vol. 1 (2008), pp. 204–217.
- [11] H. FIELD, *The power of naive truth.* *Review of Symbolic Logic*, (2020), pp. 1–34.
- [12] M. FISCHER, L. HORSTEN, and C. NICOLAI, *Hypatia's silence.* *Noûs*, vol. 55 (2021), pp. 62–85.
- [13] H. FRIEDMAN and M. SHEARD, *An axiomatic approach to self-referential truth.* *Annals of Pure and Applied Logic*, vol. 33 (1987), pp. 1–21.
- [14] K. FUJIMOTO, *Relative truth definability.* *Bulletin of Symbolic Logic*, vol. 16 (2010), pp. 305–344.
- [15] ———, *Notes on some second-order systems of iterated inductive definitions and  $\Pi_1^1$ -comprehensions and relevant subsystems of set theory.* *Annals of Pure and Applied Logic*, vol. 166 (2015), pp. 409–463.
- [16] ———, *Deflationism beyond arithmetic.* *Synthese*, vol. 196 (2019), pp. 1045–1069.
- [17] ———, *The function of truth and the conservativeness argument.* *Mind*, vol. 131 (2022), no. 521, pp. 129–157.
- [18] V. HALBACH, *A system of complete and consistent truth.* *Notre Dame Journal of Formal Logic*, vol. 35 (1994), pp. 311–327.
- [19] ———, *Disquotational truth and analyticity*, this JOURNAL, vol. 66 (2001), pp. 1959–1973.
- [20] ———, *How not to state the T-sentences.* *Analysis*, vol. 66 (2006), pp. 276–280.
- [21] ———, *Axiomatic Theories of Truth*, revised ed., Cambridge University Press, Cambridge, 2014 (first edition 2011).
- [22] ———, *Formal notes on the substitutional analysis of logical consequence.* *Notre Dame Journal of Formal Logic*, vol. 61 (2020), pp. 317–339.
- [23] ———, *The substitutional analysis of logical consequence.* *Noûs*, vol. 54 (2020), pp. 431–450.
- [24] V. HALBACH and G. LEIGH, *The Road to Paradox: A Guide to Syntax, Truth, and Modality*, Cambridge University Press, 2024, to appear.
- [25] J. VAN HEIJENOORT (editor), *From Frege to Gödel. A Source Book in Mathematical Logic, 1879–1931*, Harvard University Press, Cambridge, 1967.
- [26] W. HODGES, *Truth in a structure.* *Proceedings of the Aristotelean Society*, vol. 86 (1986), pp. 135–151.
- [27] L. HORSTEN and G. E. LEIGH, *Truth is simple.* *Mind*, vol. 126 (2017), pp. 195–232.
- [28] S. KRIPKE, *Outline of a theory of truth.* *Journal of Philosophy*, vol. 72 (1975), pp. 690–716.
- [29] A. LACHLAN, *Full satisfaction classes and recursive saturation.* *Canadian Mathematical Bulletin*, vol. 24 (1981), pp. 295–297.
- [30] H. LEITGEB, *What truth depends on.* *Journal of Philosophical Logic*, vol. 34 (2005), pp. 155–192.
- [31] V. MCGEE, *How truthlike can a predicate be? A negative result.* *Journal of Philosophical Logic*, vol. 14 (1985), pp. 399–410.
- [32] ———, *Maximal consistent sets of instances of Tarski's schema (T).* *Journal of Philosophical Logic*, vol. 21 (1992), pp. 235–241.
- [33] L. PICOLLO, *Reference in arithmetic.* *Review of Symbolic Logic*, vol. 11 (2018), pp. 573–603.

- [34] ———, *Reference and truth*. *Journal of Philosophical Logic*, vol. 49 (2020), pp. 439–474.
- [35] ———, *Minimalism, reference, and paradoxes*, *The Logica Yearbook 2015* (P. Arazim and M. Dancak, editors), College Publications, London, 2016, pp. 163–178.
- [36] W. POHLERS, *Subsystems of set theory and second order number theory*, *Handbook of Proof Theory* (S. Buss, editor), Studies in Logic and the Foundations of Mathematics, 137, Elsevier, Amsterdam, 1998, pp. 209–335.
- [37] W. REINHARDT, *Some remarks on extending and interpreting theories with a partial predicate for truth*. *Journal of Philosophical Logic*, vol. 15 (1986), pp. 219–251.
- [38] ———, *Remarks on significance and meaningful applicability*. *Mathematical Logic and Formal Systems: A Collection of Papers in Honor of Professor Newton C.A. Da Costa* (L. P. de Alcantara, editor), Lecture Notes in Pure and Applied Mathematics, 94, Marcel Dekker, 1985, pp. 227–242.
- [39] B. RUSSELL, *Mathematical logic as based on the theory of types*. *American Journal of Mathematics*, vol. 30 (1908), pp. 222–262, Reprinted in [25, 150–182].
- [40] T. SCHINDLER, *Axioms for grounded truth*. *Review of Symbolic Logic*, vol. 7 (2014), pp. 73–83.
- [41] J. STERN, *Supervaluation-style truth without supervaluations*. *Journal of Philosophical Logic*, vol. 47 (2018), pp. 817–850.
- [42] A. TARSKI, *Der Wahrheitsbegriff in den formalisierten Sprachen*. *Studia Philosophica Commentarii Societatis Philosophicae Polonorum*, vol. 1 (1935), pp. 261–405.
- [43] A. TARSKI and R. VAUGHT, *Arithmetical extensions of relational systems*. *Compositio Mathematica*, vol. 13 (1956), pp. 81–102.

NEW COLLEGE  
UNIVERSITY OF OXFORD  
OXFORD OX1 3BN  
UK

*E-mail:* [volker.halbach@new.ox.ac.uk](mailto:volker.halbach@new.ox.ac.uk)

SCHOOL OF MATHEMATICS  
UNIVERSITY OF BRISTOL  
WOODLAND ROAD  
BRISTOL BS8 1UG  
UK

*E-mail:* [kentaro.fujimoto@bristol.ac.uk](mailto:kentaro.fujimoto@bristol.ac.uk)