

# Abstract

**Title:** Axiomatic Studies of Truth

**Name:** Kentaro Fujimoto

**College:** Merton College

**Degree:** DPhil in Philosophy

**Date of Submission:** December 22, 2010.

**Approximate Number of Words:** 72000 words (both including and excluding citations)

In contemporary formal theory of truth, model-theoretic and non-classical approaches have been dominant. I rather pursue the so-called classical axiomatic approaches toward truth and my dissertation begins by arguing for the classical axiomatic approach and against the others. The classical axiomatic approach inevitably leads to abandonment of the naïve conception of truth and revision of the basic principles of truth derived from that naïve conception such as the full T-schema. In the absence of the general guiding principles based on that naïve conception, we need to conduct tedious but down-to-earth ‘field works’ of various theories of truth by examining and comparing them from various points of view in searching for satisfactory theories of truth. As such attempt, I raise two new criteria for comparison of truth theories, make a proof-theoretic study of them in connection to the foundation of mathematics.

# Contents

<b>1</b>	<b>The Classical Axiomatic Approach</b>	<b>5</b>
1.1	Definitional Approach vs Non-definitional Approach . . . . .	5
1.2	Semantic Approach vs Axiomatic Approach . . . . .	16
1.3	Revision of Logic . . . . .	22
1.3.1	Moh Shaw-Kwei's Paradox . . . . .	26
1.3.2	Hajek-Paris-Shepherdson Paradox . . . . .	35
1.4	The Classical Axiomatic Approach . . . . .	40
<b>2</b>	<b>List of Systems of Truth</b>	<b>43</b>
2.1	Notational Preliminaries . . . . .	43
2.2	Disquotational systems of truth I . . . . .	47
2.3	Systems of Typed Truth . . . . .	50
2.4	Systems of iterative compositional truth I . . . . .	54
2.5	Disquotational systems of Truth II . . . . .	61
2.6	Systems of iterative compositional truth 2 . . . . .	62
2.7	Systems of Determinate Truth . . . . .	65
2.8	Systems of Symmetric Truth . . . . .	68
2.9	Systems of iterative non-compositional truth . . . . .	75

2.10 Schematic Reflective Closures . . . . .	76
<b>3 Relative Truth Definability</b>	<b>78</b>
3.1 Comparing Conceptual Aspects of Axiomatic Truth . . . . .	78
3.2 Relative Truth Definability . . . . .	83
3.2.1 Mathematical Properties of Truth-Definability . . . . .	88
3.3 Digression – Proof-Theoretic Ordinal – . . . . .	90
3.4 Several Results on Truth Definability . . . . .	92
3.5 Some Applications . . . . .	102
3.6 Concluding Remarks of Chapter 3 . . . . .	106
<b>4 The Inner Theory of Truth Systems</b>	<b>109</b>
4.1 The Notion of Inner Theory . . . . .	109
4.2 Systems of Finitely Iterated Self-Applicable Truth . . . . .	114
4.3 Semantics of iterative compositional theories . . . . .	118
4.4 Inner Theory Comparison via Semantics . . . . .	125
4.5 Proof theory for finitely iterated self-applicable truths . . . . .	129
4.5.1 Iterated fixed-points and self-applicable truths . . . . .	130
4.5.2 Semi-Formal Systems for $\widehat{\text{ID}}_n$ . . . . .	132
4.5.3 Asymmetric Interpretation . . . . .	136
4.6 Inner theory comparison via semi-formal systems . . . . .	143
4.6.1 Systems with and without Cons . . . . .	143
4.6.2 Comparison of different self-applicable truths . . . . .	150
4.7 Summary of Chapter 4 . . . . .	151
<b>5 Truth Progressions</b>	<b>153</b>
5.1 Background . . . . .	153

5.2	Notation for Transfinitely Iterated Truths . . . . .	156
5.3	Kripke-Feferman Truth . . . . .	159
5.4	Consistency and Completeness Axioms . . . . .	162
5.5	Positive Uniform T-biconditionals . . . . .	175
5.6	Kripke-Feferman Truth with Weak Kleene Logic . . . . .	177
5.7	Determinate Truth and Feferman Logic . . . . .	184
5.8	Cantini's VF-Truth . . . . .	187
5.8.1	The Systems $\text{VF}_{\prec}$ and $\text{Aut}(\text{VF})$ . . . . .	188
5.8.2	Lower Bound . . . . .	190
5.8.3	Upper Bound . . . . .	195
5.9	Summary and Concluding Remarks of Chapter 5 . . . . .	202

**6 Further Prospects**

**204**

# Chapter 1

## The Classical Axiomatic Approach

In my thesis, I concentrate on so-called formal theories of truth. Among various types of formal theories of truth, I choose one type of theories that I call ‘classical axiomatic’ theories of truth. In this opening chapter, I will classify formal theories of truth into some categories, compare the advantages and disadvantages of them, and then explain why I prefer the classical axiomatic approach to the others.

### 1.1 Definitional Approach vs Non-definitional Approach

Over a long period of time, a huge variety of theoretical analyses of the notion of truth have been attempted. In the contemporary debates on truth, after the development of modern symbolic logic, a distinction is particularly made between philosophical and formal theories of truth. If this distinction is granted, my principal interest of the present thesis lies in the latter, although they bear a quite close and reciprocal connection to each other.

Notoriously, the notion of truth bears a paradoxical aspect. Due to this aspect, we are always faced with the danger of inconsistency in any kind of analyses of truth. The most famous and basic

example is the Liar paradox. One formulation of this paradox is:

The statement (1.1) is false. (1.1)

This example might look unnatural. Alternatively, we may consider the following situation: a stranger comes to you on the street, he mistakes you for somebody else, and then suddenly says to you ‘I never tell you the truth’; but he is a stranger to you and you have never met him before; namely, the above is the first utterance made by him toward you. Can this utterance either true or false? This ‘Liar’ problem is not easily solved at all. Admitting statements neither true nor false cannot be an immediate solution because of the so-called strengthened Liar paradox. Restricting self-reference doesn’t straightforwardly lead us to a solution either because of other paradoxes with no appeal to self-reference such as Yablo’s paradox [86].

This paradoxical phenomenon related to the notion of truth is indeed a serious obstacle for theoretical analyses of this notion. Paul Horwich proposed a distinction between ‘the theory of the concept of truth’ and ‘the theory of truth itself’. According to him, ‘the former purports to specify the conditions in which someone uses the word ‘true’ with a certain meaning ([39, p.252])’ and the latter purports to specify ‘principles about the *property* of truth on the basis of which all the facts about truth are to be explained ([39, p.253]).’ For the former, Horwich insists, the inconsistency arising from the Liar paradox may not be an overwhelming difficulty. McGee also pointed out ‘it is possible to use the word ‘true’ coherently without possessing a coherent theory of truth ([58, p.3])’. It might well be supposed that a certain condition constitutes our practical mastery of the word ‘true’ but leads to inconsistency at the same time. In the present thesis, however, I am principally interested in ‘the theory of truth itself’ in the terminology of Horwich, and thus inconsistency is the worst consequence to be avoided in the study of this type of theories of truth.

In the formal approach to the notion of truth, we facilitate the tools of contemporary formal

logic. One of the principal benefits of the use of formal logic concerns the aforementioned problem of inconsistency. By means of formal logic, we can investigate whether a given system of truth is consistent or not. Formal tools provides us with a framework in which we can measure how safe a given theory is and make it more explicit what may cause a contradiction and how we can avoid it.

Then, what are formal theories of truth like? There are various kinds of such theories. Let us try to classify them into some categories. They first branch out into two categories: definitional or reductionist theories, and non-definitional or non-reductionist theories. The former group aims at giving a definition of truth and thereby explaining away it in terms of other more basic and non-semantic notions. There can be considered two types of definitional approaches. The one tries to define the truth of a given language  $\mathcal{L}$  in the terminology of  $\mathcal{L}$  itself; the other aims at reducing it into some terminology possibly not possessed by  $\mathcal{L}$ . However, both types of definitional approaches are no longer a principal concern of truth theorists.

Only a few contemporary philosophers and logicians would choose the former stronger form of definitional approaches of defining the truth of  $\mathcal{L}$  *within*  $\mathcal{L}$ . This common renunciation comes from the well-known *undefinability theorem* of Tarski. Tarski [82] proved that a definition of truth is impossible under fairly general assumptions. Let  $\mathcal{L}$  be the language whose truth we want to consider. If the truth of  $\mathcal{L}$  be definable within  $\mathcal{L}$ , we have an  $\mathcal{L}$ -predicate  $\Phi$  of  $\mathcal{L}$ -sentences which defines the truth of  $\mathcal{L}$ . Since the truth predicate is a predicate of sentences, we assume  $\mathcal{L}$  contains names of  $\mathcal{L}$ -sentences and let  $\ulcorner \sigma \urcorner$  be the name of a sentence  $\sigma$ . Tarski showed that, under fairly moderate conditions on  $\mathcal{L}$  and the background theory over  $\mathcal{L}$ , there is no predicate  $\Phi$  such that, for all  $\mathcal{L}$ -sentence  $\sigma$ ,

$$\Phi(\ulcorner \sigma \urcorner) \text{ if and only if } \sigma;$$

this schema is called Tarski's *T-schema* and each instance of the T-schema is called *T-biconditional*.

The famous Tarski's convention T claims that

A predicate  $\Phi$  is materially adequate definition of the truth of  $\mathcal{L}$  iff ( $\alpha$ ) the T-schema holds and ( $\beta$ ) if  $\Phi(x)$  then  $x$  is equal to  $\ulcorner \sigma \urcorner$  for some  $\mathcal{L}$ -sentence  $\sigma$ .

Tarski seems to be claiming that these ( $\alpha$ ) and ( $\beta$ ) are *necessary* and *sufficient* conditions for a predicate  $\Phi$  to be a definition of truth. It is still arguable whether they are sufficient on the one hand; but it appears to be a common view among philosophers and logicians that they are necessary. Hence, most of philosophers and logicians would renounce any definitional approach of this stronger form; at least, they do not start their study of truth by assuming that such a definition of truth is possible.

In contrast, the weaker form of the definitional approach will raise little concern either, simply because the standard and commonly accepted definition of truth is already in our hands: i.e., most philosophers and logicians would agree that Tarski's original definition of the truth of a language  $\mathcal{L}$  by means of the mathematical resource outside of  $\mathcal{L}$  is successful and acceptable to a great degree.<sup>1</sup> Though his original setting was based on a slightly old-fashioned theory of types (or calculus of classes), his methodology can be easily carried over to contemporary settings. For example, we can carry out Tarski's construction of the truth of the language  $\mathcal{L}_{\text{PA}}$  of first-order arithmetic within the second-order system  $\text{ATR}_0$  of arithmetical transfinite recursion; we simulate his construction for the language  $\mathcal{L}_\in$  of first-order set theory within the second-order system  $\text{NBG}_\omega$  of  $\omega$  elementary recursion; under the standard translation of the language of arithmetic in that of set theory, the Tarskian definition of the truth of  $\mathcal{L}_\in$  can be given in a fairly weak set theory, say,  $\text{KP}\omega$ ; in these three examples of the Tarskian definition of truth,  $\text{ATR}_0$ ,  $\text{NBG}_\omega$  and  $\text{KP}\omega$  plays the role of the *metatheory* or *metasystem* (in what follows, though it is not standard, I will use the word 'metasystem' for technical reasons), their languages are called *metalanguage* and the languages

---

<sup>1</sup>By 'Tarski's theory', I mean his theory illustrated in the first half of §3 and in §4 of [82], and I exclude another theory of his developed in [82, pp.199–209]. It is the former one that we call 'Tarski's theory of truth' today.

$\mathcal{L}_{PA}$ ,  $\mathcal{L}_\in$  and  $\mathcal{L}_{PA}$  respectively correspond to the object-language. Hence, as for this weaker form of the definitional approach, the problem is not on whether or not definition is possible. Rather the problem is that such a definition bears little significance for our purposes.

**Remark 1.** It is sometimes wrongly thought that Tarski's definition of truth in his [82] just amounts to giving the ordinary classical model-theoretic structure  $\mathfrak{M}$ , which is also introduced and defined by Tarski himself and his student Vaught [83], of the language in question. Though their definitions are technically similar, they should be clearly distinguished.

For the sake of the subsequent arguments, we slightly digress for a while and look more closely at the Tarskian definition. We must start by fixing a metasystem  $Q$  over a certain language  $\mathcal{L}_Q$ . Let  $\mathcal{L}$  be the target language the truth of which is to be defined. In general, there need not be any particular connexion between  $\mathcal{L}_Q$  and  $\mathcal{L}$ ; one may be a proper sublanguage of the other; they may be identical; they may be entirely disjoint; or they may partially overlap. However, it is necessary for the Tarskian definition of the truth of  $\mathcal{L}$  that  $\mathcal{L}$  can be properly translated into  $\mathcal{L}_Q$ . Thereby, the definition of the truth of  $\mathcal{L}$  is carried out within  $Q$  as the definition of the truth of the translation of  $\mathcal{L}$  in  $\mathcal{L}_Q$ . Then, according to Tarski's theorem, only when  $\mathcal{L}$  (or, equivalently, its translation in  $\mathcal{L}_Q$ ) is essentially poorer than  $\mathcal{L}_Q$ , the truth of  $\mathcal{L}$  is definable in  $Q$ . Now, let us recall that Tarski began his arguments by explicitly stating that

The question how a certain concept is to be definable is correctly formulated only if a list is given of the terms by means of which the required definition is to be constructed. If the definition is to fulfil its proper task, the sense of the terms in this list must admit of no doubt. [82, p.152–153]

That is, the sense of  $\mathcal{L}_Q$ -terms must be given and fixed in advance before the definition of truth for  $\mathcal{L}$ .<sup>2</sup> Consequently, each  $\mathcal{L}$ -term has the same meaning as its translation in  $\mathcal{L}_Q$  which is already given by virtue of  $Q$ . Hence, in order to talk about the translatability of  $\mathcal{L}$  in  $\mathcal{L}_Q$  (with the background

---

<sup>2</sup>It is highly debatable how we can provide a term with its meaning. Some might only require some axioms and inference rules; other might require more. But I don't get into this issue and let us assume that the sense of  $\mathcal{L}_Q$ -terms is given by a formal axiomatic system  $Q$ .

system  $\mathcal{Q}$ ), the sense of  $\mathcal{L}$ -terms is also to be given. Now consider the case where we define the truth of the language  $\mathcal{L}_\epsilon$  within a second-order set theory such as Morse-Kelly impredicative class theory  $\text{MK}$  ( $\supset \text{NBG}_\omega$ ) over the language  $\mathcal{L}_2$  ( $\supset \mathcal{L}_\epsilon$ ) of second-order set theory. Then, the sense of each term of the target language  $\mathcal{L}_\epsilon$  is given within  $\text{MK}$ . Let  $\text{ZFA}$  be  $\text{ZF}$  plus Aczel's anti-foundation axiom.  $\text{ZFA}$  is a system over the same language  $\mathcal{L}_\epsilon$  but not a subsystem of  $\text{MK}$ . It is sometimes argued that the conception of sets is different between  $\text{ZF}$  and  $\text{ZFA}$ ; the former is based on, say, the iterative conception of sets, and the latter is based on the so-called graph conception of sets; therefore, we might well say that the sense of ' $\epsilon$ ' is different between  $\text{ZF}$  and  $\text{ZFA}$ . That is to say,  $\mathcal{L}_\epsilon$  with the background system  $\text{ZF}$  and  $\mathcal{L}_\epsilon$  with  $\text{ZFA}$  are different languages with different senses; the mere syntactical coincidence of their vocabulary does not necessarily entail that each term has the same meaning as each other. In the above example, this point is put in the following way:  $\text{MK}$  has defined the truth of the language  $\mathcal{L}_\epsilon$  with some meaning but has not defined the truth of the syntactically identical language  $\mathcal{L}_\epsilon$  with another meaning. In summary, we cannot treat the object-language  $\mathcal{L}$  simply as a syntactical collection of certain vocabulary; rather we need to treat it together with its background system by virtue of which the sense of each  $\mathcal{L}$ -term is given.

Keeping in mind the points developed in the last paragraph, we turn back to the methodological problems concerning the definitional approach. As is well-known, Tarski deduced from his aforementioned results the famous 'impossibility thesis' of truth one formulation of which says:

**(Tars)** The definition of truth of a language  $\mathcal{L}$  is only possible in an essentially richer language  $\mathcal{L}'$  (together with a sufficiently rich background system over  $\mathcal{L}'$ ) than  $\mathcal{L}$ , but impossible in  $\mathcal{L}$  itself.

One consequence of this thesis is crucial: if we stick to the definitional approach and require the truth predicate to be defined away, then no language  $\mathcal{L}$  possesses its own truth predicate and only an 'essentially richer' language  $\mathcal{L}'$  than  $\mathcal{L}$  may possibly contain the truth predicate of  $\mathcal{L}$ . Hence, in particular, if our natural language is counted within the scope of the formal study of truth,

definitional theorists of truth must conclude that our natural language does not possess its own truth predicate, because vocabulary outside or beyond our natural language is no longer a part of our natural language of course and we cannot express them within our natural language. Let us recall Tarski's notion of 'universal language'. According to Tarski, a universal language is a language with 'adequate facilities for expressing everything that can be expressed at all, in any language whatsoever ([81, p.89]).' One of the conclusions Tarski deduced from the above (Tars) and his undefinability theorem is:

**(Univ)** We cannot define the truth of universal languages.

Since a natural language must be universal in the above sense (at least from the viewpoint of the speakers of that natural language), the impossibility of the truth definition for a natural language can be regarded as a corollary of Tarski's thesis (Univ).<sup>3</sup> Hence, a universal language cannot possess its own truth predicate in the definitional approach.

Not only are natural or universal languages to be excluded from the definitional approach. Tarski's construction is carried out within mathematics, and the theory he uses in the construction belongs to mathematics. The language of mathematics (or, mathematics simpliciter) is not universal in Tarski's original sense. However, trivially speaking, mathematics is closed under mathematical objects and reasonings: anything mathematical is subsumed in mathematics. Hence, mathematics constitutes, say, a 'closed totality' and thus we must conclude that we cannot define the truth of any language which includes that of the entire mathematics, unless we compromise to define the mathematical truth by means of non-mathematical tools which (probably) lack sufficient rigour. To

---

<sup>3</sup>The thesis (Univ) is made in [82]. Tarski later gave another similar thesis in [80]. He introduced the notion of a 'semantically closed' language; a language is called semantically closed if it contains 'in addition to its expressions, also the names of these expressions, as well as semantic terms such as the term 'true' referring to sentences of this language [80, p.123]'. Then, he claimed that

**(Sem)** We cannot define truth for semantically closed languages.

Then, it is natural to think that semantical closedness follows from universality. At any rate, the most basic thesis is that if a theory includes a reasonably rich theory of its own syntax, then it cannot define its own truth; then both (Univ) and (Sem) are regarded as a corollary of this thesis.

overcome this problem, one might adopt the so-called ‘orthodox approach’, in which the language of mathematics is conceived as layered into certain hierarchy of sublanguages, and the truth of one sublanguage is possessed by other languages higher than it in that hierarchy. However, it wouldn’t help, because we couldn’t express the statement ‘there is at least one untrue mathematical theorem’; whether or not this statement is indeed true doesn’t matter, but the problem is that we can not even express this statement. If one conceives mathematics as a static closed totality, this consequence should be rejected, unless she decides not to talk about the truth of the entire mathematics at all. It is also to be noted that the same type of ‘orthodox approach’ could be adopted for natural or universal languages but it would be refuted in a similar manner; consider the statements ‘there is no sentence which is both true and false at the same time’; according to the current standpoint in question, this statement cannot bear a ‘universal’ or ‘logical’ character and the term ‘no sentence’ must be relativized to a certain initial fixed proper fragment of our language. In addition, Kripke’s ingenious example on self-applicability [50] illustrates another defect of the ‘orthodox approach’ with respect to natural languages and it seems to me decisive and undefeatable.<sup>4</sup>

One might counterargue that mathematics is open-ended and does not constitute a static ‘closed totality’; then, she might continue, whenever we try to define the truth of mathematics, the notion or realm of mathematics simply gets expanded. Although I myself have sympathy with the idea that mathematics is open-ended, this option shouldn’t be taken for the very problem we are concerned with at any rate, simply because she has already expanded the notion or realm of mathematics in her attempt to define the truth of ‘the former’ mathematics and thus the resulting definition is no longer the definition of the truth of ‘the new’ mathematics. Nonetheless, this view gives us another perspective. If mathematics is open-ended and it never forms a closed totality, then it would be of little importance and significance to try to define the truth of the entire mathematics any longer;

---

<sup>4</sup>To my knowledge, the term ‘orthodox approach’ was first used by Kripke in [50] to denote this type of standpoint particularly for natural languages. It is not so clear if his criticism could have the same effect on the case of the language of mathematics, since self-applicability of the truth predicate might be dispensable for the case of mathematics.

consequently, the impossibility of the definition of the truth of mathematics would be no problem at all for definitional theorists of truth. A parallel arguments for natural language can be made: i.e., there is no totality of a natural language and thus no need to define the truth of the entire natural language.

Indeed, I myself am opt for this view of ‘open-ended’ mathematics. At least for our current knowledge, it seems too much to assume that the closed totality of mathematics is already given in front of us and we have already known exactly what constitutes mathematics. If such a fixed closed totality were given to us, we could concentrate on the fragments of mathematics whose truth is definable within that given totality of mathematics but would have to abandon any discourse concerning the truth of the fragments whose truth is not definable therein; consequently, we would have to anyhow abandon the definitional approach toward some parts of mathematics. Motivated by the ‘open-ended’ view, one may argue, for the definitional approach, that, whether or not  $Q$  is a system provisionally adopted as the representation of the entire mathematics, we could simply expand  $Q$  to a certain sufficiently stronger system  $Q'$  and then give the Tarskian definition of the truth of the language  $\mathcal{L}_Q$  of  $Q$  within  $Q'$ .<sup>5</sup> As a matter of fact, the motivation and background of this approach is resonant with those of the approach I myself will take. However, even if one adopts the above suggested approach, the definitional approach still bears two serious problems and it looks less versatile than my approach.

The first problem of the definitional approach in this approach concerns the self-applicability of truth: that is, we cannot therein deal with the self-applicable character of the notion of truth. This is simply because the truth predicate of a language  $\mathcal{L}$  always belongs to the outside of  $\mathcal{L}$  and thus

---

<sup>5</sup>I would like to emphasize again that Tarski’s definition of truth is fairly satisfactory particularly for mathematical languages or any other languages, such as those of physics or chemistry, which do not intrinsically contain the truth predicate or other semantical predicates and that if the Tarskian definition of truth of a language  $\mathcal{L}$  of our interest can be successfully given within an initially fixed metasystem  $Q$ , this is an almost ideal situation and we need not demand more, with the proviso that  $Q$  is acceptable system. However, we probably want to ask and talk about the truth of any given mathematical language and my interest of the present thesis is principally in the truth *simpliciter* of any language given to us for whatever reason, whether  $Q$  is given as the foundation of mathematics or a subject of a particular mathematics research.

it is meaningless and impossible to apply the truth predicate to a sentence with the truth predicate itself. Notice that this problem applies to any type definitional approach in general. This fact is a crucial defect particularly when we consider the truth of our natural language or any language with reasonably rich semantical terms. Formal theories of truth usually treat only mathematical languages, but if the ultimate goal of such formal study is to analyze our ordinary notion of truth in our natural language and to give an appropriate theory of truth for it, this defect is fatal and should not be ignored by any means.

The second problem is methodological one and comes from the fact that there is no unique procedure of expanding the initial (meta) system  $\mathbf{Q}$ . It is to be observed that the Tarskian definition of truth of  $\mathcal{L}_{\mathbf{Q}}$  within a stronger system  $\mathbf{S}$  requires some redundant conditions on  $\mathbf{S}$  which are not pertinent to the truth definition in their own. For example, as I have mentioned, we can carry out essentially the same construction of the truth predicate of  $\mathcal{L}_{\epsilon}$  in  $\mathbf{NBG}_{\omega}$  (and thus in  $\mathbf{MK}$  *a fortiori*). However, as a matter of fact, we only need  $\mathbf{NBG}$  to define a predicate which merely satisfies the T-schema for  $\mathcal{L}_{\epsilon}$  (cf. [62]); furthermore, within the system  $\mathbf{NBG}$  plus  $\omega$ -induction for  $\Sigma_1^1$ -formulae, we can define a predicate which satisfies all the semantical (and syntactical) properties concerning  $\mathcal{L}_{\epsilon}$  that Tarski used and mentioned as the consequence of his definition in [82]. Then we naturally ask which system we should choose to define the truth of  $\mathcal{L}_{\epsilon}$ . There are at least three options to take:

- (i $^{\circ}$ ) Choose the weakest system among the systems in which we define a predicate satisfying the T-schema or a certain collections of the principles mentioned by Tarski.
- (ii $^{\circ}$ ) Choose the weakest system among the systems in which we can carry out the Tarskian definition.
- (iii $^{\circ}$ ) Choose whatever system as long as it achieves our purpose of carrying out the Tarskian definition of truth or satisfying a certain truth-theoretic principles.

First, the former two are faced with the problem of measuring the ‘strength’ of systems. To make the matter worse, even if we are given some appropriate measure of ‘strength’, it is not necessarily the case that the given measure well-orders all the relevant systems and thereby uniquely determines the ‘weakest’ one; in fact, most of the standard intertheoretic measures do not determine a unique such system. This problem might well be overcome simply by regarding all minimal systems according to a fixed measure as ‘equivalent’ and identifying them.

Second, as for the second option (ii<sup>o</sup>), the phrase ‘can carry out the Tarskian definition’ is fairly vague and it is difficult to give a general criterion to judge whether a given construction of a predicate is ‘Tarskian’ or not.

Third, as I have explained, in the Tarskian definition of truth, we apparently have to add redundantly strong assumption to the original system  $Q$  which has little direct connection to the notion of truth itself. If a certain supersystem  $S$  of  $Q$  is given for some reason independent of the truth definition of  $\mathcal{L}_Q$  (e.g.,  $S$  may be given as the foundation of mathematics) and if we happen to be able to define the truth of  $\mathcal{L}_Q$ , then there would be nothing to complain and our concern would shift to the question whether  $S$  is indeed acceptable for that purpose independently of truth definitions. However, we are currently considering expanding  $Q$  to some  $S$  specially for the sake of the truth definition of  $\mathcal{L}_Q$ . Not to mention Ockham’s razor, since we are now interested only in the truth of  $\mathcal{L}_Q$  and its analysis, we should focus on truth simpliciter and keep addition and expansion as small as possible. This point poses a difficulty for both (ii<sup>o</sup>) and (iii<sup>o</sup>), but it particularly concerns the latter; for, that we are allowed to choose arbitrary such system means that we are allowed to add arbitrary strong assumptions to  $Q$ .

Finally, it is to be noted that the expansion of  $Q$  may cause a change of the sense of  $\mathcal{L}_Q$ -terms. As we have seen, the sense of each  $\mathcal{L}_Q$ -term must be initially fixed by virtue of the background system  $Q$ . In order for us to give a truth definition to  $\mathcal{L}_Q$ , the approach currently considered suggests us to expand  $\mathcal{L}_Q$  to a certain language  $\mathcal{L}$ , adopt some sufficiently rich system  $S$  over  $\mathcal{L}$

and then define the truth of  $\mathcal{L}_Q$  in  $S$ . For this purpose, the language  $\mathcal{L}_Q$  with the fixed meaning provided by virtue of the initial system  $Q$  must be translatable into  $\mathcal{L}$ . The translation of  $\mathcal{L}_Q$  in  $\mathcal{L}$  need not coincide with  $\mathcal{L}_Q$  as a sublanguage of  $\mathcal{L}$  (though they are syntactically identical), but it would not be clear whether there is such a translation if the expansion of  $Q$  to  $S$  brings a drastic change of the sense of  $\mathcal{L}_Q$ -terms. In particular, the options (iii<sup>o</sup>) is implausible from this point of view, since a drastic strengthening of system is very likely to cause a drastic change of meanings.

The approach (i<sup>o</sup>) thus seems to be most plausible among the above three. However, what one is doing under (i<sup>o</sup>) is essentially the same as what one is pursuing under the axiomatic approach, which is my own choice and will be explained in detail in the next section. Under the approach (i<sup>o</sup>), one begins with selecting some principles of truth and then she examines which system can define a predicate satisfying the selected principles. However, the first step is exactly what the axiomatic approach does, and the second step simply amounts to examining whether a certain axiomatic system of truth can be suitably translated (or embedded) in systems in question. For example, the aforementioned fact that NBG defines the predicate satisfying the T-schema for  $\mathcal{L}_\epsilon$  simply means that the system TB (see §2.2) is syntactically embeddable in NBG. Hence, the axiomatic approach methodologically precedes the approach (i<sup>o</sup>). In addition, since axiomatic theories of truth are immune from the aforementioned problem concerning self-applicability of truth, it is much more versatile than the definitional approach of the type in question.

## 1.2 Semantic Approach vs Axiomatic Approach

Now, let us turn to the non-definitional (or non-reductionist) approach. As the attempt to define truth is abandoned in the non-definitional approach, truth must be conceived as a primitive notion of our language which is not further reduced to other notions. This doesn't mean that there is nothing left for the conceptual analysis of truth. We are still interested in the characteristics of

this notion, the role it plays in our linguistic or mathematical activities, its consequences for other notions or other subjects of philosophy or mathematics, and so forth. A non-definitional theory of truth usually amounts to providing a certain formal ‘theory’ or ‘structure’ to a language which is already equipped with its own truth predicate (or possibly more than one truth predicates).

Since we have left Tarski’s original framework, some terminologies of Tarski should be redefined for the sake of our current framework, i.e., the non-definitional theories of truth. In particular, we need to make it clear what Tarski’s notions of meta-language and object-language mean in the context of the non-definitional approach. Originally, Tarski defined them as follows:

[The object language] is the language which is ‘talked about’ and which is the subject matter of the whole discussion; the definition of truth which we are seeking applies to the sentences of this language. [The meta-language] is the language in which we ‘talk about’ the [object-language], and in terms of which we wish, in particular, to construct the definition of truth for the first language. [80, p.125]

In the non-definitional approach, we try to give a certain theory or mathematical structure to a language  $\mathcal{L}$  already equipped with some predicate(s) which are intended to be the truth predicate(s) of  $\mathcal{L}$ . Hence, though we no longer aim at defining the truth of  $\mathcal{L}$ , it is this  $\mathcal{L}$  that is ‘the language which is talked about and which is the subject matter of the whole discussion’, and I call this  $\mathcal{L}$  the object-language in the non-definitional theories of truth on the one hand. On the other hand, the meta-language is the language  $\mathcal{L}_Q$  in which the non-definitional theory in question of truth of the object-language  $\mathcal{L}$  is given and developed by virtue of a certain background system  $Q$  over  $\mathcal{L}_Q$ . It is to be emphasized that, in order to develop a non-definitional theory of truth, we need to work with some background system  $Q$  and I call it the metasystem as in the case of the definitional approach. For example, in the case of Kripke-Feferman theory  $KF$  over Peano arithmetic  $PA$  ([17]), the object language is  $\mathcal{L}_T$ , i.e., the language of arithmetic plus a unary truth predicate  $T$ , and the metasystem is any system (e.g.,  $I\Sigma_1$ ) rich enough to represent the syntax and derivations of primitive recursive systems therein; for another example, in the case of the Kripke’s fixed-point semantics over arithmetic ([50]), the object language is the same as the  $KF$  case and the

metasystem is a sufficiently rich set theory in which Kripke's construction is possible.

There are various different kinds of 'formal theories of truth'. According to Sheard [76], they are divided into two major categories: i.e., the one is called *semantical* theories of truth and the other is *axiomatic* theories of truth. Sheard doesn't give a detailed explanation of what they are; but the section of 'Axiomatic approaches' in [76] (put next to the section of 'Semantical approaches') start with the following passage: '[a]n alternative way of thinking about self-referential truth is to pose the problem as one of creating not models but consistent theories ([76, p.1044])'. Hence, Sheard seems to consider a semantical theory of truth to amount to a construction of a model-theoretic structure of the object-language in question, while an axiomatic theory aims at giving a consistent set of sentences of the object language which contains some desired principles of truths; therefore, we may well call the former 'a model-theoretic approach' alternatively. Still the distinction between them is no less clear than that between semantics and syntax.<sup>6</sup> Nevertheless, fortunately, as far as the so far presented formal theories of truth are concerned, we can distinguish them in most cases; for instance, Kripke's theory is to be classified as semantical on the one hand, and Kripke-Feferman's theory KF is axiomatic on the other hand.

The origin of the axiomatic approach dates back to Tarski. Faced with the thesis (Tars), Tarski himself wrote:

[I]t would be incorrect to infer the impossibility of operating consistently and in agreement with intuition with semantical concepts and especially with the concept of truth. But since one of the possible ways of constructing the scientific foundations of semantics is closed we must look for other methods. The idea naturally suggests itself of setting up semantics as a special deductive science with a system of morphology as its logical sub-structure. For this purpose it would be necessary to introduce into morphology a given

---

<sup>6</sup>First, it is not definitely clear whether giving a 'consistent' set of sentences over a non-recursively axiomatizable systems. For example, the logic of the system PALTr of Hájek et.al. [28] is infinite-valued Lukasiewicz Logic  $L\forall$  which is known to be non-effective (i.e., the set of all 1-tautology is not recursive). However, the motivation of Hájek et.al. is apparently axiomatic; they consider adding some truth-theoretic principles to a given set of arithmetical sentences and show the consistency of the resulting set of sentences with respect to  $L\forall$ . Second, there is a borderline case between syntax and semantics such as Lindenbaum algebra. Third, since Sheard classifies Barwise and Etchemendy's theories of truth (both Russellian and Austinian ones) as semantical theories, it seems that he does not restrict the sense of the word 'model' to the ordinary model-theoretic sense (with a set domain) and he includes theories which are committed to proper class size structures, such as Barwise and Etchemendy's and Tarski's own, into the category of 'semantical theories of truth'.

semantical notion as an undefined concept and to establish its fundamental properties by means of axioms. [82, p.255]

[E]ven respect to formalized languages of infinite order, the consistent and correct use of the concept of truth is rendered possible by including this concept in the systems of primitive concepts of the metalanguage and determining its fundamental properties by means of the axiomatic method (the question whether the theory of truth established in this way contains no contradiction remains for the present undecided). [82, p.266]

From the context in which the second passage is put, it is clear that Tarski mentioned ‘formalized languages of infinite order’ as an example of the languages for which a formally correct and materially adequate definition of true sentences *cannot* be constructed in the metalanguage; therefore, Tarski himself suggests the axiomatic approach as a remedy for the difficulty with which we are currently concerned. Thus, it is the axiomatic approach that Tarski himself suggested as a natural alternative to the definitional approach. The cited passage also explains what the axiomatic approach is. However, it is not clear whether Tarski endorsed the axiomatic approach, and Tarski himself did not explore this direction further. He only gave some negative comments to one specific axiomatic system TB of truth and closed his debates on the axiomatic approach; the system TB will be introduced in the next chapter and I will come back later to this issue of Tarski’s standpoint on the axiomatic approach in the next chapter. I myself prefer the axiomatic approach to the semantical approach and adopt it as the main topic of the present thesis. In the rest of present section, I will argue against the semantical approach.

In contrast to the axiomatic approach –although it is much more popular and dominant among philosophers working on truth today– I am not sure of the origin of the semantical approach. Roughly speaking, the semantical approach starts with a fixed semantical structure of the base language without the truth predicate; then it goes to add semantics of the truth predicate by giving its ‘interpretation’ or ‘extension’ over the given structure. In most typical cases such as Kripke’s theory [50] and Gupta and Belnap’s revision theory [26], the ‘semantical structure of the base language’ means the ordinary classical model-theoretic structure with a certain set domain. One advantage of this approach is that it is immune from the aforementioned problem on possible

changes of meanings in expanding systems; for, we can naturally assume that the meaning of each term is already fixed in the given *semantics* or model-theoretic structure and it will not be affected by the process of giving the interpretation of the truth predicate. Another advantage is that a semantical theory often gives us more information of the concept of truth in question and make it more perspicuous. However, ultimately, it seems crucially problematic, particularly from the foundational point of view, to require that our metalanguage is given such a model-theoretic or semantical structure from the outset; in other words, I think such an approach is inappropriate for the analysis of the notion of truth of the language which we are indeed using or speaking.

Vänäänen [85] gives a suggestive and useful discussion in this respect. In any context of formal logic, we first fix a ‘top-level’ system and start our mathematical inquiries from this ‘top-level’.

Reflecting this fact, Vänäänen points out that:

Formalization of mathematics involves defining a formal language with some intended meaning. Let us call this language *urlogic*. The idea is that *urlogic* is the most primitive formal language we use to study the process of doing mathematics. [85, p.501]

We extracted *urlogic* as a formalization of the act of doing mathematics. The semantics of *urlogic* is totally informal. [85, p.501]

The *urlogic* is the ‘top-level’ formal counterpart of our informal mathematics, and we do not ascend above *urlogic*. Let  $T$  denote this ‘top-level’ system. Then, Vänäänen proposes the following criteria for *urlogic* ([85, p.511]):

- Sentences of *urlogic* are finite strings of symbols. That a string of symbol is a sentence of *urlogic*, is a non-mathematical judgement.
- Some sentences are accepted as axioms. That a sentence is an axiom is a non-mathematical judgement.
- Derivation are made from axioms. The derivations obey certain rules of proof. That a derivation obeys the rules of proof, is a non-mathematical judgement.
- Derived sentences can be asserted as facts.

According to these criteria, for example, infinitary logic and full second order logic cannot be *urlogic*, since they need highly mathematical devices for their formulation; infinitary logic violates

the first condition, and full second order logic violates the second and third. Doing model theory *within* an initially accepted urlogic such as ZF is of course acceptable, but it is very likely to violate these criteria to assume that the metasytem is already equipped with a model-theoretic semantics and the meaning of its language is provided by virtue of it. I will not get into these issues on how metamathematics or the top-level metasytem should be formalized, but I myself strongly agree with Väänänen's argument.

One possible defense of the semantical approach might be this: in a semantical theory of truth, one is constructing a model-theoretic structure  $\mathfrak{A}$  of the language  $\mathcal{L}$  *within* a given metasytem which appropriately satisfies Väänänen's criteria so that the interpreted language  $\mathcal{L}$  by  $\mathfrak{A}$  has the same meaning of the meta-language. First of all, this argument poses a restriction to the variety of the metasytem: i.e., the metasytem has to be a system which can define and operate on model-theoretic structures. Second, it is not clear at all how they can establish this preservation of meanings between the meta-language and the interpreted language. It seems that this preservation at least requires that the structure  $\mathfrak{A}$  should be a model of the meta-system, say, ZF. However, in the ordinary settings, this is impossible due to Gödel's incompleteness and completeness theorems.

Another possible counterargument on behalf of the semantical approach might be this: although many semantical theories of truth operate on model-theoretic structure (as a set-theoretic object), they in fact aim at explaining the procedure of giving semantics to the truth predicate within a suitably expanded metasytem. For example, suppose the initial metasytem is ZF; Then, by expanding ZF to  $\Pi_1^1$ -CA (i.e., NBG plus the comprehension axiom schema for  $\Pi_1^1$  formulae), we can carry out essentially the same construction of Kripke's interpretation of the truth predicate; we can define the least fixed-point of each elementary positive operator and thereby interpret the Kripkean (self-applicable) truth predicate of  $\mathcal{L}_\epsilon$  by means of one such fixed-point. However, this interpretation of what semantical theories of truth do is too charitable and, at any rate, the same problem as I have raised against the definitional approach at the end of the last paragraph applies

to this approach as well.

The arguments against the semantical approach I have raised so far are probably not decisive for rejecting the semantical approach. However, Väänänen's argument, at least, begs semantical (or model-theoretic) theorists of truth for further justifications of their methodology; we pass the turn back to semantical theorists of truth and it is now their turn to defend their approach.

Now we have chosen the non-definitional axiomatic approach. Let  $\mathcal{L}$  be a language equipped with its truth predicate  $\mathcal{L}$ . We need to block paradoxes and avoid contradiction. Feferman [15] pointed out that, in general, the possible solutions to the paradoxes consist in restrictions of (1°) language, (2°) logic or (3°) basic principles. It is slightly vague what he means by the solution type (1°). Feferman classified Tarski's theory of truth as the solution under (1°). Thus, according to my reading, the option (1°) amounts to restrict the formal study of truth to the language which doesn't contain the truth predicate for itself (it might possibly allowed to contain the truth predicate of a proper sublanguage). As I have argued, Tarski's theory already gives us a fairly satisfactory solution under this category (1°), but we are now searching for another type of theory because of the aforementioned disadvantages of his theory. Hence, we will set aside the option (1°) for the discussion of the present thesis. Then, the remained options are (2°) and (3°). I myself would choose (3°) and will give some arguments against (2°) in the next section.

### 1.3 Revision of Logic

In the present section, I will consider the option (2°) of revising logic. As I have declared, the objective of the present section is to provide some arguments against this approach. Briefly speaking, my reason for preferring (3°) and avoiding (2°) is: (a) in order to keep the basic principles of our most naïve conception of truth, we apparently need to make too drastic and probably intolerable revisions on our logic; (b) logic is the most fundamental in our reasoning in every subject matter

and thus revising our logic merely for the sake of blocking paradoxes related to the notion of truth is too much; it is as if Einstein had revised the logic for accommodating the incoherency of Newtonian theory of gravity and special theory of relativity rather than to provided general theory of relativity. The problem caused by the Liar paradox and its variants is more naturally conceived as a problem of our way of conceiving truth rather than an intrinsic problem of our logic.

Of course, it is difficult to give a decisive argument for rejecting any approach under (2°); above all, we haven't yet reached the unanimous conclusion of which logic is 'the' logic. Instead, I will illustrate some of the costs we would necessarily have to pay in revising logic, and try to argue that the costs are too high. In the present section, I will present two fairly general paradoxes applicable to a large range of formal logics; thus, they mainly concern the point (a).

First of all, I pose the following assumptions for the following arguments:

**(Pred)** Truth is expressed by a predicate, i.e., a linguistic device which is applied to singular terms.

**(Synt)** A reasonably rich theory of syntax is available to any system of truth.

These assumptions look fairly reasonable and broadly accepted among truth theorists. Leitgeb [53, p.277] claims that 'there is almost unanimous agreement' for (Pred); for, we require truth to be applicable to quantified variables so that we can thereby properly express such sentences as 'There is at least one axiom of ZF which is true', and 'For all sentences, its negation is true if and only if it is not true'. I agree with his arguments and take the assumption (Pred) granted. Next, the assumption (Synt) partially concerns the problem of the truth-bearer. Although I am opt for syntactical objects, i.e., more specifically, declarative sentence types, it seems that we haven't yet reached a unanimous conclusion about what the bearers of truth are. However, it is almost beyond doubt that, as long as truth is conceived as a predicate, it is of course to be applied to its bearers and thus any system of truth must be equipped with a reasonably rich theory about the truth-bearers. Therefore, if one prefers another kind of truth-bearers, she may simply replace the

assumption (Synt) by a corresponding assumption on a theory of her choice of truth-bearers. Even if she does so, I do not think that her choice would bring any serious difference to my arguments in what follows. What we will need for the subsequent arguments is a certain structural features of truth-bearers which, I would like to claim, will be reasonably assumed for various other truth-bearers without loss of generality; my stance here coincides with the one Halbach [37] takes; I refer the reader to his book [37, pp.17–20] for more careful and detailed discussions.

Under the type of solution (2°), the language contains its own truth predicate and the basic principles of our naïve conception of truth are tried to be preserved as much as possible. In particular, the principal aim of this approach is to provide a formal system of truth in which we have the full T-schema, i.e.,

$$Tr(\ulcorner \sigma \urcorner) \text{ if and only if } \sigma,$$

for arbitrary sentences  $\sigma$  whether or not they may contain the truth predicate  $Tr$  itself. We first focus on this requirement of the full T-schema.

It must be determined how the word ‘if and only if’ should be formalized. First, we need to determine whether it is expressed at the meta-level or object-level.<sup>7</sup> Let  $S$  be any system of truth and let  $\mathcal{L}$  be its language. Then, the next two different formulations of the T-schema can be considered:

$$S \vdash \sigma \text{ if and only if } S \vdash Tr(\ulcorner \sigma \urcorner), \text{ for any sentence } \sigma. \quad (1.2)$$

$$S \vdash \sigma \leftrightarrow Tr(\ulcorner \sigma \urcorner), \text{ for any sentence } \sigma. \quad (1.3)$$

These two are sometimes mixed up and called a single word ‘T-schema’, but we have to clearly distinguish them; for example, when we take  $S = FS$  (the definition will be given in the next

---

<sup>7</sup>We notice that, in Tarski’s theory, this distinction makes no difference, simply because the two formulations are equivalent.

chapter), the former holds while the latter fails. Let us call the former type of equivalence *external T-schema*; i.e., the T-schema expressed by a certain relationship (between sentences  $\sigma$  and  $T\ulcorner\sigma\urcorner$ ) *outside* of the object-language (and within the meta-language). On the other hand, we call the latter type of the equivalence *internal T-schema*; i.e., the T-schema expressed by a certain relationship within the object-language.

Internal formulations of the T-schema has some advantages over external ones. One defect of external formulations consist in the fact that each instance of the T-schema is not a statement of the object-language  $\mathcal{L}_S$ ; it is a statement of the meta-language. Thus, the speakers of  $\mathcal{L}_S$  cannot assert it at all. This situation brings a difficulty particularly to Horwich's minimalist standpoint for instance. Horwich writes:

[T]he deflationist maintains that, since our commitment to these schemata [T-schemata] accounts for everything we do with the truth predicate, we can suppose that they implicitly define it. Our brute acceptance of their instances constitutes our grasp of the notion of truth. [39, p.240]

Thus, the minimalist requirement for the truth predicate is not be fulfilled solely by an external formulation of the T-schema from the viewpoint of the speakers of  $\mathcal{L}_S$ . Hence, an internal formulation seems preferable to an external one particularly from the minimalist standpoint. It seems that this point is disadvantageous not only for minimalists but also for many other standpoints in general. If each instance of the T-schema is not asserted within the object-language, the T-schema is not the characterization (or, say, 'implicit definition' in terms of Horwich) of truth for the speakers of the object language  $\mathcal{L}_S$  and it is observed only by the speakers of the meta-language that  $Tr$  is the truth predicate of  $\mathcal{L}_S$ .

Another defect of external formulations specially concerns the currently questioned approach (2°). Recall that the approach (2°) suggests to revise the logic of the object-language in order to thereby accommodate the paradoxes and preserve the full T-schema. In such an approach, the metalogic (i.e., the logic governing the metasytem) and the logic of the object-language may be

different. This gap might be still acceptable. In fact, even truth-theorists endorsing (2°) usually uses classical mathematics such as classical set theory or classical arithmetic at the meta-level in constructing their theory of truth. The resort to classical mathematics at the meta-level could possibly be justified by arguing that they only revise the logic which concerns truth and its relevant notions and keep the ordinary classical mathematics as a separate subject from their discourse about truth. However, if one adopts an external formulation of the T-schema, the biconditional ‘if and only if’ belongs to the meta-language and follows the law of the metalogic, say, classical logic. Then, this results in a theory of truth which revises the entire logic of the object-language but the target principle of truth follows classical logic; this is a fairly strange consequence and the aforementioned justification for the gap between the meta- and object-logics seems no longer in effect.

Next, we need to determine how the word ‘if and only if’ is formally expressed. In many cases, we work with systems with one and only ‘default’ (or ‘canonical’) logical conditional (and biconditional) and take it for granted that ‘if and only if’ in the T-schema is expressed in terms of it. However, there may be other formulations and indeed there are already such proposals. For instance, Feferman [15] proposed a system in which ‘ $\phi$  iff  $Tr(\ulcorner \phi \urcorner)$ ’ is expressed by a formula  $\phi \equiv Tr(\ulcorner \phi \urcorner)$  with an newly introduced special connective ‘ $\equiv$ ’, which can be interpreted as ‘ $\phi$  and  $Tr(\ulcorner \phi \urcorner)$  have the same truth value (in Strong Kleene Kripkean fixed-point model)’<sup>8</sup>; this system is called Aczel-Feferman system AF.

### 1.3.1 Moh Shaw-Kwei’s Paradox

In the present subsection, I introduce the so-called Moh Shaw-Kwei’s paradox which was found by Moh Shaw-Kwei [75]. The following formulation of mine is a slightly generalized version of the

---

<sup>8</sup>This interpretation is not due to Feferman himself but to Sheard [76]. Feferman himself seems to intend to characterize the connective ‘ $\equiv$ ’ in a purely syntactical manner.

original one.

Although we are now considering the axiomatic approach, I will take a more general setting so as to inclusively treat both axiomatic and semantical theories of truth in what follows. In spite of the division of axiomatic and semantical theories, they have one common feature: namely, they are supposed to specify a certain special class of ‘designated’ sentences among the object-language. When a theory is semantical, such a class can be defined as a class of ‘satisfied’ sentences in a given semantical structure; on the other hand, when a theory is axiomatic, the class of ‘provable’ sentences corresponds to it. This view seems to be shared among those who are working on formal theories of truth.<sup>9</sup> Throughout the present subsection, the word ‘system’ will denote any mathematical structure for a certain language that decides which sentences are designated, regardless of whether it is semantical or axiomatic. This definition of ‘system’ here is extremely general. Set-theoretically, a system  $\mathbf{Q}$  is simply represented by a triple  $\langle \tau, \mathcal{L}_{\mathbf{Q}}, \models_{\mathbf{Q}} \rangle$ , where  $\tau$  is a first-order vocabulary containing a truth predicate  $Tr$ ,  $\mathcal{L}_{\mathbf{Q}}$  is the language of  $\mathbf{Q}$  including the first-order language formed from  $\tau$  (i.e.,  $\mathcal{L}_{\omega\omega}[\tau]$  in abstract model-theoretical terminology), and  $\models_{\mathbf{Q}}$  is a subset of  $\mathcal{L}_{\mathbf{Q}}$  representing the set of designated  $\mathcal{L}_{\mathbf{Q}}$ -sentences. For instance, each Kripkean fixed-point structure  $\mathfrak{A}^*$  for a language  $\mathcal{L}$  is a system in this sense; we say that a sentence  $\phi$  of  $\mathcal{L}$  is designated in  $\mathfrak{A}^*$  when  $\mathfrak{A}^* \models \phi$ . For other examples, in any first-order axiomatic system  $\mathbb{T}$ ,  $\phi$  is designated when  $\mathbb{T} \vdash \phi$ ; in any theory based on  $n$ -valued Lukasiewicz logic,  $\phi$  is designated when the value assigned by the theory to  $\phi$  is 1.

**Definition 1.3.1.** Let  $\mathbf{Q}$  be any system of truth and let  $I$  and  $C$  be definable binary sentential connectives of  $\mathcal{L}_{\mathbf{Q}}$ .  $I$  is said to be *detachable* in  $\mathbf{Q}$ , if the following holds:

if  $I(\phi, \psi)$  and  $\phi$  are designated in  $\mathbf{Q}$ ,  $\psi$  is designated in  $\mathbf{Q}$ .

---

<sup>9</sup>For example, Sheard endorses this view. According to him, ‘formal representation’ of truth refers to ‘any formal framework, either semantic or syntactic, which results in the validation of a set of sentences, including some containing a truth predicate ([61, p.169])’.

$I$  is said to be *transitive* in  $\mathbf{Q}$ , if the following holds:

if  $I(\phi, \psi)$  and  $I(\psi, \theta)$  are designated in  $\mathbf{Q}$ ,  $I(\phi, \theta)$  is designated in  $\mathbf{Q}$ .

Sentences  $\phi$  and  $\psi$  are said to be congruent, when

if  $\theta$  is designated, the formula  $\theta'$  is designated,

for any sentence  $\theta$ , where  $\theta'$  is obtained from  $\theta$  by replacing  $\phi$  (or  $\psi$ , resp.) in arbitrary places by  $\psi$  (or  $\phi$ , resp.).  $C$  is called a congruence in  $\mathbf{Q}$ , if  $\phi$  and  $\psi$  are congruent whenever  $C(\phi, \psi)$  is designated.

**Remark 2.** In Priest's system  $T_0$  ([67]), its (default) conditional is neither detachable nor transitive and its (default) biconditional is not a congruence. In Aczel-Feferman AF, the newly introduced connective ' $\equiv$ ' is not a congruence; in addition, even when  $\phi \equiv \psi$  holds,  $\phi$  ( $\psi$ , resp.) does not necessarily implies  $\psi$  ( $\phi$ , resp.); see [15, p.98].

**Definition 1.3.2.** Let  $I$  be a definable binary sentential connective in a system  $\mathbf{Q}$ . For formulae  $\phi$  and  $\psi$ ,  $I^n(\phi, \psi)$  is defined recursively by

$$I^1(\phi, \psi) \text{ is } I(\phi, \psi)$$

$$I^{k+1}(\phi, \psi) \text{ is } I(\phi, I^k(\phi, \psi)).$$

**Definition 1.3.3.** Let  $I$  be a definable binary connective in a system  $\mathbf{Q}$ . We say that the rule of absorption of order  $n$ , denoted  $(A)_n$ , holds (or admissible) with respect to  $I$  in  $\mathbf{Q}$ , if the following holds:

If  $I^{n+1}(\phi, \psi)$  is designated in  $\mathbf{Q}$ , then  $I^n(\phi, \psi)$  is designated in  $\mathbf{Q}$ .

**Remark 3.**  $(A)_1$  amounts to the so-called Contraction Rule and it holds with respect to the canonical conditional in many familiar logics. For another example,  $(A)_{n-1}$  holds in  $n$ -valued Lukasiewicz logic.

**Definition 1.3.4.** A system  $Q$  is said to be *non-explosive*, if there exists at least one sentence  $\perp \in \mathcal{L}_Q$  which is not designated in  $Q$ . We call such a sentence  $\perp$  a *bottom* of  $Q$ .

It is fairly reasonable to assume that every system we consider is non-explosive.<sup>10</sup> I will assume the following throughout the following:

**(Bot)** Every system has a bottom.

**Definition 1.3.5.** Let  $Tr$  be a predicate of  $\mathcal{L}_Q$ . This predicate  $Tr$  is intended to be the truth predicate of  $\mathcal{L}_Q$ .

(1) Let  $I$  be a binary connective. We say that  $Q$  admits the internal T-schema w.r.t.  $I$ , iff

$$I(\phi, Tr(\ulcorner \phi \urcorner)) \text{ and } I(Tr(\ulcorner \phi \urcorner), \phi) \text{ are designated in } Q \text{ for all } \phi \in \mathcal{L}_Q.$$

(2) We say  $Q$  admits the congruent external T-schema, iff  $\phi$  and  $Tr(\ulcorner \phi \urcorner)$  are congruent for any

$$\phi \in \mathcal{L}_Q.<sup>11</sup>$$

**Definition 1.3.6.** Let  $Q$  be a theory and let  $P$  a unary predicate of  $\mathcal{L}_Q$ .

Let  $I$  be a binary connective of  $\mathcal{L}_Q$ . We say that  $P$  admits internal self-reference with respect to  $I$ , iff there is a sentence  $\lambda \in \mathcal{L}_Q$  such that

$$I(\lambda, P(\ulcorner \lambda \urcorner)) \text{ and } I(P(\ulcorner \lambda \urcorner), \lambda) \text{ are designated.} \tag{1.4}$$

We say that  $P$  admits congruent external self-reference, iff there is a sentence  $\lambda$  such that  $\lambda$  and  $P(\ulcorner \lambda \urcorner)$  are congruent in  $Q$ .

---

<sup>10</sup>Even a system in paraconsistent logic is supposed to have a bottom. Paraconsistent logic is a logic which renounces the law of non-contradiction but at the same time it does not allow systems to be trivial or explosive in the sense that all sentences are designated.

<sup>11</sup>Note that the congruence between  $\phi$  and  $Tr(\ulcorner \phi \urcorner)$  is a relation at the meta level.

**Remark 4.** Here, one specific formulation of self-reference is implicitly assumed: that is, a self-referential sentence  $\sigma$  w.r.t. a predicate  $P$  is such a sentence that  $\sigma$  iff  $P(\ulcorner\sigma\urcorner)$ . There are several reasons to prefer this formulation of self-reference: as Gödel showed, if a theory is provided with its structural descriptives, this form of self-reference is expected to be derived; also, in most cases, either congruent external or internal self-reference of this form follows from another formulation of self-reference and the T-schema. Of course, there are other formulations of self-reference. One natural alternative is to formulate a self-referential sentence  $\sigma$  w.r.t.  $P$  as a sentence  $\ulcorner P(\ulcorner\sigma\urcorner)\urcorner = \ulcorner\sigma\urcorner$ . Another is that of the revision theory [26] via “definitional equivalence”. The third is the so-called situation theoretical interpretation of self-reference.<sup>12</sup>

**Theorem 1.3.7.** Let  $\mathbf{Q}$  be a system with a bottom  $\perp$ ,  $I$  a binary connective in  $\mathbf{Q}$  and  $Tr$  a unary predicate of  $\mathcal{L}_{\mathbf{Q}}$ . Suppose that  $(A)_n$  holds in  $\mathbf{Q}$  for some  $n \in \mathbb{N}$ .

- (i) If  $I$  is transitive and detachable and the predicate  $I^n(Tr(x), \perp)$  admits internal self-reference with respect to  $I$ , then  $\mathbf{Q}$  does not admit the internal T-schema w.r.t.  $I$ .
- (ii) If  $I$  is detachable and the predicate  $I^n(Tr(x), \perp)$  admits internal self-reference w.r.t.  $I$ , then  $\mathbf{Q}$  does not admit the congruent external T-schema.
- (iii) If  $I$  is detachable and  $\mathbf{Q}$  admits the internal T-schema w.r.t.  $I$ , then the predicate  $I^n(Tr(x), \perp)$  does not admit congruent external self-reference.
- (iv) If  $I$  is detachable and  $I(\phi, \phi)$  holds for all  $\phi \in \mathcal{L}_{\mathbf{Q}}$  and if the predicate  $I^n(Tr(x), \perp)$  admits congruent external self-reference, then  $Tr$  does not admit the congruent external T-schema.

*Proof.* (i) Suppose  $Tr$  admits the internal T-schema w.r.t.  $I$ . Since  $I^n(Tr(x), \perp)$  admits internal

---

<sup>12</sup>A typical example of the situation theoretical interpretation can be found in Barwise and Moss’s theory [1, Ch.13]. The intuition behind the Barwise-Moss theory is that the utterance of a paradoxical sentence such as the Liar brings a shift of situation or context, and, as a consequence of this shift, a genuine self-reference fails to be made in such a sentence.

self-reference w.r.t.  $I$ , there exists a sentence  $\lambda$  such that

$$I(I^n(Tr(\ulcorner\lambda\urcorner), \perp), \lambda) \quad \text{and} \quad I(\lambda, I^n(Tr(\ulcorner\lambda\urcorner), \perp)).$$

By the internal T-schema w.r.t.  $I$ , we have

$$I(Tr(\ulcorner\lambda\urcorner), \lambda) \quad \text{and} \quad I(\lambda, Tr(\ulcorner\lambda\urcorner)).$$

Thus, by the transitivity of  $I$  we obtains

$$I(Tr(\ulcorner\lambda\urcorner), I^n(Tr(\ulcorner\lambda\urcorner), \perp)) \quad \text{and} \quad I(I^n(Tr(\ulcorner\lambda\urcorner), \perp), Tr(\ulcorner\lambda\urcorner)).$$

From the first, using  $(A)_n$ , we have  $I^n(Tr(\ulcorner\lambda\urcorner), \perp)$ . Then,  $Tr(\ulcorner\lambda\urcorner)$  is obtained from the second by the detachability of  $I$ . Finally, by applying  $Tr(\ulcorner\lambda\urcorner)$   $n$ -times to  $I^n(Tr(\ulcorner\lambda\urcorner), \perp)$ , we obtain  $\perp$  in  $\mathbb{Q}$ ; a contradiction.

**(ii)** Suppose  $\mathbb{Q}$  admits the congruent external T-schema. Since  $I^n(Tr(x), \perp)$  admits internal self-reference w.r.t.  $I$ , there exists a sentence  $\lambda$  such that

$$I(I^n(Tr(\ulcorner\lambda\urcorner), \perp), \lambda) \quad \text{and} \quad I(\lambda, I^n(Tr(\ulcorner\lambda\urcorner), \perp)).$$

By the congruence of  $Tr(\ulcorner\lambda\urcorner)$  and  $\lambda$ , we obtain that

$$I(I^n(Tr(\ulcorner\lambda\urcorner), \perp), Tr(\ulcorner\lambda\urcorner)) \quad \text{and} \quad I(Tr(\ulcorner\lambda\urcorner), I^n(Tr(\ulcorner\lambda\urcorner), \perp))$$

The rest of proof is the same as **(i)**.

**(iii)** Similar to **(ii)**.

(iv) By the law of identity and congruent external self-reference, we have

$$I(Tr(\ulcorner \lambda \urcorner), \lambda) \quad \text{and} \quad I(\lambda, Tr(\ulcorner \lambda \urcorner)).$$

Then, the congruent external T-schema implies that

$$I(Tr(\ulcorner \lambda \urcorner), I^n(Tr(\ulcorner \lambda \urcorner), \perp)) \quad \text{and} \quad I(I^n(Tr(\ulcorner \lambda \urcorner), \perp), Tr(\ulcorner \lambda \urcorner)).$$

The rest of proof is the same as (i). □

Several well-known logics fall among the victims of Theorem 1.3.7: more precisely, their default conditionals ‘ $\rightarrow$ ’ are detachable and transitive and they admit  $(A)_n$  with respect to ‘ $\rightarrow$ ’. In particular, among such logics are the following:

- Classical logic and intuitionistic logic;
- Finite-valued Lukasiewicz logic and Post logic<sup>13</sup>;
- Weak and strong Kleene Logics and Feferman Logic<sup>14</sup>;
- Infinite-valued and finite-valued Gödel Logic<sup>15</sup>;
- Any substructural logic with Associativity and Weak Contraction, wherever  $I$  is a left-to-right conditional.<sup>16</sup>

Consequently, whenever a system  $S$  of truth follows a logic among them and the T-schema and self-reference are formulated internally with respect to its default conditional, (i)-(iii) of Theorem

---

<sup>13</sup>For the definition and explanation of Post Logic, see [55].

<sup>14</sup>We will explain Feferman logic later in the next chapter.

<sup>15</sup>For the definition and explanation of Gödel Logic, see [27] or [55]. As its name indicates, finite-valued Gödel Logic was introduced by Gödel and used to prove that any finite-valued logic cannot characterize intuitionistic logic. Infinite-valued version is sometimes called Gödel-Dummett Logic.

<sup>16</sup>I follows the Terminology of [72] here. In such a substructural logic, we have, for instance, that if  $X \vdash A \rightarrow (A \rightarrow B)$  then  $X \vdash (A \rightarrow B)$ .

1.3.7 apply to  $S$ . Moreover, since each logic above, except for Weak and Strong Kleene Logics and Post Logic, admits  $\phi \rightarrow \phi$ , (iv) of Theorem 1.3.7 applies to any such system.

The possible solutions for this MSK paradox are:

- (a) To express the T-schema in the way that ‘if and only if’ is neither congruent nor transitive and detachable;
- (b) To deny  $(A_n)$  for any  $n$  with respect to  $I$ ;
- (c) To restrict self-reference.

I will not consider (c) here. In what follows, I will examine the solutions (a) and (b).

Solutions of the type (a) are problematic from the viewpoint of the usage and utility of truth predicate. It is probably reasonable to say that we extract the T-schema as a necessary principle of truth from our ordinary usage and utility of the notion of truth. Hence, if we call a formal framework  $S$  a theory of the concept of  $X$ ,  $S$  must reflect ordinary usage of the notion of  $X$  in particular cases. Reflecting this general requirement for theories of truth, McGee raises the following requirement.

**(O) Ordinary Usage Requirement.** A successful theory of truth ought to agree with ordinary usage about the applicability of ‘true’ in a wide range of particular cases. [58, p.159]

Our everyday life usage of the truth predicate, it seems obvious that we use and apply the T-schema with the combination of the transitivity and detachability of the conditional ‘if ... then ...’. For example, when we infer ‘Ann’s claim is true iff  $P \neq NP$ ’ from ‘Bob conjectures that  $P \neq NP$  is true’ and ‘Ann claims that Bob’s conjecture is true’, we apparently use the T-schema and the transitivity of ‘if ... then ...’ (twice) in this inference; if  $P \neq NP$  we infer ‘Ann’s claim is true’ by appealing to the detachability of ‘if ... then ...’. I dare not bluntly reject the option (a) simply appealing to the thesis (O). As we have mentioned, Feferman’s AF takes this strategy; for another example, the default conditional of Priest’s theory  $T_0$  is neither detachable nor transitive. However, at any rate, any approach with (a) is required to justify its violation of (O).

Let us turn to the option (b). As far as I know, we already have three systems of truth which satisfy everything we have so far demanded except  $(A)_n$ : PALTr of Hájek et.al. ([28]), Priest and Sylvan's  $T_2$  ([67]) and Field's theory ([22]). In these systems, the default conditional ' $\rightarrow$ ' is detachable and transitive, the default biconditional ' $\leftrightarrow$ ' is a congruence,  $\phi \rightarrow \phi$  holds for any  $\phi$ , and both the internal (and congruent external) T-schema and self-reference are met; but, of course,  $(A)_n$  fails for all  $n \in \mathbb{N}$ . Admittedly, the option (b) is the most straightforward solution to the MSK paradox. However, there remain some issues to be examined.<sup>17</sup>

First, we must ask whether the absorption rule  $(A)_n$  is indeed dispensable for theories of truth. Compared to detachability and transitivity, dropping  $(A)_n$  looks much less harmful. However, the next proposition indicates that  $(A)_n$  bears some close connection to the monotonicity of our reasoning.

**Proposition 1.3.8.** Let  $L$  be a finite-valued logic  $L$  (in the standard sense of [55]) and let  $I$  be a definable binary sentential connective. We denote the value of a sentence  $\phi$  by  $\|\phi\|$ . Suppose  $L$  has the following property: for each sentence  $\phi$  and  $\psi$ ,

$$\text{(Monotonicity)} \quad \|\psi\| \leq \|I(\phi, \psi)\|;$$

this condition is naturally interpreted as expressing that the reasoning is monotonous in the sense that the more assumptions we have the more consequences we obtain. Then, for all sentences  $\phi$  and  $\psi$ , there exists  $n \in \mathbb{N}$  such that  $\|I^{n+1}(\phi, \psi)\| = \|I^n(\phi, \psi)\|$ .

Notice that the proof of Theorem 1.3.7 indeed only requires some specific sentences  $\phi$  and  $\psi$  (i.e.,  $Tr(\ulcorner \lambda \urcorner)$  and  $\perp$ ) to meet  $\|I^{n+1}(\phi, \psi)\| = \|I^n(\phi, \psi)\|$ . Hence, if we adopt a finite-valued logic, Theorem 1.3.7 forces us to somehow restrict the monotonicity of our reasoning unless we adopt infinite-valued logic.<sup>18</sup> As a corollary, since the conditions of the last theorem are met in Field's

---

<sup>17</sup>Each theory has its particular problems. For instance, PALTr is not recursively axiomatizable and has no standard model. According to Leitgeb [53], it is required as a desideratum for a theory of truth that the theory should allow for standard interpretations. Thus, PALTr is problematic from this viewpoint.

<sup>18</sup>Going beyond finite-valued logic is not as outrageous as one may expect at a first glance. Non-finite valued logic

theory, it follows that the logic of Field's theory (though it is semantical) is not any finite-valued logic.

Second, systems based on such logics might possibly not be of sufficient use. For instance, it is to be examined whether a sufficient part of arithmetic or physics can be described with such logics. Lastly, even if we block the MSK paradox by abandoning the absorption rule, other paradoxes might arise. As a matter of fact, there *is* another paradox, which arises from some additional assumptions but without any need to resort to  $(A)_n$ ; we will see one such paradox in the next section.

### 1.3.2 Hajek-Paris-Shepherdson Paradox

Tarski raised the full T-schema as the formally correct and material adequate condition to be met by any definition of truth. It is a subtle issue whether he would still maintain this condition even for non-definitional theories of truth; he would probably agree that the derivability of the T-schema is still a desideratum of a non-definitional theory of truth, but he might reject that it is also sufficient. We will come back later to this issue in the next chapter after introducing some axiomatic systems of truth.

At any rate, some further conditions may well be demanded as desiderata of a satisfactory theory of truth. One possible such condition is the following type of 'compositional axiom':

$$Tr(\ulcorner \phi_0 \circ (\phi_1, \dots, \phi_k) \urcorner) \quad \text{if and only if} \quad \circ (Tr(\ulcorner \phi_0 \urcorner), \dots, Tr(\ulcorner \phi_k \urcorner)), \quad (1.5)$$

where  $\circ$  is  $k + 1$ -ary logical connective. Again, we are compelled to choose external or internal formulation of 'if and only if'. In what follows, I will only consider the internal formulation and assume that it is expressed by a detachable and transitive sentential connective; for the sake of readability, I will use  $\rightarrow$  instead of  $I$ . The new desideratum is then formulated as follows:

---

can be effectively axiomatizable of course; for example, it is well-known that intuitionistic logic is not a finite-valued logic.

**(Internal Strong Compositionality)** For each logical connective  $\circ$ ,

$$\forall x [Tr(\circ(\vec{x})) \leftrightarrow \circ(Tr(x_0), \dots, Tr(x_k))],$$

where  $\circ$  is the representation of the syntactical operation  $\circ(\phi_0, \dots, \phi_k)$  such that

$$\circ(\ulcorner \phi_0 \urcorner, \dots, \ulcorner \phi_k \urcorner) = \ulcorner \circ(\phi_0, \dots, \phi_k) \urcorner.$$

Let us abbreviate this condition by ISR. I will not get into a substantial discussion on whether ISR is indeed dispensable for a satisfactory theory of truth. However, there are some arguments in favor of ISR. For instance, it is sometimes argued that the T-schema alone is too weak for a theory of truth. Indeed, as we will see in the next chapter, Tarski himself disregards a system TB simply because it only derives the T-schema and lacks sufficient deductive power; for instance, some generalizations of important logical laws such as the excluded middle and the principle of contradiction are not derivable in TB.<sup>19</sup> In order to compensate this defect, as Halbach argues in [35], one might well strengthen the theory by adding ISR. According to Halbach, ‘the ‘inductive’ clauses have been proven to be natural axioms and all generalizations not provable from them seem to be better left undecided by a good theory of truth’, where by ‘inductive clauses’ he means the compositional truths in the above form (1.5); for another example, Field [21] expresses that ‘I think it is clear that without such general laws the truth [i.e., the inductive clauses] predicate would not serve its main purpose’. Of course, these arguments are basically concerning classical systems and the T-schema restricted to the base language and thus they might not be directly applicable to our current discussion; but, at any rate, ISR is certainly worth serious consideration.

---

<sup>19</sup>In general, Halbach [34] showed that any truth-theoretic equivalents of infinite conjunction, i.e.,  $\forall x(\phi(x) \rightarrow Tr(x))$  where the extension of  $\phi$  is infinite, is not provable in TB. Actually, by slightly modifying his proof, we have a strengthened version of this theorem: i.e., any truth-theoretic equivalents of infinite conjunction is not provable in any system obtained from the base system B by adding a *consistent* T-schema restricted to some fragment of  $\mathcal{L}_B \cup \{Tr\}$ .

The present subsection will reveal that ISR causes another paradox when combined with the internal T-schema and self-reference (w.r.t. ‘ $\rightarrow$ ’). I call it Hájek-Paris-Shepherdson paradox (HPS paradox, for short), because the prototype of this paradox was presented by Hájek, Paris and Shepherdson [28]. This paradox shows a quite different way of producing a contradiction; in particular, it does not need the absorption rule  $(A)_n$ .

Compared to the Liar and the MSK paradox, the HPS paradox is slightly complicated. What I will state as the ‘HPS paradox’ below is just one possible formulation, and there are many other formulations. The following formulation might look ad hoc, but the formulation itself is not important; rather, the background mechanism and reasoning in its yielding a contradiction is important. I will restrict the subsequent arguments to the axiomatic systems of truth, but similar arguments can be applied to semantical theories of truth with some straightforward modifications.

Throughout the following, I will assume that the object (axiomatic) system  $S$  (over its language  $\mathcal{L}_S$ ) includes the first-order quantificational language of arithmetic and Primitive Recursive Arithmetic PRA formulated thereover; i.e., all PRA-theorems are designated in  $S$ . Hence, a primitive recursive function  $F(x, y)$  such that  $F(n, m) = \ulcorner \phi \rightarrow^n \perp \urcorner$  is contained in  $S$ , where  $m$  is the code of  $\phi$  and  $\theta_0 \rightarrow^n \theta_1$  is defined analogously from  $I^n(\theta_0, \theta_1)$ . I will also assume in the following that the object system admits the universal instantiation (UI) with respect to ‘ $\rightarrow$ ’: i.e.,  $\forall x \phi(x) \rightarrow \phi(t)$  for any term  $t$ .

Informally, HPS paradox says that: assuming moderate and natural rules about existential quantification and conditional,

- If ISR holds in  $S$ , then  $S$  is  $\omega$ -inconsistent (in a modified sense).
- If ISR holds in  $S$  and arithmetical induction schema for every formula of  $\mathcal{L}_S$  is derivable in  $S$ , then  $S$  is inconsistent.

Let us see the formal statements of these.

**Theorem 1.3.9.** If the conditions (i)-(v) below hold in  $\mathbf{S}$  and ISR holds in  $\mathbf{S}$ , then  $\mathbf{S}$  is  $\omega$ -inconsistent; i.e., there is a (first-order)  $\mathcal{L}_{\mathbf{S}}$ -formula  $\phi$  such that  $\exists x\phi$  is designated in  $\mathbf{S}$  but  $\phi(\bar{n})$  is designated for no  $n \in \mathbb{N}$ , where  $\bar{n}$  is the numeral for  $n$ .

- (i) If  $\phi \rightarrow \exists x\psi$  is designated and  $x$  does not freely occur in  $\phi$ , then  $\exists x(\phi \rightarrow \psi)$ .
- (ii) If  $t = s$  is designated, then  $t$  and  $s$  are interchangeable (w.r.t. designation) in any formula.
- (iii)  $\phi(t) \rightarrow \exists x\phi$  is designated.
- (iv) If  $\phi \rightarrow \psi$  is designated, then so is  $\exists x\phi \rightarrow \exists x\psi$ .
- (v) If  $\exists x\phi$  is designated and  $x$  is not free in  $\phi$ , then  $\phi$  is designated.

*Proof.* By the internal T-schema, we take a sentence  $\lambda$  such that  $\lambda \rightarrow \exists xTr(F(x, \ulcorner\lambda\urcorner))$  and  $\exists xTr(F(x, \ulcorner\lambda\urcorner)) \rightarrow \lambda$ . Then, we obtain:

- (1)  $Tr(\ulcorner\lambda\urcorner) \rightarrow \exists xTr(F(x, \ulcorner\lambda\urcorner))$  by the internal T-schema and transitivity;
- (2)  $\exists x[Tr(\ulcorner\lambda\urcorner) \rightarrow Tr(F(x, \ulcorner\lambda\urcorner))]$  by (1) and (i);
- (3)  $[Tr(\ulcorner\lambda\urcorner) \rightarrow Tr(F(x, \ulcorner\lambda\urcorner))] \rightarrow [Tr(\ulcorner\lambda\urcorner \rightarrow F(x, \ulcorner\lambda\urcorner))]$  by ISR and (UI);
- (4)  $[Tr(\ulcorner\lambda\urcorner) \rightarrow Tr(F(x, \ulcorner\lambda\urcorner))] \rightarrow [Tr(F(x+1, \ulcorner\lambda\urcorner))]$  by (3) and (ii);
- (5)  $Tr(F(x+1, \ulcorner\lambda\urcorner)) \rightarrow \exists xTr(F(x, \ulcorner\lambda\urcorner))$  by (iii);
- (6)  $[Tr(\ulcorner\lambda\urcorner) \rightarrow Tr(F(x, \ulcorner\lambda\urcorner))] \rightarrow \exists xTr(F(x, \ulcorner\lambda\urcorner))$  by (4) and (5);
- (7)  $\exists x[Tr(\ulcorner\lambda\urcorner) \rightarrow Tr(F(x, \ulcorner\lambda\urcorner))] \rightarrow \exists x\exists xTr(F(x, \ulcorner\lambda\urcorner))$  by (6) and (iv);
- (8)  $\exists x\exists xTr(F(x, \ulcorner\lambda\urcorner))$  by (2) and (7);
- (9)  $\exists xTr(F(x, \ulcorner\lambda\urcorner))$  by (8) and (v);

However, if we have  $Tr(\lambda \rightarrow^n \perp)$  for some  $n \in \mathbb{N}$ , since  $\lambda$  follows from (9) by our choice of  $\lambda$ , we could obtain  $\perp$  by the detachability of ' $\rightarrow$ '. Thus, by (Bot),  $Tr(\lambda \rightarrow^n \perp)$  holds for no  $n$ .  $\square$

**Definition 1.3.10.** The induction rule IND is defined by

IND: For each formula  $\phi(x)$ , if  $\phi(0, \vec{v})$  and  $\forall x(\phi(x) \rightarrow \phi(x+1))$  are designated, then  $\forall x\phi(x, \vec{v})$  is designated.

**Theorem 1.3.11.** Suppose that the conditions (i)-(v) of the last theorem and (vi)-(vii) below are all designated in  $S$ . If ISR and IND are also designated in  $S$ , then  $S$  is inconsistent.

(vi) If  $\phi \rightarrow (\psi \rightarrow \theta)$ , then  $\psi \rightarrow (\phi \rightarrow \theta)$ .

(vii) If  $\phi \rightarrow \psi$ , then  $(\psi \rightarrow \theta) \rightarrow (\phi \rightarrow \theta)$ .

*Proof.*  $Tr(\lambda \rightarrow^0 \perp) \rightarrow \perp$ , i.e.,  $Tr(\perp) \rightarrow \perp$ , immediately follows from the T-schema. Thus, the base step is established.

$$(1) Tr(\ulcorner \lambda \urcorner \rightarrow F(t, \ulcorner \lambda \urcorner)) \rightarrow (Tr(\ulcorner \lambda \urcorner) \rightarrow Tr(F(t, \ulcorner \lambda \urcorner))) \quad \text{by ISR};$$

$$(2) Tr(F(t+1, \ulcorner \lambda \urcorner)) \rightarrow [Tr(\ulcorner \lambda \urcorner) \rightarrow Tr(F(t, \ulcorner \lambda \urcorner))] \quad \text{by (1) and (ii)};$$

$$(3) Tr(\ulcorner \lambda \urcorner) \rightarrow [Tr(F(t+1, \ulcorner \lambda \urcorner)) \rightarrow Tr(F(t, \ulcorner \lambda \urcorner))] \quad \text{by (2) and (vi)};$$

$$(4) Tr(F(t+1, \ulcorner \lambda \urcorner)) \rightarrow Tr(F(t, \ulcorner \lambda \urcorner)) \quad \text{by (3), since } \lambda \text{ is designated by the last theorem};$$

$$(5) [Tr(F(t, \ulcorner \lambda \urcorner)) \rightarrow \perp] \rightarrow [Tr(F(t+1, \ulcorner \lambda \urcorner)) \rightarrow \perp] \quad \text{by (4) and (vii)};$$

Thus, by IND, we obtain

$$\forall x [Tr(F(x, \ulcorner \lambda \urcorner)) \rightarrow \perp]$$

It follows from (iv) that

$$\exists x Tr(\lambda \rightarrow^x \perp) \rightarrow \exists x \perp$$

However, since we have  $\lambda$ ,  $\exists x Tr(\lambda \rightarrow^x \perp)$  holds; a contradiction. □

Note again that, as I have already mentioned, there must be many other sets of conditions for the HPS paradox than (i)-(v) (and (vi)-(vii)). In any case, I think that Theorems 1.3.9 and 1.3.11 together with MSK paradox already well illustrate how very much the revision of logic costs.

## 1.4 The Classical Axiomatic Approach

I have so far argued that the classical axiomatic approach is more versatile and methodologically superior to the other approaches that we have examined. Hence, as I have declared at the beginning, my choice is the classical axiomatic approach toward the conceptual and logical analyses of the notion of truth. In this approach, we examine various principles of truth in the form of axiom or inference rule over a certain axiomatic deductive system.

We have decided not to revise the logic and keep it classical.<sup>20</sup> Hence, this approach inevitably imposes us a revision of our naïve conception of truth and some basic principles of truth derived from this naïve conception, since they together lead to a contradiction. However, as long as we have decided to abandon our very naïve conception of truth, we have lost general guiding principles for constructing systems of truth. There are many possibilities, conditions and options to be considered. Put in this situation, we need, besides traditional kinds of philosophical speculations, to conduct tedious but down-to-earth ‘field work’ in systems of truth by examining and comparing them from various points of view; we may carry out experimental constructions of new systems, look at various systems in connection to the other relevant areas of philosophy and mathematics or investigate what kind of and how rich mathematical consequences would be brought to us by each truth-theoretical means. The rest of the present thesis is devoted to this task.

I will propose two new means to compare more conceptual aspects of theories of truth than

---

<sup>20</sup>I have not given any particular arguments for adopting classical logic, but I have argued that the revision of the logic should be kept as small as possible and we should rather reconsider the basic principles of truth without hastily revising the logic merely for preserving the principles in question. Although some people claim that intuitionistic or other non-classical logics should be adopted as our logic, still most people are working with classical logic and I can see no reason to start by rejecting it.

other so far suggested means. The one is *relative truth definability* discussed in Chapter 3; the other is *comparison by inner theories* in Chapter 4. Before arguing them, in the next chapter, I will introduce the systems of truth to be studied in the present thesis and present basic results of them; most of them are not due to myself but some new results are contained there as well.

It is to be noted that the aforementioned problem on possible changes of meaning in expanding a system may well apply to the axiomatic approach as well. When we formulate an axiomatic truth system over a given base system, it often happens that the resulting system becomes stronger than the base system and brings new non-truth-theoretical consequences (i.e., new theorems of the language of the base system). Some axiomatic truth system over Zermelo's set theory  $Z$  may yield the axiom schema of Replacement; some system over  $ZF$  may derive the existence of an inaccessible cardinal. If such happens, it may be the case that the system is not about a system of truth of the original base language but rather a system of truth of some other language. For this reason, I have certain sympathy with the *conservativity thesis* which claims that truth systems must be conservative over their base system (for the consequences of the base language), simply because conservative systems are most likely to preserve the meaning of the original base language; this conservativity thesis is largely accepted by deflationists of truth and it is sometimes argued that this thesis is essential to deflationism (cf. [74]). However, it is not clear that non-conservative systems always bring about certain changes on the meaning to some unacceptable degree; also, it is not necessarily assumed that adding the truth predicate to the language should never cause any changes on the meanings of the original languages. I will not go into this issue any more in the present thesis and leave it for another occasion. Instead, I mention another import and significance of the axiomatic study of truth. The above debates on conservativity concerns the question how the notion of truth is to be characterized by means of axioms. One may ask another question: what does it mean to accept a certain system as true? There must be some essential gap between one's simply possessing a truth predicate of a language  $\mathcal{L}$  with its proper meaning and

her accepting a mathematical system over  $\mathcal{L}$  as true; the latter should be committed to something more and stronger than the former. I will come back to this aspect of the axiomatic study of truth in Chapter 5, where I discuss and study what I call truth progression.

## Chapter 2

# List of Systems of Truth

In the present chapter, we will introduce various axiomatic systems of truth and survey the basic results about them. Before introducing those systems, we begin with some preliminary definitions.

### 2.1 Notational Preliminaries

We can talk of the truth of any kind of subject. There are various options for base languages and systems over which we formulate axiomatic truth. In the present thesis, I follow the traditional setting and adopt an arithmetical language and system as the base system. We start by fixing an arithmetical system  $\mathbf{B}$  as our base system to which every system of truth is obtained by adding truth predicate(s) and its (their) axioms and inference rules.

In what follows, the base system  $\mathbf{B}$  is assumed to be a first-order classical system of arithmetic. We are principally interested in the case where  $\mathbf{B}$  is Peano Arithmetic  $\mathbf{PA}$ , but we adopt a slightly more general setting particularly in order to deal with the relative truth definability among Feferman's schematic reflective closures (see §10) in the same framework as other systems (since we need to add an extra predicate to the base language in formulating schematic reflective closures); there-

fore, after Chapter 3 where we will discuss relative truth definability, we will restrict our discussion to the cases where  $\mathbf{B}$  is PA in Chapters 4 and 5. The language  $\mathcal{L}_{\text{PA}}$  of PA has equality as its only relation symbol and constant  $\bar{0}$  for zero as its only constant symbol. For some technical reasons, we also assume that  $\mathcal{L}_{\text{PA}}$  contains function symbols for each primitive recursive function and PA contains axioms defining them. Given a language  $\mathcal{L}' \supset \mathcal{L}_{\text{PA}}$ , *full induction for  $\mathcal{L}'$*  is the following schema:

$$\phi(0) \wedge \forall x(\phi(x) \rightarrow \phi(x+1)) \rightarrow \forall x\phi(x), \text{ for each } \mathcal{L}'\text{-formula } \phi.$$

Then, our base system  $\mathbf{B}$  is any extension of PA such that:

- (a) Its language  $\mathcal{L}_{\mathbf{B}}$  comprises  $\mathcal{L}_{\text{PA}}$  and possibly finitely many auxiliary relation symbols.
- (b) Its axioms may include a primitive recursive subschema of full induction for  $\mathcal{L}_{\mathbf{B}}$ .
- (c)  $\mathbf{B}$  may contain at most finitely many axioms other than the specified in (b).
- (d) It is arithmetically sound: namely,  $\mathbf{B}$  is satisfiable by some expansion of the standard model  $\mathbb{N}$ ; in other words, there is a model  $\mathfrak{M}$  of  $\mathbf{B}$  whose  $\mathcal{L}_{\text{PA}}$ -reduct coincides with  $\mathbb{N}$ .

We fix such a base system  $\mathbf{B}$  throughout the present and next chapters and  $\mathcal{L}_0$  denotes the language  $\mathcal{L}_{\mathbf{B}}$  of  $\mathbf{B}$ . All the statements in the following are valid regardless of the choice of the base system  $\mathbf{B}$  unless otherwise stated, subject to the conditions (a)-(d). For a given system  $\mathbf{Q}$ ,  $\mathcal{L}_{\mathbf{Q}}$  denotes the language of  $\mathbf{Q}$ . The language of every system of truth presented in this thesis is an expansion of  $\mathcal{L}_0$  with a primitive recursive set of auxiliary predicate symbols. Hence, for any system  $\mathbf{Q}$  of truth in the present paper, the  $\mathcal{L}_{\mathbf{Q}}$ -terms are the same as the  $\mathcal{L}_{\text{PA}}$ -terms.

In formulating a system  $\mathbf{Q}$  of truth, we have to fix some coding (Gödel numbering) of the language to which the truth predicate is to be applied. First of all, we assume that, for any language  $\mathcal{L}$ , the logical connectives are ‘ $\neg$ ’ (negation), ‘ $\wedge$ ’ (conjunction) and ‘ $\rightarrow$ ’ (conditional) and

the quantifier is ‘ $\forall$ ’ (universal quantification); the disjunction ‘ $\vee$ ’ and the existential quantifier ‘ $\exists$ ’ are defined in the usual manner in terms of ‘ $\neg$ ’, ‘ $\wedge$ ’ and ‘ $\forall$ ’.<sup>1</sup>

Given a language  $\mathcal{L}$ , we assume that each syntactical expression  $e$  from the vocabulary of  $\mathcal{L}$  is (primitive recursively) assigned a fixed Gödel number or ‘code’ denoted by  $\#e$ . Let  $\ulcorner e \urcorner$  denote the term (numeral) of  $\mathcal{L}$  representing  $\#e$ ; for a precise definition of the coding, we refer the readers to [17] for example. In formulating a system of truth, we also need to introduce the formal representations of some relations and operations which enable us to talk about the syntax of  $\mathcal{L}$  within the system in question; it can be seen from our assumptions on the base system  $\mathbf{B}$  and systems of truth that all the relations and functions (except the value function  $\text{val}(x)$ ) introduced below are primitive recursive, and we assume that they are represented by means of their corresponding primitive function symbols in  $\mathcal{L}_0$ .

For a given language  $\mathcal{L}$ , we first introduce some primitive recursive relations for its basic syntactical notions:

$\text{Tm}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of a term of $\mathcal{L}$ ’;
$\text{CT}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of a closed term of $\mathcal{L}$ ’;
$\text{Var}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of a variable of $\mathcal{L}$ ’;
$\text{AtFml}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of an atomic formula of $\mathcal{L}$ ’;
$\text{AtSt}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of an atomic sentence of $\mathcal{L}$ ’;
$\text{Fml}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of a formula of $\mathcal{L}$ ’;
$\text{St}_{\mathcal{L}}(x)$	$:\Leftrightarrow$	‘ $x$ is a code of a sentence of $\mathcal{L}$ ’.

---

<sup>1</sup>In the usual formulation of a classical first-order language, we only need to introduce the negation and conjunction (or disjunction). This is because other connectives can be defined in terms of these two and regarded as mere abbreviations; for example, ‘ $A \rightarrow B$ ’ is defined as ‘ $\neg A \vee B$ ’. However, it is not necessarily the case in other logics such as intuitionistic logic, where ‘ $\rightarrow$ ’ and ‘ $\wedge$ ’ must be introduced as independent of ‘ $\neg$ ’ and ‘ $\vee$ ’. In fact, in the case of the inner logic of Feferman’s new theory DT [19], the conditional ‘ $\rightarrow$ ’ is not definable by ‘ $\neg$ ’ and ‘ $\vee$ ’; namely,  $T(\ulcorner \phi \rightarrow \psi \urcorner)$  and  $T(\ulcorner \neg \phi \vee \psi \urcorner)$  are not equivalent in general. Thus, in order for us to treat DT and other systems altogether in the same framework, we introduce ‘ $\rightarrow$ ’ separately.

We will often suppress the subscript ‘ $\mathcal{L}$ ’ when it is clear from the context. In particular, since  $\mathcal{L}$ -terms and  $\mathcal{L}$ -variables are identical to those of  $\mathcal{L}_{\text{PA}}$ , we always drop the subscript ‘ $\mathcal{L}$ ’ from ‘ $\text{CT}_{\mathcal{L}}$ ’ and ‘ $\text{Var}_{\mathcal{L}}$ ’.

Since the base system extends PA, each language  $\mathcal{L}$  of a system of truth has the (standard) primitive recursive representations of the following syntactical operations:

$$\begin{array}{ll}
\text{nm}(n) = \ulcorner n \urcorner & \ulcorner t \urcorner = \ulcorner s \urcorner = \ulcorner t = s \urcorner \\
\neg \ulcorner \phi \urcorner = \ulcorner \neg \phi \urcorner & \ulcorner \phi \urcorner \vee \ulcorner \psi \urcorner = \ulcorner \phi \vee \psi \urcorner \\
\ulcorner \phi \urcorner \rightarrow \ulcorner \psi \urcorner = \ulcorner \phi \rightarrow \psi \urcorner & \ulcorner \forall v_k \urcorner, \ulcorner \phi \urcorner = \ulcorner \forall v_k \phi \urcorner \\
R(\ulcorner t_1 \urcorner, \dots, \ulcorner t_k \urcorner) = \ulcorner R(t_0, \dots, t_k) \urcorner & \text{sb}(\ulcorner e \urcorner, \ulcorner t \urcorner, \ulcorner v_k \urcorner) = \ulcorner e[t/v_k] \urcorner,
\end{array}$$

where  $n \in \mathbb{N}$ ,  $t$ ,  $s$  and  $t_1, \dots, t_k$  are terms,  $\phi$  and  $\psi$  are formulae,  $v_k$  is the  $k$ -th variable,  $R$  is a predicate symbol of  $\mathcal{L}$  with the arity  $k$ , and  $e[t/v_k]$  is the expression obtained from an expression  $e$  by replacing each (free) occurrence of a variable  $v_k$  by a term  $t$ . We also need the ‘value function’  $\text{val}(x)$  such that  $\text{val}(\ulcorner t \urcorner) = t$  for closed terms  $t$ : more precisely,  $\text{val}(x)$  is a function such that

$$\text{val}(\ulcorner s \urcorner) := F(\text{val}(s_0), \dots, \text{val}(s_k))$$

provably in, say, PA, uniformly for each closed term  $s$  of the form  $Fs_0 \dots s_k$  where  $s_0, \dots, s_k$  are also closed terms and  $F$  is a  $k + 1$ -ary function symbol. This value function  $\text{val}(x)$  can indeed be provably recursive in PA but not primitive recursive, since we could otherwise construct, for each  $n \in \mathbb{N}$ , a *primitive recursive* evaluation function  $\text{ev}_n(e, x_1, \dots, x_n)$  of all the  $n$ -ary primitive recursive functions. We will occasionally write  $\dot{x}$  and  $x^\circ$  respectively for  $\text{nm}(x)$  and  $\text{val}(x)$  to save space.

For the sake of simplicity, we introduce some other abbreviations. First, we write  $\forall a.b$  for  $\ulcorner \forall(a, b) \urcorner$

and  $a(b/c)$  for  $\text{sb}(a, \text{nm}(b), c)$ ; when it is clear from the context, we drop  $c$  and simply write  $a(b)$ . Second, for each formula  $\phi(v)$ , we write  $\ulcorner \phi(\dot{x}) \urcorner$  for  $\text{sb}(\ulcorner \phi(v) \urcorner, \text{nm}(x), \ulcorner v \urcorner)$ , i.e.,  $\ulcorner \phi(v) \urcorner(x/\ulcorner v \urcorner)$ . This definition is extended for multi-variable cases in an obvious manner. Third, we make use of the vector notation; e.g., we write  $\ulcorner \phi(\vec{x}) \urcorner$  for  $\ulcorner \phi(\dot{x}_0, \dots, \dot{x}_n) \urcorner$ . Fourth, given an formula  $\phi(\vec{x})$  where  $\vec{x} = x_0, \dots, x_n$  and  $\vec{b} \in \text{CT}$ , we write  $\text{sb}_{\vec{x}}(\ulcorner \phi \urcorner, \vec{b})$  for

$$\text{sb}(\text{sb}(\dots \text{sb}(\text{sb}(\ulcorner \phi \urcorner, b_0, \ulcorner v_0 \urcorner), b_1, \ulcorner v_1 \urcorner) \dots), b_n, \ulcorner v_n \urcorner).$$

In the rest of the present chapter, we introduce and briefly discuss various systems of truth. We have fixed an arbitrary base system  $\mathbf{B}$ , but when we need to mention a system with a specific base system  $\mathbf{C}$  we will write  $\mathbf{Q}[\mathbf{C}]$  for a system of truth  $\mathbf{Q}$  with a base system  $\mathbf{C}$ : for example,  $\text{TB}[\text{PA}]$  is the system of Tarskian biconditionals (introduced later) with the base system  $\text{PA}$ .

## 2.2 Disquotational systems of truth I

Let us begin with the most naïve and oldest type of systems of truth: i.e., the so-called *disquotational* systems of truth. The background idea of disquotational systems of truth is that truth is to be axiomatized by Tarskian *T-biconditionals*, that is, equivalences of the form

$$\ulcorner \phi \urcorner \text{ is true iff } \phi,$$

which Tarski [82, 80] regarded as the ‘material adequacy’ condition for a definition of truth. As is well-known, if we allow  $\phi$  to range over arbitrary sentences that may contain the truth predicate, the resulting system becomes inconsistent due to the Liar paradox. Hence, any disquotational approach to an axiomatization of truth must impose a certain restriction to the class of sentences over which  $\phi$  above can range; we call such a class in a given system the *domain* of T-biconditionals.

The most straightforward and indeed oldest disquotational system is TB, which dates back to Tarski [82, pp. 255–257], though he himself rejected it because of its too weak deductive and expressive power as we shall discuss later.

**Definition 2.2.1.** TB is a system over  $\mathcal{L}_T = \mathcal{L}_0 \cup \{T\}$  whose axioms consist of B-axioms and the TB-schema:

$$T(\ulcorner \sigma \urcorner) \leftrightarrow \sigma,$$

for each  $T$ -free sentence  $\sigma$ . Namely, the domain of T-biconditionals is restricted to  $T$ -free sentences in TB.

A natural expansion of T-biconditional is to allow parameters in  $\phi$ : i.e., a schema of the form

$$\forall \vec{x} [T\ulcorner \phi(\vec{x}) \urcorner \leftrightarrow \phi(\vec{x})].$$

This form of T-biconditional is called *uniform* T-biconditional.

**Definition 2.2.2.** UTB is a system over  $\mathcal{L}_T$  whose axioms consist of B-axioms and the uniform T-biconditionals for all the  $T$ -free formulae.

The next proposition is well-known.

**Proposition 2.2.3.** TB and UTB are both proof-theoretically equivalent to B for  $\mathcal{L}_0$ .<sup>2</sup>

Here we use Feferman’s notion of proof-theoretic reducibility and equivalence. For their precise definitions, we refer the readers to [16] or [18]; we will also discuss proof-theoretic reducibility in connection with relative truth definability in §3.2. We follow Feferman’s notations for proof-theoretical reducibility and equivalence: i.e.,  $\mathbb{Q}_0 \leq \mathbb{Q}_1 [\Phi]$  means that  $\mathbb{Q}_0$  is proof-theoretically

---

<sup>2</sup>For the proof, see [37] or [48]. As a matter of fact, we can say more. Suppose B contains full induction for  $\mathcal{L}_0$ . Then, even if we add full induction for  $\mathcal{L}_T$  to TB or UTB, the resulting system still remains conservative over B, since we can replace  $T$  by a definable partial truth predicate for a fixed complexity in each given derivation (see [33]). This idea is in fact what Tarski originally presented in [82]. However, since induction for partial truth predicates is not necessarily available in B in our current setting, this proof does not work for Proposition 2.2.3 in general.

reducible into  $Q_1$  for a class of formulae  $\Phi$  which is assumed to be a subset of  $\mathcal{L}_{Q_0} \cap \mathcal{L}_{Q_1}$ ;  $Q_0 \equiv Q_1 [\Phi]$  means that they are proof-theoretically equivalent for  $\Phi$ , i.e.,  $Q_0 \leq Q_1 [\Phi]$  and  $Q_1 \leq Q_0 [\Phi]$ . In the following, when considering proof-theoretical reducibility for the base language  $\mathcal{L}_0$ , we will simply say ‘proof-theoretically reducible’ or ‘proof-theoretically equivalent’ without mentioning  $\mathcal{L}_0$  and write  $Q_0 \leq (\equiv) Q_1$  dropping the part ‘ $[\mathcal{L}_0]$ ’.

From the standard proof of Proposition 2.2.3, we easily obtain:

**Proposition 2.2.4.** A system  $Q$  is said to be reflexive if  $Q$  proves the consistency of every finite subsystem of  $Q$ . Then, if  $B$  is reflexive, so are  $TB$  and  $UTB$ .

Tarski [82, p.257] observed that the principle of contradiction, i.e.,  $\forall x \in \text{St}_{\mathcal{L}_0} \neg [T(x) \wedge T(\neg x)]$ , and similar general principles of truth (in the form of universal statements) are not derivable in  $TB$ ; on the other hand, one can prove  $\neg(T \ulcorner \sigma \urcorner \wedge T \ulcorner \neg \sigma \urcorner)$  for each single standard  $\mathcal{L}_0$ -sentence  $\ulcorner \sigma \urcorner$ . Of course, the principle of contradiction or other general principles can be added to  $TB$  as a new truth-theoretic axiom, but Tarski ‘attach[es] little importance to this procedure’ since ‘it seems that every such enlargement of the axiom system has an accidental character’ ([82, p.258]).

This attitude of Tarski is puzzling and seems to be incoherent with his Convention T that any definition of truth is formally correct and materially adequate whenever the T-schema is derivable from the definition. One immediate concern is that, by taking  $TB$  as the metasystem, we could define the truth of  $\mathcal{L}_0$  in terms of the predicate  $T$  of  $\mathcal{L}_T$  in  $TB$ . As to this concern, Tarski would probably respond that the predicate  $T$  is a semantical term and any definition by means of semantical terms is not taken into consideration for his theory; Tarski declared ‘in this construction [of truth] I shall not make use of any semantical concept if I am not able to previously reduce it to other concepts ([82, p.153])’ and thus Tarski was pursuing a reductive programme of defining away truth in terms of non-semantical terms. Tarski might have thought that any definition of truth in terms of non-semantical terms would yield a sufficient array of truth-theoretic principles. However, unfortunately, Mostowski [62] later provided a natural example of a truth definition in

purely non-semantical terms that is formally correct and materially adequate by Tarski's standards but does not imply general claims such as the principle of contradiction.<sup>3</sup>

Another objection Tarski raised against TB is that TB is not 'categorical'; in other words, TB does not 'unambiguously determine the extension of the symbol  $[T]$  ([82, p.258])'. In view of this postulate 'categoricity', he requires systems of truth to meet the following: 'if we introduce into the [metasystem], alongside this symbol  $Tr$ , another primitive sign, e.g., the symbol ' $Tr'$ ' and set up analogous axioms for it, then the statement ' $Tr = Tr'$ ' must be provable ([82, p.258]).' I will not discuss whether this postulate is necessary, but many of the subsequent systems satisfies this desideratum of Tarski indeed; in most cases, this is derivable due to full-induction.

## 2.3 Systems of Typed Truth

In his seminal paper [82], Tarski provided us with a means to define truth of a formal languages by means of the so-called 'inductive clauses' for truth. As is desired, his definition entails T-biconditionals for the target language. However, this definition cannot be carried out within the target language itself; this is his famous undefinability theorem of truth. These facts motivate the 'orthodox approach' [50] toward truth which leads to the Tarskian hierarchical view of languages and their truth. According to the 'orthodox approach', a language forms a certain hierarchy of its sublanguages in which the truth of one sublanguage is defined in and belongs to another sublanguage at a higher level in the hierarchy. Systems of hierarchical truth are motivated by this view and formulated so as to incorporate Tarski's 'inductive clauses' ('compositional axioms' in my terminology) for the truth of lower languages in their axioms.

**Definition 2.3.1.** The system  $TC^-$  over  $\mathcal{L}_T$  consists of B-axioms and Tarski's 'inductive clauses' for truth:

---

<sup>3</sup>More precisely, Mostowski [62] defined a predicate in NBG which meets the T-schema for  $\mathcal{L}_E$ .

**T1.**  $\forall \vec{x} \in \text{CT}[T(R\vec{x}) \leftrightarrow R\vec{x}^\circ]$ , for each  $\mathcal{L}_0$ -atomic  $R$

**T2.**  $\forall x \in \text{St}_{\mathcal{L}_0}[T(\neg x) \leftrightarrow \neg Tx]$

**T3.**  $(\forall x \in \text{St}_{\mathcal{L}_0})(\forall y \in \text{St}_{\mathcal{L}_0})[T(x \wedge y) \leftrightarrow (Tx \wedge Ty)]$

**T4.**  $\forall x \in \text{St}_{\mathcal{L}_0} \forall y \in \text{St}_{\mathcal{L}_0}[T(x \rightarrow y) \leftrightarrow (Tx \rightarrow Ty)]$

**T5.**  $\forall x \forall y [\forall x.y \in \text{St}_{\mathcal{L}_0} \rightarrow (T(\forall x.y) \leftrightarrow \forall z Ty(z/x))]$ .

The system **TC** is obtained by adding full induction for  $\mathcal{L}_T$  to **TC**<sup>-</sup>.

The acronym **TC**<sup>(-)</sup> stands for Tarski's (inductive) clauses, and the above axioms are a direct axiomatization of Tarski's 'inductive clauses'. For example, the last axiom says that a sentence  $\forall z \phi(z)$  is true iff  $\phi(a)$  is true for each name  $a$  of an object in the domain. Note that the truth predicate characterized in this way is essentially the same thing as the so-called satisfaction class of models of arithmetic.<sup>4</sup>

**Theorem 2.3.2** (Halbach [33]). **TC**<sup>-</sup> is proof-theoretically equivalent to **B**.<sup>5</sup>

We say  $\mathbb{Q}_0$  is conservative over  $\mathbb{Q}_1$  for a set of formulae  $\Phi \subset \mathcal{L}_{\mathbb{Q}_0} \cap \mathcal{L}_{\mathbb{Q}_1}$ , written  $\mathbb{Q}_0 \subset_{\Phi} \mathbb{Q}_1$ , iff every  $\Phi$ -theorem of  $\mathbb{Q}_0$  is derivable in  $\mathbb{Q}_1$  as well. Then it is immediately observed that conservativity for  $\Phi$  follows from proof-theoretic reducibility for  $\Phi$ . Hence, for example, it follows from Proposition 2.2.3 that  $\text{TB}[\text{PA}], \text{UTB}[\text{PA}] \subset_{\mathcal{L}_0} \text{PA}$ : thus, they all have the same truth-free theorems. One more point is to be attended: if  $\mathbb{Q}_0 \subset_{\mathcal{L}_0} \mathbb{Q}_1$  then the consistency of  $\mathbb{Q}_1$  implies that of  $\mathbb{Q}_0$ , written  $\mathbb{Q}_0 \leq_{\text{Con}} \mathbb{Q}_1$ .

On the other hand, **TC** is no longer conservative over **B**; indeed **TC** can prove the consistency of **B** (and thus **TC**<sup>-</sup>). In the case of  $\mathbb{B} = \text{PA}$ , it is known that:

<sup>4</sup>For a detailed exposition of satisfaction classes, see [47].

<sup>5</sup>Historically, Kotlarski, Krajewski and Lachlan [49] first gave a model-theoretic proof that **TC**<sup>-</sup>**[PA]** is conservative over **PA** for  $\Pi_{\infty}^0$ . Then, Halbach provided a finitary proof for this conservativity via cut-elimination and generalized this result to arbitrary base systems.

**Theorem 2.3.3** (Feferman [23]).  $\text{TC}[\text{PA}]$  is proof-theoretically equivalent to  $\text{ACA}$ .

*Proof.* This match-up is established by the standard existence proof of a truth predicate of  $\text{PA}$  in  $\text{ACA}$  and the canonical embedding of second-order quantifications of  $\text{ACA}$  in  $\text{TC}[\text{PA}]$  using the inductive clauses; for more details see [37].  $\square$

Note that since  $\text{TC}$  contains full induction for  $\mathcal{L}_T$  and extends  $\text{PA}$ , it is reflexive (indeed, essentially reflexive) by Lemma III.3.47 of [29]. In addition, we can also see from Halbach’s proof of Theorem 2.3.2 that we can uniformly find a primitive recursive function  $F$  such that any derivation of an  $\mathcal{L}_0$ -sentence from a finite set  $\Delta \subset \text{TC}^-$  can be primitive recursively transformed into a derivation of the same sentence from a finite set  $F(\Delta) \subset \mathbf{B}$  (note that this says more than mere proof-theoretic reducibility). Hence, if  $\mathbf{B}$  is reflexive then so is  $\text{TC}^-$  (but not essentially reflexive).<sup>6</sup>

**Proposition 2.3.4.**  $\text{TC}$  is reflexive, and if  $\mathbf{B}$  is reflexive then so is  $\text{TC}^-$ .

$\text{TC}^-$  and  $\text{TC}$  present how to axiomatize Tarski’s ‘inductive clauses’. Then, we can apply this to ascend up to higher and higher languages in the hierarchy of the ‘orthodox approach’. Recall that the hierarchical view claims that no language can contain its own truth predicate and the truth predicate of a given language  $\mathcal{L}$  is only contained and defined in a richer language  $\mathcal{L}'$  than  $\mathcal{L}$ ; then we iterate this process of ascending to the richer and richer languages and systems from the base system  $\mathbf{B}$  and its language  $\mathcal{L}_0$ . The systems thus obtained form a hierarchy  $\mathcal{L}_0, \mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_m, \dots$ . This hierarchy can go up beyond transfinite levels by taking the union of the lower languages at a limit stage. To formalize this idea, we first fix a standard notation system up to the least critical ordinal  $\Gamma_0$ <sup>7</sup>; we can of course consider higher levels, but  $\Gamma_0$  is sufficient for our subsequent arguments. For the simplicity, we identify ordinals and their representations; e.g.,  $\bar{\alpha}$  denotes the numeral whose value is  $\alpha$  in this notation system; but we often abuse notation and write  $\alpha$  for  $\bar{\alpha}$ .

---

<sup>6</sup>As far as the I know, Proposition 2.3.4 was first explicitly noticed by Martin Fischer (by private communication).

<sup>7</sup>For its precise definition, see [11] or [66].

**Definition 2.3.5.** Let  $\alpha \leq \Gamma_0$  and  $\mathcal{L}_\alpha$  be a language  $\mathcal{L}_0 \cup \{T_\beta\}_{\beta < \alpha}$  where  $T_\beta$  is a new unary predicate expressing the truth of  $\beta$ -th level. Then, the ramified truth system up to  $\alpha$ , written as  $\text{RT}_{<\alpha}$ , is a system over  $\mathcal{L}_\alpha$  consisting of **B**-axioms, full-induction for  $\mathcal{L}_\alpha$  and, for each  $\gamma < \beta < \alpha$ ,

**R1.**  $\forall \vec{x} \in \text{CT}(T_\beta(R\vec{x}) \leftrightarrow R\vec{x}^\circ)$ , for each  $\mathcal{L}_0$ -atomic  $R$ .

**R2.**  $\forall x \in \text{St}_{\bar{\beta}}(T_\beta(\neg x) \leftrightarrow \neg T_\beta(x))$

**R3.**  $\forall x, y \in \text{St}_{\bar{\beta}}[(T_\beta(x \forall y) \leftrightarrow T_\beta(x) \vee T_\beta(y)) \wedge (T_\beta(x \rightarrow y) \leftrightarrow T_\beta(x) \rightarrow T_\beta(y))]$

**R4.**  $\forall z \forall x [\text{St}_{\bar{\beta}}(\forall z.x) \rightarrow (T_\beta(\forall z.x) \leftrightarrow \forall y T_\beta x(y/z))]$

**R5.**  $\forall x \in \text{CT}[\text{St}_{\bar{\gamma}}(\text{val}(x)) \rightarrow (T_\beta(T_\gamma x) \leftrightarrow T_\gamma x^\circ)]$

**R6.**  $\forall x \in \text{CT} \forall \delta < \bar{\beta} [\text{St}_{\bar{\delta}}(\text{val}(x)) \rightarrow (T_\beta(T_\delta x) \leftrightarrow T_\beta x^\circ)]$ ,

where  $\text{St}_a(b)$  is a binary primitive recursive predicate on  $a$  and  $b$  such that  $\text{St}_a(b)$  iff  $\text{St}_{\mathcal{L}_\delta}(b)$  for  $\delta = a$  and  $T_a b$  is a binary primitive recursive function on  $a$  and  $b$  whose value is  $\ulcorner T_\beta t \urcorner$  when  $a = \beta$  and  $b = \ulcorner t \urcorner$ .<sup>8</sup>

Then, the hierarchy of languages is represented by  $\text{RT}_{<\alpha}$  in the following way: a predicate  $T_\beta$  in  $\mathcal{L}_{\beta+1}$  is the truth predicate of  $\mathcal{L}_\beta$  in the sense that  $T_\beta$  satisfies the Tarski's inductive clauses **R1-R5** for  $\mathcal{L}_\beta$ ; **R5** and **R6** express that the higher truth predicates include all the lower truth predicates and thus the higher languages are richer than the lower ones in the hierarchy; **R6** also indicates that  $T_\lambda$  for a limit  $\lambda$  is the union of all the lower truth predicates.

Suppose  $\alpha < \beta$ . Then,  $\text{RT}_{<\beta}$  can prove the consistency of  $\text{RT}_{<\alpha}$ , since  $\text{RT}_{<\beta} \vdash \forall x \in \text{St}_\alpha[\text{Prv}_{\text{RT}_{<\alpha}}(x) \rightarrow T_\alpha(x)]$ . Hence, we have:

**Proposition 2.3.6.** For  $\alpha < \beta$ ,  $\text{RT}_{<\alpha} \subsetneq_{\mathcal{L}_0} \text{RT}_{<\beta}$ .

---

<sup>8</sup>We follow the formulation of  $\text{RT}_{<\alpha}$  by Halbach in [37].

Let  $RA_{<\alpha}$  denote the ramified analysis up to  $\alpha$ .<sup>9</sup> Then the ramified truth system and ramified analysis are related in the following way<sup>10</sup>:

**Theorem 2.3.7** (Feferman [9]).  $RT_{<\alpha}[[PA]] \equiv RA_{<\alpha}$  for each ordinal  $\alpha$ .

For  $\alpha < \Gamma_0$  and a formula  $\phi(x)$  of a language  $\mathcal{L}$ , let  $TI(\alpha; \phi)$  denote

$$\forall x(\forall y < x \phi(y) \rightarrow \phi(x)) \rightarrow (\forall x \leq \alpha) \phi(x),$$

and  $TI_{\mathcal{L}}(\alpha)$  denote the schema  $TI(\alpha; \phi)$  for all  $\phi \in \mathcal{L}$ . Then, for  $\alpha \leq \Gamma_0$ ,  $TI_{\mathcal{L}}(< \alpha)$  denotes transfinite induction up to  $\alpha$ : i.e.,  $TI_{\mathcal{L}}(\beta)$  for all  $\beta < \alpha$ . It is known that, for an epsilon number  $\alpha$  (i.e.,  $\alpha = \omega^\alpha$ ),  $RA_{<\alpha}$  proves transfinite induction up to  $\varphi_\alpha 0$  for arithmetical formulae, i.e.,  $TI_{\mathcal{L}_{PA}}(< \varphi_\alpha 0)$ , where  $\varphi$  is the binary Veblen function. It follows from the equivalence between  $RT_{<\alpha}[[PA]]$  and  $RA_{<\alpha}$  that  $RT_{<\alpha}[[PA]]$  proves the same transfinite induction (cf. [11, 17]). Then, the proof of  $RT_{<\alpha}[[PA]] \vdash TI_{\mathcal{L}_{PA}}(< \varphi_\alpha 0)$  induced from the proof of  $RA_{<\alpha} \vdash TI_{\mathcal{L}_{PA}}(< \varphi_\alpha 0)$  can be straightforwardly transformed to the proof of  $RT_{<\alpha}[[B]] \vdash TI_{\mathcal{L}_B}(< \varphi_\alpha 0)$  for an arbitrary base system  $B$  which meets the conditions we have posed. Thus, we have the next theorem.

**Theorem 2.3.8.** Let  $\alpha$  be an epsilon number. Then,  $RT_{<\alpha} \vdash TI_{\mathcal{L}_0}(< \varphi_\alpha 0)$ .

## 2.4 Systems of iterative compositional truth I

The iterative compositional conception of truth is best represented by the following passage of Feferman [17, p.18]: ‘truth or falsity is *grounded* in atomic facts from the base language  $\mathcal{L}$ , i.e., can be determined from such facts by evaluation according to the rules of truth for the connectives and quantifiers, and where statements of the form  $[T(\ulcorner A \urcorner)]$  are evaluated to be true (false) only when  $A$  itself has already been verified (falsified).’ It is ‘iterative’ because a sentence  $T(\ulcorner A \urcorner)$  is true only

<sup>9</sup>For the precise definition, see [10] or [17].

<sup>10</sup>Feferman gave a sketch of its proof in [9]. A more detailed proof can be found in [32].

if  $A$  is already true, and it is ‘compositional’ because a compositional sentence is determined only by its components according to the evaluation rule for connectives and quantifiers. Kripke [50] first focused on this iterative compositional conception of truth and gave a model-theoretic formulation of such truth, and then Feferman [17] later gave an axiomatic formalization  $\mathbf{KF}$  of this conception.<sup>11</sup> For simplicity, we write  $F(x)$  for  $T(\neg x)$  (its intended meaning is ‘false’) in the following.

**Definition 2.4.1.** The system  $\mathbf{KF}^-$  over  $\mathcal{L}_T$  consists of  $\mathbf{B}$  and:

$$\mathbf{K0} \quad \forall z \in \text{Var} \forall y \in \text{CT} \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow (T\text{sb}(x, y, z) \leftrightarrow T\text{sb}(x, \text{nm}(y^\circ), z))]$$

**K1** For each atomic  $R$  of  $\mathcal{L}_0$ ,

$$\forall \vec{x} \in \text{CT} (T(R\vec{x}) \leftrightarrow R\vec{x}^\circ) \wedge \forall \vec{x} \in \text{CT} (F(R\vec{x}) \leftrightarrow \neg R\vec{x}^\circ)$$

$$\mathbf{K2} \quad \forall x \in \text{CT} [T(Tx) \leftrightarrow Tx^\circ] \wedge \forall x \in \text{CT} [F(Tx) \leftrightarrow Fx^\circ]$$

$$\mathbf{K3} \quad \forall x \in \text{St}_{\mathcal{L}_T} [T(\neg\neg x) \leftrightarrow Tx]$$

$$\mathbf{K4} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [T(x \wedge y) \leftrightarrow (Tx \wedge Ty)]$$

$$\mathbf{K5} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [F(x \wedge y) \leftrightarrow (Fx \vee Fy)]$$

$$\mathbf{K6} \quad \forall z \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow (T(\forall z.x) \leftrightarrow \forall y T x(y/z))]$$

$$\mathbf{K7} \quad \forall z \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow (F(\forall z.x) \leftrightarrow \exists y F x(y/z))]$$

$$\mathbf{K8} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} \rightarrow [(T(x \rightarrow y) \leftrightarrow T(\neg x \vee y)) \wedge (F(x \rightarrow y) \leftrightarrow F(\neg x \vee y))]^{12}$$

Then,  $\mathbf{KF}$  is defined as  $\mathbf{KF}^-$  plus full induction for  $\mathcal{L}_T$ .

The axioms of  $\mathbf{KF}$  is regarded as a partial axiomatization of (closed off) Kripkean fixed-point models with Strong Kleene schema. As is known, Strong Kleene Logic is the three valued logic with the following truth table:

<sup>11</sup>The system  $\mathbf{KF}$  and the main results about it were in fact presented by Feferman in 1979. For more historical issues, see [30].

<sup>12</sup>This axiom is not usually postulated in the literature, but we need it since we have assumed that ‘ $\rightarrow$ ’ is a primitive logical connective and since ‘ $\rightarrow$ ’ is definable by ‘ $\neg$ ’ and ‘ $\vee$ ’ in Strong Kleene Logic.

$\neg$		$\wedge$	T	U	F
T	F	T	T	U	F
U	U	U	U	U	F
F	T	F	F	F	F

The conditional ‘ $\rightarrow$ ’ and conjunction ‘ $\wedge$ ’ are defined in the usual manner in terms of ‘ $\neg$ ’ and ‘ $\vee$ ’. The existential and universal quantifiers are defined as infinitary disjunction and conjunction respectively.

Now it is clear how **KF** represents the iterative compositional conception of truth with the Strong Kleene schema. For example, **K5** represents how falsity is assigned to a disjunctive sentence according to the Strong Kleene schema: that is, a disjunctive sentence  $\phi \vee \psi$  is false iff both conjuncts are false. For another example, **K2** represents the iterative character of truth in which a sentence  $T(\phi)$  is true (false) iff  $\phi$  is true (false). The axiom **K0** is postulated to assure that the truth of a given sentence  $\phi(t)$  depends only on the value of  $t$ .<sup>13</sup>

The proof-theoretic strength of **KF** and **KF<sup>-</sup>** are known:

**Theorem 2.4.2** (Feferman [17]).  $\text{KF}[\text{PA}] \equiv \text{RA}_{<\varepsilon_0}$ . Thus  $\text{KF}[\text{PA}]$  is also equivalent to, e.g.,  $\Sigma_1^1\text{-AC}$  and  $\widehat{\text{ID}}_1$ .<sup>14</sup>

It is known that  $\text{KF}^-\llbracket\text{PA}\rrbracket \equiv \text{PA}$ . To my knowledge, this equivalence was first observed by Cantini [5, Theorem 9.10]. Cantini showed that  $\text{KF}_{\text{tot}}$ , an extension of  $\text{KF}^-\llbracket\text{PA}\rrbracket$ , is reducible into  $\text{PA}$ . His proof does apply for any base system  $\mathbf{B}$  which meets the condition we have already put in the last subsection. We show an alternative proof of the equivalence  $\text{KF}^-\llbracket\text{PA}\rrbracket \equiv \text{PA}$ . The proof is a straightforward modification of Halbach’s conservativity proof of  $\text{TC}^-\llbracket\mathbf{B}\rrbracket$  over  $\text{PA}\llbracket\mathbf{B}\rrbracket$ .

**Theorem 2.4.3.**  $\text{KF}^- \equiv \mathbf{B}$ .

<sup>13</sup>**K0** was dropped in Feferman’s original formulation of **KF**, since it is derivable from the other axioms in **KF**. However, it is not derivable in **KF<sup>-</sup>** without full induction for  $\mathcal{L}_T$  and thus introduced by Cantini for **KF<sup>-</sup>**. It is sometimes called the regularity axiom.

<sup>14</sup>For the precise definitions of  $\Sigma_1^1\text{-AC}$  and  $\widehat{\text{ID}}_1$ , see [14].

*Proof.* We first reformulate  $\mathbf{KF}^-$  in a sequent calculus à la Tait. The axioms for the truth predicates (i.e., **K1-K8**) are reformulated as the following inference rules. For a technical reason, we drop **K0** and the following in fact shows that  $\mathbf{KF}^- - \mathbf{K0} \leq \mathbf{B}$ .<sup>15</sup>

$$\begin{array}{c}
\frac{\Gamma, R(\bar{s}^\circ)}{\Gamma, \neg\text{CT}(\bar{s}), s \neq R\bar{s}, Tt} \mathbf{K1}^+ \qquad \frac{\Gamma, \neg R(\bar{s}^\circ)}{\Gamma, \neg\text{CT}(\bar{s}), s \neq R\bar{s}, \neg Tt} \mathbf{K1}^- \\
\\
\frac{\Gamma, Ts^\circ}{\Gamma, \neg\text{CT}(s), t \neq Ts, Tt} \mathbf{K2-1}^+ \qquad \frac{\Gamma, \neg Ts^\circ}{\Gamma, \neg\text{CT}(s), t \neq Ts, \neg T(t)} \mathbf{K2-1}^- \\
\\
\frac{\Gamma, Fs^\circ}{\Gamma, \neg\text{CT}(s), t \neq Ts, Ft} \mathbf{K2-2}^+ \qquad \frac{\Gamma, \neg Fs^\circ}{\Gamma, \neg\text{CT}(s), t \neq Ts, \neg F(t)} \mathbf{K2-2}^- \\
\\
\frac{\Gamma, Ts}{\Gamma, \neg\text{St}_\beta(t), t \neq \neg\neg s, Tt} \mathbf{K3}^+ \qquad \frac{\Gamma, \neg Ts}{\Gamma, \neg\text{St}_\beta(t), t \neq \neg\neg s, \neg Tt} \mathbf{K3}^- \\
\\
\frac{\Gamma, Ts \wedge Tt}{\Gamma, \neg\text{St}_\beta(r), r \neq s \wedge t, Tr} \mathbf{K4}^+ \qquad \frac{\Gamma, \neg Ts \vee \neg Tt}{\Gamma, \neg\text{St}_\beta(r), r \neq s \wedge t, \neg Tr} \mathbf{K4}^- \\
\\
\frac{\Gamma, Fs \vee Ft}{\Gamma, \neg\text{St}_\beta(r), r \neq s \wedge t, Fr} \mathbf{K5}^+ \qquad \frac{\Gamma, \neg Fs \wedge \neg Ft}{\Gamma, \neg\text{St}_\beta(r), r \neq s \wedge t, \neg Fr} \mathbf{K5}^- \\
\\
\frac{\Gamma, \forall y Ts(y)}{\Gamma, \neg\text{St}_\beta(r), r \neq \forall t.s, Tr} \mathbf{K6}^+ \qquad \frac{\Gamma, \neg \forall y Ts(y)}{\Gamma, \neg\text{St}_\beta(r), r \neq \forall t.s, \neg Tr} \mathbf{K6}^- \\
\\
\frac{\Gamma, \exists y Fs(y)}{\Gamma, \neg\text{St}_\beta(r), r \neq \forall t.s, Fr} \mathbf{K7}^+ \qquad \frac{\Gamma, \neg \exists y Ts(y)}{\Gamma, \neg\text{St}_\beta(r), r \neq \forall t.s, \neg Fr} \mathbf{K7}^-
\end{array}$$

for any terms  $s$  and  $t$ . Let us call these  $T$ -rules. The logical rules and identity rules are as usual. Each derivation is represented by a finite tree.

We will show that, given a derivation  $\mathcal{D}$ , we can effectively construct a derivation  $\mathcal{D}'$  with the same conclusion which contains applications of cut only to truth-free formulae. The main induction is on the number  $n$  of cut to formulae with truth predicates. Suppose we have shown the claim until  $n = m + 1$ . Then we can focus on a subderivation  $\mathcal{D}_0$  of  $\mathcal{D}$  which contains only one application of cut to formulae with truth predicates and ends with it. If we can eliminate this cut from  $\mathcal{D}_0$ ,

<sup>15</sup>Dropping **K0** causes no loss of generality for the current purpose since we have  $\mathbf{KF}^- \leq \mathbf{KF}^- - \mathbf{K0}$ . This can be shown by interpreting  $T(\ulcorner\phi(\bar{x})\urcorner)$  by  $T(\text{sb}(\ulcorner\phi\urcorner, \text{nm} \circ \text{val}(\bar{x}), \bar{x}))$  in  $\mathbf{KF}^- - \mathbf{K0}$ .

our claim for  $n$  follows from IH.

Then, we define the *T-complexity* of this derivation  $\mathcal{D}_0$ . The *T-complexity* of a formula occurring in a sequent of  $\mathcal{D}_0$  is the numbers of *T*-rules used to obtain it: more precisely, the *T-complexity* of a truth-free formula is always 0; the first occurrence of a formula  $\phi$  (thus introduced by identity or logical axioms) containing truth predicates is assigned *T-complexity* 1; the *T-complexities* of side formulae of inferences are not changed at all; if  $\phi$  is obtained by a *T*-rule from  $\psi$  then the *T-complexity* of  $\phi$  is that of  $\psi$  plus 1; if  $\phi$  is obtained by a logical rules from  $\psi_i$  ( $i \in I$ ) then the *T-complexity* of  $\phi$  is the maximum of the *T-complexities* of  $\psi_i$ 's. For more detailed explanation, see [33]. Then, the *T-complexity* of  $\mathcal{D}_0$  is defined to be the sum of *T-complexity* of the critical formula  $\phi$  and  $\neg\phi$  of the last inference (i.e., cut) of  $\mathcal{D}_0$ . We will show that  $\mathcal{D}_0$  can be transformed into a derivation  $\mathcal{D}_1$  with the same conclusion but with less *T-complexity*; our claim is obtained thus by induction on *T-complexity* of  $\mathcal{D}_0$ .

Given  $\mathcal{D}_0$ , in the usual manner, we can transform  $\mathcal{D}_0$  without changing its *T-complexity* into an derivation  $\mathcal{D}'_0$  in which cut is only applied to atomic formulae (i.e., the cut-rank is 1) and in which the critical formulae of cut are introduced just in the premise of the critical sequents (i.e., the sum of the left-rank and the right-rank is 2 in the terminology of [79]). Hence, we can assume without loss of generality that the last inference of  $\mathcal{D}_0$  is of the form:

$$\frac{\frac{\frac{\vdots}{\Gamma', \Phi}}{\Gamma, Tt} \quad \frac{\frac{\vdots}{\Gamma'', \Psi}}{\Gamma, \neg Tt}}{\Gamma} \quad (2.1)$$

For the base step, suppose  $\Gamma, Tt$  is introduced by an axiom. We only consider the essential case in which it is obtained by the identity axiom. This can be shown by cases according to how  $\Gamma, \neg Tt$  is obtained, and we show only one case in which it is obtained by **K2·2<sup>-</sup>**. Then, the last inference

is of the form:

$$\frac{\Gamma_0, s \neq t, \neg Ts, Tt \quad \frac{\frac{\vdots}{\Gamma_1, \neg Tr^\circ}}{\Gamma_1, \neg CT(r), t \neq Tr, \neg Tt}}{\Delta}}$$

Let  $\Gamma'_1$  be  $\Gamma_1 \cup \{t \neq Tr\}$ . Then, for a suitable  $\Gamma'_0$ , we can derive

$$\frac{\Gamma'_0, s \neq t, Tr \neq t, Tr = s \quad \frac{\frac{\Gamma'_1, \neg Tr^\circ}{\Gamma'_1, \neg CT(r), Tr \neq s, \neg Ts}}{\Delta}}{\Delta}}$$

by applying cut to  $Tr = r$  and  $Tr \neq r$ . This derivation contains no cut to formulae with truth predicates.

For the induction step, suppose  $\Gamma, Tt$  and  $\Gamma, \neg Tt$  are obtained non-dual  $T$ -rules. Then, due to the unique readability,  $\Gamma$  can be obtained without cut to any formula containing truth predicates. For example, if they are obtained by  $\mathbf{K3}^+$  and  $\mathbf{K4}^-$ , then  $\Gamma$  contains  $t \neq \neg\neg t_0$  and  $t \neq t_1 \wedge t_2$  for some terms  $t_0, t_1, t_2$  but the sequence  $t \neq \neg\neg t_0, t \neq t_1 \wedge t_2$  is derivable in  $\mathbf{B}$  by the unique readability. Assume  $\Gamma$  is obtained by some dual  $T$ -rules. Then, it must be of the form

$$\frac{\frac{\frac{\vdots}{\Gamma', \Phi}}{\Gamma, Tt} \quad \frac{\frac{\vdots}{\Gamma', \neg\Phi}}{\Gamma, \neg Tt}}{\Gamma}}$$

which is transformed into a derivation with less  $T$ -complexity by applying cut to  $\Phi$  and  $\neg\Phi$ . We have completed the proof.  $\square$

As we have mentioned, Cantini showed that a supersystem  $\mathbf{KF}_{\text{tot}}$  of  $\mathbf{KF}^-$  is proof-theoretically reducible into PA. However,  $\mathbf{KF}_{\text{tot}}$  is not reflexive and indeed his proof does not directly yield the

reflexivity of  $\text{KF}^- \llbracket \text{PA} \rrbracket$ , since the bound of complexity of the reduced proof in PA depends on the length of the original proof in  $\text{KF}_{\text{tot}}$ . In contrast, it is observed that the above proof of  $\text{KF}^- \equiv \text{B}$  yields the reflexivity result in the same manner as Proposition 2.3.4.

**Proposition 2.4.4.**  $\text{KF}$  is (essentially) reflexive. If  $\text{B}$  is reflexive, then so is  $\text{KF}^-$ .

In the context of systems of iterative compositional truth, we sometimes consider the two special axioms: the consistency axiom and completeness axiom which are defined respectively by

$$\text{Cons} : \forall x \in \text{St}_{\mathcal{L}_T} \neg(Tx \wedge Fx) \qquad \text{Comp} : \forall x \in \text{St}_{\mathcal{L}_T} (Tx \vee Fx).$$

Feferman [17] argued that the character of Cons is different from the rules of evaluation incorporated in  $\text{KF}$  and his argument should apply to Comp as well.<sup>16</sup> In fact, neither of them is derivable from  $\text{KF}$  and adding both to  $\text{KF}$  leads to a contradiction.<sup>17</sup> On the other hand, Cantini [5] showed that either (but not both) of them can be consistently and even conservatively (for  $\mathcal{L}_0$ ) added to  $\text{KF}$ . We will see in §2.5 and §5 that Cons is particularly important in the connection between iterative compositional systems and determinate systems of truth.

Concluding this subsection, we have a closer look at iterativity and compositionality. We can now extract from the above observations more general characteristics of the iterative compositional systems of truth. That is, they may be characterized in the following way: for each logical connectives  $C$ , its axioms for compositional truth are of the form

$$\forall \vec{x} [\text{St}_{\mathcal{L}_T}(x_0) \wedge \dots \wedge \dots, \text{St}_{\mathcal{L}_T}(x_k) \rightarrow (T(C\vec{x}) \leftrightarrow R_C(\vec{x}))] \tag{2.2}$$

for a certain formula  $R_C$  suitable for a given evaluation rule for  $C$ ; the axioms for quantifiers are of

<sup>16</sup>Feferman denotes Cons by *Disj* there. The notations Cons and Comp are due to Cantini [5], but, to my knowledge, these two axioms first appears in [23].

<sup>17</sup>Given a least Kripkean fixed-point model  $M$  with Strong Kleene schema, if  $M$  is closed off in the way that all the undetermined sentences are pushed into the *extension* of  $T$ , then the resulting classical structure is a model of  $\text{KF} + \text{Comp}$ ; if  $M$  is closed-off in the usual manner then it is a model of  $\text{KF} + \text{Cons}$ .

the parallel forms for each quantifier  $Q$  and its evaluation formula  $R_Q$ ; the axioms for iterativity are of the form like  $T(TA) \leftrightarrow TA$  as in KF. Although it is highly debatable what kinds of evaluation formulae  $R_C$  and  $R_Q$  can be regarded as ‘compositional’, any system for iterative compositional truth should be axiomatized in this way. Now, let  $S$  be an iterative compositional system of truth with the full induction, and suppose that  $R_C$  and  $R_Q$  are  $T$ -positive for each connective  $C$  and quantifier  $Q$ . Then, we can show that the truth of  $S$  can be constructed as a fixed-point of a positive arithmetical operator. Therefore, when the base system  $B$  is PA,  $S$  is proof-theoretically reducible to  $\Sigma_1^1$ -AC via Aczel’s construction of a  $\Sigma_1^1$  fixed-point of a  $\Sigma_1^1$ -formula (see [17]).

## 2.5 Disquotational systems of Truth II

The domains of T-biconditionals are restricted to arithmetical sentences in TB and UTB, but one can consider a disquotational system of truth whose domain of T-biconditionals is larger and may contain sentences with  $T$ . However, it has turned out to be difficult to specify such a domain in a natural and consistent way among uncountably many possibilities.<sup>18</sup> As far as the I know, only one such expansion of the domain has so far been presented; the schema PUTB is the uniform T-biconditionals with the domain the set of sentences in which  $T$  occurs only positively. The acronym PUTB stands for *positive uniform T-biconditionals*.<sup>19</sup>

The schema PUTB properly extends UTB, but if we add PUTB to  $B$ , it is known that the resulting system is still conservative over  $B$ . On the other hand, if we add the full induction for  $\mathcal{L}_T$  together with PUTB, the resulting system is no longer conservative; we call this system PUTB.

---

<sup>18</sup>McGee [59] showed that there are uncountably many different and incompatible maximal consistent sets of T-biconditionals; to make the matter worse, Cieśliński [8] further showed that there are uncountably many different and incompatible maximal sets of T-biconditionals that are even *conservative over the base system*. It is also observed by Halbach (by private communication) that any system can be equivalently reformulated in a purely disquotational form by means of McGee’s trick [59]. Thus the strength of purely disquotational systems can range over all the possible ones.

<sup>19</sup>As far as I know, the PUTB schema was first considered by Cantini [5], and Halbach [36] used it as a formal system of truth. The appellation PUTB is due to Halbach.

**Definition 2.5.1.** Given a  $\mathcal{L}_T$ -formula  $\phi$ ,  $\phi$  is said to be  $T$ -positive if every occurrence of  $T$  in  $\phi$  is positive. The system PUTB over  $\mathcal{L}_T$  consists of B-axioms, full induction for  $\mathcal{L}_T$  and the schema PUTB: for each  $T$ -positive formula  $\phi(\vec{x})$ ,

$$\forall \vec{x}[T(\ulcorner \phi(\vec{x}) \urcorner) \leftrightarrow \phi(\vec{x})].$$

Then,  $\text{PUTB}^-$  denotes the system B plus the PUTB-schema only (without full induction for  $\mathcal{L}_T$ );  $\text{PUTB}^-$  is proof-theoretically equivalent to B as was noted.

We shall see later that PUTB is fairly strong system.

It is also known that PUTB is a subsystem of KF; more precisely, we have:

**Lemma 2.5.2** (Cantini [5]).  $\text{KF}^-$  proves the PUTB schema. Hence  $\text{PUTB}^-$  (or PUTB) is a subsystem of  $\text{KF}^-$  (KF, respectively).

## 2.6 Systems of iterative compositional truth 2

KF was an axiomatic formulation of Kripke's fixed-point definition of self-applicable truth with the Strong Kleene schema. It is then natural to consider 'KF-style' axiomatic formulation of Kripke's fixed-point truth definition *with other schemata*. As was suggested by Feferman [17], one natural and straightforward alternative is the Weak Kleene schema, which is also a three-valued evaluation schema named after Kleene. Its evaluation for the logical connectives is given as follows:

$\neg$		$\wedge$	T	U	F
T	F	T	T	U	F
U	U	U	U	U	U
F	T	F	F	U	F

The universal quantifier is interpreted as infinitary (Weak Kleene) conjunction: i.e.,  $\forall x\phi(x)$  is true iff  $\phi(x)$  is true for all  $x$ ;  $\forall x\phi$  is false iff  $\phi(x)$  is false for some  $x$  and  $\phi(x)$  is either true or false for

all  $x$ . The other logical primitives are defined in the usual manner from them.

Then, the system **WKF**, ‘KF-like’ system of self applicable truth with the Weak Kleene schema, is defined as follows.

**Definition 2.6.1.**  $\text{WKF}^-$  consists of all the axioms of PA, **K0–K4**, **K6**, **K8** and the following axioms (instead of **K5** and **K7**):

$$\mathbf{WK5} \text{ St}_{\mathcal{L}_T}(x) \wedge \text{St}_{\mathcal{L}_T} \rightarrow \left[ F(x \wedge y) \leftrightarrow \left( (Fx \wedge Fy) \vee (Fx \wedge Ty) \vee (Tx \wedge Fy) \right) \right]$$

$$\mathbf{WK7} \text{ Var}(z) \wedge \text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow \left[ F(\forall z.x) \leftrightarrow \left( \forall y [Tx(y/z) \vee Fx(y/z)] \wedge \exists y Fx(y/z) \right) \right]$$

Then **WKF** plus full induction for  $\mathcal{L}_T$ .

These axioms represent how compound sentences are evaluated according to the Weak Kleene schema in the form of (2.2) in §2.4; for example, **WK5** says that a compound sentence  $A \wedge B$  is evaluated to be false when both  $A$  and  $B$  has a definite truth value and either of them is false.

Feferman [19] introduced a new type of system **DT** of self-applicable truth. In my [20], I pointed out that his **DT** is in fact a ‘disguised’ **KF**-like system with a special evaluation schema. I call this schema *Feferman schema* and the logic based on this schema *Feferman Logic*. I will introduce **DT** itself in the next section, but I consider that the **KF**-like system with Feferman schema below in the present section; the connection between it and **DT** will be explained later in the next section.

Feferman schema evaluates the connectives and quantifiers in the same way as the Weak Kleene one only except for the conditional  $\rightarrow$ . The evaluation rule for the conditional ‘ $\rightarrow$ ’ is given by the following:

$\rightarrow$	T	U	F
T	T	U	F
U	U	U	U
F	T	T	T

Then, the ‘KF-like’ system **FKF** with Feferman schema is defined as follows.

**Definition 2.6.2.**  $\text{FKF}^-$  consists of all the axioms of PA,  $\text{WKF}^-$  minus **K8**, and the following:

$$\mathbf{FK9} \quad \text{St}_{\mathcal{L}}(x) \wedge \text{St}_{\mathcal{L}}(y) \rightarrow [T(x \rightarrow y) \leftrightarrow ((Tx \wedge Ty) \vee Fx)]$$

$$\mathbf{FK10} \quad \text{St}_{\mathcal{L}}(x) \wedge \text{St}_{\mathcal{L}}(y) \rightarrow [F(x \rightarrow y) \leftrightarrow (Tx \wedge Fy)].$$

We will briefly review some basic results of  $\text{KF}^{(-)}$ ,  $\text{FKF}^{(-)}$ ,  $\text{WKF}^{(-)}$  and  $\text{PUTB}^{(-)}$  below.

**Proposition 2.6.3.**  $\text{KF}^-$ ,  $\text{FKF}^-$ ,  $\text{WKF}^-$  and  $\text{PUTB}^-$  prove the following schemata:

- (i)  $\forall \vec{x} [T(\ulcorner \phi(\vec{x}) \urcorner) \leftrightarrow \phi(\vec{x})]$ , if  $\phi$  contains no occurrence of  $T$ .
- (ii)  $\text{KF}^- \vdash D^+(\ulcorner \phi(\vec{x}) \urcorner) \rightarrow [T\ulcorner \phi(\vec{x}) \urcorner \leftrightarrow \phi(\vec{x})]$ , where  $D^+(x)$  denotes  $(Tx \vee Fx) \wedge \neg(Tx \wedge Fx)$ ; note that  $D^+(x)$  is equivalent to  $F(x) \leftrightarrow \neg T(x)$ .
- (iii)  $\text{Cons} \rightarrow \forall \vec{x} [T(\ulcorner \phi(\vec{x}) \urcorner) \rightarrow \phi(\vec{x})]$ , for an arbitrary formula  $\phi$ .
- (iv)  $\text{Comp} \rightarrow \forall \vec{x} [\phi(\vec{x}) \rightarrow T(\ulcorner \phi(\vec{x}) \urcorner)]$ , for an arbitrary formula  $\phi$ .

**Remark 5.** As was seen in Lemma 2.5.2, the schema  $\text{PUTB}$  is derivable in  $\text{KF}^-$ . However, neither  $\text{FKF}$  nor  $\text{WKF}$  derives it. Let  $\lambda$  denote the Liar sentence: the sentence such that  $\ulcorner F\ulcorner \lambda \urcorner \urcorner = \ulcorner \lambda \urcorner$  provably, say, in PA. Then,  $\text{FKF}(\text{WKF}) \vdash 0 = 0 \vee \lambda$ , but  $\text{FKF}(\text{WKF}) \not\vdash T(\ulcorner 0 = 0 \vee \lambda \urcorner)$ , as will be shown later in Chapter 4.

**Proposition 2.6.4.** Let  $\lambda$  be the above liar sentence and  $\text{Q}$  be either  $\text{WKF}^-$ ,  $\text{FKF}^-$ ,  $\text{KF}^-$  or  $\text{PUTB}^-$ . Then,  $\text{Q} + \text{Cons} \vdash \neg \lambda$  and  $\text{Q} + \text{Comp} \vdash \lambda$ .

Define a translation  $^c$  as follows:  $\phi^c = \phi$  when  $\phi$  is  $\mathcal{L}_0$ -atomic,  $\phi^c = \neg F(t)$  when  $\phi$  is  $Tt$ , and  $^c$  commutes with the connectives and quantifiers.

**Proposition 2.6.5** (Cantini [5]). If  $\text{KF}^{(-)} \vdash \phi$  then  $\text{KF}^{(-)} \vdash \phi^c$ . As a consequence,  $\text{KF}^{(-)} + \text{Cons} \vdash \phi$  implies  $\text{KF}^{(-)} + \text{Comp} \vdash \phi^c$  and *vice versa*.

*Proof.* This is proved by showing that  $\text{KF}^- \vdash \sigma^c$  for each  $\text{KF}^-$ -axiom  $\sigma$ . □

By contrast, as we shall see in Chapter 4, this duality fails for  $\text{WKF}^{(-)}$ ,  $\text{FKF}^{(-)}$  and  $\text{PUTB}^{(-)}$ .

## 2.7 Systems of Determinate Truth

Feferman's system DT [19] brought about a new form of axiomatic systems of truth. The background idea of DT, which Feferman ascribes to Russell, is that every predicate  $R$  has a domain  $D$  of significance and it makes sense to apply  $R$  only to objects in  $D$  and therefore the principles which are supposed to characterize the concept expressed by  $R$  are to be applied only to objects in  $D$ . Thus the system of truth based on this view would have two predicates  $T$  and  $D$  representing truth and its domain of significance. According to Feferman, such a domain of significance in the case of truth consists of *meaningful and determinate* sentences, i.e., sentences with a definite truth value, true or false, and the T-schema and other principles characterizing truth can be restricted to such sentences. Feferman's argument for this identification of the domain with the set of determinate sentences is simple. Let  $D$  be the domain of truth predicate  $T$ . Then, we have  $[T(A) \vee F(A)] \rightarrow D(A)$ .<sup>20</sup> Conversely, since the T-biconditional,  $T(A) \leftrightarrow A$  and  $F(A) \leftrightarrow \neg A$ , obtains for every  $A$  in  $D$ , we have  $D(A) \rightarrow [T(A) \vee F(A) \leftrightarrow (A \vee \neg A)]$ ; thus, we have  $D(A) \rightarrow [T(A) \vee F(A)]$ . This argument is intuitively convincing and Feferman just identifies  $D(x)$  with  $T(x) \vee F(x)$ ; thus the new predicate  $D$  becomes redundant.<sup>21</sup> In addition, Feferman further requires that  $D$  must be *strongly compositional* in the sense that a compound sentence is in  $D$  iff all the substitution instances of its subformulae by meaningful terms belong to  $D$ .<sup>22</sup> Feferman presented a system DT to capture these views on truth.

**Definition 2.7.1.**  $\text{DT}^-$  is a system over  $\mathcal{L}_T$  whose axioms consist of B-axioms, **K0** and:

$$\mathbf{D1} \quad \forall x[\text{AtSent}_{\mathcal{L}_0}(x) \rightarrow D(x)]$$

<sup>20</sup>To be precise, in this argument, Feferman implicitly assumes that  $D(\neg A) \rightarrow D(A)$ .

<sup>21</sup>Formally speaking, there is a subtle matter about this argument. The argument to conclude  $D(A) \leftrightarrow [T(A) \vee F(A)]$  for each (standard) sentence  $A$  does not formally entail  $D(x) \leftrightarrow [T(x) \vee F(x)]$  for all  $x$ . This is because any (recursive) arithmetical system ( $\supset \text{PA}$ ) must have a nonstandard model  $M$  and thus nonstandard sentences in  $M$  by the usual overspill argument. Hence, when one introduces a predicate  $D$  separately from  $T$ , the equivalence  $D(x) \leftrightarrow T(x) \vee F(x)$  is not necessarily *provable*. It may be of some interest to consider and investigate a DT-like system in which  $D$  and  $T$  are separately introduced.

<sup>22</sup>As Feferman himself pointed out, strong compositionality in this sense is not met with respect to the conditional ' $\rightarrow$ ' due to **D11**. If we pose full strong compositionality, the resulting system becomes identical to the system  $\text{WKF} + \text{Cons}$  which will be introduced later. See §5.

$$\mathbf{D2} \quad \forall x \in \text{CT} [D(Tx) \leftrightarrow D(\text{val}(x))]$$

$$\mathbf{D3} \quad \forall x \in \text{St}_{\mathcal{L}_T} [D(\neg x) \leftrightarrow D(x)]$$

$$\mathbf{D4} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [D(x \wedge y) \leftrightarrow (Dx \wedge Dy)]$$

$$\mathbf{D5} \quad \forall z \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow (D(\forall z.x) \leftrightarrow \forall y Dx(y/z))]$$

$$\mathbf{D6} \quad \forall x \in \text{CT} [T(R\vec{x}) \leftrightarrow R(\text{val}(\vec{x}))], \text{ for each atomic } R \text{ of } \mathcal{L}_0$$

$$\mathbf{D7} \quad \forall x \in \text{CT} [D(\text{val}(x)) \rightarrow (T(Tx) \leftrightarrow T(\text{val}(x)))]$$

$$\mathbf{D8} \quad \forall x \in \text{St}_{\mathcal{L}_T} [D(x) \rightarrow (T(\neg x) \leftrightarrow \neg Tx)]$$

$$\mathbf{D9} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [D(x \vee y) \rightarrow (T(x \wedge y) \leftrightarrow (Tx \wedge Ty))]$$

$$\mathbf{D10} \quad \forall z \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \wedge D(\forall z.x) \rightarrow (T(\forall z.x) \leftrightarrow \forall y Tx(y/z))]$$

$$\mathbf{D11} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [D(x \rightarrow y) \leftrightarrow (Dx \wedge (Tx \rightarrow Dy))]$$

$$\mathbf{D12} \quad \forall x, y \in \text{St}_{\mathcal{L}_T} [D(x \rightarrow y) \rightarrow (T(x \rightarrow y) \leftrightarrow (Tx \rightarrow Ty))],$$

where  $D(x)$  abbreviates  $Tx \vee Fx$ . Then, DT is  $\text{DT}^-$  plus full induction for  $\mathcal{L}_T$ .

As is desired, it is shown by straightforward meta-induction on  $\phi(\vec{x})$  that

$$\text{DT}^- \vdash \forall \vec{x} (D(\ulcorner \phi(\vec{x}) \urcorner) \rightarrow (T(\ulcorner \phi(\vec{x}) \urcorner) \leftrightarrow \phi(\vec{x}))).$$

DT-axioms look different from KF-axioms *prima facie* and they have distinct background motivations. However, DT and KF are in fact very closely related. In order to see this, we first consider the ‘inner logic’ behind DT.

In his consistency proof of DT, Feferman constructed a Kripkean fixed-point model for DT in which the truth of sentences are evaluated according to a special logic. This logic is what I called Feferman Logic in the last section, and I have presented the ‘KF-like’ system FKF there.

Then the next theorem indicates that determinate truth and iterative compositional truth bear a certain very close connection.

**Theorem 2.7.2.** DT (DT<sup>-</sup>) and FKF+Cons (FKF<sup>-</sup>+Cons) are identical systems.

*Proof.* The proof is straightforwardly done by showing all the axioms of one are provable in the other. We will only illustrate two typical examples. For  $x, y \in \text{St}_{\mathcal{L}_T}$ , if  $Tx \wedge Fx$  then  $Dx$  and thus  $Fx \leftrightarrow \neg Tx$  in DT; thus  $\text{DT} \vdash \text{Cons}$ . For another example, in FKF+Cons, given any  $x, y \in \text{St}_{\mathcal{L}}$ , if  $T(x \rightarrow y)$  then  $(Tx \wedge Ty) \vee Fx$  by **FK9**. But it follows from Cons that  $(Tx \wedge Ty) \vee \neg Tx$ ; we have  $\neg Tx \vee Ty$ . Conversely, if  $\neg Tx \vee Ty$  and  $D(x \rightarrow y)$ , then  $\neg Tx \vee (Tx \wedge Ty)$  and  $Dx$ . Since  $\neg Tx$  implies  $Fx$  by  $Dx$ , we have  $(Tx \wedge Ty) \vee Fx$ . Finally,  $\text{FKF} + \text{Cons} \vdash \mathbf{D11}$ .  $\square$

We can also relate WKF to certain ‘DT-like’ system of determinate truth. As was noted in fn.22, DT does not meet strong compositionality. Although there is a good reason for Feferman’s restricting strong compositionality on the conditional ‘ $\rightarrow$ ’, it is then natural to consider a determinate system of truth with full strong compositionality. Such a system is obtained by replacing the axiom **D11** by the following:

$$\text{St}_{\mathcal{L}_T}(x) \wedge \text{St}_{\mathcal{L}_T} \rightarrow [D(x \rightarrow y) \leftrightarrow (Dx \wedge Dy)].$$

Then, we can show in a parallel manner to the last theorem that the resulting system is identical with WKF + Cons.

Feferman [19] raised a conjecture about the proof-theoretic strength of DT.

**Feferman’s First Conjecture:**  $\text{DT}[\text{PA}] \equiv \text{RA}_{<\varepsilon_0}$ .

Feferman also suggested its proof. As far as I see, his suggested proof does work for the lower bound of  $\text{DT}[\text{PA}]$ , but not for the upper bound as it is. He suggested to simulate the construction of Kripkean fixed-point model of  $\text{DT}[\text{PA}]$ , which is given in his consistency proof of  $\text{DT}[\text{PA}]$ , within

$\Sigma_1^1$ -AC, as in the proof for the upper bound of  $\text{KF}[\text{PA}]$  in [17]. However, as we have seen, every fixed-point model of DT must be consistent in the sense that  $Ta \wedge Fa$  for no  $a$ , and  $\Sigma_1^1$ -AC only assures the existence of a fixed-point of an  $\Sigma_1^1$  positive operator but not necessarily consistent one. In the next chapter, we will give a solution to this obstacle and positive answer to the conjecture together with his second conjecture explained below as one application of relative truth definability.

## 2.8 Systems of Symmetric Truth

Systems like KF and DT are called type-free systems of truth, since they allow us to prove the truth of sentences which contain the truth predicate itself. However, any extension Q of KF or DT do not possess the following *symmetry* property of truth:

$$Q \vdash \sigma \quad \text{if and only if} \quad Q \vdash T^\ulcorner \sigma \urcorner;$$

in my terminology of the last chapter, this symmetry property corresponds to the external T-schema (w.r.t. the default conditional at the meta-level); we will introduce the notions of ‘inner theory’ and ‘outer theory’ in Chapter 4, and the symmetry is equivalent to the coincidence the inner and outer theories of Q. As a corollary, neither of them can be consistently closed under Cons and Comp. In the present section, we present some systems of symmetric type-free truth.

We first introduce the system FS; the acronym FS comes from Friedman and Sheard who first introduced this system [23]. This appellation and the following formulation is due to Halbach [31].

**Definition 2.8.1.** The system FS over  $\mathcal{L}_T$  is defined as B with full-induction for  $\mathcal{L}_T$  plus the following axioms

$$\mathbf{FS1} \quad \forall \vec{x} \in \text{CT} \rightarrow [T(R\vec{x}) \leftrightarrow R(\text{val}(\vec{x}))]$$

$$\mathbf{FS2} \quad \forall x \in \text{St}_{\mathcal{L}_T} [T(\neg x) \leftrightarrow \neg Tx]$$

**FS3**  $\forall x, y \in \text{St}_{\mathcal{L}_T} [(T(x \forall y) \leftrightarrow (Tx \vee Ty)) \wedge (T(x \rightarrow y) \leftrightarrow (Tx \rightarrow Ty))]$

**FS4**  $\forall z \forall x [\text{St}_{\mathcal{L}_T}(\forall z.x) \rightarrow (T(\forall z.x) \leftrightarrow \forall y Tx(y/z))]$ ,

and the next two extra inference rules: for each  $\mathcal{L}_T$ -sentence  $\sigma$ ,

$$\frac{\sigma}{T(\ulcorner \sigma \urcorner)} \text{ (NEC)} \quad \frac{T(\ulcorner \sigma \urcorner)}{\sigma} \text{ (CONEC)}$$

**Proposition 2.8.2.**  $\text{FS} \vdash \text{Cons} + \text{Comp}$ .

The proof-theoretic strength of FS was determined by Halbach [31]. FS seems natural and has desirable properties, i.e., Cons and Comp, but it turned out to be  $\omega$ -inconsistent as can be shown by using McGee's trick [57].

**Theorem 2.8.3.** (1)  $\text{FS} \equiv \text{RT}_{<\omega}$ . (2) FS is  $\omega$ -inconsistent.

Hence, when  $B = \text{PA}$ , we have  $\text{FS}[\text{PA}] \equiv \text{RA}_{<\omega}$ . Rathjen [68] showed that  $\text{ACA}_0 + \text{BR}$ , where BR denotes the Bar Rule, and  $\text{ACA}_0^+$ , i.e.,  $\text{ACA}_0$  plus the axiom of the existence of  $\omega$ -th Turing Jump of each set  $X$ ,<sup>23</sup> prove the same  $\Pi_1^1$  theorems. Then we can observe from Rathjen's proof that  $\text{RA}_{<\omega}$  (and thus  $\text{FS}[\text{PA}]$ ) also has the same arithmetical theorems as they have.<sup>24</sup>

We next look at another system of symmetric truth. Let us consider natural generalizations of (NEC) and (CONEC): the new rules  $\text{NEC}^*$  and  $\text{CONEC}^*$  respectively defined as

$$\frac{T^\ulcorner \phi(\vec{x}^\urcorner)}{\phi(\vec{x})} \text{ (NEC}^*) \quad \frac{\phi(\vec{x})}{T^\ulcorner \phi(\vec{x})^\urcorner} \text{ (CONEC}^*)$$

for each  $\mathcal{L}_T$ -formula  $\phi(\vec{x})$ . We notice that they are admissible in FS due to the presence of NEC

<sup>23</sup>For the precise definition and more explanations of  $\text{ACA}_0^+$ , see [77, Ch.X].

<sup>24</sup>It is noted that Rathjen, in fact, already showed that  $\text{RA}_{<\omega} \leq \text{ACA}_0 + \text{BR} [\Pi_\infty^0]$  in [68]. The converse ( $\text{ACA}_0 + \text{BR} \subset_{\Pi_\infty^0} \text{RA}_{<\omega}$ ) is obtained by observing that the system  $\text{ACA}_0^\omega$ , which Rathjen intermediately uses to establish his result, is embeddable in  $\text{RA}_{<\omega}$ . Although the proof of the latter is model-theoretic, Theorem 3.5.6 below gives the proof-theoretical one, since  $\text{ACA}_0 + \text{BR}$  can be embedded in  $\text{TC}^*[\text{PA}] (\equiv \text{RA}_{<\omega})$  in the standard manner (the second-order variables are interpreted to range over  $\text{St}_{\mathcal{L}_0(P)}$ ).

and CONEC (use FS4).

Let  $\text{Pos} \subset \mathcal{L}_T$  be the set of  $T$ -positive  $\mathcal{L}_T$ -formulae. There are some different but equivalent ways of defining  $T$ -positive formulae and in fact the subsequent arguments in the present chapter does not depend on the choice of its definition. However, for a technical reason, let us assume that  $\text{Pos}$  is defined by the following BNF:

$$\text{Pos} ::= s(\vec{x})=t(\vec{x}) \mid Tt(\vec{x}) \mid \text{Pos} \wedge \text{Pos} \mid \text{Pos} \vee \text{Pos} \mid \forall z.\text{Pos} \mid \exists z.\text{Pos}$$

Then,  $\text{POS}$  is defined as  $\text{Pos} \cap \text{St}_T$ : i.e.,  $\text{PSt}$  is the set of  $T$ -positive  $\mathcal{L}_T$ -sentences.

We will consider another KF-style system  $\text{POSKF}$ .

**Definition 2.8.4.**  $\text{POSKF}$  consists of  $\text{PA}$ , the full induction over  $\mathcal{L}_T$  and:

$$\mathbf{P1} \quad x, y \in \text{CT}(x) \rightarrow [T(x=y) \leftrightarrow \text{val}(x) = \text{val}(y)] \wedge [T(x \neq y) \leftrightarrow \text{val}(x) \neq \text{val}(y)]$$

$$\mathbf{P2} \quad \text{CT}(x) \rightarrow [T\text{val}(x) \leftrightarrow T(Tx)]$$

$$\mathbf{P3} \quad \text{POS}_{\mathcal{L}}(x) \wedge \text{POS}_{\mathcal{L}}(y) \rightarrow [T(x \wedge y) \leftrightarrow (Tx \wedge Ty)]$$

$$\mathbf{P4} \quad \text{POS}_{\mathcal{L}}(x) \wedge \text{POS}_{\mathcal{L}}(y) \rightarrow [T(x \vee y) \leftrightarrow (Fx \vee Fy)]$$

$$\mathbf{P5} \quad \text{Var}(z) \wedge \text{POS}(\forall z.x) \rightarrow [T(\forall z.x) \leftrightarrow \forall yTx(y/z)]$$

$$\mathbf{P6} \quad \text{Var}(z) \wedge \text{POS}(\exists z.x) \rightarrow [T(\exists z.x) \leftrightarrow \exists yTx(y/z)],$$

where  $\text{POS}$  is the representation of the set  $\text{POS}$  of  $T$ -positive formulae.

**Lemma 2.8.5.**  $\text{POSKF} \vdash \text{PUTB}$ ; shown by meta-induction.

Hence,  $\text{POSKF}$  is a supersystem of  $\text{PUTB}$ . Indeed,  $\text{POSKF}$  is properly stronger than  $\text{PUTB}$ , since the former include  $\mathcal{T}(\text{PA})$  but the latter doesn't.

Our goal is to show the next theorem.

**Theorem 2.8.6.**  $\text{POSKF} + \text{NEC}^* + \text{CONEC}^*$  is conservative over  $\text{POSKF}$  for arithmetical sentences.

Then, since  $\text{POSKF}$  is truth-definable in  $\text{PUTB}$ , we have the next corollary:

**Corollary 2.8.7.**  $\text{PUTB} + \text{NEC}^* + \text{CONEC}^*$  is conservative over  $\text{PUTB}$  for arithmetical sentences.

**Remark 6.**  $\text{KF} + \text{NEC} + \text{CONEC}$  is inconsistent, and the same applies to  $\text{WKF} + \text{NEC} + \text{CONEC}$  and  $\text{DT} + \text{NEC} + \text{CONEC}$  as well. In addition,  $\text{PUTB} \subsetneq \text{POSKF} \subsetneq \text{KF}$  and  $\text{WKF} \not\subset \text{POSKF}$  and  $\text{POSKF} \not\subset \text{DT}$ .

I will give the proof of Theorem 2.8.6 in the rest of the present section. It suffices to show that

Given any derivation  $\mathcal{D}$  from  $\text{POSKF} + \text{NEC}^* + \text{CONEC}^*$ , for any  $(M, D') \models \text{POSKF}$ ,  
there exists  $D \subset M$  such that  $(M, D) \models \mathcal{D}$ ,

where  $(M, D) \models \mathcal{D}$  means that the universal closure of each  $\phi \in \mathcal{D}$  is true in  $(M, D)$ ; for, if  $\text{POSKF} + \text{NEC}^* + \text{CONEC}^* \vdash \sigma$  for some arithmetical  $\sigma$ , then  $\sigma$  is true in every  $M \models \text{PA}$  expandable to a model of  $\text{POSKF}$ , and thus  $\text{POSKF} \vdash \sigma$ .

Take any  $(M, D') \models \text{POSKF}$ . Let  $\mathcal{D} = \langle \phi_0, \dots, \phi_n \rangle$  be a derivation in  $\text{POSKF}$  in Hilbert style. We can assume without loss of generality that if  $\phi_i$  is  $Tt(\vec{x})$  then  $Tt(\vec{x})$  is not repeated later: i.e., there is no  $i < j \leq n$  such that  $Tt(\vec{x}) \equiv \phi_j$ . Then, we will construct a model  $(M, D)$  such that

$$(M, D) \models \text{POSKF} + \forall \vec{x} (\phi_0(\vec{x}) \wedge \dots \wedge \phi_n(\vec{x})).$$

For the sake of simplicity, we assume all the variables occurring in  $\mathcal{D}$  is included in  $\vec{x}$ . Let  $\theta_0(\vec{x}), \dots, \theta_i(\vec{x})$  be the sentences such that  $\theta_i$  ( $i \leq n$ ) is not  $T$ -positive,  $T(\ulcorner \theta_i \urcorner) = \phi_j$  for some  $j \leq n$  and  $T(\ulcorner \theta_i(\vec{x}) \urcorner)$  obtained by  $\text{NEC}^*$ .

Consider the following operation  $\mathbf{P}: \mathcal{P}(M) \rightarrow \mathcal{P}(M)$ :

$$\begin{aligned}
a \in \mathbf{P}(X) &\Leftrightarrow \exists b \exists c [\text{CT}^M(b) \wedge \text{CT}^M(c) \wedge a = b \dot{=} c \wedge \text{val}^M(b) = \text{val}^M(c)] \\
&\vee \exists b \exists c [\text{CT}^M(b) \wedge \text{CT}^M(c) \wedge a = b \dot{\neq} c \wedge \text{val}^M(b) \neq \text{val}^M(c)] \\
&\vee \exists \vec{d} [a = \text{sb}_{\vec{x}}^M(\ulcorner \theta_0(\vec{x}) \urcorner, \vec{d}) \vee \dots \vee a = \text{sb}_{\vec{x}}^M(\ulcorner \theta_l(\vec{x}) \urcorner, \vec{d})] \\
&\vee \exists b [\text{CT}^M(b) \wedge a = T^M b \wedge \text{val}(b) \in X] \\
&\vee \exists b \exists c [\text{POS}^M(b) \wedge \text{POS}^M(c) \wedge a = b \wedge^M c \wedge b \in X \wedge c \in X] \\
&\vee \exists b \exists c [\text{POS}^M(b) \wedge \text{POS}^M(c) \wedge a = b \vee^M c \wedge (b \in X \vee c \in X)] \\
&\vee \exists z \exists b [\text{Var}^M(z) \wedge \text{POS}^M(\forall z.b) \wedge a = \forall^M z.b \wedge \forall y [b(y) \in X]] \\
&\vee \exists z \exists b [\text{Var}^M(z) \wedge \text{POS}^M(\exists z.b) \wedge a = \exists^M z.b \wedge \exists y [b(y) \in X]],
\end{aligned}$$

where, for a predicate  $R$  or function  $f$  of  $\mathcal{L}_{\text{PA}}$ ,  $R^M$  is the interpretation of  $R$  or  $f$  in  $M$ .

Then, since  $\mathbf{P}$  is a positive arithmetical operator and  $\text{POSKF}(\supset \text{PUTB})$  can construct a fixed-point of such an operator (though not necessarily the least one) and  $(M, D')$  is a model of  $\text{POSKF}$ , there is  $D \subset M$  which is definable in  $(M, D')$  and a fixed-point of  $\mathbf{P}$ .

More precisely, set  $\xi(x, y)$  to be:

$$\begin{aligned}
& \exists b \exists c [\text{CT}(b) \wedge \text{CT}(c) \wedge x = b \dot{=} c \wedge \text{val}(b) = \text{val}(c)] \\
& \vee \exists b \exists c [\text{CT}(b) \wedge \text{CT}(c) \wedge x = b \dot{\neq} c \wedge \text{val}(b) \neq \text{val}(c)] \\
& \vee \exists \vec{d} [x = \text{sb}_{\vec{x}}(\ulcorner \theta_0(\vec{x}) \urcorner, \vec{d}) \vee \dots \vee x = \text{sb}_{\vec{x}}(\ulcorner \theta_l(\vec{x}) \urcorner, \vec{d})] \\
& \vee \exists b [\text{CT}(b) \wedge x = T b \wedge Ty(\text{val}(b))] \\
& \vee \exists b \exists c [\text{POS}(b) \wedge \text{POS}(c) \wedge x = b \wedge c \wedge Ty(b) \wedge Ty(c)] \\
& \vee \exists b \exists c [\text{POS}(b) \wedge \text{POS}(c) \wedge x = b \vee c \wedge (Ty(b) \vee Ty(c))] \\
& \vee \exists z \exists b [\text{Var}(z) \wedge \text{POS}(\forall z.b) \wedge x = \forall z.b \wedge \forall w [Ty(w)]] \\
& \vee \exists z \exists b [\text{Var}(z) \wedge \text{POS}(\exists z.b) \wedge x = \exists z.b \wedge \exists w [Ty(w)]],
\end{aligned}$$

Then diagonalize  $\xi$  and take its fixed-point formula  $\gamma$  such that

$$\gamma(x) \leftrightarrow \xi(x, \ulcorner \gamma \urcorner).$$

Note that since  $\xi$  is  $T$ -positive, we can assume that  $\gamma$  is also  $T$ -positive by the usual construction of fixed-point formulae. Then, we take  $D := \gamma^M(x)$ . We will use the same type of construction again in the next chapter.

**Lemma 2.8.8.**  $(M, D)$  is a model of POSKF.

*Proof.* First, PA-axioms are obviously satisfied since  $M \models \text{PA}$ . Second, the full induction over  $\mathcal{L}_T$  is trivially satisfied by  $(M, D)$ , since  $D$  is definable in  $(M, D')$  and  $(M, D')$  is a model of the full induction. Finally, it suffices to show that P1-P6 are all satisfied by  $(M, D)$ ; but it can be straightforwardly shown by using the fixed-pointness of  $D$ . The important point is that, since  $\theta_0, \dots, \theta_l$  are not  $T$ -positive, they have nothing to do with the POSKF-axioms and thus do not

affect the usual proof. I will only illustrate one example:  $(M, D) \models \text{P3}$ . Suppose  $a = b \wedge^M c$  and  $b, c \in \text{POS}^M$ .

$$\begin{aligned}
\gamma^M(b \wedge^M c) &\Leftrightarrow \xi(b \wedge^M c, \ulcorner \gamma \urcorner) && \text{(Since } \gamma \text{ is a fixed-point of } \xi) \\
&\Leftrightarrow T^M \gamma(b) \wedge T^M \gamma(c) \\
&\Leftrightarrow \gamma^M(b) \wedge \gamma^M(c) && \text{(By PUTB)} \quad \square
\end{aligned}$$

Now, we show that  $(M, D) \models \mathcal{D}$  by induction on  $\phi_k(\vec{x})$  ( $k \leq n$ ). By the last lemma, it suffices to show that  $(M, D) \models \phi_k(\vec{c})$  for all  $\vec{c} \in M$  if  $\phi_k$  is obtained by  $\text{NEC}^*$  or  $\text{CONEC}^*$ .

For the former case, let  $\phi_k \equiv T(\ulcorner \psi(\vec{x}) \urcorner)$  obtained by  $\text{NEC}$  from  $\psi(\vec{x})$ . Suppose  $\psi(\vec{x})$  is  $T$ -positive. Then since  $(M, D)$  is a model of  $\text{POSKF}$  and we already have  $(M, D) \models \psi(\vec{c})$  for all  $\vec{c}$  by IH, it follows from  $\text{PUTB}$ -schema (recall that  $\text{POSKF} \vdash \text{PUTB}$ ) that  $(M, D) \models \phi_k(\vec{c})$  for all  $\vec{c}$ . Otherwise,  $\psi(\vec{x})$  is  $\theta_i$  for some  $i \leq l$ , and thus, for any  $\vec{c} \in M$ ,  $\phi_k(\vec{c}) \in D$  by the third line of the definition of  $\mathbf{P}$ .

Finally, suppose  $\phi_k(\vec{x})$  is obtained by  $\text{CONEC}^*$  from  $\phi_j(\vec{x}) \equiv T(\ulcorner \phi_k(\vec{x}) \urcorner)$  ( $j < k$ ). If  $\phi_k$  is  $T$ -positive, it is satisfied by  $(M, D)$  again since it is a model of  $\text{PUTB}$ ; Suppose not. By IH, we have  $(M, D) \models \phi_j(\vec{c})$  for all  $\vec{c} \in M$  and thus  $\ulcorner \phi_k(\vec{c}) \urcorner \in D$  for all  $\vec{c}$ . Since  $\phi_k$  is not  $T$ -positive and thus  $\ulcorner \phi(\vec{c}) \urcorner \notin \text{POS}^M$  for all  $\vec{c}$  (see the above remark),  $\phi_k$  is identical to one of  $\theta_0, \dots, \theta_l$ ; for, if  $\ulcorner \phi_k(\vec{c}) \urcorner \in D$  then it must meet one of the disjuncts of the defining clauses of  $\mathbf{P}$ ; but then, only the second line of the definition of  $\mathbf{P}$  does apply. Let  $\phi_k = \theta_i$ . Then, by definition,  $T(\ulcorner \theta_i \urcorner)$  is obtained by  $\text{NEC}^*$  from some  $\phi_m$  ( $m \leq n$ ). However,  $T(\ulcorner \theta_i \urcorner)$  is just identical to  $\phi_j$  and we have already assumed that  $Tt$  is not repeated in  $\mathcal{D}$ ; therefore,  $m < j$  and thus by IH  $(M, D) \models \phi_m(\vec{c})$  for all  $\vec{c} \in M$ ; i.e.,  $(M, D) \models \phi_k(\vec{c})$ .

We have completed the proof.

Indeed, the above proof can be directly transformed into proof-theoretic proof; given a proof

$\mathcal{D}$ , then we just replace each  $T$  by a fixed-point of  $\mathbf{P}$  which is effectively defined in POSKF, since  $\text{POSKF} \supset \text{PUTB} \supset \text{ID}_1$ .

Furthermore, we can directly apply the above argument to  $\text{POSKF}^-$  and  $\text{PUTB}^-$  as it is. Now we have the next corollary.

**Corollary 2.8.9.**  $\text{POSKF}^- + \text{NEC}^* + \text{CONEC}^*$  is conservative over POSKF for arithmetical sentences, and the same applies to  $\text{PUTB}^-$  as well.

## 2.9 Systems of iterative non-compositional truth

Cantini [6] introduced the system VF of self-applicable truth, which is designed so as to embody the Kripkean fixed-point semantics with van Fraassen's supervaluation schema, and showed that it is proof-theoretically equivalent to the system  $\text{ID}_1$  of elementary inductive definitions.

Given a system  $\mathbf{S}$ ,  $\text{Ax}_\mathbf{S}$  represents the set of the axioms of  $\mathbf{S}$  and  $\text{Prv}_\mathbf{S}$  is a canonical provability predicate of  $\mathbf{S}$ . Let  $\text{PL}_\mathcal{L}$  be pure classical predicate logic for the language  $\mathcal{L}$ ; then  $\text{Ax}_{\text{PL}_\mathcal{L}}$  represents the set of logical axioms (for the language  $\mathcal{L}_T$ ). I adopt the simpler version of VF in [7] (or [45]) rather than the original one in [6].

**Definition 2.9.1.** The system  $\text{VF}^-$  over  $\mathcal{L}_T$  comprises PA plus the following truth axioms:

$$\mathbf{V1} \quad \forall \vec{x} [T(\ulcorner \phi(\vec{x}) \urcorner) \rightarrow \phi(\vec{x})], \text{ for each } \mathcal{L}_T\text{-formula } \phi(\vec{x}).$$

$$\mathbf{V2} \quad \forall \vec{x} \in \text{CT} [(T(R\vec{x}) \leftrightarrow R\vec{x}^\circ) \wedge (F(R\vec{x}) \leftrightarrow \neg R\vec{x}^\circ)], \text{ for each } R \in \mathcal{L}_0$$

$$\mathbf{V3} \quad \forall x \in \text{St}_T [\text{Ax}_{\text{PL}_\mathcal{L}_T}(x) \rightarrow Tx]$$

$$\mathbf{V4} \quad \forall x, y \in \text{St}_T [T(x \rightarrow y) \rightarrow (Tx \rightarrow Ty)]$$

$$\mathbf{V5} \quad \forall z \forall x [\text{St}_T(\forall z.x) \rightarrow (\forall y Tx(y/z) \rightarrow T\forall z.x)]$$

Then, the system VF is obtained by adding full-induction for  $\mathcal{L}_T$  to  $\text{VF}^-$ .

**Theorem 2.9.2** (Cantini [6, 7]).  $\text{VF}^- \equiv \text{B}$  and  $\text{VF} \equiv \text{ID}_1$ .

We will later have a closer look at iterative non-compositional truths in Chapter 5.

## 2.10 Schematic Reflective Closures

Feferman [17] introduced the notion of *schematic reflective closure* as an answer to the question ‘which schemata ... in the language of [a theory] ought to be accepted if one has accepted the basic (schematic) axioms and rules of [the theory] ([17, p.2])’.<sup>25</sup>

**Definition 2.10.1.** Given a base system  $\text{B}$  over  $\mathcal{L}_0$ , let  $\text{B}(P)$  denote a system over  $\mathcal{L}_0(P) := \mathcal{L}_0 \cup \{P\}$ , where  $P$  is a new unary predicate symbol, whose non-logical axioms are just the same as  $\text{B}$  (i.e.,  $P$  is a vacuous predicate in  $\text{B}(P)$  for which only logical axioms are assumed).

Then, the schematic reflective closure  $\text{Q}^*$  of a system  $\text{Q}$  of truth with a base system  $\text{B}$ , is a system over  $\mathcal{L}_{\text{Q}}(P) = \mathcal{L}_{\text{Q}} \cup \{P\}$  whose axioms consist of  $\text{Q}[\text{B}(P)]$ , full induction schema for  $\mathcal{L}_{\text{Q}}(P)$ , and a new inference rule  $P$ -Subst:

$$\frac{\phi(P)}{\phi(\hat{x}\psi(x))} \text{ (P-Subst) }, \quad (2.3)$$

for any formula  $\phi(P)$  of  $\mathcal{L}_0(P)$  (possibly interspersed with occurrences of  $P$ ) and a formula  $\psi$  of  $\mathcal{L}_{\text{Q}}(P)$ , where  $\phi(\hat{x}\psi(x))$  is the result of replacing each occurrence of  $Pt$  in  $\phi$  by  $\psi[t/x]$  and untangling bound variables.<sup>26</sup> Following Feferman’s notation, we sometimes write  $\phi(\hat{\psi})$  for  $\phi(\hat{x}\psi(x))$ .

Notice that  $\text{Q}[\text{B}(P)]$  contains truth-axioms for the new predicate symbol  $P$ . For example,

---

<sup>25</sup>The definition below is modified for our current arguments and not exactly the same as Feferman’s original one. Feferman’s original definition depends on what is regarded as ‘schematic axiom’. At any rate, if we focus on arithmetical base theories and only regard induction schema as ‘schematic axiom’, then Feferman’s definition and ours coincide.

<sup>26</sup>As Feferman pointed out, we should not allow arbitrary  $\phi$  of  $\mathcal{L}_{\text{Q}}(P)$  in the premise of  $P$ -Subst; otherwise, it may lead to a contradiction. For example, if we allow arbitrary  $\phi$  in  $\text{KF}^*$ , then we obtain  $T(Px) \leftrightarrow \psi$  for all  $\psi$ .

$\text{KF}[\mathbb{B}(P)]$  contains the axiom

$$\text{CT}(x) \rightarrow [P(\text{val}(x)) \leftrightarrow T(Px)],$$

and each ‘ $\text{St}_{\mathcal{L}_T}$ ’ occurring in **K0-K8** is replaced by  $\text{St}_{\mathcal{L}_T(P)}$ .

Schematic reflective closure can be viewed as an analogue to the Bar rule for subsystems of second-order arithmetic (cf. [17, §5]). It is not only a philosophically significant framework but also a useful and natural device of extending a given theory. It indeed increases the proof-theoretic strength of theories in many cases (though we will see some exceptional cases in §5).

**Theorem 2.10.2** (Feferman [17]).  $\text{KF}^*[\text{PA}] \equiv \text{RA}_{<\Gamma_0}$ . Thus  $\text{KF}^*[\text{PA}]$  is also proof-theoretically equivalent to, e.g.,  $\Sigma_1^1\text{-AC} + \text{BR}$ ,  $\widehat{\text{ID}}_{<\omega}$  and  $\text{ATR}_0$ .<sup>27</sup>

In addition to his first one, Feferman also conjectured that

**Feferman’s Second Conjecture:**  $\text{DT}^*[\text{PA}] \equiv \text{RA}_{<\Gamma_0}$ .

We will later show that this is indeed true.

We have now introduced all the systems to be considered in the rest of the present thesis. However, before concluding this section, we need to introduce some extra systems which are to be intermediately used for the proof-theoretic analyses of schematic reflective closures.

---

<sup>27</sup>For the definition of  $\widehat{\text{ID}}_{<\omega}$ , see [14]. For  $\text{ATR}_0$ , see [77].

## Chapter 3

# Relative Truth Definability

In the present chapter, I will propose one new tool, *relative truth definability*, to compare more conceptual aspects of systems of truth. We focus on the cases in which systems of truth share their base system. This is because we are interested in a comparison of the conceptions of truth and thus we want to keep all the other conceptions fixed thereby excluding inessential factors. For example, suppose the base system of one system is PA and that of another system is ZFC. Then, in most cases, the truth of the former system can be defined in the latter solely in terms of sets without appealing to its truth predicate. In such cases, what conception of the truth is employed in the latter system does not matter at all. Hence, we start by fixing a base system B to which every system is obtained by adding truth predicate(s) and its (their) axioms and inference rules.<sup>1</sup>

### 3.1 Comparing Conceptual Aspects of Axiomatic Truth

We have seen a large variety of axiomatic systems of truth. Different systems embody different views and conceptions of truth behind their axioms. It should be thus expected that one can gain

---

<sup>1</sup>The contents of the present section are based on my [24].

both mathematically and philosophically valuable insights about truth by investigating how they are related to each other and by comparing them from various points of view.

In the last chapter, we have also seen that much proof-theoretic work on axiomatic systems of truth has been done. Various systems of truth have been related to other systems of truth or certain subsystems of second-order arithmetic in terms of proof-theoretic equivalence, equiconsistency. For example, we have seen that some systems of truth have the same truth-free theorems, even if they apparently embody different views of truth. In many cases, these proof-theoretic works focus only on the arithmetical parts of systems, i.e., their theorems containing no truth predicate. The coincidences of truth-free theorems might indeed result from a certain close and intimate conceptual connection between two systems, but it seems that mere coincidences of truth-free theorems cannot solely imply the existence of such a connection. For example, let  $\mathbf{T}$  be a system over  $\mathcal{L}_T$  defined as  $\mathbf{B}$  plus a single additional axiom ' $T(\ulcorner 0 = 0 \urcorner)$ '. Then, we can show that  $\mathbf{T}$  and  $\mathbf{UTB}$  are proof-theoretically equivalent (and in fact mutually interpretable if  $\mathbf{B}$  is reflexive). However, it seems implausible to think that  $\mathbf{T}$  and  $\mathbf{UTB}$  share a common view and conception of truth, or, in other words, it is hardly agreed that they are conceptually 'equivalent'. Hence, proof-theoretic equivalence, conservativity or other well-known intertheoretic relations seem too coarse-grained for our purpose of comparing conceptual aspects of systems of truth.

If one adopts an instrumentalist approach to truth, in which truth is regarded merely as a practical device to increase the deductive or expressive power of systems, it might be sufficient to consider only what kind of and how rich truth-free consequences can be obtained by means of a truth predicate (or truth predicates if the language has more than one truth predicates). However, if we are interested in more conceptual aspects of truth systems, it seems no longer sufficient to consider only the derivability of truth-free sentences. For, different systems based on different and incompatible views of truth may have the same truth-free consequences just by accident. Hence, in order to compare conceptual aspects of systems of truth, we must take into consideration their

theorems which may contain truth predicate(s). We need mathematically finer-grained means to compare and distinguish their conceptual aspects.

We are thus interested in a more fine-grained means of comparing conceptual aspects of systems of truth. In the present section, we will propose such a criterion, by means of which we can represent a certain ‘equivalence’ or ‘reducibility’ among systems of truth from a more conceptual point of view, although it may be still too coarse to fully represent ‘conceptual equivalence’ or ‘conceptual reducibility’ (in some strong sense). As we have argued, any relation among systems which only focuses on truth-free parts seems insufficient for our purpose on the one hand. On the other hand, some relations are too strong. For example, subsystem relation and its induced equivalence relation, i.e., identity (as systems), should be too strong; e.g.,  $\text{KF} + \neg T(\mathbf{n})$  for some numeral  $\mathbf{n}$  which codes no sentence is not a subsystem of KF, but we usually think that the difference only consists in an inessential point irrelevant to their views of truth.<sup>2</sup>

In the present and next chapter, I will try to suggest two alternative means more suitable for comparing conceptual aspects of systems of truth than other so far presented means. Namely, what we will suggest is *relative truth definability*, which examine whether the truth predicate(s) of one system Q is definable in another system S in terms of S’s truth predicate(s). Although its precise definition and detailed explanations will be given in §3, I briefly explain it here. Given two systems Q and S of truth over languages  $\mathcal{L}_Q$  and  $\mathcal{L}_S$  which share the same base system B over  $\mathcal{L}_B (\subset \mathcal{L}_Q \cap \mathcal{L}_S)$ , we say that Q is relatively truth definable in S, if for each truth predicate  $T$  of  $\mathcal{L}_Q$  there exists a predicate  $R_T$  of  $\mathcal{L}_S$  such that  $Q \vdash \phi$  implies  $S \vdash \phi'$ , where  $\phi'$  is obtained from  $\phi$  by replacing all the occurrence of each truth predicate  $T$  by  $R_T$ ; in other words, Q is relatively truth definable in S, if there exists  $\mathcal{L}_B$ -conservative interpretation of Q in S. This notion of relative truth definability is in fact a finer-grained tool for comparing systems of truth than mere derivability of truth-free

---

<sup>2</sup>Subsystem relation as a representation of intertheoretic ‘reducibility’ was also criticized by Niebergall [63]. His argument is about ‘system reducibility’ in general and thus is applied to our current discussion as well.

sentences or other well-known means. For example, conservativity for truth-free theorems and relative interpretability follow from relative truth definability, but the converse does not necessarily hold. We will then argue that relative truth definability sheds more light on conceptual aspects of systems of truth and compare various systems in the light of it. In the end, we will see that relative truth definability can distinguish some systems with the same truth-free theorems but based on apparently different views of truth and also that it is useful in its applications.

I do not at all intend to deny other kinds of comparison of axiomatic systems of truth. What the present chapter aims at is rather to suggest one possible methodology for comparison and investigation of more conceptual aspects of truth systems, to which not enough attention has been paid, and to give a new perspective of recent works of axiomatic truth systems in the light of relative truth definability.

One probably comes to relative interpretability as a candidate for such an intertheoretic relation. Indeed, we often say, for instance, that geometry (or any other subject of mathematics) is reducible to set theory, and this usually means the existence of a relative interpretation of geometry in set theory. However, it seems still arguable whether this implies that the concept of geometrical objects such as ‘line’ or ‘triangle’ are ‘reducible’ to the concept of set. It is a deep and difficult problem whether this is a genuine reduction or not, but, at any rate, relative interpretability is not what we want to propose.<sup>3</sup>

First of all, we note that the system  $T$  above and  $UTB$  are mutually interpretable. Hence relative interpretability cannot exclude this kind of pathological cases. In general, the problems of relative interpretability as our means of comparison are:

- (A) Relative interpretation may interpret terms and predicates by totally different notions;
- (B) Relative interpretation may change the range of quantifiers, in other words, the ontological

---

<sup>3</sup>For more discussion of relative interpretability as a tool of reduction in more general settings, particularly on its role in foundational issues, see e.g., [18] and [63].

commitment of the interpreted system;

(C) Relative interpretability does not necessarily preserve truth-free theorems;

(D) Relative interpretability just follows from conservativity for truth-free theorems in some typical setting of systems of truth.

We will explain each of these problems below.

The characteristics (A) and (B) of relative interpretability are principal defects for our current objectives. It is true that relative interpretability is very useful and widely applicable for comparing diverse kinds of mathematical systems, but we are currently in a special situation in which only systems of truth are considered and we compare only their conceptions of truth. For example, relative interpretability is often regarded as a suitable tool for ontological reduction of mathematical (or other) objects posited by one mathematical subject into those posited by another subject. However, what we are aiming at is not ontological reduction but rather, so to say, ‘conceptual reduction’. Thus, for our current purpose, we would like to keep every other concept fixed in order to exclude inessential factors, although it may make our approach less versatile than relative interpretability. Therefore, since relative interpretability may require that the words, e.g., ‘for all natural numbers’ or ‘greater than’ of one system be translated by ‘for all even numbers’ or ‘less than’ in another system, relative interpretability seems unsuitable for our current purpose.

The problem (C) is caused by (A) and (B) and it is problematic since we are looking for a finer-grained means to compare systems. There are indeed two systems  $Q_0$  and  $Q_1$  such that  $Q_1$  can prove properly more truth-free theorems than  $Q_0$  but is relatively interpretable in  $Q_0$ ; e.g., consider TC and  $TC + \neg\text{Cons}(\text{TC})$ . If the truth of  $Q_1$  is ‘reducible’ to  $Q_0$ , then it is natural to require that the truth of  $Q_1$  should not bring richer truth-free consequences than that of  $Q_0$  does.

The fourth problem (D) is more precisely elaborated by the following fact.

**Theorem 3.1.1.** Let  $Q$  and  $S$  be first-order classical systems and primitive recursive extensions of

PA. Suppose that  $Q$  and  $S$  are both reflexive and that  $Q \subset_{\Pi_1^0} S$ . Then  $Q$  is interpretable in  $S$ .

*Proof.* Let  $\sigma(x)$  be a primitive recursive predicate which binumerates  $Q$ -axioms. Since  $Q$  is reflexive,  $Q \vdash \text{Con}_{\sigma|k}$  and thus  $S \vdash \text{Con}_{\sigma|k}$  for each  $k \in \mathbb{N}$ . Set  $\sigma^* := \sigma(x) \wedge \text{Con}_{\sigma|x}$ . Then,  $\sigma^*$  binumerates  $Q$  in  $S$ . Finally, the claim follows from the arithmetized completeness theorem.<sup>4</sup>  $\square$

Many of the systems of truth are reflexive; in fact, as we shall see in §4, all the systems of truth we have so far introduced are reflexive. Hence, in many cases, only the truth-free part (indeed  $\Pi_1^0$  part) of systems does matter to relative interpretability, and relative interpretability is no more fine-grained than the derivability of truth-free theorems.<sup>5</sup>

## 3.2 Relative Truth Definability

Now we present our own proposal, relative truth definability. Informally speaking, that a system  $Q$  is truth-definable in another system  $S$  means that the language  $\mathcal{L}_Q$  of  $Q$  can be translated into the language  $\mathcal{L}_S$  of  $S$  in such a way that the basic terms are kept unchanged and the concept of truth captured in  $Q$  is expressible and its function in  $Q$  can be fully simulated by a single formula in  $\mathcal{L}_S$  which translated the truth predicate(s) of  $\mathcal{L}_Q$ . The formal and precise definition is the following.

**Definition 3.2.1.** Let  $Q$  and  $S$  be (primitive recursive) systems of truth over languages  $\mathcal{L}_Q$  and  $\mathcal{L}_S$  respectively, and let  $\mathcal{L}_Q$  be  $\mathcal{L}_0 \cup \{T_i\}_{i \in I}$  where  $I$  is a certain index set; recall that  $\mathcal{L}_Q$  may have more than one truth predicates as in the case of  $\text{RT}_{<\alpha}$ . In order to keep the definition as general as possible, the base systems  $B$  of  $Q$  and  $B'$  of  $S$  may be different with the proviso that  $\mathcal{L}_B \subset \mathcal{L}_{B'}$ , although we usually consider the cases in which they are the same.

Given a formula  $\theta_i$  of  $\mathcal{L}_S$  for each  $i \in I$ , we define a translation  $\mathcal{T}_{\vec{\theta}}$  from  $\mathcal{L}_Q$  to  $\mathcal{L}_S$  in the following way:

---

<sup>4</sup>See §6 and §8 of [54].

<sup>5</sup>Some non-reflexive systems of truth are known. For example, Cantini's  $\text{KF}_{\text{int}}$  and its variant  $\text{KF}_{\text{tot}}$  are not reflexive, since they are finitely axiomatizable.

$$\mathcal{T}_{\bar{\theta}}(\phi) := \begin{cases} \phi & \text{if } \phi \text{ is an atomic formula of } \mathcal{L}_0 \\ \theta_i(x) & \text{if } \phi \text{ is of the form } T_i(x) \\ \neg\mathcal{T}_{\bar{\theta}}(\psi) & \text{if } \phi \text{ is of the form } \neg\psi \\ \mathcal{T}_{\bar{\theta}}(\psi_0) \vee \mathcal{T}_{\bar{\theta}}(\psi_1) & \text{if } \phi \text{ is of the form } \psi_0 \vee \psi_1 \\ \forall x\mathcal{T}_{\bar{\theta}}(\psi) & \text{if } \phi \text{ is of the form } \forall x\psi \end{cases}$$

We say  $\mathbf{Q}$  is *relatively truth definable* (more simply *truth-definable*) in  $\mathbf{S}$ , or  $\mathbf{S}$  (relatively) *defines the truth of*  $\mathbf{Q}$ , when there exists formulae  $\theta_i(x)$  of  $\mathcal{L}_{\mathbf{S}}$  for each  $i \in I$  such that  $\mathbf{Q} \vdash \phi \Rightarrow \mathbf{S} \vdash \mathcal{T}_{\bar{\theta}}(\phi)$  for all  $\phi \in \mathcal{L}_{\mathbf{Q}}$ . Then, we say that  $\mathcal{T}_{\bar{\theta}}$  is a *relative truth definition* (or simply *truth-definition*) of  $\mathbf{Q}$  in  $\mathbf{S}$  and we call  $\theta_i$  *truth-defining formula* of the truth predicate  $T_i$  of  $\mathbf{Q}$  in  $\mathbf{S}$ .

In other words,  $\mathbf{Q}$  is relatively truth definable in  $\mathbf{S}$  iff there exists an  $\mathcal{L}_0$ -conservative relative interpretation of  $\mathbf{Q}$  in  $\mathbf{S}$ , or, equivalently, iff  $\mathbf{Q}$  is a subsystem of a definitional expansion of  $\mathbf{S}$  (with the proviso that the truth predicates of  $\mathcal{L}_{\mathbf{Q}}$  and  $\mathcal{L}_{\mathbf{S}}$  are syntactically distinguished in the definitional expansion). In the literature, this kind of interpretation has been occasionally used for determining proof-theoretic strength of systems and the idea and technique themselves are nothing new; e.g., Feferman gave an  $\mathcal{L}_{\text{PA}}$ -conservative interpretation of  $\text{KF}[\text{PA}]$  in  $\Sigma_1^1\text{-AC}$  to establish  $\text{KF}[\text{PA}] \leq \Sigma_1^1\text{-AC}$ .

We will explain below that relative truth definability is finer-grained than other intertheoretic relations. First of all, we note that relative truth definability is at least as fine-grained as relative interpretability, consistency strength and conservativity of truth-free theorems: a relative truth definition is a relative interpretation; if a system  $\mathbf{Q}$  is relatively truth definable in  $\mathbf{S}$  then the truth-free theorems of the former are included in the latter, i.e.,  $\mathbf{Q} \subset_{\mathcal{L}_0} \mathbf{S}$ , and thus  $\mathbf{Q} \leq_{\text{Con}} \mathbf{S}$ .

On the other hand, we cannot generally obtain proof-theoretic reducibility from relative truth definability, since the latter has no condition on the system in which the definability itself is proved, although in practice, at least in all the cases in the present thesis, a proof of relative truth definability immediately entails proof-theoretic reducibility.

According to Feferman's original formulation in [16], a system  $\mathbf{Q}$  is proof-theoretically reducible

to  $S$  for a set of formulae  $\Phi \subset \mathcal{L}_Q \cap \mathcal{L}_S$ , if there exists a partial recursive function  $f$  such that

(I) for each proof  $p$  in  $Q$  of a formula  $\phi \in \Phi$ ,  $f(p) \downarrow$  and is a proof in  $S$  of  $\phi$

(II) the formalization of (I) is provable in  $S$ .

This definition itself is indeed debatable, since, as was pointed out by Niebergall [63], it forfeits transitivity of proof-theoretic reducibility. However, under a moderate assumption, this problem is resolved. Let us say that  $Q$  is proof-theoretically reducible to  $S$  for  $\Phi$  *provably in* a system  $T$  when (I) holds for  $Q$  and  $S$  and, instead of (II), (I) is provable in  $T$ . Then, the following fact due to Niebergall and Rathjen is known<sup>6</sup>:

**Theorem 3.2.2.** For any (primitive recursive) extensions  $Q$  and  $S$  of  $I\Sigma_1$ , if  $Q$  is proof-theoretically reducible to  $S$  for a primitive recursive set  $\Phi$  provably in  $S$  and the reduction function  $f$  can be taken to be primitive recursive, then  $Q$  is proof-theoretically reducible to  $S$  for  $\Phi$  provably in  $I\Sigma_1$ .

In practice, the reduction function  $f$  has been always primitive recursive in the literature: indeed, Rathjen rather defines proof-theoretic reducibility in [69] with the restriction that  $f$  be primitive recursive. In addition, we usually deal with systems stronger than  $I\Sigma_1$  in the context of foundational discussions and the present paper only considers such systems. Thus we can adopt, as our definition, proof-theoretic reducibility provably in  $I\Sigma_1$  via a primitive recursive reducing function, instead of Feferman's original definition.

Then, if we want relative truth definability to generally imply proof-theoretic reducibility, we need to add some new condition to Definition 3.2.1. For example, we can add the following condition (b) to Definition 3.2.1:

(b) The  $\theta_i$ 's are primitive recursively specified for each truth predicate  $T_i$  of  $Q$ , and it is provable

---

<sup>6</sup>The proof of Theorem 3.2.2 appeared in [18], which Feferman accredits to Niebergall, whereas Rathjen also pointed out the same fact in [69, §2.5].

in  $\mathbf{IS}_1$  that if  $\mathbf{Q} \vdash \phi$  then  $\mathbf{S} \vdash \mathcal{T}_{\bar{\theta}}(\phi)$ ; more formally,

$$\mathbf{IS}_1 \vdash \forall x \in \text{St}_{\mathbf{Q}} [\text{Prv}_{\mathbf{Q}}(x) \rightarrow \text{Prv}_{\mathbf{S}}(\mathcal{T}_{\bar{\theta}}(x))], \quad (3.1)$$

where  $\text{Prv}_{\mathbf{Q}}$  ( $\text{Prv}_{\mathbf{S}}$ ) are canonical representations of provability in  $\mathbf{Q}$  ( $\mathbf{S}$ ).

Then, the definition of truth-definability augmented by (b) implies proof-theoretic reducibility in general, since we can obtain by a well-known metamathematical property of  $\mathbf{IS}_1$  that (3.1) is equivalent to

$$\mathbf{S} \vdash \forall x \in \text{St}_{\mathbf{Q}} \forall y [\text{Proof}_{\mathbf{Q}}(y, x) \rightarrow \text{Proof}_{\mathbf{S}}(g(y), \mathcal{T}_{\bar{\theta}}(x))],$$

for some primitive recursive  $g$ , where  $\text{Proof}_{\mathbf{Q}}(y, x)$  ( $\text{Proof}_{\mathbf{S}}(y, x)$ ) canonically expresses that  $y$  is a code of a proof in  $\mathbf{Q}$  ( $\mathbf{S}$ ) of the formula coded by  $x$ . Furthermore, if  $\mathbf{Q}$  is a first-order classical system without extra inference rules, we only need: e.g.,

(b') The  $\theta_i$ 's are primitive recursively specified for each  $T_i$  of  $\mathbf{Q}$ , and we can construct a proof of  $\mathcal{T}_{\bar{\theta}}(\sigma)$  for each  $\mathbf{Q}$ -axiom  $\sigma$  primitive recursively on  $\sigma$ .

I think that these additional conditions neither lessen applicability and usability of truth-definability in practice nor harm the philosophical significance as a tool for comparing conceptual aspects of truth systems, although more careful discussions and philosophical evaluations are left for another occasion. At any rate, if one does not stick to general implication from relative truth definability to proof-theoretic reducibility, she need not worry too much about the issue argued so far, since relative truth definability immediately implies proof-theoretic reducibility in practice. In fact, in the subsequent proofs of relative truth definability in §4, we first primitive recursively assign a formula  $\theta$  of  $\mathcal{L}_{\mathbf{S}}$  to each truth predicate  $T$  of  $\mathcal{L}_{\mathbf{Q}}$  and then explicitly prove that  $\mathbf{S} \vdash \mathcal{T}_{\bar{\theta}}(\sigma)$  for all  $\mathbf{Q}$ -axioms  $\sigma$  and  $\mathcal{T}_{\bar{\theta}}$  preserves each inference rule of  $\mathbf{Q}$ . Given such a translation  $\mathcal{T}_{\bar{\theta}}$ , we can primitive recursively transform any proof of an  $\mathcal{L}_0$ -sentence in  $\mathbf{Q}$  to a proof of the same sentence in

S.

Whether we add some additional conditions like (b) to its definition or not, relative truth definability can be proved to be transitive; note that it is trivially reflexive. Hence relative truth definability induces a pre-ordering among systems of truth; then the classes of mutually truth-definable systems form equivalent classes. Niebergall [63] raised five conditions for a good axiomatization of the reducibility relation  $\rho$ : for classical first-order systems  $Q$  and  $S$ ,

**PRL1**  $Q \subset S \Rightarrow Q\rho S$ ;

**PRL2**  $\rho$  is transitive;

**PRL3** if  $Q\rho S$  then  $\text{Con}(S)$  implies  $\text{Con}(Q)$ ;

**PRL4** for each finite  $Q' \subset Q$ , if  $Q\rho S$ , there is a finite  $S' \subset S$  such that  $Q'\rho S'$ ;

**PRL5**  $|\Sigma_1 \vdash \text{Con}(S) \rightarrow \text{Con}(Q)$  for finitely axiomatized  $Q$  and  $S$  with  $Q\rho S$ .

It is easily shown that relative truth definability meets them all.

I would like to argue below that relative truth definability is a more suitable means than relative interpretability for expressing ‘conceptual reducibility’ of one system to another. First, the problems (A)-(D) (in §1) of relative interpretability no longer apply to relative truth definability. We have already observed that (C) is not the case for truth-definability; (D) doesn’t hold for relative truth definability as we shall see in §4; since the base language and quantification are kept unchanged in relative truth definition, neither (A) nor (B) is the case. Second, some pathological examples like  $T$  (i.e.,  $B + T(\ulcorner 0=0 \urcorner)$ ) above and  $UTB$  are excluded by relative truth definability.

**Proposition 3.2.3.**  $UTB$  is not truth-definable in  $T$ .

*Proof.* Let  $T'$  be  $T$  plus  $\forall x(x \neq \ulcorner 0=0 \urcorner \rightarrow \neg Tx)$ . Suppose  $UTB$  is truth-definable in  $T$  and thus  $T'$ . Note that  $T' = B + \forall x[Tx \leftrightarrow x = \ulcorner 0=0 \urcorner]$ . Then, if  $T' \vdash \phi$ , the formula  $\phi'$  obtained from  $\phi$  by replacing each occurrence of  $T(t)$  by  $t = \ulcorner 0=0 \urcorner$  is derivable in  $B$ . This implies that  $B$  can define the truth of  $UTB$ : this contradicts Tarski’s undefinability theorem.  $\square$

### 3.2.1 Mathematical Properties of Truth-Definability

We have mainly looked at the philosophical significance of relative truth definability, but it is also mathematically and technically useful. In the present subsection, we will see some of its mathematical properties and applicabilities.

First of all, as we have already observed, for systems  $Q$  and  $S$ , if  $Q$  is truth-definable in  $S$ , then  $Q \subset_{\mathcal{L}_0} S$ ,  $Q$  is interpretable in  $S$ ,  $Q \leq_{\text{Con}} S$ , and  $Q \leq S$  under some moderate condition like (b).

The next proposition immediately follows from the definition.

**Proposition 3.2.4.** (1) Suppose  $Q$  is truth-definable in  $S$  and  $M \models B$  can be expanded to a model of  $S$  (i.e., there exist  $X_i \subset |M|$  for each truth predicate  $T_i$  of  $S$  such that  $(M, \langle X_i \rangle_{i \in I}) \models S$ ). Then  $M$  is also expandable to a model of  $Q$ .

(2) Suppose  $Q$  is truth-definable in  $S$ . If  $Q$  is  $\omega$ -inconsistent, so is  $S$ .

(3) Suppose  $Q$  is truth-definable in  $S$ . Then  $Q$  plus full induction for  $\mathcal{L}_Q$  is truth-definable in  $S$  plus full induction for  $\mathcal{L}_S$ .

Thus, relative truth definability preserves useful mathematical properties. Proposition 3.2.4 tells us that relative truth definability can distinguish pathological systems like  $\omega$ -inconsistent system FS from others, since no  $\omega$ -consistent system can define the truth of an  $\omega$ -inconsistent system. It also tells that relative truth definability implies ‘conservativity in the semantic sense for  $\mathcal{L}_0$ ’; here a system  $S$  is conservative in the semantic sense over  $Q$  for  $\mathcal{L}_0$  iff the  $\mathcal{L}_0$ -reduct of every model of  $Q$  is expandable to a model of  $S$ . McGee [60] argued that conservativity in the semantic sense is to be preferred to and has more philosophical significance than conservativity in the ordinary sense (‘conservativity in the proof-theoretic sense’ in the terminology of [60]), since, he claims, there is no metaphysical cost in moving from one system of truth to another conservative system in the semantic sense and this move between conservative systems in the semantic sense makes no difference about what the world is like. On the other hand, for instance,  $RT_{<\omega}$  and FS are conservative over each

other for  $\mathcal{L}_0$  in the ordinary sense but the latter makes a stronger ontological commitment than the former since FS is  $\omega$ -inconsistent.

In contrast to relative truth definability, relative interpretability and proof-theoretic reducibility satisfy none of (1)-(3) above. It is thus observed that relative truth definability is a finer-grained tool for comparing systems of truth than relative interpretability and proof-theoretical reducibility (and thus consistency strength and conservativity).

We will show one more important feature of relative truth definability in its application to Feferman' *schematic reflective closure* already explained in Ch.2, §10. The next lemma shows that truth-definability carries over to schematic reflective closures.

**Lemma 3.2.5.** Let  $\mathbf{Q}$  be a first-order classical system of truth with a base system  $\mathbf{B}$  (without any extra inference rules such as NEC and CONEC). Suppose  $\mathbf{Q}[\mathbf{B}(P)]$  is truth definable in  $\mathbf{S}[\mathbf{B}(P)]$  and let  $\vec{\theta}$  be the truth-defining formulae. Then, the same translation  $\mathcal{T}_{\vec{\theta}}$  is a truth-definition of  $\mathbf{Q}^*$  in  $\mathbf{S}^*$ .

Even for  $\mathbf{Q}$  with extra inference rules, if  $\mathcal{T}_{\vec{\theta}}$  further preserves all the inference rules of  $\mathbf{Q}[\mathbf{B}(P)]$  in the sense that, for each inference rule  $\mathcal{R}$  of  $\mathbf{Q}$  in the form ‘if  $\phi_0, \dots, \phi_n$  are derived, then infer  $\psi_{\vec{\phi}}$ ’,  $\mathcal{T}_{\vec{\theta}}$  satisfies that

$$\text{if } \mathbf{S} \vdash \mathcal{T}_{\vec{\theta}}(\phi_0), \dots, \mathbf{S} \vdash \mathcal{T}_{\vec{\theta}}(\phi_n), \text{ then } \mathbf{S} \vdash \mathcal{T}_{\vec{\theta}}(\psi_{\vec{\phi}}),$$

then  $\mathbf{Q}^*$  is truth-definable in  $\mathbf{S}^*$  via  $\mathcal{T}_{\vec{\theta}}$ .

*Proof.* By induction on the length of derivation. The  $\mathcal{T}_{\vec{\theta}}$ -translation of each instance of full induction for  $\mathcal{L}_{\mathbf{Q}}(P)$  and each  $\mathbf{Q}^*$ -axiom trivially holds in  $\mathbf{S}^*$ . We show that the inference rules are preserved by  $\mathcal{T}_{\vec{\theta}}$ . Generalization and Modus Ponens are obviously preserved, and so are the other extra rules (if any) by the assumption. For  $P$ -Subst, suppose  $\mathbf{Q}^* \vdash \phi(P)$ . Then, we have  $\mathbf{S}^* \vdash \mathcal{T}_{\vec{\theta}}(\phi(P))$  by IH. Since  $\mathcal{T}_{\vec{\theta}}$  does not change  $P$ , we have  $\mathbf{S}^* \vdash (\mathcal{T}_{\vec{\theta}}(\phi))(P)$ . Thus, for any  $\psi \in \mathcal{L}_{\mathbf{Q}}(P)$ , we obtain  $\mathbf{S}^* \vdash (\mathcal{T}_{\vec{\theta}}(\phi))(\widehat{\mathcal{T}_{\vec{\theta}}(\psi)})$  by  $P$ -Subst. Finally, since  $\mathcal{T}_{\vec{\theta}}$  preserves the logical connectives and quantifiers,

$S^* \vdash \mathcal{T}_{\hat{\theta}}(\phi(\hat{\psi}))$  is obtained. □

This property is neither possessed by relative interpretability nor proof-theoretic reducibility.  $KF^*[[PA]]$  is neither proof-theoretically reducible nor relatively interpretable to  $(RT_{<\varepsilon_0})^*[[PA]]$  (by Lemma 3.5.6 below), whereas  $KF[[PA(P)]]$  is both proof-theoretically reducible to and interpretable (by Theorem 3.1.1) in  $RT_{<\varepsilon_0}[[PA(P)]]$ .

**Remark 7.** If a system  $B$  meets the conditions (a)-(d) in §2.1 for our base system, then so does  $B(P)$ . Hence, in the setting of the present thesis, whenever we establish the truth definability of  $Q$  in  $S$  without specifying the base system, it applies to both base systems  $B$  and  $B(P)$  and thus  $Q^*$  becomes automatically truth-definable in  $S^*$  as well by the last lemma. This is the reason why I didn't fix  $PA$  as the base system and rather adopt a slightly general setting (cf. Ch.2, §1).

### 3.3 Digression – Proof-Theoretic Ordinal –

The notion of proof-theoretic ordinal figures prominently among measures of the strength of subsystems of second-order arithmetic and various other systems, and the subject called ordinal analysis is now one of the major subdisciplines of proof system (see [69] or [66] for an overview). However, the usual definition of proof-theoretic ordinal cannot be directly applied to systems of truth. Some alternative definitions still seem implausible for systems of truth. Thus the ‘proof-theoretic ordinal’ cannot be an appropriate measure until some suitable definition is given for systems of truth. We shortly digress to discuss this issue.

Given a subsystem of second-order arithmetic  $Q$ , the proof-theoretic ordinal of  $Q$  usually denotes its  $\Pi_1^1$ -ordinal: i.e.,

$$\sup\{\alpha \mid Q \vdash TI(\prec) \wedge otp(\prec) = \alpha \wedge \prec \text{ is primitive recursive}\},$$

where  $TI(\prec)$  expresses the transfinite induction along  $\prec$ , that is,

$$\forall X [\forall x (\forall y \prec x (y \in X) \rightarrow x \in X) \rightarrow \forall x \in field(\prec) (x \in X)].$$

This definition presupposes that the language of  $\mathbf{Q}$  contains second-order parameters. Since they are not contained in the language of truth systems, the above definition should be modified in a certain way.

One such modification for systems without second-order parameters can be found in [69] for example. There Rathjen considers augmenting a system with a new unary predicate  $U$  and redefine

$$TI(\prec) := \forall x (\forall y \prec x U(y) \rightarrow U(x)) \rightarrow \forall x \in field(\prec) U(x)].$$

However, it remains to be decided which properties are to be assumed for  $U$ , and indeed the resulting ordinal depends on which properties are assumed. Suppose that we postulate all the axioms governing the truth predicate and mathematical induction even for  $U$ . Then, for instance, the resulting proof-theoretic ordinal of  $\mathbf{KF}[\mathbf{PA}]$  would be  $\varphi_{\varepsilon_0}(0)$ . On the other hand, when we only postulate induction for  $U$ , the proof theoretic ordinal would be  $\varepsilon_0$ .<sup>7</sup> If we postulate nothing for  $U$  then the ordinal would then be  $\omega$ .

Another possible definition is for it to be the  $\alpha$  such that  $\mathbf{Q} \equiv \mathbf{PA} + \bigcup_{\beta < \alpha} \mathbf{TI}_{\mathcal{L}_{\mathbf{PA}}}(\beta)$  or  $\sup\{\beta \mid$

---

<sup>7</sup>I only give a sketch of its proof here. Let us denote the system we are considering by  $\mathbf{KF}(U)$ . Take a suitable semi-formal system to which  $\mathbf{KF}(U)$  is embedded; such a system can be obtained by adding inference rules for the truth predicate  $T$  to  $\omega$ -arithmetic; cf. [5]. In this system, no sentence containing  $U$  can be critical in the rules for  $T$ . We call a derivation  $\mathcal{D}$  in this system *quasi-normal*, when cut is only applied to atomics of the forms  $Tt$  and  $\neg Tt$ . Then, we can show that, if  $\prec$  is a well-ordering and  $otyp(\prec)$  is limit,  $otyp(\prec)$  is always less than or equal to the least length  $\alpha$  of quasi-normal derivations of  $TI(\prec)$ . This is shown in a similar way to Theorem 6.7.2 (Boundedness Theorem) of [66]. It is crucial in its proof that  $U$  cannot be critical in any rule for  $T$ . For example, in showing the analogue to Lemma 6.6.8 (Boundedness Lemma) of [66] (we can take a fixed interpretation of  $T$  in  $\mathbb{N}$ , say, the least Kripkean fixed-point, since it doesn't depend on the interpretation of  $U$  at all), though  $T$  may occur in  $\Delta$  (we are using the same notation of [66]), the positiveness or negativness of  $U$ 's occurrences isn't affected by any rule for  $T$  and the induction steps work in the same way as the cited lemma; by contrast, in the case of  $\mathbf{KF}[\mathbf{PA}(U)]$ , this argument is no longer valid because  $U$  can be critical in the rules for  $T$  and the interpretation of  $T$  must depend on the interpretation of  $U$ . However, we can also show by the standard technique that each derivation  $\mathcal{D}$  in  $\mathbf{KF}(U)$  can be transformed into a quasi-normal derivation  $\mathcal{D}'$  in the semi-formal system of the length  $< \varepsilon_0$ . This entails that if  $\mathbf{KF}(U) \vdash TI(\prec)$  then  $otyp(\prec) < \varepsilon_0$ .

$\text{PA} + \text{TI}_{\mathcal{L}_{\text{PA}}}(\beta) \subset \mathbf{Q}$ . However, both are only focusing on the truth-free part, and thus they are unsatisfactory for the same reason as the reason why we reject proof-theoretic reducibility or conservativity.<sup>8</sup>

### 3.4 Several Results on Truth Definability

In this section, we will present several results on truth-definability among systems we have introduced in §2.

For simplicity and readability, let  $\preceq$  stand for the truth-definability relation: i.e., for systems  $\mathbf{Q}_0$  and  $\mathbf{Q}_1$ , we write  $\mathbf{Q}_0 \preceq \mathbf{Q}_1$  if  $\mathbf{Q}_0$  is truth-definable in  $\mathbf{Q}_1$ . Then,  $\mathbf{Q}_0 \prec \mathbf{Q}_1$  is for  $\mathbf{Q}_0 \preceq \mathbf{Q}_1$  &  $\mathbf{Q}_1 \not\preceq \mathbf{Q}_0$ .

Some results on truth-definability immediately follow from already known or shown facts. First, if  $\mathbf{Q}_0$  is a subsystem of  $\mathbf{Q}_1$ , then the former is trivially truth-definable in the latter (via the identity translation). Hence, for example,

$$\text{TB} \preceq \text{UTB} \quad \text{PUTB} \preceq \text{KF} \quad \text{KF} \preceq \text{KF}^*.$$

Second, since truth-definability implies conservativity for  $\mathcal{L}_0$ -sentences, any system not conservative (for  $\mathcal{L}_0$ -sentences) over a system  $\mathbf{Q}$  is not truth-definable in  $\mathbf{Q}$ . Hence, for example, the following negative results on truth-definability are already available:

$$\text{TB}, \text{UTB}, \text{TC}^-, \text{KF}^- \not\preceq \text{TC} \not\preceq \text{FS} \not\preceq \text{RT}_{<\varepsilon_0}$$

Third, the next proposition follows from Proposition 3.2.4, since no system but FS is  $\omega$ -inconsistent.<sup>9</sup>

---

<sup>8</sup>For another example, Beklemishev [2] introduced  $\Pi_n^0$ -ordinal of systems ( $n \in \mathbb{N}$ ) which can be defined without appealing to second-order parameters. However, this definition should be rejected for the same reason that it only looks at the arithmetical (truth-free) part.

<sup>9</sup>Recall that we have assumed that  $\mathbf{B}$  is arithmetically sound and thus has a model over  $\mathbb{N}$ . Then, it is known

**Lemma 3.4.1.** None of the systems introduced so far (except FS itself) can define the truth of FS.

Fourth, due to Tarski's undefinability theorem, no system of truth presented so far is truth-definable in PA, since every system can derive the TB-schema.

In the rest of this section, we will show less trivial results.

**Lemma 3.4.2.**  $\text{KF}^- \not\leq \text{TC}^-$ , and  $\text{TC}^- \not\leq \text{KF}^-$ . On the other hand, we have  $\text{TC}^- \prec \text{KF}^-$  since the former is a subsystem of the latter.

*Proof.* The first claim follows from the fact that KF is stronger than TC (cf. Theorem 3.4.4 below). For the second claim, let  $\mathbf{B}'$  be  $\mathbf{B}$  plus full-induction for  $\mathcal{L}_0$ . For the sake of contradiction, suppose  $\text{TC}^- \leq \text{KF}^-$ . Since a relative truth definition preserves the  $\mathcal{L}_0$ -part, this entails that  $\text{TC}^- \llbracket \mathbf{B}' \rrbracket \leq \text{KF}^- \llbracket \mathbf{B}' \rrbracket$  as well. However, every nonstandard model  $M \models \mathbf{B}'$  must be recursively saturated with respect to  $\mathcal{L}_0$  (i.e., every recursive type  $p(x)$  from  $\mathcal{L}_0$  over  $M$  is realized in  $M$ ) whenever it is expandable to a model of  $\text{TC}^- \llbracket \mathbf{B}' \rrbracket$ ; Lachlan [51] showed this for  $\mathbf{B}' = \text{PA}$ , but the proof of Lachlan's result (also see [47]) can be directly generalized to an arbitrary base system  $\mathbf{B}'$  over  $\mathcal{L}_0$  which meets the conditions we posed in §2 and contains full-induction for  $\mathcal{L}_0$ . On the other hand, since  $\mathbf{B}$  is arithmetically sound and thus  $\mathbf{B}'$  has an infinite model,  $\mathbf{B}'$  has a non-recursively saturated model; we can take the prime model of some suitable complete extension of  $\mathbf{B}'$ , which is not recursively saturated but a nonstandard model of this extension.<sup>10</sup> By Proposition 3.2.4 (1), these contradict the fact that any model  $M \models \mathbf{B}'$  can be expanded to a model of  $\text{KF}^- \llbracket \mathbf{B}' \rrbracket$  ([5, §5]).  $\square$

**Lemma 3.4.3.**  $\text{TB}, \text{UTB} \prec \text{TC}^-$

*Proof.* It is obvious that  $\text{TB}, \text{UTB} \leq \text{TC}^-$ , since the former two are subsystems of the latter. We can show by the same argument as the last lemma that the converse fails, since TB and UTB are

---

that all the systems in §2 except FS has a model over  $\mathbb{N}$  and thus  $\omega$ -consistent.

<sup>10</sup>For the details of the proof, we refer the readers to [47, §8.1, §11.2], where these results are proved for PA but they can be generalized for an arbitrary  $\mathbf{B}$  with full-induction.

conservative over  $\mathbf{B}$  in the semantic sense.<sup>11</sup> □

There still remains one open problem on the relative truth definability among conservative systems of truth over  $\mathbf{B}$ : we do not yet know whether  $\text{UTB}$  is truth-definable in  $\text{TB}$ .

**Theorem 3.4.4.**  $\text{RT}_{<\varepsilon_0}$  is truth-definable in  $\text{KF}$ .

We prove a more general fact from which this theorem easily follows. Preliminarily, define an  $\mathcal{L}_T$ -formula  $D^+(x)$  by

$$D^+(x) := \text{St}_{\mathcal{L}_T}(x) \wedge (Tx \vee Fx) \wedge \neg(Tx \wedge Fx).$$

Notice that  $D^+(x)$  is equivalent to  $(Fx \leftrightarrow \neg Tx) \wedge x \in \text{St}_{\mathcal{L}_T}$ .

**Lemma 3.4.5.** Let  $\mathbf{Q}$  be a system over  $\mathcal{L}_T$  which derives the following:

- (i) The axiom **K1** (of  $\text{KF}$ );
- (ii)  $[D^+(x) \wedge D^+(y)] \rightarrow [D^+(\neg x) \wedge D^+(x \vee y) \wedge D^+(x \rightarrow y)]$ ;
- (iii)  $[\text{St}_{\mathcal{L}_T}(\forall z.x) \wedge \forall y D^+x(y/z)] \rightarrow D^+(\forall z.x)$ ;
- (iv)  $D^+(x) \rightarrow (T\neg x \leftrightarrow \neg Tx)$ <sup>12</sup>;
- (v)  $D^+(x \vee y) \rightarrow [Tx \vee y \leftrightarrow (Tx \vee Ty)]$ ;
- (vi)  $D^+(x \rightarrow y) \rightarrow [Tx \rightarrow y \leftrightarrow (Tx \rightarrow Ty)]$ ;
- (vii)  $D^+(\forall z.x) \rightarrow [T\forall z.x \leftrightarrow \forall y Tx(y/z)]$ ;
- (viii) The axiom **K2** (of  $\text{KF}$ ).

Then, if  $\mathbf{Q} \vdash \text{TI}_{\mathcal{L}_T}(<\alpha)$ , then  $\text{RT}_{<\alpha}$  is truth definable in  $\mathbf{Q}$ .

*Proof.* Let  $h$  be a binary primitive recursive function such that

---

<sup>11</sup>Let  $\mathbf{B}'$  be  $\mathbf{B}$  plus full-induction for  $\mathcal{L}_0$ . We can show that  $\text{UTB}$  plus full-induction for  $\mathcal{L}_T$  is proof-theoretically reducible to  $\mathbf{B}'$  for  $\mathcal{L}_0$  by using partial truth predicates; cf. fn.5. This gives an alternative proof of this lemma, since  $\text{TC}^-\llbracket\mathbf{B}\rrbracket \preceq \text{UTB}\llbracket\mathbf{B}\rrbracket$  implies  $\text{TC}^-\llbracket\mathbf{B}'\rrbracket \preceq \text{UTB}\llbracket\mathbf{B}'\rrbracket$  and thus  $\text{TC}\llbracket\mathbf{B}'\rrbracket$  would become truth-definable in  $\text{UTB}\llbracket\mathbf{B}'\rrbracket$  plus full-induction for  $\mathcal{L}_T$ . In contrast, it is known that  $\text{UTB}$  plus full induction is not conservative over  $\mathbf{B}'$  in the semantic sense; only recursively saturated  $\mathbf{B}'$ -models can be expanded to its model (for the proof, see [47, §15]).

<sup>12</sup>This (iv) is actually redundant since it follows from the definition of  $D^+$  as was noted.

$$h(x, \beta) := \begin{cases} x & \text{if } x \in \text{St}_\beta \text{ and } \beta < \alpha \\ \ulcorner 0 = 1 \urcorner & \text{otherwise} \end{cases}$$

We write  $h_\beta(x)$  for  $h(x, \beta)$ . Then, by the primitive recursion theorem, we take another primitive recursive function  $k$  such that:

$$k(a) := \begin{cases} a & \text{if } a \in \text{AtFml}_{\mathcal{L}_0} \\ \ulcorner T(k \circ h_\gamma(b)) \urcorner & \text{if } a = T_\gamma b \text{ for } b \in \text{Term} \\ \neg \ulcorner Tk\dot{b} \urcorner & \text{if } a = \neg b \\ \ulcorner Tk\dot{b} \urcorner \vee \ulcorner Tk\dot{c} \urcorner \text{ (} \ulcorner Tk\dot{b} \urcorner \rightarrow \ulcorner Tk\dot{c} \urcorner \text{, resp.)} & \text{if } a = b \vee c \text{ (or } a = b \rightarrow c \text{)} \\ \forall c. \ulcorner Tk(\dot{b}(c/\ulcorner c \urcorner)) \urcorner & \text{if } a = \forall c. b \text{ for } c \in \text{Var} \\ \ulcorner 0 = 1 \urcorner & \text{otherwise;} \end{cases}$$

We slightly abuse notation here;  $\ulcorner T(k \circ h_\gamma(b)) \urcorner$  and  $\forall z. \ulcorner Tk(\dot{b}(c/\ulcorner c \urcorner)) \urcorner$  precisely mean  $\ulcorner T(k \circ h_\gamma(u)) \urcorner (b/\ulcorner u \urcorner)$  and  $\forall c. \ulcorner Tk(\dot{b}(u/\ulcorner c \urcorner)) \urcorner (c/\ulcorner u \urcorner)$  where  $u$  is a fresh variable; informally,  $k(T_\gamma t) = T(k \circ h_\gamma(t))$  and  $k(\forall z \phi(z)) = \forall z T(k^\ulcorner \phi(\dot{z}) \urcorner)$  for a term  $t$  and formula  $\phi$ . We notice that if  $x \in \text{St}_\beta$  for  $\beta < \alpha$  then  $kx \in \text{St}_{\mathcal{L}_T}$ . Then, for each  $\beta < \alpha$ , we set  $\theta_\beta(x)$  to be  $Tk(x)$  and we take  $\theta_\beta(x)$  as the truth-defining formula of  $T_\beta$ .

We will first show by  $\text{TI}_{\mathcal{L}_T}(\beta)$  for each fixed  $\beta < \alpha$  that

$$\mathbf{Q} \vdash (\forall \gamma \leq \beta) (\text{St}_\gamma(x) \rightarrow D^+(kx)). \quad (3.2)$$

Since  $\mathbf{Q} \vdash \text{TI}_{\mathcal{L}_T}(< \alpha)$ , it suffices to show that  $\text{St}_\gamma(x) \rightarrow D^+(kx)$  is progressive on  $\gamma$ . Let  $a \in \text{St}_\gamma$  and suppose the claim up to  $\gamma$ . We show the claim for  $\gamma$  by subinduction on the surface complexity of the sentence coded by  $a$ .

The case in which  $a \in \text{AtSent}_{\mathcal{L}_0}$  immediately follows from (i) since  $k(a) = a$ . Suppose  $a = T_\delta x$

for  $x \in \text{CT}$  and  $\delta < \gamma$ . Then, by (viii), it is observed that:

$$\begin{aligned} D^+(ka) &\Leftrightarrow [T^\Gamma Tk \circ h_\delta(x)^\neg \vee F^\Gamma Tk \circ h_\delta(x)^\neg] \wedge [\neg T^\Gamma Tk \circ h_\delta(x)^\neg \vee \neg F^\Gamma Tk \circ h_\delta(x)^\neg] \\ &\Leftrightarrow [Tk \circ h_\delta(\text{val}(x)) \vee Fk \circ h_\delta(\text{val}(x))] \wedge [\neg Tk \circ h_\delta(\text{val}(x)) \vee \neg Fk \circ h_\delta(\text{val}(x))]. \end{aligned}$$

If  $\text{val}(x) \in \text{St}_\delta$ , then  $h_\delta(\text{val}(x)) = \text{val}(x)$  and thus  $k \circ h_\delta(\text{val}(x)) = k(\text{val}(x))$ ; the claim follows from IH. Otherwise,  $k \circ h_\delta(\text{val}(x)) = \ulcorner 0 = 1 \urcorner$ ; the claim trivially follows from (i). For the cases in which  $a = \neg b$ ,  $a = b \vee c$  or  $a = b \rightarrow c$ , we first observe that  $D^+(kx) \Leftrightarrow D^+(\ulcorner Tk \dot{x} \urcorner)$  for all  $x \in \text{St}_\gamma$ , which is shown by (viii). Then the claim follows from (ii). Next, suppose  $a$  is of the form  $\forall c.b$ . We already have  $(\forall y)D^+(k(b(y/c)))$  by SIH, and we also have:

$$D^+(\ulcorner Tk(\dot{b}(c/\ulcorner c \urcorner))^\neg(y/c) \urcorner) \Leftrightarrow D^+(\ulcorner Tk(\dot{b}(y/\ulcorner c \urcorner))^\neg \urcorner) \Leftrightarrow D^+k(b(y/c))$$

for each  $y$  in  $\mathbf{Q}$ . This entails the claim  $D^+(\forall c.\ulcorner Tk(\dot{b}(c/\ulcorner c \urcorner))^\neg \urcorner)$  by (iii). Thus, we have shown (3.2).

Now we can straightforwardly show using (iv)-(viii) that  $\mathcal{T}_{\bar{\theta}}$  is indeed a truth definition of  $\text{RT}_{<\alpha}$  in  $\mathbf{Q}$ . We illustrate  $\mathbf{Q} \vdash \mathcal{T}_{\bar{\theta}}(\mathbf{R2})$  and  $\mathbf{Q} \vdash \mathcal{T}_{\bar{\theta}}(\mathbf{R6})$  for examples. For the former, let  $\beta < \alpha$  and  $x \in \text{St}_\beta$ . Then  $\mathcal{T}_{\bar{\theta}}(T_\beta \neg x)$  is equivalent to  $T(\ulcorner \neg Tk \dot{x} \urcorner)$ . By (3.2), we already know both  $D^+(\ulcorner \neg Tk \dot{x} \urcorner)$  and  $D^+(kx)$ . Hence,

$$T(\ulcorner \neg Tk \dot{x} \urcorner) \Leftrightarrow \neg T(\ulcorner Tk \dot{x} \urcorner) \Leftrightarrow \neg Tkx = \mathcal{T}_{\bar{\theta}}(\neg T_\beta x);$$

we thus obtain  $\mathbf{Q} \vdash \mathcal{T}_{\bar{\theta}}(\mathbf{R2})$ . Next, for the latter, let  $\beta < \alpha$  and take any  $\delta < \beta$  and  $x \in \text{CT}$  with  $\text{val}(x) \in \text{St}_\delta$ . Then, it follows from (viii) that

$$\mathcal{T}_{\bar{\theta}}((T_\beta(T_\delta x)) \Leftrightarrow T(\ulcorner Tk \circ h_\delta x \urcorner) \Leftrightarrow Tk \circ h_\delta(\text{val}(x))$$

Since we have assumed that  $\text{val}(x) \in \text{St}_\delta$ , we have  $k \circ h_\delta(\text{val}(x)) = k(x)$  and thus the last formula is equivalent to  $\mathcal{T}_{\bar{\theta}}(T_\beta \text{val}(x)) = Tk(\text{val}(x))$ . The other cases can be shown similarly.  $\square$

Since  $\text{KF}^-$  satisfies the conditions (i)-(viii)<sup>13</sup> and  $\text{KF} \vdash \text{TI}_{\mathcal{L}_T}(< \varepsilon_0)$  (by the standard proof), Theorem 3.4.4 follows from this lemma. In addition,  $\text{FKF}^-$  also meets all the conditions (i)-(viii), the next theorem follows.

**Theorem 3.4.6.**  $\text{RT}_{<\varepsilon_0}$  is truth-definable in  $\text{FKF}$  (and thus  $\text{DT}$ ).

**Corollary 3.4.7.** If  $\alpha \leq \varepsilon_0$ , then  $\text{DT}^-, \text{KF}^- \not\leq \text{RT}_{<\alpha}$ .<sup>14</sup>

*Proof.* We only show the claim for  $\text{KF}^-$ , but  $\text{DT}^-$  case is shown in the same manner. When  $\alpha < \varepsilon_0$ , the claim is obvious since  $\text{RT}_{<\varepsilon_0} \leq \text{KF}$ . Let  $\alpha = \varepsilon_0$ . For the sake of contradiction, suppose  $\text{KF}$  is truth-definable in  $\text{RT}_{<\varepsilon_0}$ . Let  $\theta(x)$  be the  $\mathcal{L}_{\varepsilon_0}$ -formula defining the truth of  $\text{KF}$ . Take a finite  $\Gamma \subset \text{RT}_{<\varepsilon_0}$  such that  $\Gamma \vdash \mathcal{T}_\theta(\mathbf{K0}) \wedge \cdots \wedge \mathcal{T}_\theta(\mathbf{K8})$ . Take the maximum  $\beta < \varepsilon_0$  such that  $T_\beta$  occurs in  $\Gamma$  or  $\theta$ . Then, we have  $\Gamma \cup \{\theta\} \subset \text{RT}_{<\beta+1}$  and thus  $\text{KF}$  would be truth-definable in  $\text{RT}_{<\beta+1}$ ; however, since  $\beta + 1 < \varepsilon_0$ , it leads to a contradiction.  $\square$

These techniques can be modified for schematic reflective closures.

**Lemma 3.4.8.** Suppose  $\mathbf{Q}$  over  $\mathcal{L}_T$  satisfies (i)-(viii) of Lemma 3.4.5. Then,  $\text{RT}_{<\Gamma_0}$  is truth-definable in  $\mathbf{Q}^*$ . We can take the same truth-definition  $\mathcal{T}_\theta$  as in the proof of Lemma 3.4.5.

*Proof.* Define ordinals  $\beta_n$  for each  $n \in \mathbb{N}$  by  $\beta_0 = \varepsilon_0$  and  $\beta_{n+1} = \varphi_{\beta_n}(0)$ . Note that  $\Gamma_0 = \lim_{n \rightarrow \infty} \beta_n$  and  $\beta_n$  is an epsilon number for all  $n \in \mathbb{N}$ . By Lemma 3.4.5, it suffices to show that  $\mathbf{Q}^* \vdash \text{TI}_{\mathcal{L}_T(P)}(< \beta_n)$  for all  $n \in \mathbb{N}$ . This is shown by induction on  $n$ .

Recall that  $\mathbf{Q}^*$  contains full induction for  $\mathcal{L}_T(P)$ . Thus  $\mathbf{Q}^* \vdash \text{TI}_{\mathcal{L}_T(P)}(< \varepsilon_0)$ : we have shown the base step. For the induction step, suppose  $\mathbf{Q}^* \vdash \text{TI}_{\mathcal{L}_T(P)}(< \beta_n)$ . Then, by Lemma 3.4.5,

<sup>13</sup>See Lemma 3.2.4 and 3.2.5 of [17], where  $D^+(x)$  is denoted by  $D(x)$ .

<sup>14</sup>We might expect that  $\varepsilon_0$  can be replaced by  $\Gamma_0$  in this corollary. I do not yet know whether it is true or false, but we can at least show that there is no truth definition of  $\text{KF}^- \llbracket \mathbf{B}(P) \rrbracket$  (and  $\text{DT}^- \llbracket \mathbf{B}(P) \rrbracket$ ) in  $\text{RT}_{<\Gamma_0} \llbracket \mathbf{B}(P) \rrbracket$ . For there would otherwise be a truth definition  $\text{KF}^*$  in  $(\text{RT}_{<\Gamma_0})^*$  and then we could derive a contradiction in a similar way to the below proof by using Lemmata 3.4.8 and 3.5.6.

$\text{RT}_{<\beta_n} \llbracket \mathbb{B}(P) \rrbracket$  is truth-definable in  $\mathbf{Q}^*$ . It follows from Theorem 2.3.8 that  $\text{RT}_{<\beta_n} \llbracket \mathbb{B}(P) \rrbracket \vdash \text{TI}_{\mathcal{L}_0(P)}(<\beta_{n+1})$  and thus  $\mathbf{Q}^* \vdash \text{TI}_{\mathcal{L}_0(P)}(<\beta_{n+1})$ . In particular  $\mathbf{Q}^* \vdash \text{TI}(<\beta_{n+1}; P(x))$ . Finally, we apply  $P$ -Subst and obtain  $\text{TI}_{\mathcal{L}_T(P)}(<\beta_{n+1})$ .  $\square$

**Theorem 3.4.9.**  $\text{RT}_{<\Gamma_0}$  is truth-definable in  $\text{KF}^*$  and  $\text{FKF}^*$  (and thus in  $\text{DT}^*$ ).

We now turn to the relative truth definability among type-free systems of truth. Cantini [5] showed that  $\text{KF}$  can construct a fixed-point of a positive arithmetical operator. Looking closely into his proof, it is observed that only the PUTB schema is needed for this construction. Let  $\mathcal{L}_0^{(2)}$  be the second-order language which expands  $\mathcal{L}_0$  with second-order variables. We have the following.

**Lemma 3.4.10** (Cantini [5]). Let  $\Phi(X, x, \vec{v})$  be a  $\mathcal{L}_0^{(2)}$ -formula with no second-order quantifier and with only displayed variables free in which  $X$  occurs only positively. Then, there exists a  $\mathcal{L}_T$ -formula  $\phi(x, \vec{v})$  such that

$$\text{PUTB}^- \vdash \phi(x, \vec{v}) \leftrightarrow \Phi(\hat{u}\phi(u, \vec{v}), x, \vec{v}),$$

where  $\Phi(\hat{u}\phi(u, \vec{v}), x, \vec{v})$  is the result of substituting  $\phi(t, \vec{v})$  for each occurrence of  $t \in X$  in  $\Phi$ .

*Proof.* By parametrized diagonalization (applied to a  $\mathcal{L}_T$ -formula  $\Phi(\hat{u}T a(u/\ulcorner x \urcorner, \vec{v}/\ulcorner \vec{v} \urcorner), x, \vec{v})$  w.r.t. the variable  $a$ ), we can find a  $\gamma$  such that

$$\text{PUTB}^- \vdash \gamma(x, \vec{v}) \leftrightarrow \Phi(\hat{u}T^\ulcorner \gamma(\hat{u}, \vec{v}) \urcorner, x, \vec{v})$$

From the construction of the diagonal formula  $\gamma$ , we can assume that  $\gamma$  is  $T$ -positive. Hence, the PUTB schema implies

$$\text{PUTB}^- \vdash \gamma(x, \vec{v}) \leftrightarrow \Phi(\hat{u}\gamma(u, \vec{v}), x, \vec{v})$$

We can take this  $\gamma$  as the desired  $\phi$ . □

It is observed by Feferman [17] that the truth predicate of KF can be defined by a fixed-point of a certain positive arithmetical operator; Feferman made use of this fact to show that  $\Sigma_1^1\text{-AC}$ , within which a fixed-point (but not necessarily the least) of any positive  $\Sigma_1^1$ -operator can be constructed, can define the truth of KF. Since full induction isn't used in Feferman's construction, we indeed carry out the same argument in  $\text{PUTB}^-$ :

**Theorem 3.4.11.**  $\text{KF}^-$  is truth-definable in  $\text{PUTB}^-$ .

Recall that, since a base system  $\mathbf{B}$  is arbitrary and thus we can replace  $\mathbf{B}$  by  $\mathbf{B}(P)$ , truth-definability carries over to schematic reflective closures (also see the Remark in §3.3).

**Corollary 3.4.12.**  $\text{KF}$  and  $\text{KF}^*$  are truth-definable in  $\text{PUTB}$  and  $\text{PUTB}^*$  respectively.

**Theorem 3.4.13.**  $\text{PUTB}^- \equiv \mathbf{B}$ ,  $\text{PUTB} \equiv \text{KF}$  and  $\text{PUTB}^* \equiv \text{KF}^*$ . In particular, we have  $\text{PUTB}^*[[\text{PA}]] \equiv \text{RA}_{<\Gamma_0}$ .

It follows from Feferman's consistency proof of DT in [19] that the truth of  $\text{FKF}^-$  is also defined as an arbitrary fixed-point of a certain positive arithmetical operator. Thus,  $\text{FKF}$  and  $\text{FKF}^-$  are truth definable in  $\text{PUTB}$  and  $\text{PUTB}^-$  respectively. However, as we have noted in §2.5, fixed-points constructible in  $\text{PUTB}^{(-)}$  (or  $\Sigma_1^1\text{-AC}$ ) are not necessarily consistent. Thus, these techniques are not directly applied to obtain that systems with Cons are truth-definable in  $\text{PUTB}^{(-)}$ . Fortunately, we can nonetheless modify these techniques so that  $\text{PUTB}^-$  can define the truth of  $\text{FKF}^- + \text{Cons}$ , i.e.,  $\text{DT}^-$ .<sup>15</sup>

**Theorem 3.4.14.**  $\text{DT}^-$  is truth-definable in  $\text{KF}^-$  (and thus in  $\text{PUTB}^-$ ).

*Proof.* First, using the recursion theorem, we define (in PA) a preliminary translation  $I$  (on the codes of formulae) such that:

---

<sup>15</sup>In contrast, it is still open whether  $\text{PUTB}^-$  can define the truth of  $\text{KF}^- + \text{Cons}$ .

$$I(\phi) = \begin{cases} \phi & \text{if } \phi \text{ is } t = s \text{ or } t \neq s, \\ TI\dot{t} & \text{if } \phi \text{ is } Tt, \\ TI\dot{\neg}t & \text{if } \phi \text{ is } \neg Tt, \\ I(\psi) & \text{if } \phi \text{ is } \neg\neg\psi \\ [I\psi_0 \wedge I\psi_1 \wedge \neg I\neg\psi_0 \wedge \neg I\neg\psi_1] & \\ \quad \vee [I\neg\psi_0 \wedge I\psi_1 \wedge \neg I\psi_0 \wedge \neg I\neg\psi_1] & \\ \quad \vee [I\psi_0 \wedge I\neg\psi_1 \wedge \neg I\neg\psi_0 \wedge \neg I\psi_1] & \text{if } \phi \text{ is } \psi_0 \vee \psi_1 \\ I\neg\psi_0 \wedge I\neg\psi_1 \wedge \neg I\psi_0 \wedge \neg I\psi_1 & \text{if } \phi \text{ is } \neg(\psi_0 \vee \psi_1) \\ [I\psi_0 \wedge I\psi_1 \wedge \neg I\neg\psi_0 \wedge \neg I\neg\psi_1] \vee [I\neg\psi_0 \wedge \neg I\psi_0] & \text{if } \phi \text{ is } \psi_0 \rightarrow \psi_1, \\ I\psi_0 \wedge I\neg\psi_1 \wedge \neg I\neg\psi_0 \wedge \neg I\psi_1 & \text{if } \phi \text{ is } \neg(\psi_0 \rightarrow \psi_1), \\ \forall x(I\psi \wedge \neg I\neg\psi) & \text{if } \phi \text{ is } \forall x\psi, \\ \forall x((I\psi \wedge \neg I\neg\psi) \vee (I\neg\psi \wedge \neg I\psi)) \wedge \exists x I\neg\psi & \text{if } \phi \text{ is } \neg\forall x\psi \\ 0 & \text{otherwise,} \end{cases}$$

where  $I$  is the representation of  $I$ . Let  $\theta(x)$  be

$$TIx \wedge \neg TI\neg x \wedge FI\neg x \wedge \neg FIx.$$

Then we will show that this  $\theta$  defines the truth of  $\text{FKF}^- + \text{Cons}$  in  $\text{KF}^-$ . This can be done in a straightforward but somewhat tedious manner; we will show some typical cases.

First, note that  $\mathcal{T}_\theta(Ft)$  is  $TI\neg t \wedge \neg TI\dot{t} \wedge FI\dot{t} \wedge \neg FI\neg t$ . Then, it immediately follows that  $\text{KF}^- \vdash \mathcal{T}_\theta(\text{Cons}) \wedge \mathcal{T}_\theta(\mathbf{K2})$ . By definition of  $I$ , we obviously have  $\text{KF}^- \vdash \mathcal{T}_\theta(\mathbf{K1}) \wedge \mathcal{T}_\theta(\mathbf{K3})$  as well.

For **FK4**, let  $x, y \in \text{St}_{\mathcal{L}_T}$ . Then,  $\mathcal{T}_\theta(Tx \vee y)$  is equivalent to

$$\begin{aligned}
& [(TIx \wedge TIy \wedge FI\bar{x} \wedge FI\bar{y}) \vee (TI\bar{x} \wedge TIy \wedge FIx \wedge FI\bar{y}) \\
& \quad \vee (TIx \wedge TI\bar{y} \wedge FI\bar{x} \wedge FIy)] \\
& \wedge [\neg TI\bar{x} \vee \neg TI\bar{y} \vee \neg FIx \vee \neg FIy] \\
& \wedge [FI\bar{x} \vee FI\bar{y} \vee TIx \vee TIy] \\
& \wedge [(\neg FIx \wedge \neg FIy \wedge \neg TI\bar{x} \wedge \neg TI\bar{y}) \vee (\neg FI\bar{x} \wedge \neg FIy \wedge \neg TIx \wedge \neg TI\bar{y}) \\
& \quad \vee (\neg FIx \wedge \neg FI\bar{y} \wedge \neg TI\bar{x} \wedge \neg TIy)]
\end{aligned}$$

Since the second and third clauses are entailed from the fourth and first respectively, they can be dropped. Then, the resulting formula is equivalent to:

$$\begin{aligned}
& [TIx \wedge TIy \wedge FI\bar{x} \wedge FI\bar{y} \wedge \neg FIx \wedge \neg FIy \wedge \neg TI\bar{x} \wedge \neg TI\bar{y}] \\
& \quad \vee [TI\bar{x} \wedge TIy \wedge FIx \wedge FI\bar{y} \wedge \neg FI\bar{x} \wedge \neg FIy \wedge \neg TIx \wedge \neg TI\bar{y}] \\
& \quad \vee [TIx \wedge TI\bar{y} \wedge FI\bar{x} \wedge FIy \wedge \neg FIx \wedge \neg FI\bar{y} \wedge \neg TI\bar{x} \wedge \neg TIy)].
\end{aligned}$$

This formula is equivalent to  $\mathcal{T}_\theta((Tx \wedge Ty) \vee (Tx \wedge Fy) \vee (Fx \wedge Ty))$ . We finally obtain  $\text{KF}^- \vdash \mathcal{T}_\theta(\mathbf{FK4})$ . □

**Theorem 3.4.15.** DT and DT\* are truth-definable in KF and KF\* respectively (and thus in PUTB and PUTB\* respectively).

From the truth-definability of PUTB<sup>-</sup> and DT<sup>-</sup> in KF<sup>-</sup>, we know:

**Corollary 3.4.16.** If B is reflexive, then PUTB<sup>-</sup> and DT<sup>-</sup> are reflexive. PUTB and DT are also (essentially) reflexive since they contain full induction.

However, it still remains open whether the converses of Theorem 3.4.15 hold. Concluding this

section, we restate the two open problems for future work.

**Open Problem 1:** Is  $\text{PUTB}^-$  (or  $\text{PUTB}$ ) truth-definable in  $\text{DT}^-$  ( $\text{DT}$ , resp.)?

**Open Problem 2:** Is  $\text{UTB}$  truth-definable in  $\text{TB}$ ?

### 3.5 Some Applications

In the present section, we exhibit some applications of the results and techniques on relative truth definability developed so far. First of all, as we have announced, we give positive answers to both the first and the second conjectures of Feferman.

**Theorem 3.5.1.**  $\text{DT}[\text{PA}] \equiv \text{RA}_{<\varepsilon_0}$  and  $\text{DT}^*[\text{PA}] \equiv \text{RA}_{<\Gamma_0}$ .

*Proof.* The upper bound of  $\text{DT}[\text{PA}]$  immediately follows from Theorem 3.4.15. The converse follows from Theorems 2.3.7 and 3.4.6.

Similarly, the upper bound of  $\text{DT}^*[\text{PA}]$  is also obtained by Theorem 3.4.15, and it follows from Theorem 2.3.7 and 3.4.9 that this bound is exact.  $\square$

As the next application, we will show some proof-theoretic equivalences among systems of iterative compositional truth.

Feferman [17, p.28] suggested that the proof-theoretic strength of  $\text{WKF}[\text{PA}]$  is at least as strong as  $\text{RA}_{<\varepsilon_0}$ . This is indeed correct since  $\text{WKF}$  meets the conditions (i)-(viii) of Lemma 3.4.5 and proves  $\text{TI}_{\mathcal{L}_T}(< \varepsilon_0)$ . For the converse, we can show that a simple interpretation which translates ‘ $T(x \rightarrow y)$ ’ by ‘ $T(\neg x \vee y)$ ’ is a truth-definition of  $\text{WKF}^- (+\text{Cons})$  in  $\text{FKF}^- (+\text{Cons}, \text{resp.})$ ; note that we use the recursion theorem to define this translation. Hence, by Theorems 2.4.2 and 3.4.15, we know the bound is exact. Consequently, we have the following:

**Theorem 3.5.2.** (1)  $\text{WKF}^-$  and  $\text{WKF}^- + \text{Cons}$  are truth-definable in  $\text{KF}^-$ .

(2)  $\text{WKF}^- \equiv \text{WKF}^- + \text{Cons} \equiv \text{B}$ .

(3)  $\text{WKF}[\text{PA}] \equiv \text{WKF}[\text{PA}] + \text{Cons} \equiv \text{RA}_{<\varepsilon_0}$ .

(4)  $\text{WKF}^*[\text{PA}] \equiv \text{WKF}^*[\text{PA}] + \text{Cons} \equiv \text{RA}_{<\Gamma_0}$ .

We have so far focused on Cons. Then what would happen when we add Comp? As a matter of fact, Comp collapses all the differences among KF, FKF and WKF: that is, we have the following.

**Lemma 3.5.3.**  $\text{KF}^- + \text{Comp}$ ,  $\text{FKF}^- + \text{Comp}$  and  $\text{WKF}^- + \text{Comp}$  are all identical.

*Proof.* The proof is straightforward. For example, observe that  $Tx \vee Ty$  is equivalent to  $(Tx \wedge Fx) \vee (Fx \wedge Ty) \vee (Tx \wedge Ty)$  under Comp.  $\square$

In the rest of the present section, we will determine the proof-theoretic strength of the schematic reflective closures of  $\text{RT}_{<\alpha}$  and FS by using the results and techniques on truth-definability developed so far. As we shall see, schematic reflective closure does not increase proof-theoretic strength in the cases of FS and  $\text{RT}_{<\alpha}$  with  $\alpha = \omega^\xi$  for some  $\xi > 0$ .

For each  $k \in \mathbb{N}$  and ordinal  $\alpha$ ,  $(\text{RT}_{<\alpha})^* \upharpoonright_k$  denotes the system obtained from  $(\text{RT}_{<\alpha})^*$  by restricting the number of applications of  $P$ -Subst in each derivation to at most  $k$  times; therefore  $(\text{RT}_{<\alpha})^* \upharpoonright_0$  is just identical with  $\text{RT}_{<\alpha}[\mathbf{B}(P)]$ . In the following, we assume that derivations are formulated in Hilbert style in which inference rules are Modus Ponens and  $P$ -Subst only.

**Lemma 3.5.4.**  $(\text{RT}_{<\alpha})^* \upharpoonright_{k+1}$  is proof-theoretically reducible to  $(\text{RT}_{<\alpha+\alpha})^* \upharpoonright_k$  for  $\mathcal{L}_\alpha(P)$ .

*Proof.* Take an arbitrary derivation  $\mathcal{D} = \langle \theta_0, \dots, \theta_n \rangle$  in  $(\text{RT}_{<\alpha})^* \upharpoonright_{k+1}$ . Let  $\mathcal{D}_0 = \langle \theta_0, \dots, \theta_m \rangle$  ( $m \leq n$ ) be the initial segment of  $\mathcal{D}$  which ends with the first application of  $P$ -Subst  $\frac{\theta_i(P)}{\theta_j(\psi)}$  for some  $\psi \in \mathcal{L}_\alpha(P)$ : thus  $\theta_m$  is  $\theta_j(\hat{\psi})$  ( $j < m$ ). Then, there exist  $\beta, \gamma \leq \alpha$  such that  $\theta_0, \dots, \theta_m \in \mathcal{L}_\beta$  and  $\psi \in \mathcal{L}_\gamma$ . We define a translation  $\mathcal{S}_\psi^\beta$  from  $\mathcal{L}_\beta(P)$  to  $\mathcal{L}_{\gamma+\beta}(P)$ , which ‘lifts up’ the level of truth predicates in each  $\theta_i$  ( $i < m$ ) so that  $\psi$  can be treated as if it were at the lowest level; note that

$\gamma + \beta \leq \alpha + \alpha$ . First we define a preliminary translation  $g_\psi^\beta$  (on the codes of formulae) by the recursion theorem such that:

$$g_\psi^\beta(\phi) := \begin{cases} \phi & \text{if } \phi \text{ is an } \mathcal{L}_0\text{-atomic} \\ \psi(t) & \text{if } \phi \text{ is } P(t) \\ T_{\gamma+\delta}(g_\psi^\beta(t)) & \text{if } \phi \text{ is } T_\delta(t) \\ \neg g_\psi^\beta(\theta) & \text{if } \phi \text{ is } \neg\theta \\ g_\psi^\beta(\theta_0) \vee g_\psi^\beta(\theta_1) \text{ (} g_\psi^\beta(\theta_0) \rightarrow g_\psi^\beta(\theta_1), \text{ resp.)} & \text{if } \phi \text{ is } \theta_0 \vee \theta_1 \text{ (or } \theta_0 \rightarrow \theta_1) \\ \forall z.g_\psi^\beta(\theta) & \text{if } \phi \text{ is } \forall z\theta(z), \end{cases}$$

where  $g_\psi^\beta$  is the representation of  $g_\psi^\beta$ . We notice that  $\phi \in \mathcal{L}_\delta$  implies  $g_\psi^\beta(\phi) \in \mathcal{L}_{\gamma+\delta}$ . The translation  $\mathcal{S}_\psi^\beta$  is thereby defined as:  $\mathcal{S}_\psi^\beta(R\vec{x})$  is  $R\vec{x}$  for each atomic formula of  $\mathcal{L}_0$ ,  $\mathcal{S}_\psi^\beta(T_\delta x)$  is  $T_{\gamma+\delta}(g_\psi^\beta(x))$ ,  $\mathcal{S}_\psi^\beta(Px)$  is  $\psi(x)$ , and  $\mathcal{S}_\psi^\beta$  commutes with connectives and quantifiers.

Then, we can straightforwardly show that  $\text{RT}_{<\gamma+\beta}[\mathbb{B}(P)] \vdash \mathcal{S}_\psi^\beta(\theta_i)$  for all  $i < m$ . Since  $\theta_j(P) \in \mathcal{L}_0(P)$  and  $\mathcal{S}_\psi^\beta$  keeps the  $\mathcal{L}_0$ -part unchanged, we have  $\mathcal{S}_\psi^\beta(\theta_j) = \theta_j(\hat{\psi})$  and thus  $\text{RT}_{<\gamma+\beta}[\mathbb{B}(P)] \vdash \theta_m$ . Let  $\mathcal{D}_1$  be this derivation of  $\theta_m$  in  $\text{RT}_{<\gamma+\beta}[\mathbb{B}(P)]$ . Set  $\mathcal{D}' := \langle \theta_0, \dots, \theta_{m-1} \rangle * \mathcal{D}_1 * \langle \theta_{m+1}, \dots, \theta_n \rangle$ , where  $*$  denotes the concatenation of sequences. Then  $\mathcal{D}'$  is a derivation of  $\theta_n$  in  $(\text{RT}_{<\alpha+\alpha})^* \upharpoonright_k$ .  $\square$

**Lemma 3.5.5.** Let  $\alpha$  be a limit principal ordinal, i.e.,  $\alpha = \omega^\xi$  for some  $\xi > 0$ . Then,  $(\text{RT}_{<\alpha})^*[\mathbb{B}]$  and  $\text{RT}_{<\alpha}[\mathbb{B}(P)]$  are identical (as systems).

*Proof.* Suppose  $(\text{RT}_{<\alpha})^* \upharpoonright_k \vdash \phi$ . Since  $\alpha$  is limit, there exist  $\beta < \alpha$  and  $k \in \mathbb{N}$  such that  $\phi \in \mathcal{L}_\beta$  and  $(\text{RT}_{<\beta})^* \upharpoonright_k \vdash \phi$ . Then, it follows from the last lemma that  $(\text{RT}_{<\beta, 2^k})^* \upharpoonright_0 \vdash \phi$ . Since  $\alpha$  is principal, we have  $\text{RT}_{<\alpha}[\mathbb{B}(P)] \vdash \phi$ .  $\square$

**Theorem 3.5.6.**  $(\text{RT}_{<\alpha})^* \upharpoonright_k$  is proof-theoretically reducible to  $\text{RT}_{<\alpha \cdot 2^k}$  for  $\mathcal{L}_0$ -formulae. Hence, in particular,  $\text{TC}^* \leq \text{RT}_{<\omega}[\mathcal{L}_0]$ .

*Proof.* Suppose  $(\text{RT}_{<\alpha})^* \upharpoonright_k \vdash \phi$  for  $\mathcal{L}_0$ -formula  $\phi$ . By Lemma 3.5.4, we have  $\text{RT}_{<\alpha \cdot 2^k}[\mathbb{B}(P)] \vdash \phi$  by some derivation  $\mathcal{D} = \langle \theta_0, \dots, \theta_n \rangle$  ( $\theta_n$  is equal to  $\phi$ ). Let  $\beta \leq \alpha \cdot 2^k$  be an ordinal such that

$\theta_0, \dots, \theta_n \in \mathcal{L}_\beta$  and let  $\psi$  be an arbitrary  $\mathcal{L}_0$ -formula (e.g.,  $(x)_0 = (x)_1$ ). Take a translation  $\mathcal{S}_\psi^\beta$  from  $\mathcal{L}_\beta(P)$  to  $\mathcal{L}_\beta$  in the same way as Lemma 3.5.4. Then, we can show that  $\text{RT}_{<\beta}[\mathbf{B}] \vdash \mathcal{S}_\psi^\beta(\theta_n)$ . Since  $\theta_n \in \mathcal{L}_0$  and  $\mathcal{S}_\psi^\beta$  keeps the  $\mathcal{L}_0$ -part unchanged, we finally have  $\text{RT}_{<\alpha \cdot 2^k}[\mathbf{B}] \vdash \theta_n$ .  $\square$

**Theorem 3.5.7.** Let  $\alpha$  be a limit principal ordinal. Then,  $(\text{RT}_{<\alpha})^*$  is proof-theoretically equivalent to  $\text{RT}_{<\alpha}$  for  $\mathcal{L}_0$ -formulae.

*Proof.* Suppose  $(\text{RT}_{<\alpha})^* \vdash \phi$  for  $\mathcal{L}_0$ -formula  $\phi$ . Since  $\alpha$  is limit, there exist  $\beta < \alpha$  and  $k \in \mathbb{N}$  such that  $(\text{RT}_{<\beta})^* \upharpoonright_k \vdash \phi$ . Then, by the last theorem, we have  $\text{RT}_{<\beta \cdot 2^k} \vdash \phi$ . Since  $\alpha$  is principal, we obtain  $\text{RT}_{<\alpha} \vdash \phi$ .  $\square$

In the same way as Corollary 3.4.7, we can show that:

**Corollary 3.5.8.** If  $\alpha \leq \Gamma_0$ , then  $\text{DT}^*, \text{KF}^* \not\leq (\text{RT}_{<\alpha})^*$ .

Now we can determine the proof-theoretic strength of  $\text{FS}^*$ .

**Theorem 3.5.9.**  $\text{FS}^*$  is proof-theoretically equivalent to  $\text{RT}_{<\omega}$ . Thus, in particular,  $\text{FS}^*[\mathbf{PA}] \equiv \text{RA}_{<\omega}$  (and thus  $\equiv \text{ACA}_0 + \text{BR}$ ; cf. fn.26 and Theorem 3.5.6).

*Proof.* Halbach [31, §5] already showed that  $\text{RT}_{<\omega}$  is truth-definable in  $\text{FS}$ , and thus  $(\text{RT}_{<\omega})^*$  is truth-definable in  $\text{FS}^*$  by Lemma 3.2.5.

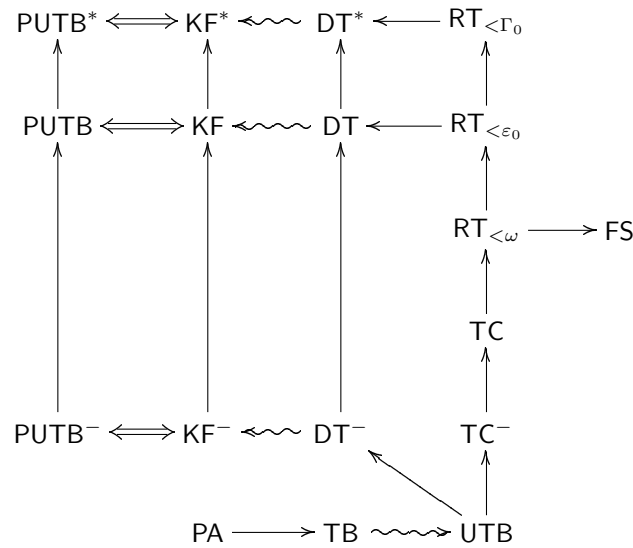
For the converse, Halbach [31, §5] also gave a truth-definition  $g_n$  of  $\text{FS}_n$  in  $\text{RT}_{<2n}$ , where  $\text{FS}_n$  is obtained from  $\text{FS}$  by restricting the numbers of applications of  $\text{NEC}$  and  $\text{CONEC}$  in each derivation to at most  $(n-1)$ -times respectively. It is observed from his proof that if  $(\text{FS}_n)^* \vdash \phi$  then  $(\text{RT}_{<2n})^* \vdash g_n(\phi)$ ; for, it can be shown in the same way as [31] (the extra cases for  $P$ -Subst are obvious) that  $(\text{FS}_i)^* \vdash \phi$  implies  $(\text{RT}_{<2n})^* \vdash g_i(\phi) \wedge \dots \wedge g_{2n-i}(\phi)$  for each  $i \leq n$ . This implies that  $\text{FS}^* \leq (\text{RT}_{<\omega})^*$  and thus  $\text{RT}_{<\omega}$  by Theorem 3.5.7.<sup>16</sup>  $\square$

---

<sup>16</sup>Halbach's truth-definitions were for  $\mathbf{B} = \text{PA}$ , but they can be straightforwardly generalized for an arbitrary  $\mathbf{B}$  which meets our conditions for base systems.

### 3.6 Concluding Remarks of Chapter 3

The results on relative truth definability obtained so far are partially summarized by the following diagram: for  $B = PA$ ,<sup>17</sup>



In this diagram, for each system  $Q$  and  $S$ , ' $Q \leftarrow S$ ' means  $Q \succ S$ , ' $Q \Leftrightarrow S$ ' means  $Q \preceq S$  and  $Q \succeq S$ , and the curly arrows indicate that it is still open whether the converse holds. The systems in the same horizontal line are all proof-theoretically equivalent (for  $\mathcal{L}_0$ -sentences); but the systems in the 5th and 6th lines are equivalent as well.

Most of the results on relative truth definability so far are one-way and mutual truth definability fails in many cases. This is what we expected, since we introduced the notion in order to distinguish

<sup>17</sup>As a matter of fact, the diagram applies to an arbitrary base system  $B$  (subject to the conditions (a)-(d) in §2.1). However, in order to establish this, we need to show  $KF^* \not\preceq KF$  and  $DT^* \not\preceq DT$  for arbitrary  $B$  in general, but we have not proved them in the present paper. Its proof needs extra technical arguments which go beyond the scope of the present paper. I only give a sketch of the proof here. We first formulate  $KF$  (or  $DT$ ) in a certain suitable semi-formal system with  $\omega$ -rule. Then, we can show that if  $KF$  ( $DT$ )  $\vdash \phi$  for an arithmetical  $\phi$  then  $\phi$  is provable in the semi-formal system by a derivation with the length  $< \varphi_{\varepsilon_0} 0$  which contains no application of cut to formulae with  $T$  (but may contain cuts to arithmetical formulae). Since  $KF^*$  ( $DT^*$ ) can prove  $TI_{\mathcal{L}_T}(\varphi_{\varepsilon_0} 0)$ , it can model those infinitary derivations of arithmetical formulae by a formula  $Tx$ : i.e.,  $KF^*$  ( $DT^*$ ) can prove that if an arithmetical  $\phi$  is provable in  $KF$  ( $DT$ ) then  $T^\top \phi^\top$  holds. Thus,  $KF^*$  ( $DT^*$ ) can prove the consistency of  $KF$  ( $DT$ , resp.).

the systems of truth which have the same truth-free theorems but apparently embody different views of truth. Concluding the present paper, we make some comments on the formal results obtained so far.

First, it is observed that ramified truth systems and type-free truth systems are clearly differentiated in view of relative truth definability. Type-free truth systems are not relatively truth definable in ramified truth systems even with much stronger proof-theoretic strength; e.g.,  $\text{KF}^- \not\leq \text{RT}_{<\varepsilon_0}$ .

Second, relative truth definability can distinguish disquotational systems like UTB and compositional systems like  $\text{TC}^-$ . However, we have also seen that an iterative compositional system KF and a disquotational system PUTB are equivalent in view of relative truth definability. On the one hand, this fact may be regarded as a defect of relative truth definability as a tool for comparing conceptual aspects of systems of truth, since KF and PUTB look different in some respects. For example, although PUTB can define the truth of KF, it can neither prove the compositionality nor iterativity of its own truth predicate in contrast to KF: i.e., PUTB does not prove statements like

$$x, y \in \text{St}_{\mathcal{L}_T} \rightarrow [T(x \forall y) \leftrightarrow Tx \vee Ty] \quad \text{or} \quad x \in \text{CT} \rightarrow [T \neg Tx \leftrightarrow T \neg \text{val}(x)].$$

For another example, PUTB can be consistently closed by NEC and CONEC, whereas it is not the case for KF; indeed, the resulting system is even conservative over PUTB for truth-free theorems, and, furthermore, the same holds even when we add uniform versions of NEC and CONEC, in which the premises and conclusions of the rules can contain parameters. On the other hand, their mutual truth definability may still be interpreted to indicate that KF and PUTB are indeed ‘conceptually equivalent’ despite these ‘superficial’ differences.

Third, it is suggested that iterative compositional systems of truth with weaker evaluation schemata are ‘reducible’ to those with stronger schemata in a stronger sense than mere relative interpretability, proof-theoretic reducibility and others. But it is still open whether this reduction is

one-way. The solution to the Open Problem 1 would answer the question whether iterative compositional systems of truth with different schemata indeed give different or incompatible conceptions of truth.

Fourth, as we have seen in §5, relative truth definability has useful applications. In particular, relative truth definability makes it easier to compare the schematic reflective closures of systems and mutual truth definability gives more information on model-theoretic features of systems of truth, since it entails conservativity in the semantic sense.

## Chapter 4

# The Inner Theory of Truth

## Systems

### 4.1 The Notion of Inner Theory

Let us start by repeating some historical issues to explain and clarify the motivation of the present chapter. Tarski's theory of truth [82] showed us how the notoriously difficult concept of truth could be formalized in the modern formal framework. As a consequence of his theory, he suggested a negative thesis about formalized truth: i.e., the thesis that no language can possess its own truth predicate and it could be contained only in other essentially richer languages. This thesis naturally leads to the 'orthodox approach' or the 'hierarchical view' of truth [50].

It took four decades since Tarski before a new proposal of formalizing another conception of truth was submitted. Among others (e.g., [56] and [84]), Kripke [50] challenged the Tarskian hierarchical view of truth and presented a method for formalizing a new conception of truth in which the truth predicate is properly applicable to sentences containing the truth predicate itself.

Kripke gave a model-theoretic semantics of a language which contain its own truth predicate and thus his theory counts a semantical theory of truth. Soon after that , an axiomatic treatment of self-applicable truth was presented by Feferman [9, 17]; his system is called Kripke-Feferman system KF today. Since KF, various axiomatic systems of self-applicable truth have been presented and much work has been done on proof-theoretic the analysis of those systems. As I have discussed in Chapter 1, the self-applicable character of the notion of truth can be represented in no definitional theories of truth and one of the major advantages of the non-definitional approach consists in the fact they can properly deal with it. The present chapter particularly focuses on axiomatic systems of self-applicable truth and proposes one means of comparison among them.

Since a system of self-applicable truth has its own truth predicate, it is naturally motivated to distinguish and compare a statement ‘ $A$ ’ and the assertion of its truthhood ‘the statement  $A$  is true’ within the same system; both types of statements are expressible in such a system of self-applicable truth. This distinction and comparison were first explicitly made by Reinhardt [70, 71]. He studied the class of the statements  $\sigma$  which are *asserted to be true* or *provably true* in KF, i.e.,  $\{\sigma \in \mathcal{L}_T \mid \text{KF} \vdash T^\top \sigma^\top\}$ . This class is called the *inner theory* of KF. This approach of Reinhardt and the notion of inner theory can be naturally generalized to any other system  $\mathbf{Q}$  of self-applicable truth: i.e., the inner theory of  $\mathbf{Q}$  is  $\{\sigma \mid \sigma \text{ is a sentence of } \mathbf{Q} \text{ and } \mathbf{Q} \vdash T^\top \sigma^\top\}$ . Let us denote it by  $\mathcal{I}(\mathbf{Q})$ . In contrast against the inner theory, we call the set of theorems of  $\mathbf{Q}$  the *outer theory* of  $\mathbf{Q}$  and denote it by  $\mathcal{O}(\mathbf{Q})$ : i.e.,  $\mathcal{O}(\mathbf{Q}) = \{\sigma \mid \mathbf{Q} \vdash \sigma\}$ . As we will argue in the following, the notion of inner theory has philosophical and mathematical significance in some respects. The objective of the present chapter is to investigate the inner theories of various systems of self-applicable truth and correlate these systems in terms of their inner theories.

Now let us look at Reinhardt’s original motivation of his study of inner theories. Reinhardt interpreted KF as a system of *meaningful applicability*, the notion of which he ascribed to Gödel [25]. This idea and its motivation are briefly summarized by the following passage: ‘[T]he paradoxes can

perhaps be resolved by noting that some properties and relations ... may not always be meaningfully applicable ... [70, p.227]'; in other words, Reinhardt suggests that the paradoxes such as Russell's paradox and the Liar paradox may be resolved by restricting the application of the principles which leads to these paradoxes (i.e., the full comprehension axiom of set theory and the full T-schema in truth theory) only to a certain domain the objects of which those principles are *meaningfully applicable* to. First, Reinhardt argued that a sentence  $\sigma$  is meaningfully applicable or *significant* (they are equivalent notions for Reinhardt) if and only if  $\sigma$  has determinate truth value, i.e., truth or falsity.<sup>1</sup> Then, Reinhardt proposed to use KF as a tool of producing true (and thus significant) statements and make the following thesis:

**(R1)** If  $A$  is a sentence such that  $\text{KF} \vdash T[A]$ , then  $A$  is true.

**(R1)** If  $A$  is a sentence such that  $\text{KF} \vdash T[A] \vee F[A]$ , then  $A$  is significant.

Hence, for Reinhardt, the inner theory of KF comprises true and significant sentences. Reinhardt makes an analogy between his project and Hilbert's programme; in this analogy of his,  $\mathcal{O}(\text{KF})$  and  $\mathcal{I}(\text{KF})$  respectively corresponds to Hilbert's 'ideal' mathematics and 'real' (or 'contentual') mathematics. Then, Reinhardt raised two questions of KF ([70, p.239]):

**(1)** Is there a natural axiomatization of  $\mathcal{I}(\text{KF})$  which consists of only significant sentences?

**(2)** If  $\text{KF} \vdash T^\top \sigma^\top$ , is there a proof  $\mathcal{D} = \langle \theta_0, \dots, \theta_n \rangle$  of  $\sigma$  such that  $\text{KF} \vdash T^\top \theta_i^\top$  for each  $i \leq n$ ?

The question (1) is important for Reinhardt since a positive answer to it implies that we have natural axiomatization of the set of true statements by true sentences. A positive answer to (2) 'justifies the use of insignificant sentences' as a mere tool for deducing true sentences 'entirely within the framework of the significant sentences'. Here we can again see an analogy between Reinhardt's project and Hilbert's programme; (2) corresponds to the conservativity requirement of Hilbert's

---

<sup>1</sup>This view is also taken by Feferman and motivated his theory DT, though Feferman's motivation is probably independent of Reinhardt.

programme. Consequently, since we can show that the desired truth-theoretic principles such as the T-schema and compositional axioms (or inductive clauses) hold of significant sentences within KF, KF would then be regarded as a fairly satisfactory system of truth for Reinhardt.

Of course, the answer to (1) depends on what ‘natural’ means there; but, at any rate, it is easily seen that the answer to (1) would be negative if the required axiomatization need to be based on classical logic, since  $\lambda \vee \neg\lambda \notin \mathcal{I}(\text{KF})$  where  $\lambda$  is the liar sentence. To make the matter worse, (2) was already negatively answered by Halbach and Horsten [38]. However, at the time of Reinhardt’s writing [70] and [71], no alternative axiomatic system of self-applicable truth to KF was known, whereas we now have much more variety of such systems. Thus, we may propose the thesis like **(R1)** and **(R2)** with respect to other systems of self-applicable truth and then ask similar questions to (1) and (2) for them. We have already seen that the three systems of symmetric truths, i.e., FS, POSKF + NEC + CONEC and PUTB + NEC + CONEC (as well as their variants), and they all meet both (1) and (2); in addition, more systems such as KF + Comp meet (2) although they fail to meet (1). It is to be noted, however, that the positive answers to (1) and (2) were of importance to Reinhard because he held the theses **(R1)** and **(R2)** particularly with respect to KF. Hence, it is still debatable whether other systems can be a substitute for KF for Reinhardt’s project. At any rate, Reinhardt’s project can be generalized and has not yet been put an end.

The notion of inner theory thus arose from Reinhardt’s specific project and philosophical motivation. However, besides Reinhardt’s project, it seems to be philosophically and mathematically significant in other respects. In what follows, we will explain some of them.

From the mathematical point of view, the notion of inner theory could be a useful tool for comparing systems of self-applicable truth. As in the case of relative truth definability, comparison by inner theory brings more information about systems of self-applicable truth than comparison by the mere derivability of truth-free theorems. I repeatedly emphasize that the proof-theoretic analyses so far given only focus on the truth-free part of systems of truth. However, we often

compare two subsystems of second-order arithmetic by their derivability of a certain larger class of formulae; this type of comparison is often found in the subject of reverse mathematics (cf. [77]). For example,  $\Sigma_1^1$ -DC and  $(\Pi_1^0\text{-CA})_{<\varepsilon_0}$ , which are both proof-theoretically equivalent to **KF**, prove the same  $\Pi_2^1$ -sentences but the former has properly more  $\Pi_3^1$  theorems than the latter. In view of such a practice in mathematical logic, it seems natural to consider and compare by larger classes of theorems of truth systems than that of arithmetical ones.

One natural generalization of the existing approach is perhaps to expand an arithmetical hierarchy by incorporating the truth predicate and then to compare systems in the light of this new hierarchy. We can define classes  $\Pi_0^T$  and  $\Sigma_0^T$  of formulae as the least set containing all arithmetical atomics and those of the form  $Tt(\vec{x})$  closed under Boolean combinations and bounded quantifications; then  $\Sigma_n^T$  and  $\Pi_n^T$  can be defined in parallel manners to the ordinary arithmetical hierarchy. This direction should be of some interest. However, in many cases, systems of truth are incomparable even for  $\Pi_1^T$  or  $\Sigma_1^T$  class: that is, the class of  $\Pi_1^T$  ( $\Sigma_1^T$ ) theorems of one system neither includes nor is included in that of another system in many cases. As we shall see in §3, as far as the systems considered in this chapter concern, they are all incomparable by means of thus defined expanded arithmetical hierarchy except for the trivial cases in which one is a subtheory of another.

In contrast to the comparison of the above expanded arithmetical hierarchy, many theories are suitably comparable to one another by means of their inner theories. Hence, I suggest the comparison of inner theories as a better behaved tool to compare systems of self-applicable truth. I also notice that, since every system considered in the present thesis derives T-biconditionals for arithmetical (truth-free) formulae, the inner theory comparison is properly finer grained than the mere comparison by arithmetical (truth-free) consequences.

Comparison of systems of self-applicable truth by their inner theories also has a philosophical motivation. Different systems of truth embody different views or conceptions of truth. However, as we have extensively discussed in the last chapter, comparison only by the arithmetical part

often bring no difference among systems of truth based on different views on truth. Now, let us consider KF and WKF for instance. They are proof-theoretically equivalent, but we probably have the intuition (or at least the feeling) that weak Kleene logic is weaker than strong Kleene logic in some sense (judging not simply from their names ‘weak Kleene’ and ‘strong Kleene’ but from their truth tables) and thus WKF should be weaker than KF in some sense. This intuition is correctly reflected by their inner theories; we will show that  $\mathcal{I}(\text{DT}) \subsetneq \mathcal{I}(\text{KF})$ , and thus WKF is indeed ‘weaker’ than KF in view of their inner theories.

The present chapter is basically a technical one but motivated by the aforementioned rather philosophical considerations. In what follows, we always assume that the base system  $\mathbf{B}$  is PA; as I have explained, we no longer need to discuss within the former general setting after obtaining some results concerning relative truth definability in Chapter 3.

## 4.2 Systems of Finitely Iterated Self-Applicable Truth

For the discussion that follows, we technically need to introduce the systems of finitely iterated self-applicable truth; they will be used for the analysis of the inner theories of schematic reflective closures of KF, WKF and FKF (we will later consider the systems of transfinitely iterated self-applicable truth in the next chapter).

**Definition 4.2.1.** For each natural number  $n > 0$ , we introduce a primitive recursive extension  $\text{KF}_n$  of KF. The language of  $\text{KF}_n$  consists of  $\mathcal{L}_n = \mathcal{L}_0 \cup \{T_k\}_{1 \leq k \leq n}$ , where  $T_k$ ’s are new unary predicate. We set  $\text{KF}_0$  to be PA. For  $n > 0$ , the system  $\text{KF}_n$  over  $\mathcal{L}_n$  consists of  $\text{KF}_{n-1}$ , the full induction for  $\mathcal{L}_n$  and the following axioms for the predicate  $T_n$ :

$$\mathbf{K1}_n \quad \forall x, y \in \text{CT} [(T_n(x \doteq y) \leftrightarrow x^\circ = y^\circ) \wedge (F_n(x \doteq y) \leftrightarrow x^\circ \neq y^\circ)].$$

$$\mathbf{K2}_n \quad \forall \text{CT}(x) [(T_n(T_n x) \leftrightarrow T_n x^\circ) \wedge (F_n(T_n x) \leftrightarrow F_n x^\circ)].$$

$$\mathbf{K3}_n \quad \forall x \in \text{St}_n [T_n(\neg \neg x) \leftrightarrow T_n x]$$

**K4<sub>n</sub>**  $\forall x, y \in \text{St}_n [(T_n(x \wedge y) \leftrightarrow (T_n x \wedge T_n y)) \wedge (F_n(x \wedge y) \leftrightarrow (F_n \vee F_n y))]$

**K5<sub>n</sub>**  $\forall x, y \in \text{St}_n [(T_n(x \rightarrow y) \leftrightarrow (F_n x \vee T_n y)) \wedge (F_n(x \rightarrow y) \leftrightarrow (T_n \wedge F_n y))]$

**K6<sub>n</sub>**  $\forall x \forall y [\text{St}_n(\forall x. y) \rightarrow (T_n(\forall x. y) \leftrightarrow \forall z T_n y(z)) \wedge (F_n(\forall x. y) \leftrightarrow \exists z F_n y(z))]$

**K7<sub>n</sub>**  $\forall x \in \text{CT} [(T_n(T_m x) \leftrightarrow T_m x^\circ) \wedge (F_n(T_m x) \leftrightarrow \neg T_m x^\circ)]$ , for each  $m < n$ .

Then, we can identify KF and KF<sub>1</sub>. We also define the system KF<sub><ω</sub> as  $\bigcup_n \text{KF}_n$ .

WKF<sub>n</sub> and WKF<sub><ω</sub> (FKF<sup>n</sup> and WKF<sup><ω</sup>) are analogously defined.

**Lemma 4.2.2.** Let  $n > m$ . For every  $\mathcal{L}_m$ -formula  $\phi(\vec{x})$ ,  $\text{KF}_n \vdash \forall \vec{x} [T_n(\ulcorner \phi(\vec{x}) \urcorner) \leftrightarrow \phi(\vec{x})]$ . WKF<sub>n</sub> and FKF<sub>n</sub> derive the same.

**Corollary 4.2.3.** Let  $n > m$  and  $\phi(\vec{x})$  be an  $\mathcal{L}_m$ -formula. If  $\text{KF}_m \vdash \phi(\vec{x})$ , then  $\text{KF}_n \vdash T_n(\ulcorner \phi(\vec{x}) \urcorner)$ .

The same also holds in the cases of the iterated WKF and FKF.

**Lemma 4.2.4.** Let  $n \geq m \geq l$ . We have  $\text{KF}_n \vdash \forall x \in \text{St}_m (T_m x \leftrightarrow T_l x)$ . WKF<sub>n</sub> and FKF<sub>n</sub> derive the same.

For the reader's convenience, we repeat the definition of  $\mathbf{Q}^* \upharpoonright_n$  made in §5 of the last chapter: given  $n \in \mathbb{N}$  and a system  $\mathbf{Q}$ , the system  $\mathbf{Q}^* \upharpoonright_n$  is obtained from the schematic reflective closure  $\mathbf{Q}^*$  by restricting the number of application of the  $P$ -Subst rule to at most  $n$ -times. Recall that the language of  $(\text{KF}_n)^*$  is  $\mathcal{L}_n(P)$  which contains a new (vacuous) predicate  $P$ . In the rest of the present section, we assume that each derivation of  $(\text{KF}_n)^*$  is given in Hilbert style in which the inference rules are either Modus Ponens or  $P$ -Subst.

The proof of the next lemma is parallel to that of Lemma 3.5.4 .

**Lemma 4.2.5.**  $(\text{KF}_n)^* \upharpoonright_{k+1}$  is a subtheory of  $(\text{KF}_{2n})^* \upharpoonright_k$ .

*Proof.* Given a derivation  $\mathcal{D} = \langle \theta_0, \dots, \theta_l \rangle$  in  $(\text{KF}_n)^* \upharpoonright_{k+1}$ , let  $\mathcal{D}_0 = \langle \theta_0, \dots, \theta_{l'} \rangle$  be the initial segment of  $\mathcal{D}$  which ends with the first application of  $P$ -Subst  $\frac{\theta_j(P)}{\theta_j(\hat{\psi})}$  for some  $j < l'$ ; thus,  $l' \leq l$  and  $\theta_{l'} \equiv \theta_j(\hat{\psi})$ .

Let  $\psi \in \mathcal{L}_m(P)$  ( $m \leq n$ ). We define a translation  $\mathcal{S}_\psi^n$  from  $\mathcal{L}_n(P)$  to  $\mathcal{L}_{m+n}(P)$ , which ‘lifts up’ the level of truth predicates in each  $\theta_i$  ( $i < m$ ) so that  $\psi$  can be treated as if it were at the lowest level. First we preliminarily define a primitive recursive function  $g_\psi^n$  (by the primitive recursion theorem) by:

$$g_\psi^n(a) := \begin{cases} a & \text{if } a \in \text{AtSt}_0 \\ \ulcorner \psi(t) \urcorner & \text{if } a \text{ is } \ulcorner P(t) \urcorner \\ \ulcorner T_{m+i}(g_\psi^\beta(t)) \urcorner & \text{if } a \text{ is } \ulcorner T_i t \urcorner \ (1 \leq i \leq n) \\ \neg g_\psi^n(b) & \text{if } a = \neg b \\ g_\psi^n(b) \wedge g_\psi^n(c) \ (g_\psi^n(b) \rightarrow g_\psi^n(c), \text{ resp.}) & \text{if } a = b \wedge c \ (\text{or } a = b \rightarrow c) \\ \forall z. g_\psi^n(b) & \text{if } a = \forall z. b. \end{cases}$$

$\mathcal{S}_\psi^n$  is thereby defined as:  $\mathcal{S}_\psi^n(R\vec{x})$  is  $R\vec{x}$  for each atomic formula of  $\mathcal{L}_0$ ,  $\mathcal{S}_\psi^n(T_i x)$  is  $T_{m+i}(g_\psi^n(x))$ ,  $\mathcal{S}_\psi^n(Px)$  is  $\psi(x)$ , and  $\mathcal{S}_\psi^n$  commutes with the connectives and quantifiers.

We can straightforwardly show that  $(\text{KF}_{m+n})^* \upharpoonright_0 \vdash \mathcal{S}_\psi^n(\theta_q)$  for each  $q < l'$ . We thus obtain  $(\text{KF}_{m+n})^* \upharpoonright_0 \vdash \mathcal{S}_\psi^n(\theta_j(P))$ . Since  $\mathcal{S}_\psi^n$  preserves connectives and quantifiers and keeps the  $\mathcal{L}_0$ -part unchanged, we have  $\mathcal{S}_\psi^n(\theta_j(P)) \equiv \theta_j(\hat{\psi}) (\equiv \theta_{l'})$ . Let  $\mathcal{D}_1$  be this derivation of  $\theta_{l'}$  in  $(\text{KF}_{m+n})^* \upharpoonright_0$ . Set  $\mathcal{D}' := \langle \theta_0, \dots, \theta_{l'-1} \rangle * \mathcal{D}_1 * \langle \theta_{l'+1}, \dots, \theta_l \rangle$ , where  $*$  denotes the concatenation of sequences. Then  $\mathcal{D}'$  is a derivation of  $\theta_l$  in  $(\text{KF}_{m+n})^* \upharpoonright_k$ .  $\square$

**Corollary 4.2.6.**  $(\text{KF}_n)^* \upharpoonright_k$  is a subtheory of  $(\text{KF}_{2^k \cdot n})^* \upharpoonright_0$ .

**Corollary 4.2.7.**  $(\text{KF}_n)^*$  is a subtheory of  $(\text{KF}_{<\omega})^* \upharpoonright_0$ .

**Corollary 4.2.8.**  $(\text{KF}_{<\omega})^*$  is a subtheory of  $(\text{KF}_{<\omega})^* \upharpoonright_0$ .

**Lemma 4.2.9.**  $(\text{KF}_n)^*$  is proof-theoretically reducible to  $\text{KF}_{<\omega}$  for  $\mathcal{L}_0$ -formulae.

*Proof.* Suppose  $(\text{KF}_n)^* \vdash \phi$  for  $\phi \in \mathcal{L}_0$ . For some  $k \in \mathbb{N}$ , we have  $(\text{KF}_k)^* \upharpoonright_0 \vdash \phi$  by some derivation  $\mathcal{D} = \langle \theta_0, \dots, \theta_m \rangle$  ( $\theta_m$  is equal to  $\phi$ ). Let  $\psi$  be an arbitrary  $\mathcal{L}_0$ -formula (e.g.,  $(x)_0 = (x)_1$ ). Take a translation  $\mathcal{S}_\psi^k$  from  $\mathcal{L}_k(P)$  to  $\mathcal{L}_k$  in the same way as the last lemma. Then, we can show that

$\text{KF}_k \vdash \mathcal{S}_\psi^k(\theta_m)$ . Since  $\theta_m \in \mathcal{L}_0$  and  $\mathcal{S}_\psi^k$  keeps the  $\mathcal{L}_0$ -part unchanged, we finally have  $\text{KF}_k \vdash \theta_m$ .  $\square$

One might expect that lemma 4.2.9 would hold even for the case where  $\phi \in \mathcal{L}_n$  and it would hold that  $(\text{KF}_n)^* \leq \text{KF}_{<\omega}[\mathcal{L}_n]$ . However, it is not the case since  $T^\Gamma P\dot{x} \vee \neg P\dot{x}^\neg$  is provable in  $(\text{KF}_m)^*$  but not in  $\text{KF}_{<\omega}$ .

Next, we turn to the schematic reflective closures with Cons. For a given  $n \in \mathbb{N}$ , let  $\text{Cons}_n$  and  $\text{Comp}_n$  respectively denote the axioms of consistency and completeness with respect to the  $n$ -iterated truth: i.e.,

$$\text{Cons}_n \equiv \forall x \in \text{St}_n \neg [T_n(x) \wedge T_n(\neg x)] \quad \text{and} \quad \text{Comp}_n \equiv \forall x \in \text{St}_n [T_n(x) \vee T_n(\neg x)];$$

we define  $\text{Cons}_{<\omega} := \bigcup_{n \in \mathbb{N}} \text{Cons}_n$  and  $\text{Comp}_{<\omega} := \bigcup_{n \in \mathbb{N}} \text{Comp}_n$ .

**Proposition 4.2.10.** Let  $m \leq n$ . Then,  $\text{KF}_n \vdash (\text{Cons}_n \rightarrow \text{Cons}_m) \wedge (\text{Comp}_n \rightarrow \text{Comp}_m)$ . The same for  $\text{WKF}_n$  and  $\text{FKF}_n$ .

**Lemma 4.2.11.**  $\text{KF}_n + \text{Cons}_n$  and  $\text{KF}_n + \text{Comp}_n$  are mutually truth definable. As in the case of non-iterated truth, this doesn't hold for  $\text{WKF}_n$  and  $\text{FKF}_n$ .

*Proof.* We define in PA (by the recursion theorem or diagonalization) a function  $f$  by:

$$f(x) := \begin{cases} x & \text{if } x \in \text{AtFml}_0 \\ \neg F_m(fy)^\neg \text{ (i.e., sb}(\neg F_m f(\cdot)^\neg; y)) & \text{if } x = T_m y \text{ for } y \in \text{Tm} \text{ and } m \leq n \\ \neg f(y) \text{ (} f(y) \wedge f(z), \text{ resp.)} & \text{if } x = \neg y \text{ (or } x = y \wedge z) \\ \forall z. f(y) & \text{if } x = \forall z. y; \end{cases}$$

otherwise,  $f(x) = 0$ . The truth definition  $\mathcal{T}$  for both claims is given by  $T_z^\prec x \mapsto \neg F_z f(x)$ . We only

demonstrate that  $\mathbf{K2}_n$  is preserved by  $\mathcal{T}$ : for  $x \in \text{CT}$  and  $m < n$ ,

$$\mathcal{T}(T_n T_m x) \equiv \neg F_n \ulcorner \neg F_m(fx) \urcorner \leftrightarrow \neg T_n \ulcorner F_m(fx) \urcorner \leftrightarrow \neg F_m f(x^\circ) \equiv \mathcal{T}(T_m x^\circ)$$

$$\mathcal{T}(F_n T_m x) \equiv \neg F_n \ulcorner \neg F_m(fx) \urcorner \leftrightarrow \neg F_n \ulcorner F_m(fx) \urcorner \leftrightarrow \neg \neg F_n f(x^\circ) \equiv \mathcal{T}(\neg T_m x^\circ).$$

The other cases are straightforward. □

**Lemma 4.2.12.**  $(\text{KF}_n + \text{Cons}_n)^* \upharpoonright_{k+1}$  is a subtheory of  $(\text{KF}_{2n} + \text{Cons}_{2n})^* \upharpoonright_k$ .

**Corollary 4.2.13.**  $(\text{KF}_n + \text{Cons}_n)^* \upharpoonright_k$  is a subtheory of  $(\text{KF}_{2^k \cdot n} + \text{Cons}_{2^k \cdot n})^* \upharpoonright_0$ .

**Corollary 4.2.14.**  $(\text{KF}_n + \text{Cons}_n)^*$  is a subtheory of  $(\text{KF}_{<\omega} + \text{Cons}_{<\omega})^* \upharpoonright_0$ .

**Corollary 4.2.15.**  $(\text{KF}_{<\omega} + \text{Cons}_{<\omega})^*$  is a subtheory of  $(\text{KF}_{<\omega} + \text{Cons}_{<\omega})^* \upharpoonright_0$ .

**Lemma 4.2.16.**  $(\text{KF}_n + \text{Cons}_n)^*$  is proof-theoretically reducible to  $\text{KF}_{<\omega} + \text{Cons}_{<\omega}$  for  $\mathcal{L}_0$ -formulae.

### 4.3 Semantics of iterative compositional theories

In the present section, we generalize the semantics of  $\text{KF}^-$  developed by Cantini [5] to the other systems  $\text{WKF}^-$ ,  $\text{FKF}^-$  and  $\text{PUTB}^-$ . It will be used to determine how their inner theories of truth systems are correlated.

Let  $M$  be an arbitrary model of PA, and fix this  $M$  throughout this section.  $|M|$  denote the domain of  $M$ . Given a formula  $\phi$  and a function  $f$  of  $\mathcal{L}$ ,  $\phi^M$  and  $f^M$  denote the interpretation of  $\phi$  and  $f$  in  $M$  respectively. Given a set  $\mathcal{C}$  of subsets of  $|M|$  (i.e.,  $\mathcal{C} \subset \mathcal{P}(|M|)$ ),  $\text{sup } \mathcal{C}$  and  $\text{inf } \mathcal{C}$  denote  $\bigcup \mathcal{C}$  ( $\subset |M|$ ) and  $\bigcap \mathcal{C}$  ( $\subset |M|$ ) respectively. For each  $Y \subset |M|$ , we write  $\neg Y$  for  $|M| \setminus Y$ . For an  $\mathcal{L}_T$ -formula  $\phi$  and  $\vec{a} \in |M|$ , we write  $X \models \phi(\vec{a})$  iff  $(M, X) \models \phi(\vec{a})$ , where  $T$  is interpreted by  $X$ .

**Definition 4.3.1.** (1)  $X \subset |M|$  is regular iff, for any  $a \in \text{For}^M$  and  $b \in \text{CT}^M$ ,

$$\text{sb}^M(a, b) \in X \Leftrightarrow \text{sb}^M(a, \text{nm}(\text{val}(b))) \in X.$$

We denote the set of all regular  $X \subset |M|$  by  $\text{REG}_M$ . Note that  $X \in \text{REG}_M$  assures that  $X$  is a model of the axiom K0.

(2) Given  $X \in \text{REG}_M$ ,

- we say that  $X$  is consistent iff there is no  $a \in |M| \cap \text{St}_{\mathcal{L}_T}^M$  such that both  $a$  and  $\neg^M a$  are in  $X$ , and  $\text{CONS}_M$  denotes the set of all consistent  $X \in \text{REG}_M$ ;
- we say that  $X$  is complete iff either  $a$  or  $\neg^M a$  is in  $X$  for every  $a \in |M| \cap \text{St}_{\mathcal{L}_T}^M$ , and  $\text{COMP}_M$  denotes the set of all complete  $X \in \text{REG}_M$

**Definition 4.3.2.** We will define operators **S**, **F** and **W** from  $\text{REG}_M$  to  $\text{REG}_M$  according to Strong Kleene, Feferman, and Weak Kleene schemata, respectively.

(I)  $a \in \mathbf{W}(X)$  iff  $a \in \text{St}_{\mathcal{L}_T}^M$  and

- (i)  $M \models a = R(\vec{b}) \wedge \text{CT}(\vec{b}) \wedge R(\vec{b}^\circ)$ , or
- (ii)  $M \models a = \neg R(\vec{b}) \wedge \text{CT}(\vec{b}) \wedge \neg R(\vec{b}^\circ)$ , or
- (iii)  $M \models \text{CT}(b) \wedge a = Tb$  and  $X \models T(b^\circ)$ , or
- (iv)  $M \models \text{CT}(b) \wedge a = \neg Tb$  and  $X \models F(b^\circ)$ , or
- (v)  $M \models \text{CT}(b) \wedge a = Fb$  and  $X \models F(b^\circ)$ , or
- (vi)  $M \models \text{CT}(b) \wedge a = \neg Fb$  and  $X \models T(b^\circ)$ , or
- (vii)  $M \models a = \neg \neg b$  and  $X \models Tb$ , or
- (viii)  $M \models a = b \wedge c$  and  $X \models Tb \wedge Tc$ , or
- (ix)  $M \models a = \neg(b \wedge c)$  and  $X \models (Fb \wedge Fc) \vee (Tb \wedge Fc) \vee (Fb \wedge Tc)$ , or

(**x**)  $M \models a = \forall z.a$  and  $X \models \forall zTa(z)$ , or

(**xi**)  $M \models a = \neg(\forall z.a)$  and  $X \models [\forall z(Ta(z) \vee Fa(z))] \wedge \exists zFa(z)$ .

(**II**)  $a \in \mathbf{F}(X)$  iff  $a \in \text{St}_{\mathcal{L}_T}^M$  and either of (i)-(xi) and the following (xii)–(xiii) holds for  $a$ :

(**xii**)  $M \models a = b \rightarrow c$  and  $X \models (Ta \wedge Tc) \vee Fb$

(**xiii**)  $M \models a = \neg(b \rightarrow c)$  and  $X \models (Ta \wedge Fc)$ .

(**III**)  $a \in \mathbf{S}(X)$  iff  $a \in \text{St}_{\mathcal{L}_T}^M$ , and one of (i)-(vii), (xii)-(xiii) and the following holds

(**viii'**)  $M \models a = b \wedge c$  and  $X \models Tb \wedge Tc$  or

(**ix'**)  $M \models a = \neg(b \wedge c)$  and  $X \models T\neg b \vee T\neg c$  or

(**x'**)  $M \models \text{Var}(z) \wedge a = \forall z.a$  and  $X \models \forall zTa(z)$  or

(**xi'**)  $M \models \text{Var}(z) \wedge a = \neg\forall z.a$  and  $X \models \exists zT\neg a(z)$ .

**Proposition 4.3.3.**  $\mathbf{S}$ ,  $\mathbf{F}$  and  $\mathbf{W}$  are operators from  $\text{REG}_M$  to  $\text{REG}_M$ .

*Proof.* The proof for  $\mathbf{S}$  is already given in [5]. The other cases are parallel. We remark that the assumption that  $\text{sb}$  commutes with (the representations of) logical operations is needed to prove this proposition.  $\square$

For an operator  $\mathbf{G}$  from  $\mathcal{P}(|M|)$  to  $\mathcal{P}(|M|)$ , we call it monotone if  $\mathbf{G}(X) \subset \mathbf{G}(Y)$  whenever  $X \subset Y$ . Then it is observed that  $\mathbf{S}$ ,  $\mathbf{F}$  and  $\mathbf{W}$  are all monotone.

**Definition 4.3.4.** Let  $\mathbf{G}$  be any monotone operator from  $\mathcal{P}(|M|)$  to  $\mathcal{P}(|M|)$  and let  $X \subset |M|$ .

1. For each ordinal  $\alpha$ ,  $X_\alpha^{\mathbf{G}}$  is recursively defined by:

$$X_0^{\mathbf{G}} = X, \quad X_{\beta+1}^{\mathbf{G}} = X_\beta^{\mathbf{G}} \cup \mathbf{G}(X_\beta), \quad X_\lambda^{\mathbf{G}} = \bigcup_{\beta < \lambda} X_\beta^{\mathbf{G}},$$

where  $\lambda$  is a limit ordinal. Dually,  $X_{\mathbf{G}}^{\alpha}$  is defined by

$$X_{\mathbf{G}}^0 = X, \quad X_{\mathbf{G}}^{\beta+1} = X_{\mathbf{G}}^{\beta} \cap \mathbf{G}(X_{\mathbf{G}}^{\beta}), \quad X_{\mathbf{G}}^{\lambda} = \bigcap_{\beta < \lambda} X_{\mathbf{G}}^{\beta}.$$

2.  $\text{UP}^{\mathbf{G}}(X) \stackrel{\text{def}}{=} X_{\infty}^{\mathbf{G}} = \bigcup_{\alpha \in \mathcal{O}_n} X_{\alpha}^{\mathbf{G}}$ .  $\text{DOWN}^{\mathbf{G}}(X) \stackrel{\text{def}}{=} X_{\mathbf{G}}^{\infty} = \bigcap_{\alpha \in \mathcal{O}_n} X_{\mathbf{G}}^{\alpha}$ .
3.  $X_{\text{up}}^{\mathbf{G}} \stackrel{\text{def}}{=} \text{DOWN}^{\mathbf{G}}(\text{UP}^{\mathbf{G}}(X))$ , and  $X_{\text{down}}^{\mathbf{G}} \stackrel{\text{def}}{=} \text{UP}^{\mathbf{G}}(\text{DOWN}^{\mathbf{G}}(X))$ .
4.  $\text{FIX}_M^{\mathbf{G}} \stackrel{\text{def}}{=} \{X \in \text{REG}_M \mid \mathbf{G}(X) = X\}$ .

Given  $X \subset |M|$ , we say  $X$  is  $\mathbf{G}$ -dense if  $X \subset \mathbf{G}(X)$ , and  $X$  is  $\mathbf{G}$ -closed if  $\mathbf{G}(X) \subset X$ .

**Proposition 4.3.5.** For any monotone operator  $\mathbf{G}: \mathcal{P}(|M|) \rightarrow \mathcal{P}(|M|)$ , the following hold:

- (1)  $\text{UP}^{\mathbf{G}}(X)$  is the the  $\subset$ -least  $\mathbf{G}$ -closed set including  $X$  among  $\text{REG}_M$ , and  $\text{DOWN}^{\mathbf{G}}(X)$  is the  $\subset$ -largest  $\mathbf{G}$ -dense set included in  $X$  among  $\text{REG}_M$ .
- (2)  $X_{\text{up}}^{\mathbf{G}}, X_{\text{down}}^{\mathbf{G}} \in \text{FIX}_M^{\mathbf{G}}$  and  $X_{\text{down}}^{\mathbf{G}} \subset X_{\text{up}}^{\mathbf{G}}$ .
- (3)  $\text{UP}^{\mathbf{G}}(\emptyset) = \emptyset_{\text{up}}^{\mathbf{G}} = \emptyset_{\text{down}}^{\mathbf{G}}$ , and  $\text{DOWN}^{\mathbf{G}}(|M|) = |M|_{\text{down}}^{\mathbf{G}} = |M|_{\text{up}}^{\mathbf{G}}$ . Therefore,  $\emptyset_{\text{up}}^{\mathbf{G}}$  is minimal fixed-point and  $|M|_{\text{down}}^{\mathbf{G}}$  is maximal fixed-point with respect to  $\mathbf{G}$ .
- (4) If  $\mathbf{G}$  preserves regularity (i.e.,  $\mathbf{G}(X) \in \text{REG}_M$  whenever  $X \in \text{REG}_M$ ), then  $X_{\alpha}^{\mathbf{G}}, X_{\mathbf{G}}^{\alpha} \in \text{REG}_M$  for all  $X \in \text{REG}_M$ .

**Proposition 4.3.6.** Let  $\mathbf{G}$  be  $\mathbf{W}$ ,  $\mathbf{F}$  or  $\mathbf{S}$ . If  $X \in \text{CONS}_M$  then  $X_{\text{down}}^{\mathbf{G}} \in \text{CONS}_M$ , and if  $X \in \text{COMP}_M$  then  $X_{\text{up}}^{\mathbf{G}} \in \text{COMP}_M$ .

**Proposition 4.3.7.** Let  $\mathbf{G}$  be  $\mathbf{S}$ ,  $\mathbf{F}$  or  $\mathbf{W}$ . Then,  $\emptyset_{\text{up}}^{\mathbf{G}} \in \text{CONS}_M$  and  $|M|_{\text{down}}^{\mathbf{G}} \in \text{COMP}_M$ .

**Lemma 4.3.8.** Let  $X \subset |M|$ .

- (1)  $X \in \text{FIX}_M^{\mathbf{W}}$  iff  $X \models \text{WKF}^-$ .  $X \in \text{FIX}_M^{\mathbf{W}} \cap \text{CONS}_M$  iff  $X \models \text{WKF}^- + \text{Cons}$ .  $X \in \text{FIX}_M^{\mathbf{W}} \cap \text{COMP}_M$  iff  $X \models \text{WKF}^- + \text{Comp}$ .

(2)  $X \in \text{FIX}_M^{\mathbf{F}}$  iff  $X \models \text{FKF}^-$ .  $X \in \text{FIX}_M^{\mathbf{F}} \cap \text{CONS}_M$  iff  $X \models \text{FKF}^- + \text{Cons}$ .  $X \in \text{FIX}_M^{\mathbf{F}} \cap \text{COMP}_M$  iff  $X \models \text{FKF}^- + \text{Comp}$ .

(3)  $X \in \text{FIX}_M^{\mathbf{S}}$  iff  $X \models \text{KF}^-$ .  $X \in \text{FIX}_M^{\mathbf{S}} \cap \text{CONS}_M$  iff  $X \models \text{KF}^- + \text{Cons}$ .  $X \in \text{FIX}_M^{\mathbf{S}} \cap \text{COMP}_M$  iff  $X \models \text{KF}^- + \text{Comp}$ .

(4) If  $M$  is standard, (1)-(3) hold even for WKF, FKF and KF respectively.

**Definition 4.3.9.** We also consider the semantics of PUTB. In the following  $\text{SPOS}_M$  and  $\text{SNEG}_M$  denotes the set

$$\begin{aligned} & \{\text{sb}_{\vec{x}}^M(\ulcorner \phi(\vec{x}) \urcorner, \vec{b}) \mid \phi(\vec{x}) \text{ is a (standard) } T\text{-positive } \mathcal{L}_T\text{-formula and } \vec{b} \in \text{CT}^M\} \\ & \{\text{sb}_{\vec{x}}^M(\ulcorner \phi(\vec{x}) \urcorner, \vec{b}) \mid \phi(\vec{x}) \text{ is a (standard) } T\text{-negative } \mathcal{L}_T\text{-formula and } \vec{b} \in \text{CT}^M\}; \end{aligned}$$

note that  $\text{SPOS}_M$  is different from  $\text{POS}^M$  (defined in CH.2-8); the extra ‘S’ stands for ‘standard’.

We define operators  $\mathbf{P}_0, \mathbf{P}_1, \mathbf{P}_2$  from  $\text{REG}_M$  to  $\text{REG}_M$  for PUTB is defined as follows.

(a)  $a \in \mathbf{P}_0(X)$  iff  $a \in \text{SPOS}_M$  and  $a = \text{sb}_{\vec{x}}(\ulcorner \phi \urcorner, \vec{b})$  for some positive (standard)  $\mathcal{L}_T$ -formula  $\phi(\vec{x})$  and  $\vec{b} \in \text{CT}^M$  such that  $X \models \phi(\text{val}(\vec{b}))$ ; note that  $\mathbf{P}_0(X) \subset \text{SPOS}_M$  for all  $X \subset |M|$ .

(b)  $\mathbf{P}_1$  is defined by

$$\mathbf{P}_1(Z) := \mathbf{P}_0(Z) \cup ((Z \cap \text{St}_{\mathcal{L}_T}^M) \setminus \text{SPOS}_M).$$

(c)  $\mathbf{P}_2$  is defined by

$$\mathbf{P}_2(Z) := \mathbf{P}_0(Z) \cup (\text{St}_{\mathcal{L}_T}^M \setminus \text{SPOS}_M).$$

It is easy to see that  $\mathbf{P}_0, \mathbf{P}_1$  and  $\mathbf{P}_2$  are all monotone; for, if  $X \subset Y$  and  $X \subset \phi(\vec{a})$  then  $Y \models \phi(\vec{a})$

for all positive formulae  $\phi$ .

**Lemma 4.3.10.**  $\mathbf{P}_0$ ,  $\mathbf{P}_1$  and  $\mathbf{P}_2$  are operators from  $\text{REG}_M$  to  $\text{REG}_M$ .

*Proof.* Let  $\phi(x)$  be a  $T$ -positive  $\mathcal{L}_T$ -formula. For each  $a \in \text{CT}^M$  and  $X \subset |M|$ ,  $X \models \phi(\text{val}(a))$  iff  $X \models \phi(\text{val} \circ \text{nm} \circ \text{val}(c))$ . The claim for  $\mathbf{P}_0$  follows from this. The other cases trivially follow.  $\square$

**Lemma 4.3.11.**  $\mathbf{P}_0$ ,  $\mathbf{P}_1$  and  $\mathbf{P}_2$  are all monotone.

*Proof.* It suffices to show that  $\mathbf{P}_0$  is monotone. The monotonicity of  $\mathbf{P}_0$  follows from the  $T$ -positivity:  $X \models \phi$  implies  $Y \models \phi$  for all  $T$ -positive formula  $\phi$  and  $X \subset Y \subset |M|$ .  $\square$

**Lemma 4.3.12.** Let  $X \in \text{REG}_M$ . Then  $\mathbf{P}_0(X) \in \text{CONS}_M$  and  $\mathbf{P}_2(X) \in \text{COMP}_M$

**Proposition 4.3.13.** If  $X \in \text{CONS}_M$  then  $X_{\text{down}}^{\mathbf{P}_0} \in \text{CONS}_M$ , and if  $X \in \text{COMP}_M$  then  $X_{\text{up}}^{\mathbf{P}_2} \in \text{COMP}_M$ .

*Proof.* Since  $\mathbf{P}_0(X) \subset \text{POS}_M$  for all  $X$ , we have  $\text{DOWN}_M^{\mathbf{P}_0} \subset \text{POS}_M$ . For simplicity, let us write  $Y$  for  $\text{DOWN}^{\mathbf{P}_0}(X)$ . Then we can show by induction on  $\alpha$  that  $Y_\alpha^{\mathbf{P}_0} \subset \text{POS}_M \in \text{CONS}_M$  for all ordinals  $\alpha$ . For the second claim, we first observe that  $\text{UP}^{\mathbf{P}_2}(X) \in \text{COMP}_M$ , since  $\mathbf{P}_2(X) \in \text{COMP}_M$  for all  $X \in \text{REG}_M$  by the last lemma. We write  $Y$  for  $\text{UP}^{\mathbf{P}_2}(X)$ . For the sake of contradiction, suppose  $\text{DOWN}^{\mathbf{P}_2}(Y) = X_{\text{up}}^{\mathbf{P}_2} \notin \text{COMP}_M$ . Take the least  $\alpha$  such that there exists  $a \in \text{St}_{\mathcal{L}_T}^M$  with  $a, \neg a \notin Y_{\mathbf{P}_2}^{\alpha+1}$ . Then  $a$  or  $\neg a$  is in  $Y_{\mathbf{P}_2}^\alpha$ . Assume  $a \in Y_{\mathbf{P}_2}^\alpha$ . If  $a \notin \text{SPOS}_M$  then  $a \in Y_{\mathbf{P}_2}^{\alpha+1}$  by definition. Thus,  $a \in \text{SPOS}_M$ . But then  $\neg a \in (\text{St}_{\mathcal{L}_T}^M \setminus \text{SPOS}_M)$  and thus  $\neg a \in Y_{\mathbf{P}_2}^{\alpha+1}$ .  $\square$

**Lemma 4.3.14.**  $X \models \text{PUTB}^-$  iff  $X \in \text{FIX}_M^{\mathbf{P}_1}$ . Hence  $X \models \text{PUTB}^- + \text{Cons}$  (or  $\text{Comp}$ ) iff  $X \in \text{FIX}_M^{\mathbf{P}_1} \cap \text{CONS}_M$  ( $\text{FIX}_M^{\mathbf{P}_1} \cap \text{COMP}_M$  resp.).

*Proof.* The direction from the latter to the former is immediate. Conversely, suppose the former. Then, for every  $T$ -positive  $\phi(\vec{x})$  and  $\vec{a} \in \text{St}_{\mathcal{L}_T}^M$ ,

$$\text{sb}_{\vec{x}}(\ulcorner \phi^\urcorner, \vec{a}) \in X \Leftrightarrow X \models \phi(\vec{a}^\circ) \Leftrightarrow \ulcorner \phi(\vec{a})^\urcorner \in \mathbf{P}_1(X).$$

Finally, for each  $a \in X \setminus \text{SPOS}_M$ , we trivially have  $a \in X \Leftrightarrow a \in \mathbf{P}_1(X)$ .  $\square$

**Proposition 4.3.15.**  $\emptyset_{\text{up}}^{\mathbf{P}_0} = \emptyset_{\text{up}}^{\mathbf{P}_1}$  is the smallest model of  $\text{PUTB}^-$  and  $\text{PUTB}^- + \text{Cons}$ , and  $(\text{St}_{\mathcal{L}_T}^M)^{\mathbf{P}_1}_{\text{down}} = (\text{St}_{\mathcal{L}_T}^M)^{\mathbf{P}_2}_{\text{down}}$  is the largest model of  $\text{PUTB}^-$  and  $\text{PUTB}^- + \text{Comp}$ .

*Proof.* We can show the claimed identity by transfinite induction on the constructions of  $\emptyset_{\text{up}}^{\mathbf{P}_0}$  and  $\emptyset_{\text{up}}^{\mathbf{P}_1}$  (or,  $(\text{St}_{\mathcal{L}_T}^M)^{\mathbf{P}_1}_{\text{down}}$  and  $(\text{St}_{\mathcal{L}_T}^M)^{\mathbf{P}_2}_{\text{down}}$ , resp.) The minimality and maximality follows from the last lemma and Proposition 4.3.5.  $\square$

Next we consider the semantics of schematic reflective closures. The idea is originally due to Feferman [17, pp.22–23]. Let  $\mathbf{G}$  be any of  $\mathbf{S}$ ,  $\mathbf{F}$ ,  $\mathbf{W}$ ,  $\mathbf{P}_i$  ( $i = 0, 1, 2$ ). Given  $Z \subset |M|$ , a new operator  $\mathbf{G}_Z^* : \text{REG}_M \rightarrow \text{REG}_M$  is obtained by (i) adding the following clause to its defining clauses (in Definition 4.3.2 or 4.3.9)

$$M \models a = P(b) \wedge \text{CT}^M(b) \wedge b \in Z;$$

and (ii) by replacing the occurrence of ‘ $\text{St}_{\mathcal{L}_T}^M$ ’ in its definition by ‘ $\text{St}_{\mathcal{L}_T(P)}^M$ ’.

**Lemma 4.3.16.** (1) If  $\text{WKF}^* \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{W}_Z^*}) \models \phi$  and  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{W}_Z^*}) \models \phi$  for all  $X, Z \subset \mathbb{N}$ ,

where  $P$  is interpreted by  $Z$  and  $T$  is interpreted by  $X_{\text{down}}^{\mathbf{W}_Z^*}$ . If  $\text{WKF}^* + \text{Cons} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{W}_Z^*}) \models \phi$  for all  $X \in \text{CONS}_{\mathbb{N}}$ . If  $\text{WKF}^* + \text{Comp} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{W}_Z^*}) \models \phi$  for all  $X \in \text{COMP}_{\mathbb{N}}$ .

(2) If  $\text{FKF}^* \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{F}_Z^*}) \models \phi$  and  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{F}_Z^*}) \models \phi$  for all  $X, Z \subset \mathbb{N}$ . If  $\text{FKF}^* + \text{Cons} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{F}_Z^*}) \models \phi$  for all  $X \in \text{CONS}_{\mathbb{N}}$ . If  $\text{FKF}^* + \text{Comp} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{F}_Z^*}) \models \phi$  for all  $X \in \text{COMP}_{\mathbb{N}}$ .

(3) If  $\text{KF}^* \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{S}_Z^*}) \models \phi$  and  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{S}_Z^*}) \models \phi$  for all  $X, Z \subset \mathbb{N}$ . If  $\text{KF}^* + \text{Cons} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{down}}^{\mathbf{S}_Z^*}) \models \phi$  for all  $X \in \text{CONS}_{\mathbb{N}}$ . If  $\text{KF}^* + \text{Comp} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{S}_Z^*}) \models \phi$  for all  $X \in \text{COMP}_{\mathbb{N}}$ .

**Lemma 4.3.17.** If  $\text{PUTB}^* \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{\mathbf{G}^*}) \models \phi$  for all  $X \subset \mathbb{N}$ , where  $\mathbf{G}$  can be either of  $\mathbf{P}_0$ ,  $\mathbf{P}_1$  and  $\mathbf{P}_2$ . If  $\text{PUTB}^* + \text{Cons} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{down}}^{(\mathbf{P}_0)^*}) \models \phi$  for all  $X \in \text{CONS}_{\mathbb{N}}$ . If  $\text{PUTB}^* + \text{Comp} \vdash \phi$  then  $(\mathbb{N}, Z, X_{\text{up}}^{(\mathbf{P}_2)^*}) \models \phi$  for all  $X \in \text{COMP}_{\mathbb{N}}$ .

**Lemma 4.3.18.** If  $\text{PUTB}^* + \text{Cons} \vdash T^\Gamma \sigma^\neg$ , then  $\sigma$  is  $T$ -positive. Thus, the same can be said to  $\text{PUTB}^- (+\text{Cons})$  and  $\text{PUTB} (+\text{Cons})$ .

*Proof.* This is shown by induction on  $\alpha$  that, for any non  $T$ -positive  $\phi(\vec{x})$ ,  $\emptyset_\alpha^{(\mathbf{P}_0)^*}$  (with respect to the standard model  $\mathbb{N}$ ) contains no  $\text{sb}_{\vec{x}}(\ulcorner \phi \urcorner, \vec{a})$  for any  $\vec{a} \in \mathbb{N}$ .  $\square$

**Theorem 4.3.19.** If  $\text{PUTB}^* \vdash T^\Gamma \sigma^\neg$  then  $\text{PUTB}^* \vdash \sigma$ . Thus,  $\mathcal{I}(\text{PUTB}^*) \subset \text{PUTB}^*$  and  $\text{PUTB}^* + \text{CONEC}$  is identical with  $\text{PUTB}^*$ . The same applies even to  $\text{PUTB}$  and  $\text{PUTB}^-$ .

*Proof.* It immediately follows from the last lemma, since  $\text{PUTB}^* \vdash T^\Gamma \sigma^\neg$  then  $\sigma$  is  $T$ -positive and thus  $\sigma$  is derivable in  $\text{PUTB}^*$  by the  $\text{PUTB}$ -schema.  $\square$

## 4.4 Inner Theory Comparison via Semantics

In the present section, we will compare the systems of truth by their inner theories, in other words, by their assertable truths. As we have mentioned in §1, the distinction of the outer and inner theories was first introduced by Reinhardt [70]. Given a system  $\mathbf{Q}$  over  $\mathcal{L}_T$ , the outer theory of  $\mathbf{Q}$  denotes  $\{\sigma \in \mathcal{L}_T \mid \mathbf{Q} \vdash \sigma\}$ , i.e., the ‘theory’ in the ordinary sense; on the other hand, the inner theory of  $\mathbf{Q}$  denotes  $\{\sigma \in \mathcal{L}_T \mid \mathbf{Q} \vdash T^\Gamma(\sigma^\neg)\}$ , i.e., the set of assertable truths in  $\mathbf{Q}$  in other words. Let  $\mathcal{I}(\mathbf{Q})$  denote the inner theory of a given  $\mathbf{Q}$ . Then, note that if  $\mathcal{I}(\mathbf{Q}_0) \subset \mathcal{I}(\mathbf{Q}_1)$ , and if  $\mathbf{Q}_0$  and  $\mathbf{Q}_1$  prove the  $\mathbf{T}$ -schema for truth-free sentences, then  $\mathbf{Q}_1$  proves as many truth-free theorems as  $\mathbf{Q}_0$  does. However, as we shall see, the converse does not hold. Thus, since all the theories we have considered so far prove the  $\mathbf{T}$ -schema for truth-free sentences, comparison by the inner theories gives properly stronger and finer-grained criterion than the mere derivability of truth-free theorems does.

For the sake of simplicity of the writing, we introduce some notation; for theories  $Q_0$  and  $Q_1$  over  $\mathcal{L}_T$ , we write

- $Q_0 \leq_{\mathcal{I}} Q_1$ , when if  $Q_0 \vdash Tt$  then  $Q_1 \vdash Tt$  for all closed term  $t$  (thus, if  $Q_0 \leq_{\mathcal{I}} Q_1$  then  $\mathcal{I}(Q_0) \subset \mathcal{I}(Q_1)$ ),
- $Q_0 \equiv_{\mathcal{I}} Q_1$ , if  $Q_0 \leq_{\mathcal{I}} Q_1$  and  $Q_1 \leq_{\mathcal{I}} Q_0$ ,
- $Q_0 <_{\mathcal{I}} Q_1$ , if  $Q_0 \leq_{\mathcal{I}} Q_1$  and  $Q_1 \not\leq_{\mathcal{I}} Q_0$ , and
- $Q_0 \not\leq_{\mathcal{I}} Q_1$ , if  $Q_0 \not\leq_{\mathcal{I}} Q_1$  and  $Q_1 \not\leq_{\mathcal{I}} Q_0$ .

**Proposition 4.4.1.** (1)  $WKF^- <_{\mathcal{I}} FKF^-$ . (2)  $FKF^- <_{\mathcal{I}} KF^-$ . (3)  $KF^- \equiv_{\mathcal{I}} KF^- + \text{Cons.}$  (4)  $FKF^- \equiv_{\mathcal{I}} FKF^- + \text{Cons.}$  (5)  $WKF^- \equiv_{\mathcal{I}} WKF^- + \text{Cons.}$  (6)  $KF^- <_{\mathcal{I}} KF^- + \text{Comp.}$

*Proof.* First, fix an arbitrary model  $M$  of PA and let  $Q^-$  be one of the above theories and  $\mathbf{G}$  be the operator on  $M$  for  $Q^-$ . The crucial point for the proofs is the fact that if  $X \models Q^-$  then  $\emptyset_{\text{up}}^{\mathbf{G}} \subset X$ . Thus, (3)-(5) immediately follows from this fact, since  $\emptyset_{\text{up}}^{\mathbf{G}} \models \text{Cons.}$

For (1) and (2), it then suffices to show that  $\emptyset_{\text{up}}^{\mathbf{W}} \subsetneq \emptyset_{\text{up}}^{\mathbf{F}} \subsetneq \emptyset_{\text{up}}^{\mathbf{S}}$ . We can show by transfinite induction that  $\emptyset_{\alpha}^{\mathbf{W}} \subset \emptyset_{\alpha}^{\mathbf{F}} \subset \emptyset_{\alpha}^{\mathbf{S}}$  for each  $\alpha$ ; thus  $\emptyset_{\text{up}}^{\mathbf{W}} \subset \emptyset_{\text{up}}^{\mathbf{F}} \subset \emptyset_{\text{up}}^{\mathbf{S}}$ . In order to show that these inclusions are proper, consider the following two sentences:

$$\tau_0 := 0 = 0 \vee \lambda, \text{ where } \lambda \text{ is the liar sentence (defined in Ch.2.6);}$$

$$\tau_1 := T\forall x [\exists y, z(\text{CT}(y) \wedge \text{CT}(z) \wedge x = (y=z) \wedge \text{val}(y) = \text{val}(z)) \rightarrow Tx]$$

Then, we can show by induction on  $\alpha$  that  $\tau_0 \in \emptyset_{\alpha}^{\mathbf{S}}$  but  $\tau_0 \notin \emptyset_{\alpha}^{\mathbf{F}}$ , and  $\tau_1 \in \emptyset_{\alpha}^{\mathbf{F}}$  but  $\tau_1 \notin \emptyset_{\alpha}^{\mathbf{W}}$  for a sufficiently large  $\alpha$ ; cf. Remark 5 in Ch.2.6. We have obtained the claim.

Finally for (6), since  $KF^-$  is a subtheory of  $KF^- + \text{Comp}$ , it suffices to show that  $KF^- \not\leq_{\mathcal{I}} KF^- + \text{Comp}$ . This follows from (3) and Prop 2.6.4.  $\square$

**Proposition 4.4.2.** (1)  $\text{PUTB}^- <_{\mathcal{I}} \text{KF}^-$ . (2)  $\text{PUTB}^- \equiv_{\mathcal{I}} \text{PUTB}^- + \text{Cons}$ . (3)  $\text{PUTB}^- <_{\mathcal{I}} \text{PUTB}^- + \text{Comp}$ . (4)  $\text{PUTB}^- + \text{Comp} \not\prec_{\mathcal{I}} \text{KF}^-$ . (5)  $\text{PUTB}^- \not\prec_{\mathcal{I}} \text{WKF}^-$ . (6)  $\text{PUTB}^- \not\prec_{\mathcal{I}} \text{FKF}^-$ . (7)  $\text{PUTB}^- + \text{Comp} <_{\mathcal{I}} \text{KF}^-$ . (8)  $\text{WKF}^- \not\prec_{\mathcal{I}} \text{PUTB}^- + \text{Comp}$ . (9)  $\text{FKF}^- \not\prec_{\mathcal{I}} \text{PUTB}^- + \text{Comp}$ .

*Proof.* (1). It is trivial that  $\text{PUTB} \leq_{\mathcal{I}} \text{KF}^-$ , since  $\text{PUTB}^-$  is a subtheory of  $\text{KF}^-$ . To show that it

is a proper subtheory, we note that  $\ulcorner \neg T(\ulcorner 0 \neq 0 \urcorner) \urcorner \notin \emptyset_{\text{up}}^{\mathbf{P}_0}$  but it is provable in  $\text{KF}^-$ .

(2). It is because  $\emptyset_{\text{up}}^{\mathbf{P}_0}$  is the least model for both  $\text{PUTB}^-$  and  $\text{PUTB}^- + \text{Cons}$ .

(3) & (4).  $\text{PUTB}^- + \text{Comp} \vdash \lambda$  but  $\text{KF}^- \not\vdash \lambda$ .

(5) & (6). It suffices to show that  $\text{PUTB}^- \not\prec_{\mathcal{I}} \text{FKF}^-$  and  $\text{WKF}^- \not\prec_{\mathcal{I}} \text{PUTB}^-$ . The former follows from the fact that a sentence  $T(\ulcorner 0 = 0 \vee \lambda \urcorner)$  is derivable from  $\text{PUTB}^-$  but not contained in  $\emptyset_{\text{up}}^{\mathbf{F}}$ . The latter obtains by observing that the sentence  $T(\ulcorner \neg T(\ulcorner 0 \neq 0 \urcorner) \urcorner)$  is derivable from  $\text{WKF}^-$  but not contained in  $\emptyset_{\text{up}}^{\mathbf{P}_0}$ .

(7). Let  $\tau$  be a sentence  $T(\ulcorner 0 = 0 \urcorner) \wedge \neg T(\ulcorner 0 \neq 0 \urcorner)$ . We have  $\text{KF}^- \vdash T^{\ulcorner \tau \urcorner}$ . Hence, it suffices to show that  $\text{PUTB}^- + \text{COMP} \not\vdash T^{\ulcorner \tau \urcorner}$ . Define  $X := \text{St}_{\mathcal{L}_T}^M \setminus \{\ulcorner \tau \urcorner\}$ . Note that neither  $\tau$  nor  $\neg\tau$  is  $T$ -positive. First, we can show by induction on  $\alpha$  that  $\ulcorner \tau \urcorner \notin X_{\mathbf{P}_1}^{\alpha}$  for all  $\alpha$ ; thus  $\ulcorner \tau \urcorner \notin \text{DOWN}^{\mathbf{P}_1}(X)$ . Second, we can similarly show that  $\ulcorner \tau \urcorner \notin \text{UP}^{\mathbf{P}_1}(\text{DOWN}^{\mathbf{P}_1}(X)) = X_{\text{down}}^{\mathbf{P}_1}$ . Finally, we will show that  $X_{\text{down}}^{\mathbf{P}_1} \in \text{COMP}_M$ . For any  $a \in \text{St}_{\mathcal{L}_T}^M \setminus \text{SPOS}_M$  and  $Y \subset \text{St}_{\mathcal{L}_T}^M$ ,  $a \in Y$  entails  $a \in \mathbf{P}_1(Y)$ ; therefore, we have

$$\text{DOWN}^{\mathbf{P}_1}(Y) \setminus \text{SPOS}_M = \text{UP}^{\mathbf{P}_1}(Y) \setminus \text{SPOS}_M = Y \setminus \text{SPOS}_M,$$

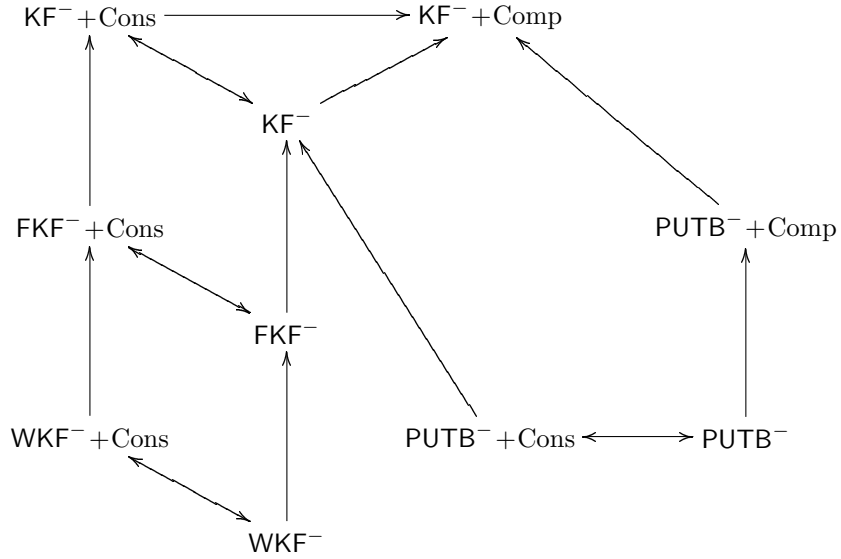
in general. Hence,  $\text{St}_{\mathcal{L}_T}^M \setminus (\text{SPOS}_M \cup \{\ulcorner \tau \urcorner\}) \subset X_{\text{down}}^{\mathbf{P}_1}$ . Then, since  $\neg\tau$  is not  $T$ -positive, we have  $\ulcorner \neg\tau \urcorner \in X_{\text{down}}^{\mathbf{P}_1}$ ; thus we obtain  $X_{\text{down}}^{\mathbf{P}_1} \in \text{COMP}_M$ .

(8) & (9). Immediate from the proof of (7), since  $T^{\ulcorner \tau \urcorner}$  is derivable both in  $\text{WKF}^-$  and  $\text{FKF}^-$ .  $\square$

It follows from the arguments so far that all the systems of truth are distinct except the cases which we have already remarked.

**Proposition 4.4.3.** The following are all distinct theories:  $\text{KF}^-$ ,  $\text{FKF}^-$ ,  $\text{WKF}^-$ ,  $\text{KF}^- + \text{Cons}$ ,  $\text{KF}^- + \text{Comp}$ ,  $\text{FKF}^- + \text{Cons}$ ,  $\text{WKF}^- + \text{Cons}$ ,  $\text{PUTB}^-$ ,  $\text{PUTB}^- + \text{Cons}$ , and  $\text{PUTB}^- + \text{Comp}$ . They all are still distinct even when we add the full induction to them.

We have thus completely correlated all the systems without full-induction in terms of  $<_{\mathcal{I}}$ ,  $\leq_{\mathcal{I}}$ ,  $\equiv_{\mathcal{I}}$  and  $\not\sim_{\mathcal{I}}$ . The results are summarized in the next diagram.



This diagram is read as follows: regarding the diagram as a directed graph,

- if there exists a path from a system  $Q_0$  to a system  $Q_1$ , then  $Q_0 \leq_{\mathcal{I}} Q_1$ , and
- if there is no path from a system  $Q_0$  to a system  $Q_1$ , then  $Q_0 \not\leq_{\mathcal{I}} Q_1$ .

We can see from this diagram that even when the outer theories are different, the inner theories may coincide; e.g.,  $\mathcal{O}(\text{KF} + \text{Cons}) \neq \mathcal{O}(\text{KF})$  but  $\mathcal{I}(\text{KF} + \text{Cons}) = \mathcal{I}(\text{KF})$ .

We have considered systems without full induction (for  $\mathcal{L}_T$ ). When we turn our eyes to systems with it, the situation will become harder to deal with. We cannot make use of the arguments

developed in the present section. For example, even when  $M \models \text{PA}$  can be expanded to a model of  $\text{WKF}$ , it is still open whether  $M$  is also expandable to a model of  $\text{KF}$ , though the converse does hold as was shown in the last chapter. For another example, even when  $M \models \text{PA}$  is expandable to a model of  $\text{KF}^-$ , it is not yet known whether  $\emptyset_{\text{up}}^{\mathbf{S}}$  is a model of  $\text{KF}$ . However, only some partial results on their inner theories can be shown using the so far developed semantical arguments. The other cases will be shown in the subsequent sections by totally different techniques.

**Proposition 4.4.4.** The following hold:

- (1)  $\text{KF}^{(*)} <_{\mathcal{I}} \text{KF}^{(*)} + \text{Comp}$ ,  $\text{FKF}^{(*)} <_{\mathcal{I}} \text{FKF}^{(*)} + \text{Comp}$ ,  $\text{WKF}^{(*)} <_{\mathcal{I}} \text{WKF}^{(*)} + \text{Comp}$  and  $\text{PUTB}^{(*)} <_{\mathcal{I}} \text{PUTB}^{(*)} + \text{Comp}$
- (2)  $\text{PUTB}^{(*)} <_{\mathcal{I}} \text{KF}^{(*)}$
- (3)  $\text{PUTB}^{(*)} \not\sim_{\mathcal{I}} \text{WKF}^{(*)}$ ,  $\text{PUTB}^{(*)} \not\sim_{\mathcal{I}} \text{FKF}^{(*)}$ , and  $\text{KF}^{(*)} \not\sim_{\mathcal{I}} \text{PUTB}^{(*)} + \text{Comp}$ .

*Proof.* By Theorems 4.3.8, 4.3.14 and 4.3.16, we can construct the same unprovable sentences within (a certain expansion of) the standard model  $\mathbb{N}$  even for the cases of systems with full-induction and their schematic reflective closures. □

## 4.5 Proof theory for finitely iterated self-applicable truths

The proof-theoretic strength of  $\text{KF} + \text{Cons}$  was determined by Cantini [5]. In the present section, we will generalize Cantini's technique for other systems  $\text{WKF}$  and  $\text{FKF}$  of self-applicable truth and also for systems of iterated self-applicable truth.

As we have seen, the truth predicates of  $\text{KF}$ ,  $\text{FKF}$  and  $\text{WKF}$  can be defined as a fixed-point of a certain positive arithmetical operator. Hence, they are all syntactically embeddable in the system  $\widehat{\text{ID}}_1$  of fixed-points. With an easy generalization, we can syntactically embed  $\text{KF}_n$ ,  $\text{FKF}_n$  and  $\text{WKF}_n$  in the system  $\widehat{\text{ID}}_n$  of  $n$  iterated fixed-points.

### 4.5.1 Iterated fixed-points and self-applicable truths

In the present subsection, we introduce the systems  $\widehat{\text{ID}}_n$  of  $n$ -iterated fixed-points for  $n > 0$ . The proof-theoretic analyses of  $\widehat{\text{ID}}_n$  are well-known and we will make use of them to correlate the inner theories of systems of self-applicable truth.

We recursively define the language  $\mathcal{L}_n^{\text{fix}}$  of  $\widehat{\text{ID}}_n$  in the following manner. First we set  $\mathcal{L}_0^{\text{fix}}$  to be  $\mathcal{L}_0$ . For each  $n \in \mathbb{N}$ , let  $\mathcal{L}_n^{\text{fix}}(R)$  be  $\mathcal{L}_n^{\text{fix}} \cup \{R\}$  where  $R$  is a fresh unary predicate symbol. An  $\mathcal{L}_n^{\text{fix}}(R)$ -formula  $\mathcal{A}(R, x)$  is called  $(n+1)$ -inductive operator form, if  $\mathcal{A}$  contains at most  $x$  free and  $R$  occurs only positively in  $\mathcal{A}$ . The language  $\mathcal{L}_{n+1}^{\text{fix}}$  is defined as  $\mathcal{L}_n^{\text{fix}}$  plus *unary* predicates  $P_{n+1}^{\mathcal{A}}$  associated to each  $(n+1)$ -inductive operator form  $\mathcal{A}(R, x)$ . Then, the system  $\widehat{\text{ID}}_{n+1}$  comprises PA, full-induction for  $\mathcal{L}_{n+1}^{\text{fix}}$  and the fixed-point axiom schema: for each  $(n+1)$ -inductive operator form  $\mathcal{A}(R, x)$ ,

$$\forall x [P_{n+1}^{\mathcal{A}}(x) \leftrightarrow \mathcal{A}(P_{n+1}^{\mathcal{A}}, x)].$$

Next, we briefly explain how we can embed systems of  $n$ -iterated self-applicable truth in  $\widehat{\text{ID}}_n$ .

First, for the sake of the embedding of  $\mathsf{KF}_n$ , we recursively define  $\mathcal{A}_n(R, x)$  as follows:

$$\begin{aligned}
& x \in \mathsf{St}_a \wedge \forall y, z \in \mathsf{CT}[x = y \doteq z \rightarrow y^\circ = z^\circ] \wedge \forall y, z \in \mathsf{CT}[x = y \neq z \rightarrow y^\circ \neq z^\circ] \\
& \wedge \bigwedge_{m < n} \forall y \in \mathsf{CT}[x = T_m y \rightarrow P_m^{\mathcal{A}_m}(y^\circ)] \wedge \bigwedge_{m < n} \forall y \in \mathsf{CT}[x = \neg T_m y \rightarrow \neg P_m^{\mathcal{A}_m}(y^\circ)] \\
& \wedge \forall y \in \mathsf{CT}[x = T_n y \rightarrow y^\circ \in R] \wedge \forall y \in \mathsf{CT}[x = \neg T_n y \rightarrow \neg y^\circ \in R] \\
& \wedge \forall y \in \mathsf{St}_n[(x = \neg \neg y) \rightarrow y \in R] \\
& \wedge \forall y, z \in \mathsf{St}_n[(x = y \wedge z) \rightarrow Ry \wedge Rz] \wedge \forall y, z \in \mathsf{St}_n[(x = \neg(y \wedge z)) \rightarrow (\neg y \in R \vee \neg z \in R)] \\
& \wedge \forall y, z \in \mathsf{St}_n[(x = y \rightarrow z) \rightarrow (R(\neg y) \vee Rz)] \wedge \forall y, z \in \mathsf{St}_n[(x = \neg(y \rightarrow z)) \rightarrow (Ry \wedge R(\neg z))] \\
& \wedge \forall y, z[\mathsf{St}_n(\forall y.z) \wedge x = \forall y.z \rightarrow \forall w(z(w) \in R)] \\
& \wedge \forall y, z[\mathsf{St}_n(\forall y.z) \wedge x = \neg \forall y.z \rightarrow \exists w(\neg z(w) \in R)].
\end{aligned}$$

Then, we can straightforwardly show that, for each  $n \in \mathbb{N}$ ,

$$\mathsf{KF}_n \vdash \Psi \quad \Rightarrow \quad \widehat{\mathsf{ID}}_n \vdash \Psi[P_1^{\mathcal{A}_1}/T_1, \dots, P_n^{\mathcal{A}_n}/T_n]. \quad (4.1)$$

We can assume without loss of generality that this derivation of  $\Psi[P_1^{\mathcal{A}_1}/T_1, \dots, P_n^{\mathcal{A}_n}/T_n]$  within  $\widehat{\mathsf{ID}}_n$  only uses the vocabulary from  $\mathcal{L}_0 \cup \{P^{\mathcal{A}_1}, \dots, P^{\mathcal{A}_n}\}$ . Let  $\widehat{\mathsf{ID}}_n(\mathcal{A})$  be the subsystem of  $\widehat{\mathsf{ID}}_n$  obtained by restricting the vocabulary to  $\mathcal{L}_n(\mathcal{A}) = \mathcal{L}_0 \cup \{P^{\mathcal{A}_1}, \dots, P^{\mathcal{A}_n}\}$ . Thus, the above (4.1) is strengthened to the following:

$$\mathsf{KF}_n \vdash \Psi \quad \Rightarrow \quad \widehat{\mathsf{ID}}_n(\mathcal{A}) \vdash \Psi[P_1^{\mathcal{A}_1}/T_1, \dots, P_n^{\mathcal{A}_n}/T_n]. \quad (4.2)$$

Similarly, we can define  $\mathcal{B}_m$  and  $\mathcal{C}_m$  ( $m \in \mathbb{N}$ ) such that, for each  $n \in \mathbb{N}$ ,

$$\mathsf{FKF}_n \vdash \Psi \quad \Rightarrow \quad \widehat{\mathsf{ID}}_n(\mathcal{B}) \vdash \Psi[P_1^{\mathcal{B}_1}/T_1, \dots, P_n^{\mathcal{B}_n}/T_n] \quad (4.3)$$

$$\text{WKF}_n \vdash \Psi \quad \Rightarrow \quad \widehat{\text{ID}}_n(\mathcal{C}) \vdash \Psi[P_1^{C_1}/T_1, \dots, P_n^{C_n}/T_n], \quad (4.4)$$

where  $\widehat{\text{ID}}_n(\mathcal{B})$  and  $\widehat{\text{ID}}_n(\mathcal{C})$  are defined analogously to  $\widehat{\text{ID}}_n(\mathcal{A})$  with restricted vocabulary  $\mathcal{B}_1, \dots, \mathcal{B}_n$  and  $\mathcal{C}_1, \dots, \mathcal{C}_n$  respectively; their languages are denoted by  $\mathcal{L}_n(\mathcal{B})$  and  $\mathcal{L}_n(\mathcal{C})$  respectively. For saving space, I explicitly give only the definition of  $\mathcal{C}_n$ :

$$\begin{aligned} & x \in \text{St}_a \wedge \forall y, z \in \text{CT}[x = y \overset{\circ}{=} z \rightarrow y^\circ = z^\circ] \wedge \forall y, z \in \text{CT}[x = y \overset{\neq}{=} z \rightarrow y^\circ \neq z^\circ] \\ & \wedge \bigwedge_{m < n} \forall y \in \text{CT}[x = T_m y \rightarrow P_m^{A_m}(y^\circ)] \wedge \bigwedge_{m < n} \forall y \in \text{CT}[x = \neg T_m y \rightarrow \neg P_m^{A_m}(y^\circ)] \\ & \wedge \forall y \in \text{CT}[x = T_n y \rightarrow y^\circ \in R] \wedge \forall y \in \text{CT}[x = \neg T_n y \rightarrow \neg y^\circ \in R] \\ & \wedge \forall y \in \text{St}_n[x = \neg \neg y \rightarrow y \in R] \wedge \forall y, z \in \text{St}_n[x = y \wedge z \rightarrow Ry \wedge Rz] \\ & \wedge \forall y, z \in \text{St}_n[x = \neg(y \wedge z) \rightarrow ((\neg y \in R \wedge \neg z \in R) \vee (\neg y \in R \wedge z \in R) \vee (y \in R \wedge \neg z \in R))] \\ & \wedge \forall y, z \in \text{St}_n[(x = y \rightarrow z) \rightarrow (R(\neg y) \vee Rz)] \wedge \forall y, z \in \text{St}_n[(x = \neg(y \rightarrow z)) \rightarrow (Ry \wedge R(\neg z))] \\ & \wedge \forall y, z[\text{St}_n(\forall y.z) \wedge x = \forall y.z \rightarrow \forall w(z(w) \in R)] \\ & \wedge \forall y, z[\text{St}_n(\forall y.z) \wedge x = \neg \forall y.z \rightarrow \exists w(\neg z(w) \in R)]; \end{aligned}$$

$\mathcal{B}_n$  can be straightforwardly obtained from  $\mathcal{C}_n$  simply by making suitable changes on the clauses for the conditional ‘ $\rightarrow$ ’.

#### 4.5.2 Semi-Formal Systems for $\widehat{\text{ID}}_n$

In order to give proof-theoretic analyses to  $\widehat{\text{ID}}_n$ , we need to make use of a semi-formal system  $\text{H}_n$  for it in which infinitary derivations are allowed. By (4.2)–(4.4), we can restrict our proof-theoretic analyses to certain subsystems of  $\widehat{\text{ID}}_n$  with restricted vocabulary.

In the present subsection, we will formulate the semi-formal systems  $\text{H}_n(\mathcal{A})$  over  $\mathcal{L}_n(\mathcal{A})$ ,  $\text{H}_n(\mathcal{B})$  over  $\mathcal{L}_n(\mathcal{B})$  and  $\text{H}_n(\mathcal{C})$  over  $\mathcal{L}_n(\mathcal{C})$  for each  $n > 0$ . They are formulated in Tait-calculus and thus we modify the language by assuming that the formulae are all given in their negative normal forms: i.e.,

formulae are assumed to be constructed from literals by  $\wedge$ ,  $\vee$ ,  $\exists$  and  $\forall$ . As is usual in  $\omega$ -arithmetic, we also assume that only closed formulae (i.e., sentences) are allowed throughout this section. For more details of the formulation of  $\omega$ -arithmetic in Tait Calculus, see [64] or [66].

For each  $n \in \mathbb{N}$ , we first define the semi-formal system  $H_n(\mathcal{A})$  for  $\widehat{ID}_n(\mathcal{A})$  as follows.

**Axioms:**

Ax0.  $\Gamma, A$ , for a true  $\mathcal{L}_{PA}$ -literal  $A$ .

Ax1.  $\Gamma, \phi(t), \neg\phi(s)$ , for an atomic  $\phi$  of the form  $P_m(t)$  and closed terms  $s$  and  $t$  with  $s^{\mathbb{N}} = t^{\mathbb{N}}$ .

**Logical rules and Cut:** the same as the ordinary  $\omega$ -arithmetic; more explicitly,

$$\frac{\Gamma, A}{\Gamma, A \vee B} \quad \frac{\Gamma, B}{\Gamma, A \vee B} \quad \frac{\Gamma, A \quad \Gamma, B}{\Gamma, A \wedge B} \quad \frac{\Gamma, A(t)}{\Gamma, \exists x A(x)} \quad \frac{\Gamma, A(\bar{n}) \text{ for all numerals } \bar{n}}{\Gamma, \forall x A(x)}$$

**Fixed-point rules:** for each  $1 \leq m \leq n$ ,

$$\frac{\Gamma, \mathcal{A}(P_m^{\mathcal{A}m}, s)}{\Gamma, P_m^{\mathcal{A}m}(s)} \quad \frac{\Gamma, \neg\mathcal{A}(P_m^{\mathcal{A}m}, s)}{\Gamma, \neg P_m^{\mathcal{A}m}(s)}$$

For each derivation  $\mathcal{D}$  of  $H_n(\mathcal{A})$ , we can assign an ordinal to  $\mathcal{D}$  as its length in the usual manner. We also assign its cut-rank to  $\mathcal{D}$  as the maximum surface complexity of the critical formulae of cut. Then, we write  $H_n(\mathcal{A})|_{\rho}^{\delta} \Gamma$ , if  $\Gamma$  is obtained by a derivation in  $H_n(\mathcal{A})$  with the length  $\beta$  and cut-rank  $\rho$ . We also write  $H_n(\mathcal{A})|_{<\rho}^{\leq\delta} \Gamma$ , when  $H_n|_{\rho'}^{\delta'} \Gamma$  for some  $\delta' < \delta$  and  $\rho' < \rho$ . We write  $H_n|_{\star}^{\delta} \Gamma$  (or  $H_n|_{\star}^{\leq\delta} \Gamma$ ) when there is a derivation of  $\Gamma$  in  $H_n$  with the length  $\delta$  ( $< \delta$ , resp.) in which all its cut formulae are of the forms  $P_{\beta}^{\mathcal{A}}t$  and  $\neg P_{\beta}^{\mathcal{A}}$  ( $\beta < \alpha$ ).

Next, we analogously define the semi-formal systems  $H_n(\mathcal{B})$  and  $H_n(\mathcal{C})$  for  $\widehat{ID}_n(\mathcal{B})$  and  $\widehat{ID}_n(\mathcal{C})$ ;

they are obtained from  $H_n(\mathcal{A})$  by replacing the two Fixed-point rules respectively by

$$\frac{\Gamma, \mathcal{B}(P_m^{\mathcal{B}}, s)}{\Gamma, P_m^{\mathcal{B}}(s)} \qquad \frac{\Gamma, \neg \mathcal{B}(P_m^{\mathcal{B}}, s)}{\Gamma, \neg P_m^{\mathcal{B}}(s)}$$

and

$$\frac{\Gamma, \mathcal{C}(P_m^{\mathcal{C}}, s)}{\Gamma, P_m^{\mathcal{C}}(s)} \qquad \frac{\Gamma, \neg \mathcal{C}(P_m^{\mathcal{C}}, s)}{\Gamma, \neg P_m^{\mathcal{C}}(s)} .$$

The aforementioned standard properties are also parallelly provable in  $\widehat{ID}_n(\mathcal{B})$  and  $\widehat{ID}_n(\mathcal{C})$ .

The desired basic properties such as Structural Lemma, Inversion Lemma for  $\wedge$  and  $\forall$ , Exportation Lemma for  $\vee$  and Numerical Equivalence (cf. [66]), can be shown in the usual manners for those systems. Since Tautology Lemma and Induction Lemma can also be shown in the usual manner, we standardly obtain the following two lemmata; cf. [4, 43, 66].

**Lemma 4.5.1.** If  $\widehat{ID}_n(\mathcal{A}) \vdash \phi$ , then  $H_n(\mathcal{A}) \vdash_{\rho}^{\omega+\omega} \phi$  for some finite  $\rho \in \mathbb{N}$ . The parallel statement holds between  $\widehat{ID}_n(\mathcal{B})$  and  $H_n(\mathcal{B})$  and between  $\widehat{ID}_n(\mathcal{C})$  and  $H_n(\mathcal{C})$ .

**Lemma 4.5.2.** If  $H_n(\mathcal{A}) \vdash_{\rho}^{\beta} \Gamma$  then  $H_n(\mathcal{A}) \vdash_{\frac{1}{1}}^{\leq \varepsilon(\beta)} \Gamma$ , where  $\varepsilon(\beta)$  is the least  $\varepsilon$ -number greater than  $\beta$ . The analogous statements hold for  $H_n(\mathcal{B})$  and  $H_n(\mathcal{C})$ .

Now, we define a ‘ramified’ semi-formal system  $AH_n(\mathcal{A})$ . The language  $\mathcal{L}_n^A$  of  $AH_n(\mathcal{A})$  consists of the sentences (in Tait-style) of  $\mathcal{L}_0 \cup \{P_m \mid 1 \leq m < n\} \cup \{P_n^\eta \mid \eta\}$ , where  $\eta$  ranges over ordinals up to a sufficiently large ordinal, say,  $\Gamma_0$ ; in other words, the predicate  $P_n$  of the highest level is further typed by ordinals. We drop the superscripts  $\mathcal{A}_m$ ’s ( $m \leq n$ ); but this will cause no ambiguity.

For each  $n \in \mathbb{N}$ , the system  $AH_n(\mathcal{A})$  over  $\mathcal{L}_n^A$  is defined as follows.

**Axioms:**

Ax0.  $\Gamma, A$ , for a true  $\mathcal{L}_{PA}$ -literal  $A$ .

Ax1.  $\Gamma, \phi(t), \neg\phi(s)$ , for an atomic  $\phi$  of the form  $P_m(t)$  ( $m < n$ ) and closed terms  $s$  and  $t$  with  $s^{\mathbb{N}} = t^{\mathbb{N}}$ .

Ax2.  $\Gamma, \neg P_n^0(t)$ , for any closed term  $t$ .

**Logical rules and Cut** the same as  $H_n$  (i.e., those of  $\omega$ -arithmetic).

**Fixed-point rules for  $P_m$  ( $m < n$ ):** the same as  $H_{n-1}$ .

**Fixed-point rules for  $P_n^\eta$ :** for closed terms  $s$  and  $t$  with  $|t| = \alpha$ ,

$$\frac{\Gamma, \mathcal{A}_m(P_m^\eta, s)}{\Gamma, P_m^{\eta+1}(s)} \quad \frac{\Gamma, \neg\mathcal{A}_m(P_m^\eta, s)}{\Gamma, \neg P_m^{\eta+1}(s)}$$

**Limit rules for  $P_n^\lambda$ :** for a limit  $\lambda$  and a closed term  $s$ ,

$$\frac{\Gamma, P_n^\eta(s) \quad (\text{for some } \eta < \lambda)}{\Gamma, P_n^\lambda(s)} \quad \frac{\Gamma, \neg P_n^\eta(s) \quad (\text{for all } \eta < \lambda)}{\Gamma, \neg P_n^\lambda(s)}$$

For each  $\mathcal{L}_n^A$ -formula  $\phi$ , we define its  $n$ -rank,  $rk_n(\phi)$ , as follows: If  $\phi$  is a literal of the form  $P_n^\eta t$  or  $\neg P_n^\eta t$  then  $rk_n(\phi) = \omega \cdot \eta$ ; if  $\phi$  is a literal of any other form,  $rk_n(\phi) = 0$ ;  $rk_n(\phi)$  for a compositional formula  $\phi$  is inductively defined in the standard manner; e.g.,  $rk_n(\psi_0 \wedge \psi_1) := \max\{rk_n(\psi_0), rk_n(\psi_1)\} + 1$ . We write  $\text{AH}_n(\mathcal{A})|_{\rho}^{\delta} \Gamma$  (or  $\text{AH}_n(\mathcal{A})|_{<\rho}^{\leq\delta} \Gamma$ ) when  $\Gamma$  is derived in  $\text{AH}_n(\mathcal{A})$  with the length  $\delta$  ( $< \delta$  resp.) and cut-rank (w.r.t.  $rk_n$ )  $\rho$  ( $< \rho$ , resp.).

The systems  $\text{AH}_n(\mathcal{B})$  and  $\text{AH}_n(\mathcal{C})$  are analogously defined: i.e., they are obtained from  $\text{AH}_n(\mathcal{A})$  by replacing the Fixed-point rules for  $P_n^\eta$  respectively by

$$\frac{\Gamma, \mathcal{B}_m(P_m^\eta, s)}{\Gamma, P_m^{\eta+1}(s)} \quad \frac{\Gamma, \neg\mathcal{B}_m(P_m^\eta, s)}{\Gamma, \neg P_m^{\eta+1}(s)}$$

and

$$\frac{\Gamma, \mathcal{C}_m(P_m^\eta, s)}{\Gamma, P_m^{\eta+1}(s)} \quad \frac{\Gamma, \neg \mathcal{C}_m(P_m^\eta, s)}{\Gamma, \neg P_m^{\eta+1}(s)} .$$

Also,  $\text{AH}_n(\mathcal{B})|_{\frac{\delta}{\rho}} \Gamma$  (or  $\text{AH}_n(\mathcal{B})|_{\frac{\leq \delta}{< \rho}} \Gamma$ ) and  $\text{AH}_n(\mathcal{C})|_{\frac{\delta}{\rho}} \Gamma$  (or  $\text{AH}_n(\mathcal{C})|_{\frac{\leq \delta}{< \rho}} \Gamma$ ) are defined in analogous manners.

### 4.5.3 Asymmetric Interpretation

Let us fix  $n > 0$  throughout the present subsection. In the present subsection, for the sake of simplicity and readability, we will focus on systems with the restricted vocabulary  $\mathcal{A}_1, \dots, \mathcal{A}_n$ , i.e.,  $\widehat{\text{ID}}_n(\mathcal{A})$ ,  $\text{H}_n(\mathcal{A})$  and  $\text{AH}_n(\mathcal{A})$  ( $n \in \mathbb{N}$ ). However, all the following arguments can be directly applied to the other two cases. We also suppress ‘ $\mathcal{A}$ ’ from the names of these systems and their predicates.

Let  $A$  be an  $\mathcal{L}_n^{\text{fix}}$ -formula. Then, an  $\mathcal{L}_n^\infty$ -formula  $A[\delta, \gamma]$  is the result of replacing each negative and positive occurrences of  $T_n$  respectively by  $P_n^\delta$  and  $P_n^\gamma$ . We particularly notice that  $(\neg A)[\delta, \gamma]$  is identical with  $\neg(A[\gamma, \delta])$ .

**Lemma 4.5.3.** Let  $\delta \leq \gamma$ . Then,  $\text{AH}_n|_{\frac{\omega \cdot \delta}{\omega \cdot \delta}} \neg P_n^\delta s, P_n^\gamma t$ .

*Proof.* Induction on  $\delta$ . The case in which  $\delta = 0$  is trivial from the Ax2. Let us assume  $\delta > 0$ .

First, suppose  $\gamma = \gamma' + 1$  and  $\delta = \delta' + 1$ . Since  $\delta' \leq \gamma'$ , we obtain by IH

$$\text{AH}_n|_{\frac{\omega \cdot \delta'}{\omega \cdot \delta'}} \neg P_n^{\delta'} r_0, P_n^{\gamma'} r_1,$$

for arbitrary terms  $r_0$  and  $r_1$ . Using this, we can show by the argument of the so-called Tautology Lemma (cf. [66]) that

$$\text{AH}_n|_{\frac{\omega \cdot \delta' + k}{\omega \cdot \delta'}} \neg \mathcal{A}_n(P_n^{\delta'}, s), \mathcal{A}_n(P_n^{\gamma'}, t),$$

for some finite  $k$  depending on  $sc(\mathcal{A})$  (i.e., the surface complexity of  $\mathcal{A}$ ). We obtain

$$\left| \frac{\omega \cdot \delta' + k + 2}{\omega \cdot \delta'} \right| \neg P_n^\delta s, P_n^\gamma t.$$

Second, suppose  $\gamma = \gamma' + 1$  and  $\delta$  is limit. Take an arbitrary  $0 < \delta' < \delta$ . Since  $\delta' < \gamma'$ , it follows from IH that

$$\text{AH}_n \left| \frac{\omega \cdot \delta'}{\omega \cdot \delta'} \right| \neg P_n^{\delta'} r_0, P_n^{\gamma'} r_1,$$

for arbitrary terms  $r_0$  and  $r_1$ . Then, again by the Tautology Lemma argument, we have

$$\text{AH}_n \left| \frac{\omega \cdot \delta' + k}{\omega \cdot \delta'} \right| \neg \mathcal{A}_n(P_n^{\delta'}, s), \mathcal{A}_n(P_n^{\gamma'}, t),$$

for some finite  $k$ . Hence, we obtain

$$\text{AH}_n \left| \frac{\omega \cdot \delta' + k + 2}{\omega \cdot \delta'} \right| \neg P_n^{\delta'+1} s, P_n^\gamma.$$

Since  $\delta$  is limit, we have  $\delta' + 1 < \delta$  and we already have by IH that

$$\text{AH}_n \left| \frac{\omega \cdot \delta'}{\omega \cdot \delta'} \right| \neg P_n^{\delta'} s, P_n^{\delta'+1}.$$

Since  $\delta' + 1 < \delta$ , we obtain by cut that

$$\text{AH}_n \left| \frac{\omega \cdot \delta' + k + 3}{\omega \cdot (\delta' + 1) + 1} \right| \neg P_n^{\delta'} s, P_n^\gamma.$$

Since  $\delta' < \delta$  is an arbitrary and  $\omega \cdot \delta' + k' + 3, \omega \cdot (\delta' + 1) + 1 < \omega \cdot \delta$ , we finally obtain by the limit

rule that

$$\text{AH}_n \mid \frac{\omega \cdot \delta}{\omega \cdot \delta} \neg P_n^\delta s, P_n^\gamma.$$

Third, suppose  $\gamma$  is limit and  $\delta = \delta' + 1$ . Similarly to the previous cases, we obtain by IH and the Tautology Lemma argument, we have

$$\text{AH}_n \mid \frac{\omega \cdot \delta' + k}{\omega \cdot \delta'} \neg \mathcal{A}_n(P_n^{\delta'}, s), \mathcal{A}_n(P_n^{\delta'}, t),$$

for some finite  $k$ , and we then have

$$\text{AH}_n \mid \frac{\omega \cdot \delta' + k + 2}{\omega \cdot \delta'} \neg P_n^\delta s, P_n^\delta.$$

By the limit rule, we obtain

$$\text{AH}_n \mid \frac{\omega \cdot \delta' + k + 3}{\omega \cdot \delta'} \neg P_n^\delta s, P_n^\gamma t.$$

Finally suppose  $\gamma$  and  $\delta$  are both limit. For each  $\delta' < \delta$ , by IH we have

$$\text{AH}_n \mid \frac{\omega \cdot \delta'}{\omega \cdot \delta'} \neg P_n^{\delta'} s, P_n^{\delta'} t.$$

By limit rule we have

$$\text{AH}_n \mid \frac{\omega \cdot \delta' + 1}{\omega \cdot \delta'} \neg P_n^{\delta'} s, P_n^\gamma t.$$

Since  $\delta' < \gamma$  is arbitrary, we obtain by (the other) limit rule again that

$$\text{AH}_n \left| \frac{\omega \cdot \delta}{\omega \cdot \delta} \right| \neg P_n^\delta s, P_n^\gamma t. \quad \square$$

The following three lemmata can be shown standardly using the last lemma.

**Lemma 4.5.4.** Let  $\delta_0 \leq \delta_1 \leq \gamma_1 \leq \gamma_0$ . Then,  $\text{AH}_n \left| \frac{\omega \cdot \gamma_1 + 2 \cdot \text{sc}(A)}{\omega \cdot \delta_1} \right| A[\delta_0, \gamma_0], \neg A[\gamma_1, \delta_1]$ .

**Lemma 4.5.5.** Let  $\delta_0 \leq \delta_1 \leq \gamma_1 \leq \gamma_0$ . Suppose  $\left| \frac{\beta}{\rho} \right| A[\delta_1, \gamma_1], \Gamma$ . Then,

$$\text{AH}_n \left| \frac{\max\{\omega \cdot \gamma_1 + 2 \cdot \text{sc}(A), \beta\} + 1}{\max\{\omega \cdot \gamma_1 + \text{sc}(A) + 1, \rho\}} \right| A[\delta_0, \gamma_0], \Gamma.$$

**Lemma 4.5.6** (Persistency Lemma). Let  $\Gamma = \{\theta_1, \dots, \theta_k\}$  and let  $\vec{\delta}, \vec{\delta}', \vec{\gamma}, \vec{\gamma}'$  be such that  $\delta'_i \leq \delta_i \leq \gamma_i \leq \gamma'_i$  for each  $1 \leq i \leq k$ . Then, if  $\left| \frac{\beta}{\rho} \right| \Gamma[\vec{\delta}, \vec{\gamma}], \Delta$ , it holds that

$$\text{AH}_n \left| \frac{\max\{\omega \cdot |\vec{\gamma}'| + 2 \cdot |\Gamma|, \beta\} + k}{\max\{\omega \cdot |\vec{\gamma}'| + |\Gamma| + 1, \rho\}} \right| \Gamma[\vec{\delta}', \vec{\gamma}'], \Delta,$$

where  $|\vec{\gamma}'|$  is  $\max\{\gamma'_1, \dots, \gamma'_k\}$  and  $|\Gamma|$  is  $\max\{\text{sc}(\theta_1), \dots, \text{sc}(\theta_k)\}$ .

The next lemma can also be standardly shown by using the previous four lemmata, but I will give some details since it is the main lemma of the present subsection.

**Lemma 4.5.7** (Embedding Lemma). Suppose  $\text{H}_n \left| \frac{\delta}{\star} \right| \Gamma$ . Then for each ordinal  $\eta$ ,

$$\text{AH}_n \left| \frac{\omega(\eta + \omega^\delta)}{\omega(\eta + \omega^\delta)} \right| \Gamma[\eta, \eta + \omega^\delta].$$

*Proof.* The claim is shown by induction on  $\delta$ .

For the base step, suppose  $\text{H}_n \left| \frac{0}{\star} \right| \Gamma$  is obtained by an axiom. The case where it is obtained by (Ax0) is trivial. Suppose it is obtained by (Ax1). If the critical formulae are in  $\mathcal{L}_m$  ( $m < n$ ), then

the claim obtains by the same axiom. Otherwise, the critical formulae are  $P_n t$  and  $\neg P_n s$  for some  $t^{\mathbb{N}} = s^{\mathbb{N}}$ ; then, by the previous Lemma and Structural Lemma, we obtain

$$\text{AH}_n \mid \frac{\omega \cdot (\eta + \omega^0)}{\omega \cdot (\eta + \omega^0)} \Gamma[\eta, \eta + \omega^0].$$

In the rest of proof, we assume that  $\delta > 0$ . We only illustrate some typical cases.

First, suppose the last inference is by the positive fixed-point rule. Then, the premise is of the form:

$$\text{H}_n \mid \frac{\delta'}{\star} \mathcal{A}_n(P_n, s), \Gamma;$$

we thus obtain by IH that

$$\text{AH}_n \mid \frac{\omega(\eta + \omega^{\delta'})}{\omega(\eta + \omega^{\delta'})} \Gamma[\eta, \eta + \omega^{\delta'}], \mathcal{A}_n(P_n^{\eta + \omega^{\delta'}}, s).$$

By Persistency Lemma,

$$\text{AH}_n \mid \frac{\omega(\eta + \omega^{\delta'}) + k}{\omega(\eta + \omega^{\delta'}) + k} \Gamma[\eta, \eta + \omega^{\delta'}], \mathcal{A}_n(P_n^{\omega^{\delta'}}, s),$$

for some finite  $k$ . Then, by fixed-point rule, we have

$$\text{AH}_n \mid \frac{\omega(\eta + \omega^{\delta'}) + k + 1}{\omega(\eta + \omega^{\delta'}) + k} \Gamma[\eta, \eta + \omega^{\delta'}], P_n^{\eta + \omega^{\delta'} + 1}(s).$$

By Persistency Lemma again, we finally obtain that, for some finite  $k'$ ,

$$\mid \frac{\omega \cdot (\eta + \omega^{\delta'} + 1) + k'}{\omega \cdot (\eta + \omega^{\delta'} + 1) + k'} \Gamma[\eta, \eta + \omega^{\delta'}], P_n^{\eta + \omega^{\delta'}}(s).$$

Next, suppose the last inference is the negative fixed-point rule: that is, the premise is of the form

$$\mathsf{H}_n \left| \frac{\delta'}{\star} \right. \neg \mathcal{A}_n(P_n, s), \Gamma,$$

for some  $\delta' < \delta$ . We obtain by IH that

$$\mathsf{AH}_n \left| \frac{\omega(\eta + \omega^{\delta'})}{\omega(\eta + \omega^{\delta'})} \right. \Gamma[\eta, \eta + \omega^{\delta'}], \neg \mathcal{A}_n(P_n^\eta, s).$$

Then, by Persistency Lemma, we have

$$\mathsf{AH}_n \left| \frac{\omega(\eta + \omega^{\delta'}) + k}{\omega(\eta + \omega^{\delta'}) + k} \right. \Gamma[\eta, \eta + \omega^{\delta'}], \neg \mathcal{A}(P_n^\eta, s)$$

for some finite  $k$ . By fixed-point rule, we have

$$\mathsf{AH}_n \left| \frac{\omega(\eta + \omega^{\delta'}) + k + 1}{\omega(\eta + \omega^{\delta'}) + k} \right. \Gamma[\eta, \eta + \omega^{\delta'}], \neg P_\beta^{\omega^\eta + 1}(s).$$

By Persistency Lemma again, we finally obtain

$$\mathsf{AH}_n \left| \frac{\omega(\eta + \omega^\delta) + k + 2}{\omega(\eta + \omega^\delta) + k} \right. \Gamma[\eta, \eta + \omega^\delta], \neg P_\beta^\eta(s).$$

As the final and crucial example, suppose the last inference is cut. Thus the premises are of the forms

$$\mathsf{H}_n \left| \frac{\delta_0}{\star} \right. \Gamma, P_n t \quad \mathsf{H}_n \left| \frac{\delta_1}{\star} \right. \Gamma, \neg P_n t,$$

for some  $\delta_0, \delta_1 < \delta$ . By IH, we have

$$\begin{aligned} & \text{AH}_n \left| \frac{\omega(\eta + \omega^{\delta_0})}{\omega(\eta + \omega^{\delta_0})} \Gamma[\eta, \eta + \omega^{\delta_0}], P_n^{\eta + \omega^{\delta_0}}(t) \right. \\ & \left. \text{AH}_n \left| \frac{\omega(\eta + \omega^{\delta_0} + \omega^{\delta_1})}{\omega(\eta + \omega^{\delta_0} + \omega^{\delta_1})} \Gamma[\eta + \omega^{\delta_0}, \eta + \omega^{\delta_0} + \omega^{\delta_1}], \neg P_n^{\eta + \omega^{\delta_0}}(t) \right. \right. \end{aligned}$$

Since  $\omega^{\delta_0} + \omega^{\delta_1} < \omega^\delta$  and thus by Persistency Lemma,

$$\begin{aligned} & \text{AH}_n \left| \frac{\omega(\eta + \omega^{\delta'}) + k_0}{\omega(\eta + \omega^{\delta_0}) + k_0} P_n^{\eta + \omega^{\delta_0}}(t), \Gamma[\eta, \eta + \omega^\delta] \right. \\ & \left. \text{AH}_n \left| \frac{\omega(\eta + \omega^{\delta_0} + \omega^{\delta_1}) + k_1}{\omega(\eta + \omega^{\delta_0} + \omega^{\delta_1}) + k_1} \neg P_n^{\eta + \omega^{\delta_0}}(t), \Gamma[\eta, \eta + \omega^\delta], \right. \right. \end{aligned}$$

for some finite  $k_0$  and  $k_1$ . We finally obtain by applying cut that

$$\text{AH}_n \left| \frac{\omega(\eta + \omega^\delta)}{\omega(\eta + \omega^\delta)} \Gamma[\eta, \eta + \omega^\delta] \right.$$

The other cases can be shown in similar ways and we complete the proof.  $\square$

The next two lemmata are all standardly shown; cf. [66].

**Lemma 4.5.8** (Reduction Lemma). Suppose  $\left| \frac{\beta}{1+\rho} \phi, \Gamma \right.$  and  $\left| \frac{\gamma}{1+\rho} \neg\phi, \Gamma \right.$ . Then  $\left| \frac{\beta\#\gamma}{1+\rho} \Gamma \right.$

**Lemma 4.5.9** (Predicative Cut Elimination). If  $\text{AH}_n \left| \frac{\delta}{1+\rho+\omega^\beta} \Gamma \right.$ , then  $\text{AH}_n \left| \frac{\varphi_\beta(\delta)}{1+\rho} \Gamma \right.$

The next immediately follows from this last lemma.

**Lemma 4.5.10.** If  $\widehat{\text{ID}}_n \vdash \phi$  then  $\text{AH}_n \left| \frac{<\varphi_{\varepsilon_0}^0}{1} \phi[0, \xi] \right.$  for some  $\xi < \varepsilon_0$ .

**Lemma 4.5.11.** If  $\text{AH}_n \left| \frac{\delta}{\rho} P_n^0, \Gamma \right.$ , then  $\text{AH}_n \left| \frac{\delta}{\rho} \Gamma \right.$

*Proof.* This is shown by induction on  $\delta$ ; use the fact that an atomic of the form  $P_n^0$  cannot be critical in any derivation in  $\text{AH}_n$ .  $\square$

**Lemma 4.5.12** (Partial Cut Elimination). If  $H_n \mid_{\rho}^{\delta} \Gamma$  (for any  $\rho < \omega$ ), then  $H_n \mid_{\star}^{\varepsilon(\delta)}$ .

**Lemma 4.5.13.** Let  $n > 1$ . If  $\widehat{ID}_n \vdash \phi$  for  $\phi \in \mathcal{L}_m^{\text{fix}}(\mathcal{A})$  then  $H_{n-1} \mid_{\star}^{\leq \varphi_{\varepsilon_0} 0} \phi$ , for some  $\xi < \varepsilon_0$ .

*Proof.* Suppose  $\widehat{ID}_n \vdash \phi$ . By Lemma 4.5.10, we have  $AH_n \mid_1^{\delta} \phi$  for some  $\xi < \varepsilon_0$  and  $\delta < \varphi_{\varepsilon_0} 0$ . By Lemma 4.5.11, we can assume without loss of generality that this derivation contains cuts applied only to  $\mathcal{L}_{n-1}^{\text{fix}}$ -literals. Hence, we have  $H_n \mid_1^{\delta} \phi$ . Finally, it follows from the last lemma that  $H_n \mid_{\star}^{\varepsilon(\delta)} \phi$  and  $\varepsilon(\delta) < \varphi_{\varepsilon_0} 0$ .  $\square$

**Lemma 4.5.14.** Let us recursively define  $\alpha_0 = \varepsilon_0$  and  $\alpha_{k+1} = \varphi_{\alpha_k} 0$ . Suppose  $\widehat{ID}_n \vdash \phi$  for an  $\mathcal{L}_m^{\text{fix}}$ -formula  $\phi$  ( $m \leq n$ ). Then  $AH_m \mid_1^{\leq \alpha_{n-m+1}} \phi[0, \xi]$  for some  $\xi < \alpha_{n-m}$ .

## 4.6 Inner theory comparison via semi-formal systems

In the present section, we will present how the inner theories of systems of iterative compositional truth with full-induction are compared.

### 4.6.1 Systems with and without Cons

The main theorem of the present subsection is:

**Theorem 4.6.1.** (1) If  $\text{KF} + \text{Cons} \vdash Tt$  then  $\text{KF} \vdash Tt$ . The parallel statements holds between

$\text{FKF} + \text{Cons}$  and  $\text{FKF}$  and between  $\text{WKF} + \text{Cons}$  and  $\text{WKF}$ .

(2) If  $(\text{KF} + \text{Cons})^* \vdash Tt$  then  $\text{KF}^* \vdash Tt$ . The parallel statements holds between  $(\text{FKF} + \text{Cons})^*$

and  $\text{FKF}^*$  and between  $(\text{WKF} + \text{Cons})^*$  and  $\text{WKF}^*$ .

We will show this theorem in the rest of the present subsection.

For the sake of readability, let us simply write  $\text{Cons}_n$  for  $\text{Cons}_n[P_1^{\mathcal{A}_1}/T_1, \dots, P_n^{\mathcal{A}_n}/T_n]$  in what follows. Again, we will focus on  $\text{KF}_n$  but the following arguments parallely apply to  $\text{WKF}_n$  and  $\text{FKF}_n$  as well.

First, suppose  $\text{KF}_n + \text{Cons} \vdash Tt$ . By (4.2), we have  $\widehat{\text{ID}}_1 \vdash \text{Cons} \rightarrow P_1^{A_1}t$ . Then, it follows from Lemma 4.5.10 that  $\text{AH}_1(\mathcal{A}) \upharpoonright_{\frac{\leq \alpha_1}{1}} \neg \text{Cons}[0, \eta], P_1^\eta t$  for some  $\eta < \varepsilon_0$ ; recall that  $\alpha_1 := \varphi_{\varepsilon_0} 0$ .

Second, for the reflective closure, if  $\text{KF}^* + \text{Cons} \vdash Tt$ , we have  $\text{KF}_n + \text{Cons}_n \vdash Tt$  for some  $n$ ; then we similarly obtain  $\text{AH}_n(\mathcal{A}) \upharpoonright_{\frac{\leq \alpha_1}{1}} \neg \text{Cons}_n, P_1 t$ : that is, we have

$$\text{KF}^* + \text{Cons} \vdash Tt \quad \Rightarrow \quad \text{KF}_n \vdash \text{Cons}_n \rightarrow Tt \quad \Rightarrow \quad \text{AH}_n(\mathcal{A}) \upharpoonright_{\frac{\leq \alpha_1}{1}} \neg \text{Cons}_n, P_1 t.$$

The next is the main lemma of the present subsection.

**Lemma 4.6.2.** Let  $\rho > 0$ . For arbitrary  $\eta_0$  and  $\eta_1$ , Suppose that (a)  $\text{AH}_n(\mathcal{A}) \upharpoonright_{\frac{\delta_0}{\rho}} P_n^{\eta_0} s, \Gamma$  and (b)  $\text{AH}_n(\mathcal{A}) \upharpoonright_{\frac{\delta_1}{\rho}} P_n^{\eta_1} \neg s, \Gamma$ . Then, we have  $\text{AH}_n(\mathcal{A}) \upharpoonright_{\frac{\delta_0 \# \delta_1}{\rho}} \Gamma$ , where  $\delta_0 \# \delta_1$  is the natural sum of  $\delta_0$  and  $\delta_1$ .

*Proof.* This lemma informally corresponds to the fact that the minimal Kripkean fixed-point is consistent; in  $\text{AH}_n(\mathcal{A})$ , each  $P_n^\eta$  corresponds to the  $\eta$ -th stage of the Kripkean construction (of the  $n$ -iterated truth) from the empty set, since we start with  $\neg P_n^0 t$  for all  $t$ . In what follows, we suppress ‘ $\mathcal{A}$ ’.

We show the claim by induction on  $\delta_0 \# \delta_1$ . The base case is easy, since an atomic of the form  $P_n^\eta t$  cannot be critical in any axiom. For the induction step, the cases where either  $P_n^{\eta_0} s$  or  $P_n^{\eta_1} \neg s$  is not critical in the last inference follow from IH. Let us assume  $P_n^{\eta_0} s$  and  $P_n^{\eta_1} \neg s$  are both critical in the last inferences. When either of the last inferences is by the limit rule, the premises contain  $\text{AH}_n \upharpoonright_{\frac{\delta_2}{\rho}} P_n^{\eta_2} (\neg) s, \Gamma$  for  $\delta_2 < \delta_0$  (or  $< \delta_1$ ) and  $\eta_2 < \eta_0$  (or  $< \eta_1$ ); then, the claim follows from IH. Hence, we focus on the essential case in which the last inferences are both the fixed-point rule. Then, we can assume that the premises are

$$(c) \text{AH}_n \upharpoonright_{\frac{\delta'_0}{\rho}} \Gamma, \mathcal{A}_n(P_\alpha^{\eta_0-1}, s) \quad \text{and} \quad (d) \text{AH}_n \upharpoonright_{\frac{\delta'_1}{\rho}} \Gamma, \mathcal{A}_n(P_\alpha^{\eta_1-1}, \neg s),$$

for some  $\delta'_0 < \delta_0$  and  $\delta'_1 < \delta_1$ . By the form of  $\mathcal{A}_n$  (and inversion), we can assume that  $s \in \text{St}_n$ ; let

$s = \ulcorner \sigma \urcorner$ . The claim is shown by cases according to the form of  $\sigma$ .

Suppose no negation ‘ $\neg$ ’ is attached to  $\sigma$ . We illustrate two cases. First, assume  $\sigma$  is  $\phi \wedge \psi$ . By  $\wedge$ -inversion,  $\forall$ -inversion and  $\vee$ -exportation, it follows from (c) that

$$\text{AH}_n \mid_{\rho}^{\delta'_0} \Gamma, \neg \text{St}_n(\ulcorner \phi \urcorner), \neg \text{St}_n(\ulcorner \psi \urcorner), s \neq \ulcorner \phi \wedge \psi \urcorner, P_n^{\eta_0-1}(\ulcorner \phi \urcorner) \wedge P_n^{\eta_0-1}(\ulcorner \psi \urcorner);$$

check with the form of  $\mathcal{A}_n$ . Since  $\neg \text{St}_n(\ulcorner \phi \urcorner), \neg \text{St}_n(\ulcorner \psi \urcorner), s \neq \ulcorner \phi \wedge \psi \urcorner$  are all false  $\mathcal{L}_{\text{PA}}$ -literals, we have  $\text{AH}_n \mid_{\rho}^{\delta'_0} \Gamma, P_n^{\eta_0-1}(\ulcorner \phi \urcorner) \wedge P_n^{\eta_0-1}(\ulcorner \psi \urcorner)$ . Then, by  $\wedge$ -inversion,

$$\text{AH}_n \mid_{\rho}^{\delta'_0} \Gamma, P_n^{\eta_0-1}(\ulcorner \phi \urcorner) \quad \text{and} \quad \text{AH}_n \mid_{\rho}^{\delta'_0} \Gamma, P_n^{\eta_0-1}(\ulcorner \psi \urcorner). \quad (4.5)$$

It also follows in the same way from (d) that

$$\text{AH}_n \mid_{\rho}^{\delta'_1} \Gamma, P_n^{\eta_1-1}(\ulcorner \neg \phi \urcorner) \vee P_n^{\eta_1-1}(\ulcorner \neg \psi \urcorner). \quad (4.6)$$

Now, we will show by side induction on  $\delta'_1$  that, under the (main) induction hypothesis up to  $< \delta_0 \# \delta_1$ , (4.5) and (4.6) (for  $\delta'_0 < \delta_1$  and  $\delta'_1 < \delta_1$ ) implies that  $\mid_{\rho}^{\delta'_0 \# \delta'_1} \Gamma$ . The base step in which (4.6) is obtained by axiom is trivial. For the induction steps, the essential case is that  $P_n^{\eta_1-1}(\ulcorner \neg \phi \urcorner) \vee P_n^{\eta_1-1}(\ulcorner \neg \psi \urcorner)$  is critical in the last inference; otherwise the claim follows from SIH.

In this case, the premise is

$$\text{AH}_n \mid_{\rho}^{\delta''_1} \Gamma, P_n^{\eta_1-1} \ulcorner \neg \phi \urcorner \quad \text{or} \quad \text{AH}_n \mid_{\rho}^{\delta''_1} \Gamma, P_n^{\eta_1-1} \ulcorner \neg \psi \urcorner,$$

where  $\delta''_1 < \delta'_1$ . Then our claim follows from IH and (4.5).

Second, assume  $s = \ulcorner \sigma \urcorner = \ulcorner T_m r \urcorner$  for  $b, r \in \text{CT}$  with  $m < n$ . By  $\wedge$ -inversion,  $\forall$ -inversion and

$\vee$ -exportation, it follows from (c) that

$$\text{AH}_n \Big|_{\rho}^{\delta'_0} \Gamma, \neg\text{CT}(b), \neg\text{CT}(r), s \neq T_m r, P_m(r^\circ).$$

Thus we obtain  $\text{AH}_n \Big|_{\rho}^{\delta'_0} \Gamma, P_m(r^\circ)$ , since  $\neg\text{CT}(b)$ ,  $\neg\text{CT}(r)$  and  $s \neq T_m r$  are all false  $\mathcal{L}_0$ -literals. It also follows in the same way from (d) that  $\text{AH}_n \Big|_{\rho}^{\delta'_1} \Gamma, \neg P_m(r^\circ)$ . Finally, since  $\rho > 0$ , we obtain by cut that  $\text{AH}_n \Big|_{\rho}^{\delta_0 \# \delta_1} \Gamma$ . The cases where  $\sigma$  is of another form can be shown similarly; but note that the case where  $s = T_n r$  needs neither the side-induction step nor cut with rank 0 and the claim directly follows from IH.

Next, assume  $\sigma$  is of the form  $\neg\sigma'$ ; then  $\neg s = \ulcorner \neg\neg\sigma' \urcorner$ . In the same manner as above, we obtain from (d) that  $\text{AH}_n \Big|_{\rho}^{\delta'_1} P_n^{\eta_1-1}(\ulcorner \sigma' \urcorner), \Gamma$ . By IH, the claim  $\text{AH}_n \Big|_{\rho}^{\delta_0 \# \delta_1} \Gamma$  follows from this and (c); the proof is completed.  $\square$

**Lemma 4.6.3.** (1) Let  $\rho > 0$ . If  $\text{AH}_n(\mathcal{A}) \Big|_{\rho}^{\delta} P_n^{\eta_0} s \wedge P_n^{\eta_1} \neg s, \Gamma$ , then  $\text{AH}_n(\mathcal{A}) \Big|_{\rho}^{\delta \# \delta} \Gamma$ .

(2) Let  $\rho > 0$ . If  $\text{AH}_n(\mathcal{A}) \Big|_{\rho}^{\delta} \neg\text{Cons}_n, \Gamma$ , then  $\text{AH}_n(\mathcal{A}) \Big|_{\rho}^{\delta \# \delta} \Gamma$ .

Now we have so far obtained the following:

$$\text{KF} + \text{Cons} \vdash Tt \quad \Rightarrow \quad \text{AH}_1(\mathcal{A}) \Big|_{1}^{\leq \alpha_1} \neg\text{Cons}[0, \eta], P_1^\eta t \quad \Rightarrow \quad \text{AH}_1(\mathcal{A}) \Big|_{1}^{\leq \alpha_1} P_1^\eta t,$$

for some  $\eta < \varepsilon_0$ . For the reflective closures, we now have:

$$\text{KF}^* + \text{Cons} \vdash Tt \quad \Rightarrow \quad \text{AH}_n(\mathcal{A}) \Big|_{1}^{\leq \alpha_1} P_1 t \quad \Rightarrow \quad \text{H}_{n-1}(\mathcal{A}) \Big|_{\star}^{\leq \alpha_1} P_1 t \quad \Rightarrow \quad \text{AH}_{n-1}(\mathcal{A}) \Big|_{1}^{\leq \alpha_2} P_1 t.$$

By repeating this procedure  $n$ -times, we obtain:

$$\text{KF}^* + \text{Cons} \vdash Tt \quad \Rightarrow \quad \text{AH}_1 \Big|_{1}^{\leq \alpha_{n+1}} P_1^\eta t,$$

for some  $\eta < \alpha_n (< \Gamma_0)$ .

Hence, it suffices for our purpose to show that

(I) KF proves that  $\text{AH}_{n-1} \upharpoonright_1^{\delta} P_1^{\eta} t$  implies  $Tt$  for  $\delta < \alpha_1$  and  $\eta < \alpha_0 (= \varepsilon_0)$ .

(II) KF\* proves that  $\text{AH}_1 \upharpoonright_1^{\delta} P_1^{\eta} t$  implies  $Tt$  for  $\delta < \alpha_{n+1}$  and  $\eta < \alpha_n$ ;

these are shown by suitably modeling the infinitary derivations in question within KF and KF\*;

Recall that given a derivation of  $\text{KF} + \text{Cons} \vdash Tt$  or  $(\text{KF} + \text{Cons})^* \vdash Tt$ , we can primitive recursively compute the  $\delta$  and  $\eta$  below the above bounds.

The next is well-known.

**Theorem 4.6.4.**  $\text{KF} \vdash \forall a \in tc\text{TI}(\beta; Ta(x))$  for any  $\beta < \alpha_1 (= \varphi_{\varepsilon_0} 0)$ , and  $\text{KF}^* \vdash \forall a \in tc\text{TI}(\beta; Ta(x))$  for any  $\beta < \Gamma_0$ .

By diagonalization, we define a predicate  $H_{\gamma}(x)$  (with the parameter  $\gamma$  and  $x$ ) represents the

truth of  $\gamma$ -th stage in its Kripke construction. For  $\gamma = 0$ , set  $H_0 = \emptyset$ ; then,

$$\begin{aligned}
H_{\gamma+1}(a) : &\Leftrightarrow \exists x, y \in \text{CT}[a = x=y \wedge x^\circ = y^\circ] \vee \exists x, y \in \text{CT}[a = x \neq y \wedge x^\circ \neq y^\circ] \\
&\vee \exists b, c \in \text{St}_{\mathcal{L}_T}[a = b \wedge c \wedge (T^\Gamma H_\gamma(\dot{b})^\neg \wedge T^\Gamma H_\gamma(\dot{c})^\neg)] \\
&\vee \exists b, c \in \text{St}_{\mathcal{L}_T}[a = \neg(b \wedge c) \wedge (T^\Gamma H_\gamma(\neg \dot{b})^\neg \vee T^\Gamma H_\gamma(\neg \dot{c})^\neg)] \\
&\vee \exists b, c \in \text{St}_{\mathcal{L}_T}[a = b \rightarrow c \wedge (T^\Gamma H_\gamma(\neg \dot{b})^\neg \vee T^\Gamma H_\gamma(\dot{c})^\neg)] \\
&\vee \exists b, c \in \text{St}_{\mathcal{L}_T}[a = \neg(b \rightarrow c) \wedge (T^\Gamma H_\gamma(\dot{b})^\neg \wedge T^\Gamma H_\gamma(\neg \dot{c})^\neg)] \\
&\vee \exists b \in \text{St}_{\mathcal{L}_T}[a = \neg \neg b \wedge T^\Gamma H_\gamma(\dot{b})^\neg] \\
&\vee \exists x \in \text{CT}[a = T x \wedge T^\Gamma H_\gamma(x^\circ)^\neg] \\
&\vee \exists x \in \text{CT}[a = \neg T x \wedge T^\Gamma H_\gamma(\neg x^\circ)^\neg] \\
&\vee \exists b \exists c [a \in \text{St}_{\mathcal{L}_T} \wedge a = \forall c. b \wedge \forall y T^\Gamma H_\gamma(\dot{b}(\dot{y})^\neg)] \\
&\vee \exists b \exists c [a \in \text{St}_{\mathcal{L}_T} \wedge a = \neg \forall c. b \wedge \exists y T^\Gamma H_\gamma(\neg \dot{b}(\dot{y})^\neg)] \\
&\vee a \in H_\gamma;
\end{aligned}$$

for a limit  $\lambda$ ,  $H_\lambda(a)$  iff  $a \in \bigcup_{\xi < \lambda} H_\xi$ ; more precisely,  $(\exists n) T^\Gamma H_{\lambda[n]}(\dot{a})^\neg$ , where  $\lambda[n]$  is the fundamental sequence of  $\lambda$ .

**Lemma 4.6.5.** We can show by transfinite induction up to  $\varepsilon_0$  for  $\mathcal{L}_T$  that, for each given  $\beta < \varepsilon_0$ ,

$$\begin{aligned}
\text{KF} \vdash \forall \gamma < \bar{\beta} \forall x [(T^\Gamma H_\gamma(\dot{x})^\neg \rightarrow T x) \wedge (H_\gamma(x) \rightarrow T x)] \\
\text{KF} \vdash \forall \gamma < \bar{\beta} \forall x [(T^\Gamma H_\gamma(\dot{x})^\neg \rightarrow D^+ x) \wedge (H_\gamma(x) \rightarrow D^+ x)] \\
\text{KF} \vdash \forall \gamma < \bar{\beta} \forall x [\ulcorner H_\gamma(\dot{x})^\neg \in D^+ \wedge \ulcorner T^\Gamma H_\gamma(\dot{x})^\neg \in D^+ \urcorner]: \text{ thus, } \ulcorner H_\gamma(x)^\neg \in tc \text{ for all } \gamma < \bar{\beta} \\
\text{KF} \vdash \forall \gamma < \bar{\beta} \forall x [T^\Gamma H_\gamma(\dot{x})^\neg \leftrightarrow H_\gamma(x)].
\end{aligned}$$

Given  $\alpha < \varepsilon_0$ , we define a translation  $\mathcal{H}_\alpha$  from  $\mathcal{L}_1^\Delta$  into  $\mathcal{L}_T$  by  $T_\beta x \mapsto H_\beta x$  for  $\beta < \alpha$  and

$T_\gamma x \rightarrow \ulcorner 0 = 0 \urcorner$  for  $\gamma \not\prec \alpha$ ; all the other vocabulary and logical symbols are preserved. Let  $h_\alpha$  be a (primitive recursive) representation of  $\mathcal{H}_\alpha$  (in, say, PA). For a technical reason, we stipulate that  $h_\alpha(x) = \ulcorner 0 = 0 \urcorner$  for all  $x \notin \mathcal{L}_1^A$ .

**Corollary 4.6.6.** Let  $\alpha < \varepsilon_0$ . The following holds within KF: for any  $\mathcal{L}_1^A$ -sentence  $\sigma$ ,  $\mathcal{H}_\alpha(\sigma)$  is in  $D^+$ . More precisely, we have

$$\text{KF} \vdash (\forall x \in \text{St}_{\mathcal{L}_1^A})(h_{\bar{\alpha}}(\ulcorner \sigma \urcorner) \in D^+);$$

therefore, we have  $h_{\bar{\alpha}}(x) \in D^+$  for all  $x$  and thus  $\ulcorner Th_{\bar{\alpha}}(x) \urcorner \in tc$ .

*Proof.* This is shown by induction on  $rk_1(\sigma)$  using the last lemma.  $\square$

For  $\alpha < \varepsilon_0$ , let  $\mathcal{L}_1^\alpha := \mathcal{L}_0 \cup \{P_1^\eta \mid \eta < \alpha\} (\subset \mathcal{L}_1^A)$ . Suppose  $\text{AH}_1 \upharpoonright_{\mathcal{L}_1^\alpha} \Gamma$  for  $\Gamma \subset \mathcal{L}_1^\alpha$ . Then, we can assume by Lemma 4.5.11 that all sequents appearing in its derivation are from  $\mathcal{L}_1^\alpha$ . Let  $\text{ah}(\delta, C)$  is a standard representation (in PA) of  $\text{AH}_1 \upharpoonright_{\mathcal{L}_1^\alpha} \Gamma$ , where  $C$  is the code of the sequent  $\Gamma$ ; let  $C \subset \text{St}_{\mathcal{L}_1^A}$  express that  $C$  is a code of a sequent from  $\mathcal{L}_1^A$ . We set  $\Phi_\alpha(\delta)$  to be the following  $\mathcal{L}_T$ -formula:

$$\forall C \subset \text{St}_{\mathcal{L}_1^A} [\text{ah}_1(\delta, C) \rightarrow T(h_{\bar{\alpha}}(\bigwedge C))]$$

By the last corollary, we have  $\ulcorner \Phi_\alpha(\delta) \urcorner \in tc$ ; we can also show that  $\text{KF} \vdash \forall \delta [\ulcorner T \ulcorner \Phi_\alpha(\delta) \urcorner \urcorner \leftrightarrow \ulcorner \Phi_\alpha(\delta) \urcorner]$ .

Then, by applying Lemma 4.6.4 to the case where  $a = \ulcorner \Phi_\alpha \urcorner$ , we obtain the next lemma.

**Lemma 4.6.7.** Let  $\alpha < \varepsilon_0$  and  $\beta < \varphi_{\varepsilon_0} 0$ . Then  $\text{KF} \vdash \forall \delta < \bar{\beta} \Phi(\delta)$ : i.e.,

$$\text{KF} \vdash \forall \delta < \bar{\beta} \forall C \subset \text{St}_{\mathcal{L}_1^A} [\text{ah}_1(\delta, C) \rightarrow T(h_{\bar{\alpha}}(\bigwedge C))].$$

Finally, Lemmata 4.6.5 and 4.6.7 yield the above (I). We can similarly obtain (II) as well; in

fact, we need not use the translation like  $\mathcal{H}_\alpha$  and can directly show that, for each given  $\beta < \Gamma_0$ ,

$$\text{KF}^* \vdash \forall \delta < \bar{\beta} \forall C \subset \text{St}_{\mathcal{L}_1^A} [\text{ah}_1(\delta, C) \rightarrow T(\bigwedge C)];$$

this yields (II).

All the arguments so far applies to the cases of the other iterative compositional truths with no serious modifications. Hence, we finally obtain:

**Theorem 4.6.8.**  $\text{KF}^{(*)} + \text{Cons} \leq_{\mathcal{I}} \text{KF}^{(*)}$ ,  $\text{FKF}^{(*)} + \text{Cons} \leq_{\mathcal{I}} \text{FKF}^{(*)}$ , and  $\text{WKF}^{(*)} + \text{Cons} \leq_{\mathcal{I}} \text{WKF}^{(*)}$ .

## 4.6.2 Comparison of different self-applicable truths

The main theorem of the present subsection is:

**Theorem 4.6.9.**  $\text{WKF}^{(*)} \leq_{\mathcal{I}} \text{FKF}^{(*)}$  and  $\text{FKF}^{(*)} \leq_{\mathcal{I}} \text{KF}^{(*)}$ .

This theorem follows from the next two lemmata by similar arguments to the last subsection.

**Lemma 4.6.10.** Let  $\Gamma$  be a finite set of  $\mathcal{L}_1^A$ -sentences in which every occurrence of  $P_1^\xi$  (for some  $\xi$ ) is positive. Suppose  $\text{AH}_1(\mathcal{C})|_{\mathbb{1}}^{\delta} \Gamma, P_1^\eta(t)$ , for some ordinal  $\eta$  and a (closed) term  $t$ . Then, we have  $\text{AH}_n(\mathcal{B})|_{\mathbb{1}}^{\delta} \Gamma, P_1^\eta(t)$ .

**Lemma 4.6.11.** Let  $\Gamma$  be a finite set of  $\mathcal{L}_1^A$ -sentences in which every occurrence of  $P_1^\xi$  (for some  $\xi$ ) is positive. Suppose  $\text{AH}_1(\mathcal{B})|_{\mathbb{1}}^{\delta} \Gamma, P_1^\eta(t)$ , for some ordinal  $\eta$  and a (closed) term  $t$ . Then, we have  $\text{AH}_n(\mathcal{A})|_{\mathbb{1}}^{\delta} \Gamma, P_1^\eta(t)$ .

Both are proved by straightforward induction on  $\delta$  by using the fact that if  $\text{AH}_1(\mathcal{C})|_{\mathbb{1}}^{\delta} \Delta$  (or  $\text{AH}_1^{\mathcal{B}}|_{\mathbb{1}}^{\delta}$ ) for a  $\Delta$  in which every occurrence of  $P_1^\xi$  (for some  $\xi$ ) is positive, then we can assume that the derivation contains no occurrence of  $\neg P_1^{\xi'}$  (for any  $\xi'$ ) and thus no negative fixed-point rule is used therein; therefore, the proofs reduces to showing that, for each  $\Gamma$  satisfying the conditions of

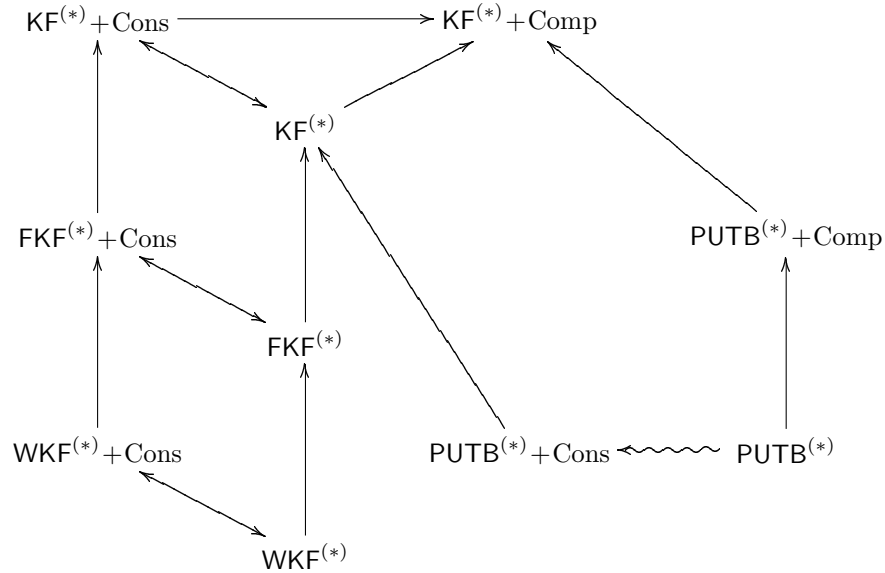
the above lemmata,

$$\begin{aligned} \text{AH}_1(\mathcal{C}) \Big|_1^{\delta} \Gamma, \mathcal{C}_1(P_1^\beta, t) &\Rightarrow \text{AH}_1(\mathcal{B}) \Big|_1^{\delta} \Gamma, \mathcal{B}_1(P_1^\beta, t) \\ \text{AH}_1(\mathcal{C}) \Big|_1^{\delta} \Gamma, \mathcal{B}_1(P_1^\beta, t) &\Rightarrow \text{AH}_1(\mathcal{B}) \Big|_1^{\delta} \Gamma, \mathcal{A}_1(P_1^\beta, t); \end{aligned}$$

the detail is somehow tedious but straightforward, since, roughly speaking,  $\mathcal{C}_1$  implies  $\mathcal{B}_1$  and  $\mathcal{B}_1$  implies  $\mathcal{A}_1$  in PA without using no fixed-point axiom.

## 4.7 Summary of Chapter 4

The results obtained in §§5–6 is summarized in the next diagram:



This diagram is to be read in the same way as the last diagram except that the curly arrow means that the converse is not yet known to hold.

Unfortunately, as the diagram indicates, two open problems are still left to us: i.e., it isn't yet known whether  $\text{PUTB} + \text{Cons} \leq_{\mathcal{I}} \text{PUTB}$  and whether  $(\text{PUTB} + \text{Cons})^* \leq_{\mathcal{I}} \text{PUTB}^*$ ; I conjecture

that both indeed hold but don't know their proofs.

One obvious direction one can take for further studies on inner theory is to investigate the inner theories of other systems of self-applicable truth: e.g.,  $\mathbf{VF}$  or the systems of Friedman and Sheard [23] (their analyses are already given by Cantini [6], Halbach [31] and Leigh and Rathjen [52]). However, it is particularly when the target systems all have the same truth-free consequences that inner theory comparison of systems are of the most significance, since the principal objective of my proposing inner theory comparison consists in the aforementioned motivation of differentiating systems of truth with distinct conceptions by means of a finer-grained criterion than the mere comparison by their truth-free consequences.

# Chapter 5

## Truth Progressions

### 5.1 Background

The present chapter studies several systems of the transfinite iteration and autonomous progression of self-applicable truth and determines their proof-theoretic strength.

To my knowledge, Jäger et al. [43] first introduced and studied a transfinitely iterated self-applicable truth. They presented the system of transfinitely iterated Kripke-Feferman truth and gave its proof-theoretic analysis. Then, Strahm [78] determined the proof-theoretic strength of the autonomous progression of Kripke-Feferman truth. the present chapter extends these studies to systems of other kinds of self-applicable truths.

Besides Kripke-Feferman truth, we can consider the autonomous progressions and transfinite iterations of various other conceptions and axiomatizations of self-applicable truth which have already been introduced so far. In the present chapter, we formulate the iterations and autonomous progressions of them and give their proof-theoretic analyses. Specifically, we study those of: Kripke-Feferman truth based on weak Kleene schema axiomatized by *WKF*; a purely disquotational truth represented by *PUTB*; Feferman's determinate truth represented by *DT*; their variants with Fried-

man and Sheard's consistency or completeness axiom; and iterative but non-compositional self-applicable truth of Cantini's system VF. In particular, the system  $\text{Aut}(\text{VF})$  of autonomously iterated VF truth has a fairly high proof-theoretic strength compared to the other systems of truth presented so far; indeed, as far as I know, it is the strongest system of truth so far.

The results of the present chapter, particularly on autonomous progressions, bear a close connection to the foundational questions in mathematics. In the rest of the present introductory section, we explain this point.

The notion of autonomous progression of theories was proposed by Kreisel and Feferman. It arose from the following questions: What is implicit in accepting a certain mathematical system  $S$  and what ought we to accept, on the same fundamental grounds, once we have made a commitment to  $S$ ? This, in turn, also tells us what we *can* accept or be justified to accept, once we become justified in accepting a starting system.

One accepts a mathematical system  $S$  only when she believes that  $S$  possesses a certain favorable property for her interests. For example, if one considers  $S$  inconsistent, she doesn't accept  $S$ , and, in the usual setting, consistency of  $S$  is implicitly accepted in her accepting the system  $S$  itself. For another example, one might argue that by accepting  $S$  we mean to accept that  $S$  is sound, and thus a certain form of reflection principle for  $S$  would be implicit in our acceptance of  $S$ . Gödel's theorems tell us that neither consistency nor reflection principle of  $S$  is derivable within  $S$  itself. Hence, according to these views of 'acceptance of system', one implicitly accepts something beyond  $S$  in accepting  $S$  itself, and, as Feferman wrote, '[Gödel's theorems] point to the possibility of systematically generating larger and larger systems whose acceptability is implicit in acceptance of the starting theory ([17, p.1])'. Autonomous progression is one of the most natural formulations of such a systematic generation for any given view of 'acceptance of system'.

In contrast to this consideration on what we *ought* to accept, the notion of autonomous progression also suggests an answer to the question of what we *can* accept, on the same fundamental

grounds as the starting system  $S$ , once we become justified in accepting  $S$ . Feferman's approach to predicativism is a typical example of this idea: once a system  $S$  is accepted, all the sets definable in  $S$  can be justifiably posited and quantification over them are allowed; then we can autonomously iterate this procedure and achieve the system of predicative analysis.

Now, let us focus on one specific but very natural view of 'acceptance of system': namely, the view that by accepting a system  $S$  we mean that we accept  $S$  as *true*. Then, what does it mean to accept a system as true? We consider the following interpretation: once we accept a system  $S_0$ , we implicitly accept it as true; we add a new predicate  $T$  to  $S_0$  for expressing its truth (a new predicate is needed due to Tarski's undefinability theorem); since we have accepted  $S_0$  as true, the axioms expressing ' $T$  represents the truth predicate for the language of  $S_0$ ' and ' $S_0$  is true (in terms of the truth predicate  $T$ )' ought to and can be accepted. In this way, we obtain and accept a new system  $S_1$ . Then, we can iterate this procedure of generating (and accepting)  $S_0, S_1, \dots, S_n, \dots$  even for transfinitely many steps (by taking the union at a limit level) along autonomously obtained well-orderings.

This progression depends on how truth is axiomatized. It is a highly profound and difficult question: By what axioms and rules the mathematical truth is to be characterized? The best known proposal is the Tarskian conception of truth, and the ordinal of its autonomous progression is known to be the Feferman-Schütte ordinal  $\Gamma_0$ ; see [17] or [24]. However, the Tarskian truth is typed and has been criticized by philosophers and logicians as not capturing the essential self-applicable character of the notion of truth, and various approaches have been attempted for axiomatizing self-applicable truth. It is of course not easy at all to pick out the 'genuine' axiomatization of self-applicable truth. the present chapter will examine the autonomous progressions (and iterations) of some possible options. We expect that our results, together with the cited works of Jäger et al. and Strahm, shed more light on the study of self-applicable truth.

## 5.2 Notation for Transfinitely Iterated Truths

For the subsequent arguments, we need slightly different notations from what has been used so far. First of all, for technical simplicity, I will exclude ‘ $\rightarrow$ ’ from the primitive logical connective and assume that they are only ‘ $\neg$ ’ and ‘ $\wedge$ ’ throughout the present chapter.

Let  $\mathcal{L}_{\text{PA}}$  denote the language of first order arithmetic with symbols for all primitive recursive (p.r.) functions; we assume that Peano Arithmetic PA contains all the axioms defining them. Since truth systems lack a second-order parameter, we add a new unary predicate  $U$ , which acts like a free set variable, so as to deal with autonomy rules and proof-theoretic ordinals. Our base language  $\mathcal{L}_0$  is  $\mathcal{L}_{\text{PA}} \cup \{U\}$ . We use ‘ $\equiv$ ’ for (meta-) identity between syntactical objects such as formulae (instead of ‘ $=$ ’).

Given any language  $\mathcal{L}' (\supset \mathcal{L}_{\text{PA}})$ , *full-induction* for  $\mathcal{L}'$  is the schema:

$$\phi(0) \wedge \forall x(\phi(x) \rightarrow \phi(x+1)) \rightarrow \forall x\phi(x), \text{ for each } \mathcal{L}'\text{-formula } \phi.$$

Our base system is Peano Arithmetic over  $\mathcal{L}_0$  and we denote it by  $\text{PA}(U)$ : i.e., PA plus full-induction for  $\mathcal{L}_0$ .<sup>1</sup> Given a formula  $\phi(x)$  of any language  $\mathcal{L}'$  and a p.r. ordering  $\prec$ , *transfinite induction* along  $\prec$  with respect to  $\phi(x)$ , written  $\text{TI}(\prec; \phi)$ , is defined by

$$\forall z[(\forall y \prec z)\phi(y) \rightarrow \phi(z)] \rightarrow \forall x \in \text{fd}(\prec)\phi(x),$$

where  $\text{fd}(\prec)$  denotes the field of the relation  $\prec$ . Then, we define the *proof-theoretic ordinal*  $|\mathbb{S}|$  of

---

<sup>1</sup>We adopt PA in our framework as the base system. One may choose another setting. For example, Cantini [7] and Kahle [46] rather take applicative theories with natural numbers as their base systems; a certain transfinite iteration of self-applicable truth in such a framework is studied in [46].

a system  $S$  over any language including  $\mathcal{L}_0$  as

$$\sup\{otyp(\prec) \mid \prec \text{ is a p.r. ordering and } S \vdash \text{TI}(\prec; U)\},$$

where  $otyp(\prec)$  denotes the order-type of  $\prec$ . For example,  $|\text{PA}(U)| = \varepsilon_0$ . Given any language  $\mathcal{L}'$ ,  $\text{TI}_{\mathcal{L}'}(\prec)$  denotes the schema  $\text{TI}_{\mathcal{L}'}(\prec; \phi)$  for all  $\phi \in \mathcal{L}'$ . In all the systems we consider in this chapter,  $\text{TI}(\prec; U)$  entails  $\text{TI}_{\mathcal{L}_0}(\prec)$ .

Given a p.r. ordering  $\prec$ ,  $\mathcal{L}^\prec$  denotes the language  $\mathcal{L}_0 \cup \{T^\prec\}$ , where  $T^\prec$  is a new *binary* predicate symbol. We write  $T_a^\prec(b)$  for  $T^\prec(b, a)$ . For  $a \in \text{fd}(\prec)$ ,  $T_a^\prec$  is meant to express the  $a$ -times iterated truth (or ‘true at the  $a$ -th level’) along  $\prec$ . Given  $a \in \text{fd}(\prec)$ ,  $\mathcal{L}_a^\prec$  denotes the sublanguage of  $\mathcal{L}^\prec$  obtained by restricting atoms of the form  $T_s^\prec t$  to closed terms  $s$  with  $s^{\mathbb{N}} \preceq a^{\mathbb{N}}$  ( $t$  can be an arbitrary term). The language  $\mathcal{L}$  for systems of autonomous progression is defined as  $\bigcup\{\mathcal{L}^\prec \mid \prec \text{ is a p.r. ordering}\}$ .

In order to formulate axiomatic systems of truth, we presuppose a standard Gödelization of  $\mathcal{L}$ . Each syntactical expression  $e$  from the vocabulary of  $\mathcal{L}$  is assigned a fixed Gödel number denoted by  $\#e$ . Let  $\ulcorner e \urcorner$  denote the numeral whose value is  $\#e$ .

First, we need the representations of some basic syntactical relations.  $\text{Tm}(x)$ ,  $\text{CT}(x)$  and  $\text{Var}(x)$  are representations of the sets of all terms, closed terms and variables respectively.  $\text{AtFml}_0(x)$ ,  $\text{Fml}_0(x)$  ( $\text{St}_0(x)$ ) and  $\text{For}_0(x)$  respectively represent the sets of all atomic formulae (atomic sentences), formulae (sentences) and formulae with at most one free variable of  $\mathcal{L}_0$ ; for each p.r. ordering  $\prec$ ,  $\text{AtFml}^\prec$  ( $\text{AtSt}^\prec$ ),  $\text{Fml}^\prec$  ( $\text{St}^\prec$ ) and  $\text{For}^\prec$  represent those of  $\mathcal{L}^\prec$ ; given  $a \in \text{fd}(\prec)$ ,  $\text{AtFml}_a^\prec$  ( $\text{AtSt}_a^\prec$ ),  $\text{Fml}_a^\prec$  ( $\text{St}_a^\prec$ ) and  $\text{For}_a^\prec$  are, respectively, the representations of those of  $\mathcal{L}_a^\prec$ .

Next, we introduce the representations of the following syntactical operations:  $\text{nm}(x)$  (‘the  $x$ -th numeral’);  $\text{val}(x)$  (‘the value of the closed term  $x$ ’);  $R\vec{x}$  (‘the atomic formula obtained by applying terms represented by  $\vec{x}$  to  $R$ ’) for each predicate  $R$ ;  $\neg x$  (‘the negation of the formula  $x$ ’);

$x \wedge y$  ('the conjunction of the formulae  $x$  and  $y$ ');  $\vee$  for disjunction defined analogously;  $\forall z.x$  ('the universal quantification of the formula  $x$  w.r.t. the variable  $z$ ');  $\exists$  for existential quantification defined analogously. We will usually write  $\dot{x}$  and  $x^\circ$  respectively for  $\text{nm}(x)$  and  $\text{val}(x)$  for saving space. For a code  $a$  of a term or formula and  $c_1, \dots, c_n \in \text{Tm}$ ,  $\text{sb}(a; c_1, \dots, c_n)$  denotes (the code of) the term or formula obtained from  $a$  by substituting the term  $c_i$  for the  $i$ -th free variable (in some fixed enumeration of variables) in  $a$  for each  $i \leq n$ ; e.g.,  $\text{sb}(\ulcorner \phi(x_0, x_1) \urcorner; \ulcorner s \urcorner, \ulcorner t \urcorner) = \ulcorner \phi(s, t) \urcorner$  for terms  $s$  and  $t$ . Then, we set  $a(\vec{x}) := \text{sb}(a; \vec{x})$ , and, given a term or formula  $e$ , we write  $\ulcorner e(\vec{x}) \urcorner$  for  $\ulcorner e \urcorner(\vec{x}) (= \text{sb}(\ulcorner e \urcorner; \vec{x}))$ ; we also write  $\ulcorner e(\vec{a}) \urcorner$  for  $\text{sb}(\ulcorner e(\vec{v}) \urcorner; \vec{a})$ .<sup>2</sup> When we need not explicitly specify a variable, we often omit it; e.g.,  $\ulcorner T_a \dot{b}(\cdot) \urcorner$  means  $\ulcorner T_a \dot{b}(w) \urcorner$  for a fresh variable  $w$ . We define the *falsity predicate*  $F_a \prec x$  (at the  $a$ -th level) as  $T_a \prec (\neg x)$ ; then  $F$  has its obvious meaning.

All the above syntactical relations and operations are primitive recursive except for  $\text{val}(x)$ ,  $\text{AtFml}_a \prec$ ,  $\text{AtSt}_a \prec$ ,  $\text{Fml}_a \prec$ ,  $\text{St}_a \prec$  and  $\text{For}_a \prec$ . Hence, we assume that they are all represented by their corresponding function symbols in  $\mathcal{L}_0$ , since  $\mathcal{L}_0$  possesses symbols for all p.r. functions. The value function  $\text{val}(x)$  for  $\mathcal{L}_0$ -terms is recursive and definable in PA but not p.r., since we could otherwise construct a p.r. evaluation function of all the  $n$ -ary p.r. functions (for each  $n$ ) by means of  $\text{val}(x)$ ;  $\text{AtFml}_a \prec$ ,  $\text{AtSt}_a \prec$ ,  $\text{Fml}_a \prec$ ,  $\text{St}_a \prec$  and  $\text{For}_a \prec$  are also recursive and PA-definable but not p.r., since they need  $\text{val}(x)$  in their standard definitions. Hence,  $\mathcal{L}_0$  does not contain primitive symbols for these six notions. Let  $\text{PA}^+(U)$  be the definitional expansion of  $\text{PA}(U)$  augmented by the symbols and defining axioms for them and let  $\mathcal{L}_0^+$  be its language. We can either work within this definitional expansion or simply regard them as abbreviations of suitable  $\mathcal{L}_0$ -formulae representing them.

**Remark 8.** There are some subtle issues to be noted concerning coding particularly when we work within definitional expansions. First, our coding was fixed before taking definitional expansions. Hence, definitionally introduced symbols are not assigned their own codings. Of course, we can

---

<sup>2</sup>To avoid possible confusion, we notice that  $\text{sb}(a; \vec{c})$  is denoted by  $a(\vec{c})$  in [43] and  $a(\vec{c})$  in our notation corresponds to their  $a(\vec{c})$ , but  $\ulcorner \phi(\vec{x}) \urcorner$  has the same meaning for a formula  $\phi$ .

redefine coding so as to accommodate the new symbols, but then  $\text{val}(x)$  is no longer the value function in this new coding schema. Thus, even when we work within  $\mathcal{L}_0^+$  or any further definitional expansion  $\mathcal{L}'$ , the codes of  $\mathcal{L}'$ -expressions should be regarded as the codes of corresponding  $\mathcal{L}_0$ -expressions with the defining formulae of the new symbols. Second, truth systems sometimes behave in unexpected ways on definitionally introduced functions. For example, let  $f$  be a definable function and  $\phi(x, y)$  be its defining formula. Consider a sentence  $T_s^{\prec} \ulcorner T_s^{\prec} f(\dot{x}) \urcorner$  for a closed term  $s$  with  $s^{\mathbb{N}} \in \text{fd}(\prec)$ . Then, one might expect that it is equivalent to  $T_s(f(x))$  due to self-applicability of truth. However,  $T_s^{\prec} \ulcorner T_s^{\prec} f(\dot{x}) \urcorner$  indeed means, say,  $T_s^{\prec} \ulcorner \forall y [\phi(\dot{x}, y) \rightarrow T_s^{\prec}(y)] \urcorner$ . Then, whether they are equivalent depends on the inner logic (or the compositional axioms) of the system in question. As a matter of fact, this equivalence fails in  $\text{WKF}_{\prec}$  (+Cons); nonetheless, as we shall see in §6, this is no obstacle to well-ordering proofs and has no essential effect on the proof-theoretic strength. In contrast, they are equivalent and definable functions can generally be treated as if they were primitive symbols within all the other systems in the present chapter; this fact is implicitly made use of particularly in Lemmata 5.4.2, 5.5.2 and 5.7.4.

### 5.3 Kripke-Feferman Truth

In the present section, we introduce the transfinite iteration and autonomous progression of Kripke-Feferman truth and review the results of Jäger et al. [43] and Strahm [78] on them.

**Definition 5.3.1.** Given a p.r. ordering  $\prec$ , the system  $\text{KF}_{\prec}$  over  $\mathcal{L}^{\prec}$  consists of  $\text{PA}(U)$  plus full-induction for  $\mathcal{L}^{\prec}$ ,  $\text{TI}_{\mathcal{L}^{\prec}}(\prec)$  and the following axioms:

$$\mathbf{K1}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall \vec{x} \in \text{CT} \left[ \left( T_a^{\prec}(R\vec{x}) \leftrightarrow R\vec{x}^{\circ} \right) \wedge \left( F_a^{\prec}(R\vec{x}) \leftrightarrow \neg R\vec{x}^{\circ} \right) \right],$$

for each predicate symbol  $R$  of  $\mathcal{L}_0$

$$\mathbf{K2}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, b \in \text{CT}$$

$$\left[ b^{\circ} \prec a \rightarrow \left( T_a^{\prec}(T_b^{\prec} x) \leftrightarrow T_{b^{\circ}}^{\prec}(x^{\circ}) \right) \wedge \left( F_a^{\prec}(T_b^{\prec} x) \leftrightarrow \neg T_{b^{\circ}}^{\prec}(x^{\circ}) \right) \right]$$

$$\mathbf{K3}_{\prec} \quad \forall x, a \in \text{CT} [a^\circ \in \text{fd}(\prec) \rightarrow (T_a^\prec(T_a^\prec x) \leftrightarrow T_a^\circ x^\circ) \wedge (F_a^\prec(T_a^\prec x) \leftrightarrow F_a^\circ x^\circ)]$$

$$\mathbf{K4}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x \in \text{St}_a^\prec [T_a^\prec(\neg \neg x) \leftrightarrow T_a^\prec x]$$

$$\mathbf{K5}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, y \in \text{St}_a^\prec [(T_a^\prec(x \wedge y) \leftrightarrow T_a^\prec x \wedge T_a^\prec y) \wedge (F_a^\prec(x \wedge y) \leftrightarrow F_a^\prec x \vee F_a^\prec y)]$$

$$\mathbf{K6}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall z, x$$

$$[\forall z. x \in \text{St}_a^\prec \rightarrow (T_a^\prec(\forall z. x) \leftrightarrow \forall y T_a^\prec x(y)) \wedge (F_a^\prec(\forall z. x) \leftrightarrow \exists y F_a^\prec x(y))]$$

The system  $\text{Aut}(\text{KF})$  over  $\mathcal{L}$  consists of  $\text{PA}(U)$ , full-induction for  $\mathcal{L}$ , and the following two rules:  
for each p.r. linear ordering (provably in, say,  $\text{PA}(U)$ ),

$$\frac{\text{TI}(\prec; U)}{\mathbf{K1}_{\prec} \wedge \dots \wedge \mathbf{K6}_{\prec}} \quad \text{and} \quad \frac{\text{TI}(\prec; U)}{\text{TI}(\prec; \psi)} \quad \text{for each } \mathcal{L}\text{-formula } \psi.$$

The right-hand rule corresponds to the Bar rule, BR. Indeed,  $\text{KF}_{\prec}$  is inconsistent for some orderings  $\prec$ , but this autonomy rule blocks inconsistency by restricting  $\prec$  to well-orderings.

Ordinals up to the Feferman-Schütte ordinal  $\Gamma_0$  are not sufficient for ordinal analyses of the systems we will consider. Recall that  $\Gamma_0$  is the least ordinal ( $> 0$ ) closed under the binary Veblen function  $\varphi$ . We introduce the *ternary Veblen function*  $\varphi$  by a simple generalization of the binary one:  $\varphi 0 \alpha \beta = \varphi_\alpha \beta$ ;  $\varphi \alpha 0 \gamma$  for  $\alpha > 0$  is the  $\gamma$ -th strongly critical ordinal with respect to  $\lambda \xi, \eta. \varphi \alpha' \xi \eta$  for  $\alpha' < \alpha$ ;  $\varphi \alpha \beta \gamma$  for  $\alpha, \beta > 0$  is the  $\gamma$ -th common fixed-point of  $\lambda \xi. \varphi \alpha \beta' \xi$  for  $\beta' < \beta$ . Let  $\Lambda_3$  be the least ordinal  $> 0$  closed under the ternary Veblen function. Throughout §§3-7, we fix a natural ordinal notation system up to a sufficiently large ordinal, say,  $\Lambda_3$ , and let  $\triangleleft$  denote the corresponding well-ordering. Given a closed term  $t \in \text{fd}(\triangleleft)$ ,  $|t|$  denotes the unique ordinal associated to  $t$  by this notation system. For  $\alpha \in \text{fd}(\triangleleft)$ , let  $\triangleleft \upharpoonright_\alpha$  be  $\{\langle \gamma, \beta \rangle \mid \gamma \triangleleft \beta \triangleleft \alpha\}$ . We write  $\text{KF}_\alpha$  for  $\text{KF}_{\triangleleft \upharpoonright_\alpha}$  for  $\alpha \in \text{fd}(\triangleleft)$  (and similar for other systems). For a language  $\mathcal{L}'$ , we write  $\text{TI}_{\mathcal{L}'}(\alpha)$  for  $\text{TI}_{\mathcal{L}'}(\triangleleft \upharpoonright_\alpha)$ .

Jäger et al. [43] determined the proof-theoretic ordinal of the system  $\widehat{\text{ID}}_\alpha$  of  $\alpha$ -times iterated fixed-points. Let  $\mathcal{L}_0^{(+)}(R, Q)$  be  $\mathcal{L}_0^{(+)} \cup \{R, Q\}$  where  $R$  and  $Q$  are fresh unary predicate symbols.

An  $\mathcal{L}_0^{(+)}$ ( $R, Q$ )-formula  $\mathcal{A}(R, Q, x, y)$  with at most  $x$  and  $y$  free is called an *inductive operator form*, when  $R$  occurs only positively in  $\mathcal{A}$ . The language  $\mathcal{L}_{\text{fix}}$  of  $\widehat{\text{ID}}_{\prec}$  for a given  $\prec$  is defined as  $\mathcal{L}_0$  plus *unary* predicates  $P^{\mathcal{A}}$  associated to each inductive operator form  $\mathcal{A}(R, Q, x, y) \in \mathcal{L}_0(R, Q)$ . We write  $P_x^{\mathcal{A}}(y)$  for  $P^{\mathcal{A}}(\langle y, x \rangle)$  and  $P_{\prec x}^{\mathcal{A}}(z)$  for  $(z)_1 \prec x \wedge P^{\mathcal{A}}(z)$ , where  $\langle \cdot, \cdot \rangle$  is a p.r. paring function. We assume that the paring function is bijective from  $\mathbb{N} \times \mathbb{N}$  to  $\mathbb{N}$  and we let  $(\cdot)_0$  and  $(\cdot)_1$  be the corresponding projections; thus, in contrast to [43], we do not need the clause  $z \in \text{Pair}$  (' $z$  is a code of a pair') in the definition of  $P_{\prec x}^{\mathcal{A}}(z)$ . Then, the system  $\widehat{\text{ID}}_{\prec}$  comprises  $\text{PA}(U)$ , full-induction for  $\mathcal{L}_{\text{fix}}$ ,  $\text{TI}_{\mathcal{L}_{\text{fix}}}(\prec)$ , and the fixed-point axiom schema: for each inductive operator form  $\mathcal{A}$ ,

$$\forall a \in \text{fd}(\prec) \forall x [P_a^{\mathcal{A}}(x) \leftrightarrow \mathcal{A}(P_a^{\mathcal{A}}, P_{\prec a}^{\mathcal{A}}, x, a)].$$

In giving the lower bound of  $|\widehat{\text{ID}}_{\alpha}| (= |\widehat{\text{ID}}_{\triangleleft_{\uparrow \alpha}}|)$ , Jäger et al. used an intermediate system  $\text{SRT}_{\alpha}$ , which is embeddable in  $\text{ID}_{\alpha}$ . This  $\text{SRT}_{\alpha}$  is essentially the same system as  $\text{KF}_{\alpha}$ .<sup>3</sup> It is observed from their proof that  $\text{KF}_{\alpha}$  and  $\widehat{\text{ID}}_{\alpha}$  have the same proof-theoretic ordinal:

**Theorem 5.3.2** (Jäger et al. [43]). Given an ordinal  $\delta$ ,  $\varepsilon(\delta)$  denotes the least epsilon number greater than  $\delta$ , and we define  $(\delta|m)$  inductively by  $(\delta|0) = \varepsilon(\delta)$  and  $(\delta|k+1) = \varphi_{(\delta|k)}0$ . Let an ordinal  $\alpha$  be such that

$$\alpha = \omega^{1+\alpha_n} + \omega^{1+\alpha_{n-1}} + \dots + \omega^{1+\alpha_1} + m \quad (\alpha_n \geq \dots \geq \alpha_1). \quad (5.1)$$

Then,  $|\widehat{\text{ID}}_{\alpha}| = |\text{KF}_{\alpha}| = \varphi 1 \alpha_n (\varphi 1 \alpha_{n-1} (\dots \varphi 1 \alpha_1 (\alpha|m) \dots))$ .

Strahm [78] introduced the system  $\text{Aut}(\widehat{\text{ID}})$  for autonomous fixed-point progression, and showed that its proof-theoretic ordinal is  $\varphi 200$ . His proof can be directly applied to  $\text{Aut}(\text{KF})$  and thus yields the following theorem:

---

<sup>3</sup>The acronym  $\text{SRT}$  stands for 'self-reflective truth'. Because all the truths we consider in the present chapter are self-reflective, we rename it in order to contrast it against other conceptions of truth.

**Theorem 5.3.3** (Strahm [78]).  $|\text{Aut}(\text{KF})| = |\text{Aut}(\widehat{\text{ID}})| = \varphi 200$ .

**Remark 9.** Feferman’s original motivation for developing the notion of *reflective closure*, which is what truth theorists now call KF, might be incompatible with its (autonomous) iteration. Feferman’s principal motivation behind it was to suggest a ‘more realistic and perspicuous’ (finite) alternative to autonomous progression for approximating or finding out what is implicit in accepting a mathematical system. In the present chapter, we regard KF not as means of reflective closure but as one axiomatization of (self-applicable) truth.

## 5.4 Consistency and Completeness Axioms

Friedman and Sheard [23] introduced two special axioms for self-applicable truth: the consistency axiom  $\text{Cons}$  (expressing ‘no sentence is both true and false’) and completeness axiom  $\text{Comp}$  (‘every sentence is either true or false’). Cantini [5] then considered adding them to KF and studied the resulting systems.<sup>4</sup> Although neither of them is derivable in KF and adding both to KF leads to a contradiction, either (but not both) of them can be consistently and even conservatively (for  $\mathcal{L}_0$ ) added to KF.

In this section, we study consistency and completeness axioms for iterated truth: for a given p.r. ordering  $\prec$ ,

$$\text{Cons}_{\prec} \equiv \forall a \in \text{fd}(\prec) \forall x \in \text{St}_a^{\prec} \neg(T_a^{\prec}x \wedge F_a^{\prec}x)$$

$$\text{Comp}_{\prec} \equiv \forall a \in \text{fd}(\prec) \forall x \in \text{St}_a^{\prec} (T_a^{\prec}x \vee F_a^{\prec}x).$$

As in the last section, for  $\alpha \in \text{fd}(\triangleleft)$ , we write  $\text{Cons}_{\alpha}$  and  $\text{Comp}_{\alpha}$  for  $\text{Cons}_{\triangleleft|_{\alpha}}$  and  $\text{Comp}_{\triangleleft|_{\alpha}}$  respectively. Throughout this section, we write  $\gamma < \beta$  for  $\gamma \triangleleft \beta$  for readability and drop the

---

<sup>4</sup>Feferman [17] also considered the consistency axiom which is denoted by *Disj* there.

superscript  $\triangleleft \upharpoonright_\alpha$  whenever there is no worry of confusion. We shall see that  $\text{Cons}_\alpha$  or  $\text{Comp}_\alpha$  can also be conservatively added to  $\text{KF}_\alpha$ .

**Definition 5.4.1.**  $\text{Aut}(\text{KF} + \text{Cons})$  and  $\text{Aut}(\text{KF} + \text{Comp})$  over  $\mathcal{L}$  are defined as  $\text{Aut}(\text{KF})$  plus

$$\frac{\text{TI}(\prec; U)}{\text{Cons}_\prec} \quad \text{or} \quad \frac{\text{TI}(\prec; U)}{\text{Comp}_\prec},$$

for each p.r. linear ordering  $\prec$ , respectively.

The goal of the present section is to determine the proof-theoretic ordinals of  $\text{KF}_\alpha + \text{Cons}_\alpha$ ,  $\text{KF}_\alpha + \text{Comp}_\alpha$ ,  $\text{Aut}(\text{KF} + \text{Cons})$ , and  $\text{Aut}(\text{KF} + \text{Comp})$ . Cantini [5] showed that  $\text{KF} + \text{Cons}$  and  $\text{KF} + \text{Comp}$  are syntactically embeddable in each other; by *syntactical embedding*, we mean relative interpretation which keeps the vocabulary of the base language unchanged. As the next lemma shows, the same applies to  $\text{KF}_\prec$  and thus it suffices for our goal to show  $|\text{KF}_\alpha + \text{Cons}_\alpha| \leq |\text{KF}_\alpha|$  via a suitable cut-elimination procedure.

**Lemma 5.4.2.**  $\text{KF}_\prec + \text{Cons}_\prec$  and  $\text{KF}_\prec + \text{Comp}_\prec$  are mutually syntactically embeddable.

*Proof.* We define in  $\text{PA}(U)$  (by the recursion theorem or diagonalization) a function  $f$  by:

$$f(x) := \begin{cases} x & \text{if } x \in \text{AtFml}_0 \\ \ulcorner \neg F_z^\prec(fy) \urcorner \text{ (i.e., } \text{sb}(\ulcorner \neg F(\cdot, \cdot) \urcorner; y, z)) & \text{if } x = T_z^\prec y \text{ for } y, z \in \text{Tm} \\ \neg f(y) \text{ (} f(y) \wedge f(z), \text{ resp.)} & \text{if } x = \neg y \text{ (or } x = y \wedge z) \\ \forall z. f(y) & \text{if } x = \forall z. y; \end{cases}$$

otherwise,  $f(x) = 0$ . The embedding  $\mathcal{F}$  for both directions is given by  $T_z^\prec x \mapsto \neg F_z^\prec f(x)$ . We only demonstrate that  $\mathbf{K2}_\prec$  is preserved by  $\mathcal{F}$ : for  $b, x \in \text{CT}$  with  $b^\circ \prec a \in \text{fd}(\prec)$ :

$$\mathcal{F}(T_a^\prec T_b^\prec x) \equiv \neg F_a^\prec \ulcorner \neg F_b^\prec(fx) \urcorner \leftrightarrow \neg T_a^\prec \ulcorner F_b^\prec(fx) \urcorner \leftrightarrow \neg F_b^\prec f(x^\circ) \equiv \mathcal{F}(T_b^\prec x^\circ)$$

$$\mathcal{F}(F_a^\prec T_b^\prec x) \equiv \neg F_a^\prec \ulcorner \neg F_b^\prec(fx) \urcorner \leftrightarrow \neg F_a^\prec \ulcorner F_b^\prec(fx) \urcorner \leftrightarrow \neg \neg F_b^\prec f(x^\circ) \equiv \mathcal{F}(\neg T_b^\prec x^\circ).$$

The other cases are straightforward.  $\square$

We assume familiarity with the cut-elimination procedure of  $\widehat{\text{ID}}_\alpha$  in [43] and use the same notations as there. For technical reasons, we explicitly work with the definitional expansion  $\text{PA}^+(U)$ ; thus, in particular,  $\widehat{\text{ID}}_\alpha$  (over  $\text{PA}^+(U)$ ) possesses  $\text{val}(x)$  and  $\text{St}_z^\prec$  as primitive symbols. In general,  $\text{KF}_\prec$  can be embedded in  $\widehat{\text{ID}}_\prec$  for any  $\prec$ ; there exists an inductive operator form  $\mathcal{A}$  such that if  $\text{KF}_\alpha (= \text{SRT}_\alpha) \vdash \phi$  then  $\widehat{\text{ID}}_\alpha \vdash \phi[P^\mathcal{A}/T]$ , where  $\phi[P^\mathcal{A}/T]$  is the result of substituting  $P_s^\mathcal{A}t$  for  $T_s t$  in  $\phi$  for each terms  $s$  and  $t$ . We can assume that  $\widehat{\text{ID}}_\alpha \vdash \phi[P^\mathcal{A}/T]$  above is obtained by a derivation in which only the vocabulary from  $\mathcal{L}_0^+ \cup \{P^\mathcal{A}\}$  is used and no other atomic of the form  $P^\mathcal{B}$  ( $\mathcal{B} \neq \mathcal{A}$ ) occurs. For the subsequent arguments, we precisely describe one such  $\mathcal{A}$  and fix it. We take the following  $\mathcal{A}(R, Q, x, a)$ :

$$\begin{aligned}
& x \in \text{St}_a \wedge \forall y, z \in \text{CT}[x = y \doteq z \rightarrow y^\circ = z^\circ] \wedge \forall y, z \in \text{CT}[x = y \neq z \rightarrow y^\circ \neq z^\circ] \\
& \wedge \forall y \in \text{CT}[(x = U y) \rightarrow U(y^\circ)] \wedge \forall y \in \text{CT}[(x = \neg U y) \rightarrow \neg U(y^\circ)] \\
& \wedge \forall b, y \in \text{CT}[(b^\circ < a \wedge x = T_b y) \rightarrow \langle y^\circ, b^\circ \rangle \in Q] \\
& \wedge \forall b, y \in \text{CT}[(b^\circ < a \wedge x = \neg T_b y) \rightarrow \langle y^\circ, b^\circ \rangle \notin Q] \\
& \wedge \forall b, y \in \text{CT}[(b^\circ = a \wedge x = T_b y) \rightarrow y^\circ \in R] \\
& \wedge \forall b, y \in \text{CT}[(b^\circ = a \wedge x = \neg T_b y) \rightarrow \neg y^\circ \in R] \\
& \wedge \forall y \in \text{St}_a[(x = \neg \neg y) \rightarrow y \in R] \wedge \forall y, z \in \text{St}_a[(x = y \wedge z) \rightarrow R y \wedge R z] \\
& \wedge \forall y, z \in \text{St}_a[(x = \neg(y \wedge z)) \rightarrow (\neg y \in R \vee \neg z \in R)] \\
& \wedge \forall y, z[\text{St}_a(\forall y.z) \wedge x = \forall y.z \rightarrow \forall w(z(w) \in R)] \\
& \wedge \forall y, z[\text{St}_a(\forall y.z) \wedge x = \neg \forall y.z \rightarrow \exists w(\neg z(w) \in R)].
\end{aligned}$$

We use the same semi-formal system  $\text{H}_\alpha$  for  $\widehat{\text{ID}}_\alpha$  as in [43, p.62-3] and follow the same notations; but we need to make slight modifications to accommodate the new symbols in  $\mathcal{L}_0^+$  and we postulate

all the true  $\mathcal{L}_0^+$ -literals (e.g.,  $\text{val}(\ulcorner 1 + 1 \urcorner) = 2$ ) as axioms; cf. the definition of  $\text{AH}_{\alpha'}$  below.<sup>5</sup> Since we focus only on the above specific  $\mathcal{A}$ , we can assume that the language  $\mathcal{L}_{\text{fix}}^\alpha$  (written as  $\mathcal{L}_\alpha$  in [43]) of  $\text{H}_\alpha$  extends  $\mathcal{L}_0^+$  only by unary symbols  $P_\beta^{\mathcal{A}}$  ( $\beta < \alpha$ ) and  $P_{<\gamma}^{\mathcal{A}}$  ( $\gamma \leq \alpha$ ). Recall that  $\text{H}_\alpha \upharpoonright_{\star}^{\delta} \Gamma$  (or  $\text{H}_\alpha \upharpoonright_{\star}^{\leq \delta} \Gamma$ ) means that there is a derivation of  $\Gamma$  in  $\text{H}_\alpha$  with the length  $\delta$  ( $< \delta$ , resp.) in which all its cut formulae are of the forms  $P_\beta^{\mathcal{A}}t$  and  $\neg P_\beta^{\mathcal{A}}$  ( $\beta < \alpha$ ). As was shown in [43] (Proposition 15), if  $\widehat{\text{ID}}_\alpha \vdash A$  for an  $\mathcal{L}_{\text{fix}}$ -sentence  $A$  then  $\text{H}_\alpha \upharpoonright_{\star}^{\leq \varepsilon(\alpha)} A^\alpha$ , where the translation  $A^\alpha$  is obtained from  $A$  by replacing  $P^{\mathcal{A}}$  by  $P_{<\alpha}^{\mathcal{A}}$ . In particular, if  $\text{KF}_\alpha + \text{Cons}_\alpha \vdash \sigma$  for some  $\mathcal{L}_0^+$ -sentence  $\sigma$ , then we have  $\widehat{\text{ID}}_\alpha \vdash \text{Cons}_\alpha[P^{\mathcal{A}}/T] \rightarrow \sigma$  and thus  $\text{H}_\alpha \upharpoonright_{\star}^{\leq \varepsilon(\alpha)} (\neg \text{Cons}_\alpha[P^{\mathcal{A}}/T])^\alpha, \sigma$ . In what follows, we just identify  $\text{Cons}_\alpha$  with  $(\text{Cons}_\alpha[P^{\mathcal{A}}/T])^\alpha$  and omit the superscript  $\mathcal{A}$  for simplicity.

**Lemma 5.4.3.** Let  $\gamma < \beta \leq \alpha$  and  $s$  and  $t$  be closed terms with  $(t)_0^{\mathbb{N}} = s^{\mathbb{N}}$  and  $|(t)_1| = \gamma$ . If  $\text{H}_\alpha \upharpoonright_{\star}^{\delta} \phi, \Gamma$  with arbitrary cut-rank, then  $\text{H}_\alpha \upharpoonright_{\star}^{1+\delta} \phi[P_\gamma(s)/P_{<\beta}(t)], \Gamma$  with the same cut-rank.

*Proof.* Induction on  $\delta$ . There are two essential cases. One is the base step where  $\phi$  is obtained by the Axiom I ([43, p.62]) and  $\phi \equiv P_{<\beta}(t)$  or  $\phi \equiv \neg P_{<\beta}(t)$ . We illustrate the former case; then,  $\Gamma = \Gamma' \cup \{\neg P_{<\beta}(r)\}$  for some  $\Gamma'$  and  $r^{\mathbb{N}} = t^{\mathbb{N}}$ . The same axiom yields  $\text{H}_\alpha \upharpoonright_{\star}^0 \Gamma', P_\gamma(s), \neg P_\gamma((r)_0)$ . Then we obtain  $\text{H}_\alpha \upharpoonright_{\star}^1 \Gamma', P_\gamma(s), \neg P_{<\beta}(r)$  by the Fixed-point rule III ([43, p.63]). The other essential case is where  $\phi$  is obtained by the Fixed-point rule III. In this case, the premise is numerically equivalent to the claimed sequent itself.  $\square$

**Lemma 5.4.4.** Let  $\gamma \leq \beta \leq \alpha$  and  $t$  be a closed term with  $|(t)_1| < \gamma$ . Suppose  $\text{H}_\alpha \upharpoonright_{\star}^{\delta} \Gamma, \phi$ . Then,  $\text{H}_\alpha \upharpoonright_{\star}^{2+\delta} \Gamma, \phi[P_{<\gamma}(t)/P_{<\beta}(t)]$  (with the same cut-rank).

*Proof.* Induction on  $\delta$ . The essential cases are the same. For the first case, we apply the Fixed-point rule III twice to a suitable instance of the Axiom I. For the second, we obtain  $(\neg)P_{<\gamma}(t)$  instead of  $(\neg)P_{<\beta}(t)$  from the premise by the same rule.  $\square$

<sup>5</sup>This change results in a slight strengthening of the ordinary  $\omega$ -arithmetic, but, even with this change, all the arguments in [43, §6] can be carried out in the same way and cut-elimination of  $\text{H}_\alpha$  properly gives the upper bound of  $\widehat{\text{ID}}_\alpha$ ; cf. the proof of Theorem 1.3.6 of [64].

We make use of the technique of asymmetric interpretation of  $H_\alpha$ . Given a successor  $\alpha' = \alpha + 1$ , we introduce a Tait-style semi-formal system  $AH_{\alpha'}$ . The language  $\mathcal{L}_A^{\alpha'}$  ( $\supset \mathcal{L}_{\text{fix}}^\alpha$ ) of  $AH_{\alpha'}$  consists of the sentences (in Tait-style) of  $\mathcal{L}_0^+ \cup \{P_\xi \mid \xi < \alpha\} \cup \{P_{<\xi} \mid \xi \leq \alpha\} \cup \{P_\alpha^\eta \mid \eta\}$ , where  $\eta$  ranges over ordinals up to a sufficiently large ordinal, e.g.,  $\Lambda_3$ ; namely, the predicate  $P_\alpha$  of the highest level is further typed by ordinals. Note that, in contrast to  $\mathcal{L}_{\text{fix}}^{\alpha'}$ ,  $\mathcal{L}_A^{\alpha'}$  does not contain the predicate  $P_{<\alpha}$ . The axioms and rules of  $AH_{\alpha'}$  are given as follows.<sup>6</sup>

### Axioms

Ax0.  $\Gamma, A$ , for a true  $\mathcal{L}_{PA}^+$ -literal (i.e.,  $(\mathcal{L}_0^+ \setminus \{U\})$ -literal)  $A$ .

Ax1.  $\Gamma, \phi(t), \neg\phi(s)$ , for an atomic  $\phi$  of the form  $P_\xi(t)$  ( $\xi < \alpha$ ),  $P_{<\xi'}(t)$  ( $\xi' \leq \alpha$ ) or  $U(t)$  and closed terms  $s$  and  $t$  with  $s^{\mathbb{N}} = t^{\mathbb{N}}$ .

Ax2.  $\Gamma, \neg P_{<\xi}(t)$ , for an ordinal  $\xi \leq \alpha$  and a closed term  $t$  with  $|(t)_1| \not\leq \xi$ .

Ax3.  $\Gamma, \neg P_\alpha^0(t)$ , for any closed term  $t$ .

**Logical rules and Cut** the same as  $H_{\alpha'}$  (i.e., those of  $\omega$ -arithmetic).

**Fixed-point rules for  $P_\xi$  and  $P_{<\xi'}$  ( $\xi < \alpha$ ,  $\xi' \leq \alpha$ )** the same as  $H_{\alpha'}$  (i.e., the Fixed-point rules III and IV in [43, p.63] for  $\xi < \alpha$  and  $\xi' \leq \alpha$ )

**Fixed-point rules for  $P_\alpha^\eta$**  for closed terms  $s$  and  $t$  with  $|t| = \alpha$ ,

$$\frac{\Gamma, \mathcal{A}(P_\alpha^\eta, P_{<\alpha}, s, t)}{\Gamma, P_\alpha^{\eta+1}(s)} \quad \frac{\Gamma, \neg \mathcal{A}(P_\alpha^\eta, P_{<\alpha}, s, t)}{\Gamma, \neg P_\alpha^{\eta+1}(s)}$$

**Limit rules for  $P_\alpha^\lambda$**  for a limit  $\lambda$  and a closed term  $s$ ,

$$\frac{\Gamma, P_\alpha^\eta(s) \quad (\text{for some } \eta < \lambda)}{\Gamma, P_\alpha^\lambda(s)} \quad \frac{\Gamma, \neg P_\alpha^\eta(s) \quad (\text{for all } \eta < \lambda)}{\Gamma, \neg P_\alpha^\lambda(s)}$$

<sup>6</sup>Since we have assumed that  $\langle \cdot, \cdot \rangle$  is bijective, we can drop the conditions like  $t \in \text{Pair}$  from the original formulation in [43].

For each  $\mathcal{L}_A^{\alpha'}$ -formula  $\phi$ , we define its AH-rank,  $rk_A(\phi)$ , as follows: If  $\phi$  is a literal of the form  $P_\alpha^\eta t$  or  $\neg P_\alpha^\eta t$  then  $rk_A(\phi) = \omega \cdot \eta$ ; if  $\phi$  is a literal of any other form,  $rk_A(\phi) = 0$ ;  $rk_A(\phi)$  for a compositional formula  $\phi$  is inductively defined in the standard manner; e.g.,  $rk_A(\psi_0 \wedge \psi_1) := \max\{rk_A(\psi_0), rk_A(\psi_1)\} + 1$ . We write  $AH_{\alpha'} \upharpoonright_{\frac{\delta}{\rho}} \Gamma$  (or  $AH_{\alpha'} \upharpoonright_{< \frac{\delta}{\rho}} \Gamma$ ) when  $\Gamma$  is derived in  $AH_{\alpha'}$  with the length  $\delta$  ( $< \delta$  resp.) and cut-rank (w.r.t.  $rk_A$ )  $\rho$  ( $< \rho$ , resp.).

Suppose  $H_{\alpha'} \upharpoonright_{\star} \Gamma$  and  $\Gamma$  contains no  $P_{< \alpha'}$ ; then its derivation contains no occurrence of  $P_{< \alpha'}$ . Given arbitrary  $\eta$ , by a straightforward modification of Cantini's asymmetric interpretation in [4], we obtain  $AH_{\alpha'} \upharpoonright_{< \frac{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}} \Gamma[\eta, \eta + \omega^\delta]$ , where  $\Gamma[\xi_0, \xi_1]$  is obtained from  $\Gamma$  by replacing each occurrence of  $\neg P_\alpha^\xi$  and  $P_\alpha^\xi$  respectively by  $\neg P_\alpha^{\xi_0}$  and  $P_\alpha^{\xi_1}$ .<sup>7</sup> Let  $\beta$  be  $\max\{\varepsilon(\delta), \varepsilon(\eta)\}$ . Then, by predicative cut elimination, we can obtain  $AH_{\alpha'} \upharpoonright_{\frac{\leq \varphi_0 \beta 0}{1}} \Gamma[\eta, \eta + \omega^\delta]$ . We omit the details of this embedding of  $H_{\alpha'}$  in  $AH_{\alpha'}$  via asymmetric interpretation, but we will carry out a very similar embedding in Lemma 5.4.7 and the reader may refer to it. We state the so-called *Persistency Lemma* in  $AH_{\alpha'}$  for later use:

$$\text{If } \xi \leq \xi' \leq \eta' \leq \eta \text{ and } AH_{\alpha'} \upharpoonright_{\frac{\delta}{\rho}} \Gamma[\xi', \eta'], \Delta, \text{ then } AH_{\alpha'} \upharpoonright_{< \frac{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}{\max\{\varepsilon(\eta), \varepsilon(\rho)\}}} \Gamma[\xi, \eta], \Delta; \quad (5.2)$$

this property is indeed heavily used in the above embedding of  $H_{\alpha'}$  in  $AH_{\alpha'}$ .<sup>8</sup>

When  $\Gamma \subset \mathcal{L}_{\text{fix}}^\alpha$  (and thus  $\Gamma[\xi_0, \xi_1] \equiv \Gamma$  for any  $\xi_0$  and  $\xi_1$ ), we take  $\eta = 0$  and obtain  $AH_{\alpha'} \upharpoonright_{\frac{\leq \varphi_0 \varepsilon(\delta) 0}{1}} \Gamma$ . We can assume without loss of generality that its derivation contains no rules for  $P_\alpha^\xi$  and thus can be regarded as a derivation of  $H_{\alpha'} \upharpoonright_{< \frac{\varphi_0 \varepsilon(\delta) 0}{1}} \Gamma$  in which cut is only applied to literals.<sup>9</sup> Then, we eliminate all the cuts applied to the literals of the forms  $P_{< \beta}$  and  $\neg P_{< \beta}$  ( $\beta \leq \alpha$ ) or arithmetical literals without changing the upperbound of the length of derivation<sup>10</sup>; we have

<sup>7</sup>In fact, we can primitive recursively compute, from given  $\delta$  and  $\eta$ , the exact length and cut-rank of  $AH_{\alpha'} \upharpoonright_{< \frac{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}} \Gamma[\eta, \eta + \omega^\delta]$ : that is, we have  $\zeta_0(\eta, \delta), \zeta_1(\eta, \delta) < \max\{\varepsilon(\eta), \varepsilon(\delta)\}$  such that  $AH_{\alpha'} \upharpoonright_{\frac{\zeta_0(\eta, \delta)}{\zeta_1(\eta, \delta)}} \Gamma[\eta, \eta + \omega^\delta]$ . For example, we can take  $\zeta_0(\eta, \delta) = \zeta_1(\eta, \delta) = \omega \cdot (\eta + \omega^\delta)$ .

<sup>8</sup>The length and cut-rank in (5.2) are both primitive recursively computed as well.

<sup>9</sup>In general, if  $AH_{\alpha'} \upharpoonright_{\frac{\delta}{\rho}} P_\alpha^0 t, \Gamma$ , then  $AH_{\alpha'} \upharpoonright_{\frac{\delta}{\rho}} \Gamma$ ; for, an atomic of the form  $P_\alpha^0 t$  cannot be critical in any derivation.

<sup>10</sup>As a matter of fact, we can assume that the above derivation of  $AH_{\alpha'} \upharpoonright_{\frac{\leq \varphi_0 \varepsilon(\delta) 0}{1}} \Gamma$  contains no cut to literals of

thus obtained  $H_\alpha \mid_{\star}^{\leq \varphi 0 \varepsilon(\delta) 0} \Gamma$ . In summary,

$$\text{If } H_{\alpha'} \mid_{\star}^{\delta} \Gamma \text{ for } \Gamma \text{ without } P_{<\alpha'}, \text{ then } AH_{\alpha'} \mid_{\star}^{\leq \max\{\varepsilon(\delta), \varepsilon(\eta)\}} \Gamma[\eta, \eta + \omega^\delta] \quad (5.3)$$

$$\text{If } H_{\alpha'} \mid_{\star}^{\delta} \Gamma \text{ for } \Gamma \subset \mathcal{L}_{\text{fix}}^\alpha, \text{ then } H_\alpha \mid_{\star}^{\delta'} \Gamma \text{ for some } \delta' < \varphi 0 \varepsilon(\delta) 0. \quad (5.4)$$

In fact,  $\delta'$  in (5.4) can be primitive recursively computed from  $\delta$  (see fn.7). Let  $\mathcal{G}$  be a function such that  $\mathcal{G}(\delta)$  is such  $\delta'$ . We write  $\mathcal{G}_k(\delta)$  for  $\mathcal{G} \circ \dots \circ \mathcal{G}(\delta)$  ( $k$ -times) and stipulate that  $\mathcal{G}_0(\delta) = \delta$ . We can assume that  $\mathcal{G}(\delta') < \mathcal{G}(\delta)$  for  $\delta' < \delta$ ; note that  $\mathcal{G}_k(\delta) < (\delta|k)$ . By iterating the procedure (5.4) for finite times, we have

$$\text{If } H_{\alpha+m} \mid_{\star}^{\delta} \Gamma \text{ for } \Gamma \subset \mathcal{L}_{\text{fix}}^{\alpha+n} \text{ and } n < m, \text{ then } H_{\alpha+n} \mid_{\star}^{\mathcal{G}_{m-n}(\delta)} \Gamma.$$

The proof of Main Lemma II of [43] is essentially the transfinite repetition of the procedure (5.4). What we will do for our corresponding lemma (Lemma 5.4.10 below) essentially amounts to transfinitely iterating the procedure of Lemma 5.4.8 below instead of (5.4): i.e., the procedure of transforming  $H_{\alpha'} \mid_{\star}^{\delta} \neg \text{Cons}_{\alpha'}, \Gamma$  into  $H_\alpha \mid_{\star}^{\mathcal{G}'(\delta)} \neg \text{Cons}_\alpha, \Gamma$  ( $\mathcal{G}'$  will be given later).

**Lemma 5.4.5.** Let  $\rho > 0$ . For arbitrary  $\eta_0$  and  $\eta_1$ , if (a)  $AH_{\alpha'} \mid_{\rho}^{\delta_0} P_\alpha^{\eta_0} s, \Gamma$  and (b)  $AH_{\alpha'} \mid_{\rho}^{\delta_1} P_\alpha^{\eta_1} \neg s, \Gamma$ , then  $AH_{\alpha'} \mid_{\rho}^{\delta_0 \# \delta_1} \Gamma$ .

*Proof.* This lemma informally corresponds to the fact that the minimal Kripkean fixed-point is consistent; in  $AH_{\alpha'}$ ,  $P_\alpha^\eta$  corresponds to the  $\eta$ -th stage of the Kripkean construction (of the  $\alpha$ -th truth) from the empty set, since we start with  $\neg P_\alpha^0 t$  for all  $t$ .

We show the claim by induction on  $\delta_0 \# \delta_1$ . The base case is easy, since an atomic of the form  $P_\alpha^\eta t$  cannot be critical in any axiom. For the induction step, the cases where either  $P_\alpha^{\eta_0} s$  or  $P_\alpha^{\eta_1} \neg s$  is not critical in the last inference follow from IH. Let us assume  $P_\alpha^{\eta_0} s$  and  $P_\alpha^{\eta_1} \neg s$  are both critical

---

the forms  $P_{<\beta}$  and  $\neg P_{<\beta}$  ( $\beta \leq \alpha$ ) and thus this elimination procedure is in fact redundant here.

in the last inferences. When either of the last inferences is by the limit rule, the premises contain  $\text{AH}_{\alpha'} \mid_{\rho}^{\delta_2} P_{\alpha}^{\eta_2}(\neg)s, \Gamma$  for  $\delta_2 < \delta_0$  (or  $< \delta_1$ ) and  $\eta_2 < \eta_0$  (or  $< \eta_1$ ); then, the claim follows from IH. Hence, we focus on the essential case in which the last inferences are both the fixed-point rule. Then, we can assume that the premises are

$$(c) \text{AH}_{\alpha'} \mid_{\rho}^{\delta'_0} \Gamma, \mathcal{A}(P_{\alpha}^{\eta_0-1}, P_{<\alpha}, s, t) \quad \text{and} \quad (d) \text{AH}_{\alpha'} \mid_{\rho}^{\delta'_1} \Gamma, \mathcal{A}(P_{\alpha}^{\eta_1-1}, P_{<\alpha}, \neg s, t),$$

for some  $|t| = \alpha$ ,  $\delta'_0 < \delta_0$  and  $\delta'_1 < \delta_1$ . By the form of  $\mathcal{A}$  (and inversion), we can assume that  $s \in \text{St}_t$ ; let  $s = \ulcorner \sigma \urcorner$ . The claim is shown by cases according to the form of  $\sigma$ .

Suppose no negation ‘ $\neg$ ’ is attached to  $\sigma$ . We illustrate two cases. First, assume  $\sigma$  is  $\phi \wedge \psi$ . By  $\wedge$ -inversion,  $\forall$ -inversion and  $\forall$ -exportation, it follows from (c) that

$$\text{AH}_{\alpha'} \mid_{\rho}^{\delta'_0} \Gamma, \neg \text{St}_t(\ulcorner \phi \urcorner), \neg \text{St}_t(\ulcorner \psi \urcorner), s \neq \ulcorner \phi \wedge \psi \urcorner, P_{\alpha}^{\eta_0-1}(\ulcorner \phi \urcorner) \wedge P_{\alpha}^{\eta_0-1}(\ulcorner \psi \urcorner);$$

check with the form of  $\mathcal{A}$ . Since  $\neg \text{St}_t(\ulcorner \phi \urcorner), \neg \text{St}_t(\ulcorner \psi \urcorner), s \neq \ulcorner \phi \wedge \psi \urcorner$  are all false  $\mathcal{L}_{\text{PA}}^+$ -literals, we have  $\text{AH}_{\alpha'} \mid_{\rho}^{\delta'_0} \Gamma, P_{\alpha}^{\eta_0-1}(\ulcorner \phi \urcorner) \wedge P_{\alpha}^{\eta_0-1}(\ulcorner \psi \urcorner)$ . Then, by  $\wedge$ -inversion,

$$\text{AH}_{\alpha'} \mid_{\rho}^{\delta'_0} \Gamma, P_{\alpha}^{\eta_0-1}(\ulcorner \phi \urcorner) \quad \text{and} \quad \text{AH}_{\alpha'} \mid_{\rho}^{\delta'_0} \Gamma, P_{\alpha}^{\eta_0-1}(\ulcorner \psi \urcorner). \quad (5.5)$$

It also follows in the same way from (d) that

$$\text{AH}_{\alpha'} \mid_{\rho}^{\delta'_1} \Gamma, P_{\alpha}^{\eta_1-1}(\ulcorner \neg \phi \urcorner) \vee P_{\alpha}^{\eta_1-1}(\ulcorner \neg \psi \urcorner). \quad (5.6)$$

Now, we will show by side induction on  $\delta'_1$  that, under the (main) induction hypothesis up to  $< \delta_0 \# \delta_1$ , (5.5) and (5.6) (for  $\delta'_0 < \delta_1$  and  $\delta'_1 < \delta_1$ ) implies that  $\mid_{\rho}^{\delta'_0 \# \delta'_1} \Gamma$ . The base step in which (5.6) is obtained by axiom is trivial. For the induction steps, the essential case is that

$P_\alpha^{\eta_1-1}(\ulcorner\neg\phi\urcorner) \vee P_\alpha^{\eta_1-1}(\ulcorner\neg\psi\urcorner)$  is critical in the last inference; otherwise the claim follows from SIH.

In this case, the premise is

$$\text{AH}_{\alpha'} \left| \frac{\delta_1''}{\rho} \Gamma, P_\alpha^{\eta_1-1} \ulcorner\neg\phi\urcorner \right. \quad \text{or} \quad \text{AH}_{\alpha'} \left| \frac{\delta_1''}{\rho} \Gamma, P_\alpha^{\eta_1-1} \ulcorner\neg\psi\urcorner \right.,$$

where  $\delta_1'' < \delta_1'$ . Then our claim follows from IH and (5.5).

Second, assume  $s = \ulcorner\sigma\urcorner = T_b r$  for  $b, r \in \text{CT}$  with  $b^\circ < \alpha$ . By  $\wedge$ -inversion,  $\forall$ -inversion and  $\vee$ -exportation, it follows from (c) that

$$\text{AH}_{\alpha'} \left| \frac{\delta_0'}{\rho} \Gamma, \neg\text{CT}(b), \neg\text{CT}(r), b^\circ \not\prec t, s \neq T_b r, P_{<\alpha}(\langle r^\circ, b^\circ \rangle) \right.$$

Thus we obtain  $\text{AH}_{\alpha'} \left| \frac{\delta_0'}{\rho} \Gamma, P_{<\alpha}(\langle r^\circ, b^\circ \rangle) \right.$ , since  $\neg\text{CT}(b)$ ,  $\neg\text{CT}(r)$ ,  $b^\circ \not\prec t$  and  $s \neq T_b r$  are all false  $\mathcal{L}_0^+$ -literals. It also follows in the same way from (d) that  $\text{AH}_{\alpha'} \left| \frac{\delta_1'}{\rho} \Gamma, \neg P_{<\alpha}(\langle r^\circ, b^\circ \rangle) \right.$ . Finally, we obtain, by cut (with rank 0), that  $\text{AH}_{\alpha'} \left| \frac{\delta_0 \# \delta_1}{\rho} \Gamma \right.$ . The cases where  $\sigma$  is of another form can be shown similarly; but note that the case where  $s = T_b r$  for  $b^\circ = \alpha$  needs neither the side-induction step nor cut with rank 0 and the claim directly follows from IH.

Next, assume  $\sigma$  is of the form  $\neg\sigma'$ ; then  $\ulcorner\neg\sigma\urcorner = \ulcorner\neg\neg\sigma'\urcorner$ . In the same manner as above, we obtain from (d) that  $\text{AH}_{\alpha'} \left| \frac{\delta_1'}{\rho} P_\alpha^{\eta_1-1}(\ulcorner\neg\sigma\urcorner), \Gamma \right.$ . By IH, the claim  $\text{AH}_{\alpha'} \left| \frac{\delta_0 \# \delta_1}{\rho} \Gamma \right.$  follows from this and (c); the proof is completed.  $\square$

**Lemma 5.4.6.** (1) If  $\text{AH}_{\alpha'} \left| \frac{\delta}{\rho} P_\alpha^{\eta_0} s \wedge P_\alpha^{\eta_1} \ulcorner\neg s\urcorner, \Gamma \right.$ , for  $\rho > 0$ , then  $\text{AH}_{\alpha'} \left| \frac{\delta \cdot 2}{\rho} \Gamma \right.$ ; by  $\wedge$ -inversion and the last lemma.

(2) For any closed term  $t$ , if  $\text{AH}_{\alpha'} \left| \frac{\delta}{\rho} \exists x \in \text{St}_t[P_\alpha^{\eta_0}(x) \wedge P_\alpha^{\eta_1}(\ulcorner\neg x\urcorner)], \Gamma \right.$ , for  $\rho > 0$ , then  $\text{AH}_{\alpha'} \left| \frac{\delta \cdot 2}{\rho} \Gamma \right.$ ; by induction on  $\delta$ .

**Lemma 5.4.7.** Let  $\Gamma$  contain no  $P_{<\alpha'}$ . Suppose  $\text{H}_{\alpha'} \left| \frac{\delta}{*} \neg\text{Cons}_{\alpha'}, \Gamma \right.$ . Then, for each  $\eta$ , we have  $\text{AH}_{\alpha'} \left| \frac{\max\{\varepsilon(\delta), \varepsilon(\eta)\}}{\max\{\varepsilon(\delta), \varepsilon(\eta)\}} \neg\text{Cons}_\alpha, \Gamma[\eta, \eta + \omega^{1+\delta}] \right.$ ; in fact, the length and cut-rank can be computed to

be  $\omega \cdot (\eta + \omega^{1+\delta}) \cdot 2$  (cf. fn.7 and fn.8).

*Proof.* Suppose  $H_{\alpha'} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha'}, \Gamma$ . The claim is shown by induction on  $\delta$ . The proof is essentially parallel to that of the embedding (5.3) except that  $P_{<\alpha'}$  occurs in  $\neg \text{Cons}_{\alpha'}$ . Thus, each case where  $\Gamma$  is critical in the last inference is dealt with in a parallel manner; we illustrate some typical cases below and the rest are similarly shown. For simplicity, we write  $\|\xi_0, \xi_1\|$  for  $\max\{\varepsilon(\xi_0), \varepsilon(\xi_1)\}$ .

Suppose  $\Gamma$  is obtained from the premises  $H_{\alpha'} \upharpoonright_{\star}^{\delta_0} \neg \text{Cons}_{\alpha'}, \Gamma, P_{\alpha} t$  and  $H_{\alpha'} \upharpoonright_{\star}^{\delta_1} \neg \text{Cons}_{\alpha'}, \Gamma, \neg P_{\alpha} t$  for some  $\delta_0, \delta_1 < \delta$  by cut. By IH, we have

$$\begin{aligned} \text{AH}_{\alpha'} \upharpoonright_{\star} & \frac{< \|\delta_0, \eta\|}{< \|\delta_0, \eta\|} \neg \text{Cons}_{\alpha}, \Gamma[\eta, \eta + \omega^{1+\delta_0}], P_{\alpha}^{\eta + \omega^{1+\delta_0}} t \\ \text{AH}_{\alpha'} \upharpoonright_{\star} & \frac{< \|\delta_1, \eta + \omega^{1+\delta_0}\|}{< \|\delta_1, \eta + \omega^{1+\delta_0}\|} \neg \text{Cons}_{\alpha}, \Gamma[\eta + \omega^{1+\delta_0}, \eta + \omega^{1+\delta_0} + \omega^{1+\delta_1}], \neg P_{\alpha}^{\eta + \omega^{1+\delta_0}} t. \end{aligned}$$

Since  $\varepsilon(\eta + \omega^{1+\delta}), \|\delta_0, \eta\|, \|\delta_1, \eta + \omega^{1+\delta_0}\| \leq \|\delta, \eta\|$ , we obtain

$$\begin{aligned} \text{AH}_{\alpha'} \upharpoonright_{\star} & \frac{< \|\delta, \eta\|}{< \|\delta, \eta\|} \neg \text{Cons}_{\alpha}, \Gamma[\eta, \eta + \omega^{1+\delta}], P_{\alpha}^{\eta + \omega^{1+\delta}} t \\ \text{AH}_{\alpha'} \upharpoonright_{\star} & \frac{< \|\delta, \eta\|}{< \|\delta, \eta\|} \neg \text{Cons}_{\alpha}, \Gamma[\eta, \eta + \omega^{1+\delta}], \neg P_{\alpha}^{\eta + \omega^{1+\delta}} t, \end{aligned}$$

by persistency (5.2). Since  $\omega \cdot (\eta + \omega^{1+\delta_0}) < \|\delta, \eta\|$ , the claim follows by cut.

Suppose  $\Gamma$  is obtained from the premise  $H_{\alpha'} \upharpoonright_{\star}^{\delta'} \neg \text{Cons}_{\alpha'}, \Gamma', \mathcal{A}(P_{\alpha}, P_{<\alpha}, t, s)$ , for  $\delta' < \delta$ ,  $|s| = \alpha$  and  $\Gamma = \Gamma' \cup \{P_{\alpha} t\}$ , by Fixed-point rule. Then, by applying (5.2) to IH in the same manner as above, we obtain

$$\text{AH}_{\alpha'} \upharpoonright_{\star} \frac{< \|\delta, \eta\|}{< \|\delta, \eta\|} \neg \text{Cons}_{\alpha}, \Gamma[\eta, \eta + \omega^{1+\delta}], \mathcal{A}(P_{\alpha}^{\eta + \omega^{1+\delta}}, P_{<\alpha}, t, s).$$

By the same Fixed-point rule, we obtain

$$\text{AH}_{\alpha'} \left| \frac{\leq \|\delta, \eta\|}{< \|\delta, \eta\|} \right. \neg \text{Cons}_{\alpha}, \Gamma[\eta, \eta + \omega^{1+\delta}], P_{\alpha}^{\eta + \omega^{1+\delta'} + 1}(t);$$

the claim follows by (5.2), since  $\varepsilon(\eta + \omega^{1+\delta'} + 1) \leq \|\delta, \eta\|$ .

Now we turn to the essential case in which  $\neg \text{Cons}_{\alpha'}$  is critical in the last inference. Then the premise is: for some  $\delta' < \delta$  and closed term  $t$ ,

$$\text{H}_{\alpha'} \left| \frac{\delta'}{\star} \right. t < \alpha' \wedge \exists x \in \text{St}_t [P_{<\alpha'}(\langle x, t \rangle) \wedge P_{<\alpha'}(\langle \neg x, t \rangle)], \Gamma.$$

If  $|t| \not< \alpha'$ , the claim trivially follows. Let  $|t| < \alpha'$ . By  $\wedge$ -inversion, we have

$$\text{H}_{\alpha'} \left| \frac{\delta'}{\star} \right. (\exists x \in \text{St}_t) [P_{<\alpha'}(\langle x, t \rangle) \wedge P_{<\alpha'}(\langle \neg x, t \rangle)], \Gamma.$$

If  $|t| < \alpha$ , then it follows from Lemma 5.4.4 that

$$\text{H}_{\alpha'} \left| \frac{2 + \delta'}{\star} \right. (\exists x \in \text{St}_t) [P_{<\alpha}^{\eta}(\langle x, t \rangle) \wedge P_{<\alpha}^{\eta}(\langle \neg x, t \rangle)], \Gamma.$$

Hence, we have  $\text{AH}_{\alpha'} \left| \frac{\leq \|\delta, \eta\|}{< \|\delta, \eta\|} \right. \exists x \in \text{St}_t [P_{<\alpha}^{\eta}(\langle x, t \rangle) \wedge P_{<\alpha}^{\eta}(\langle \neg x, t \rangle)], \Gamma[\eta, \eta + \omega^{1+\delta}]$  by the procedure (5.3); we obtain the claim by Ax0 and logical rules. Finally, assume  $|t| = \alpha$ . By Lemma 5.4.3, we have  $\text{H}_{\alpha'} \left| \frac{1+\delta'}{\star} \right. (\exists x \in \text{St}_t) [P_{\alpha}(x) \wedge P_{\alpha}(\neg x)], \Gamma$ ; therefore, since  $1 + \delta' \leq \delta$ , we have

$$\text{AH}_{\alpha'} \left| \frac{\beta_0}{\beta_1} \right. (\exists x \in \text{St}_t) [P_{\alpha}^{\eta + \omega^{1+\delta}}(x) \wedge P_{\alpha}^{\eta + \omega^{1+\delta}}(\neg x)], \Gamma[\eta, \eta + \omega^{1+\delta}],$$

for some  $\beta_0, \beta_1 < \|\delta, \eta\|$ , by the procedure (5.3). By the last lemma, we obtain  $\text{AH}_{\alpha'} \left| \frac{\beta_0 \cdot 2}{\beta_1} \right. \Gamma[\eta, \eta + \omega^{1+\delta}]$ ; the claim follows from this since  $\beta_0 \cdot 2 < \|\delta, \eta\|$ .  $\square$

**Lemma 5.4.8.** Let  $\Gamma \subset \mathcal{L}_{\text{fix}}^\alpha$ . If  $H_{\alpha'} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha'}, \Gamma$ , then  $H_{\alpha} \upharpoonright_{\star}^{\leq \varphi 0 \varepsilon(\delta) 0} \neg \text{Cons}_{\alpha}, \Gamma$ .

*Proof.* Suppose  $H_{\alpha'} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha'}, \Gamma$ . Taking  $\eta = 0$  in the last lemma, we obtain  $AH_{\alpha'} \upharpoonright_{\star}^{\frac{\leq \varepsilon(\delta)}{\leq \varepsilon(\delta)}} \neg \text{Cons}_{\alpha}, \Gamma$ . By predicative cut elimination,  $AH_{\alpha'} \upharpoonright_{\star}^{\frac{\leq \varphi 0 \varepsilon(\delta) 0}{1}} \neg \text{Cons}_{\alpha}, \Gamma$ . Since this sequent is of  $\mathcal{L}_{\text{fix}}^\alpha$ , we can replace  $AH_{\alpha'} \upharpoonright_{\star}^{\frac{\leq \varphi 0 \varepsilon(\delta) 0}{1}}$  by  $H_{\alpha} \upharpoonright_{\star}^{\leq \varphi 0 \varepsilon(\delta) 0}$ .  $\square$

It is observed from this proof (and Lemma 5.4.7) that we can primitive recursively compute the length of  $H_{\alpha} \upharpoonright_{\star}^{\leq \varphi 0 \varepsilon(\delta) 0} \neg \text{Cons}_{\alpha}, \Gamma$ , from given  $\delta$ ; in fact, we can take  $\varphi 0(\omega^{1+\delta} \cdot 2)(\omega^{1+\delta} \cdot 2)$  as such for instance. Let  $\mathcal{G}'(\delta)$  is such a length: i.e.,  $\mathcal{G}'(\delta) < \varphi 0 \varepsilon(\delta) 0$  and  $H_{\alpha} \upharpoonright_{\star}^{\mathcal{G}'(\delta)} \neg \text{Cons}_{\alpha}, \Gamma$ , for  $\Gamma \subset \mathcal{L}_{\text{fix}}^\alpha$ , when  $H_{\alpha'} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha'}, \Gamma$ . For each  $k \in \mathbb{N}$ ,  $\mathcal{G}'_k(\delta)$  is defined in the analogous way to  $\mathcal{G}_k(\delta)$ ; thus  $\mathcal{G}'_k(\delta) < (\delta|k)$ .

**Corollary 5.4.9.** Let  $m, k \in \mathbb{N}$  and  $\Gamma \subset \mathcal{L}_{\text{fix}}^{\alpha+m}$ . If  $H_{\alpha+m+k} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha+m+k}, \Gamma$ , then we have  $H_{\alpha+m} \upharpoonright_{\star}^{\mathcal{G}'_k(\delta)} \neg \text{Cons}_{\alpha+m}, \Gamma$ .

**Lemma 5.4.10.** Suppose  $H_{\alpha+\omega^{1+\rho}} \upharpoonright_{\star}^{\delta} \neg \text{Cons}_{\alpha+\omega^{1+\rho}}, \Gamma$  for  $\Gamma \subset \mathcal{L}_{\text{fix}}^{\alpha+\xi}$  and  $\xi < \omega^{1+\rho}$ . Then we have  $H_{\alpha+\xi} \upharpoonright_{\star}^{\frac{\varphi 1 \rho(2+\delta)}{\star}} \neg \text{Cons}_{\alpha+\xi}, \Gamma$ .

*Proof.* This corresponds to Main Lemma II of [43]. We prove the claim by induction on  $\rho$  with side induction on  $\delta$ . The proof is divided into three cases:  $\rho = 0$ ;  $\rho$  is successor;  $\rho$  is limit. We only illustrate the successor case; the other cases can be shown in the same pattern; but note that we use the last corollary for the base case ( $\rho = 0$ ) instead of IH.

Let  $\rho = \rho' + 1$ . Since  $\xi < \omega^{1+\rho'+1}$ , there is a natural number  $n$  such that  $\xi = \omega^{1+\rho'} \cdot n + \eta$  for some  $\eta < \omega^{1+\rho'}$ . If  $\delta = 0$  then  $\Gamma$  is critical and the claim trivially obtains. Let us assume  $\delta > 0$ . If  $\Gamma$  is critical in the last inference and it is not cut, the claim follows from SIH. The crucial cases are where  $\neg \text{Cons}_{\alpha+\omega^{1+\rho}}$  is critical and where the last inference is cut.

Suppose the first case. By  $\wedge$ -inversion,  $H_{\alpha+\omega^{1+\rho}} \upharpoonright_{\star}^{\delta'} t < \alpha + \omega^{1+\rho}, \Gamma$  and

$$H_{\alpha+\omega^{1+\rho}} \upharpoonright_{\star}^{\delta'} \exists x \in \text{St}_t [P_{<\alpha+\omega^{1+\rho}}(\langle x, t \rangle) \wedge P_{<\alpha+\omega^{1+\rho}}(\langle \neg x, t \rangle)], \Gamma,$$

for some  $t$  and  $\delta' < \delta$ . The case where  $|t| \not\leq \alpha + \omega^{1+\rho}$  is trivially obtained and we assume  $|t| < \alpha + \omega^{1+\rho}$ . The case in which  $|t| < \alpha + \xi$  follows from Lemma 5.4.4. Let  $\alpha + \xi \leq |t| < \alpha + \omega^{1+\rho'}$  for some natural number  $m > n$ . It follows from Lemma 5.4.4 that

$$\mathsf{H}_{\alpha+\omega^{1+\rho}} \Big| \frac{2+\delta'}{\star} \exists x \in \text{St}_t [P_{<\alpha+\omega^{1+\rho'}.m}(\langle x, t \rangle) \wedge P_{<\alpha+\omega^{1+\rho'}.m}(\langle \neg x, t \rangle)], \Gamma. \quad (5.7)$$

Since this sequent is of  $\mathcal{L}_{\text{fix}}^{\alpha+\omega^{1+\rho'}.m}$ , we can apply Main Lemma II of [43] to (5.7) and obtain

$$\mathsf{H}_{\alpha+\omega^{1+\rho'}.m} \Big| \frac{\varphi 1\rho(2+\delta')}{\star} \exists x \in \text{St}_t [P_{<\alpha+\omega^{1+\rho'}.m}(\langle x, t \rangle) \wedge P_{<\alpha+\omega^{1+\rho'}.m}(\langle \neg x, t \rangle)], \Gamma;$$

thus we have  $\mathsf{H}_{\alpha+\omega^{1+\rho'}.m} \Big| \frac{\varphi 1\rho(2+\delta')+2}{\star} \neg\text{Cons}_{\alpha+\omega^{1+\rho'}.m}, \Gamma$ . Since we have assumed that  $\xi < \omega^{1+\rho'}$ .  $(n+1) \leq \omega^{1+\rho'}$ , we obtain by applying IH for  $(m-n)$ -times (cf. [43]) that

$$\mathsf{H}_{\alpha+\xi} \Big| \frac{\varphi 1\rho'(\dots(\varphi 1\rho'(\varphi 1\rho(2+\delta')+2))\dots)}{\star} \neg\text{Cons}_{\alpha+\xi}, \Gamma,$$

where  $\varphi 1\rho'(\dots(\varphi 1\rho'(\varphi 1\rho(2+\delta')+2))\dots)$  is the result of operating  $\varphi 1\rho'(\cdot)$  for  $(m-n)$ -times. Since  $\rho' < \rho$  and  $\varphi 1\rho(2+\delta')+2 < \varphi 1\rho(2+\delta)$ , it is  $< \varphi 1\rho(2+\delta)$ ; we are done.

The case in which the last inference is cut remains to be shown; but this can be shown in a completely parallel manner to the one in the cited lemma of [43] (use SIH first, apply cut, and then use IH finitely many times). □

**Corollary 5.4.11.** Let  $\Gamma \subset \mathcal{L}_0^+$  and  $\alpha = \omega^{1+\alpha_n} + \dots + \omega^{1+\alpha_1} + m$ . If  $\mathsf{H}_\alpha \Big| \frac{\delta}{\star} \neg\text{Cons}_\alpha, \Gamma$  then  $\mathsf{H}_0 \Big| \frac{\leq \varphi 1\alpha_n(\dots(\varphi 1\alpha_1((\alpha|m)))\dots)}{\star} \neg\text{Cons}_0, \Gamma$ .

**Lemma 5.4.12.** If  $\mathsf{H}_\alpha \Big| \frac{\delta}{\star} \neg\text{Cons}_0, \Gamma$ , then  $\mathsf{H}_\alpha \Big| \frac{\delta}{\star} \Gamma$ ; by induction on  $\delta$  (cf.  $\neg\text{Cons}_0$  is simply false).

**Theorem 5.4.13.**  $|\text{KF}_\alpha + \text{Cons}_\alpha| = |\text{KF}_\alpha| = |\text{KF}_\alpha + \text{Comp}_\alpha|$ . It also follows, in the same way as Theorem 2 of [78], that  $|\text{Aut}(\text{KF} + \text{Cons})| = |\text{Aut}(\text{KF} + \text{Comp})| = \varphi 200$ .

## 5.5 Positive Uniform T-biconditionals

Recall that the purely disquotationalist system PUTB comprises the schema of *positive uniform T-biconditionals*, PUTB: namely,  $\phi(\vec{x})$  is true iff  $\phi(\vec{x})$ , for each formula  $\phi(\vec{x})$  in which the truth predicate occurs only positively. In the present section, we consider the iteration and autonomous progression of this schema.

**Definition 5.5.1.** Let  $\prec$  be a p.r. ordering and  $T_{\prec z}^{\prec}(x, y)$  denote  $y \prec z \wedge T_y^{\prec}(x)$ . The system  $\text{PUTB}_{\prec}$  over  $\mathcal{L}^{\prec}$  is defined as  $\text{PA}(U)$  plus full-induction for  $\mathcal{L}^{\prec}$ ,  $\text{TI}_{\mathcal{L}^{\prec}}(\prec)$  and the axiom schema  $\text{PUTB}_{\prec}$ : for each  $R$ -positive formula  $\mathcal{A}(R, Q, \vec{x}, y) \in \mathcal{L}_0(R, Q)$ ,

$$\forall y \in \text{fd}(\prec) \forall \vec{x} [\mathcal{A}(T_y^{\prec}, T_{\prec y}^{\prec}, \vec{x}, y) \leftrightarrow T_y^{\prec}(\ulcorner \mathcal{A}(T_y^{\prec}, T_{\prec y}^{\prec}, \vec{x}, y) \urcorner)].$$

The system  $\text{Aut}(\text{PUTB})$  denotes  $\text{PA}(U)$  plus full-induction for  $\mathcal{L}$  and the following rules:

$$\frac{\text{TI}(\prec; U)}{\phi} \quad \text{and} \quad \frac{\text{TI}(\prec; U)}{\text{TI}(\prec; \psi)},$$

for each p.r. linear ordering  $\prec$ , any instance  $\phi$  of  $\text{PUTB}_{\prec}$  and arbitrary  $\mathcal{L}$ -formula  $\psi$ .

**Lemma 5.5.2.** There exists a syntactical embedding of  $\text{PUTB}_{\prec}$  in  $\text{KF}_{\prec}$ .

*Proof.* It is known that  $\text{KF} \vdash \text{PUTB}$  and its proof is a straightforward meta-induction; see [5].

However, the same proof doesn't work for iterated truth and we will construct a suitable translation of  $\text{PUTB}_{\prec}$  in  $\text{KF}_{\prec}$ . By the recursion theorem (or diagonalization), we define a binary function  $f$  such that

$$f(x, y) = \begin{cases} x & \text{if } x \in \text{AtFml}_0 \\ T_y^{\prec} \text{ nm}(\ulcorner T_z^{\prec} f(w, z) \urcorner) & \text{if } x = T_z^{\prec} w \text{ for } z, w \in \text{Tm} \\ \neg f(z, y) \quad (f(z, y) \wedge f(w, y), \text{ resp.}) & \text{if } x = \neg z \text{ (or } x = z \wedge w) \\ \forall w. f(z, y) & \text{if } x = \forall w. z; \end{cases}$$

otherwise  $f(x, y) = 0$ ; recall that  $\ulcorner T_z^{\prec} f(w, z) \urcorner$  means  $\text{sb}(\ulcorner T^{\prec}(f(\cdot, \cdot), \cdot) \urcorner; w, z)$ . Now, let  $\mathcal{T}$  be a syntactical translation which preserves  $\mathcal{L}_0$ -atomics, translates  $T_z^{\prec} x$  by  $T_z^{\prec} f(x, z)$  and commutes with the logical connectives and quantifiers. We will show  $\mathcal{T}$  is an interpretation of  $\text{PUTB}_{\prec}$  in  $\text{KF}_{\prec}$  by meta-induction on  $\mathcal{A}$ . First, we trivially have  $\mathcal{T}(T_z^{\prec}(\ulcorner \phi \urcorner)) \leftrightarrow \mathcal{T}(\phi)$  for all  $\mathcal{L}_0$ -atomics  $\phi$ . It is also easily observed that  $\mathcal{T}(T_z^{\prec}(\ulcorner T_{\vec{z}}^{\prec} t(\vec{x}) \urcorner))$  is equivalent to  $\mathcal{T}(T_z^{\prec} t(\vec{x}))$  in  $\text{KF}_{\prec}$ . Next, we show that  $\mathcal{T}(T_{\vec{z}}^{\prec}(s(\vec{x}), t(\vec{x})))$  and  $\mathcal{T}(\neg T_{\vec{z}}^{\prec}(s(\vec{x}), t(\vec{x})))$  are respectively equivalent to  $\mathcal{T}(T_z^{\prec}(\ulcorner T_{\vec{z}}^{\prec}(s(\vec{x}), t(\vec{x})) \urcorner))$  and  $\mathcal{T}(F_z^{\prec}(\ulcorner T_{\vec{z}}^{\prec}(s(\vec{x}), t(\vec{x})) \urcorner))$ . We only demonstrate the latter:

$$\begin{aligned} \mathcal{T}(F_z^{\prec}(\ulcorner t(\vec{x}) \prec z \wedge T_{t(\vec{x})}^{\prec} s(\vec{x}) \urcorner)) &\leftrightarrow F_z^{\prec}(\ulcorner t(\vec{x}) \prec z \wedge T_{\vec{z}}^{\prec} \text{nm}(\ulcorner T_{t(\vec{x})}^{\prec} f(s(\vec{x}), t(\vec{x})) \urcorner) \urcorner) \\ &\leftrightarrow t(\vec{x}) \not\prec z \vee F_z^{\prec}(T_{\vec{z}}^{\prec} \text{nm}(\ulcorner T_{t(\vec{x})}^{\prec} f(s(\vec{x}), t(\vec{x})) \urcorner)) \\ &\leftrightarrow t(\vec{x}) \not\prec z \vee \neg T_{t(\vec{x})}^{\prec} f(s(\vec{x}), t(\vec{x})) \equiv \mathcal{T}(\neg T_{\vec{z}}^{\prec}(s(\vec{x}), t(\vec{x}))); \end{aligned}$$

it is crucial here that  $f(T_w^{\prec} u, z)$  is always in  $\text{St}_z^{\prec}$  for any  $u, w \in \text{CT}$ . We have thus established the base step. Finally, by observing that  $f(\ulcorner \mathcal{A}(T_{\vec{y}}^{\prec}, T_{\vec{z}}^{\prec}, \vec{x}, \vec{y}) \urcorner, y) \in \text{St}_y^{\prec}$  for each  $y \in \text{fd}(\prec)$  and  $\mathcal{A}$ , the induction steps are straightforwardly obtained.  $\square$

**Corollary 5.5.3.** There exists a syntactical embedding of  $\text{PUTB}_{\prec} + \text{Cons}_{\prec}$  ( $\text{PUTB}_{\prec} + \text{Comp}_{\prec}$ ) in  $\text{KF}_{\prec} + \text{Cons}_{\prec}$  ( $\text{KF}_{\prec} + \text{Comp}_{\prec}$ , resp.).

*Proof.* The above interpretation  $\mathcal{T}$  works as well. Observe that, for any  $z \in \text{fd}(\prec)$  and  $x \in \text{St}_z^{\prec}$ ,  $\mathcal{T}(T_z^{\prec} x \wedge F_z^{\prec} x)$  is equivalent to  $T_z^{\prec} f(x, z) \wedge T_z^{\prec} \neg f(x, z)$ ; the first claim follows from this. The second claim is shown similarly.  $\square$

**Lemma 5.5.4.** There is a syntactical embedding of  $\widehat{\text{ID}}_{\prec}$  in  $\text{PUTB}_{\prec}$  (and thus in  $\text{KF}_{\prec}$  as well).

*Proof.* Given an inductive operator  $\mathcal{A}(R, Q, x, y)$ , set  $\mathcal{A}'(x, y, z) \in \mathcal{L}^{\prec}$  be:

$$\mathcal{A}(\lambda u.T_y^{\prec} z(u, y), \lambda w.[(w)_1 \prec y \wedge T_{(w)_1}^{\prec} z((w)_0, (w)_1)], x, y).$$

Recall that  $z(u, v)$  is the result of respectively substituting  $u$ -th and  $v$ -th numerals for the first and second free variables in (a code of formula)  $z$ ; we assume here that  $x$  comes earlier than  $y$  in the fixed variable enumeration. Then, by parametrized diagonalization (for  $z$ ), there is a formula  $\psi_{\mathcal{A}}(x, y)$  such that

$$\psi_{\mathcal{A}}(x, y) \leftrightarrow \mathcal{A}(\lambda u.T_y^{\prec} \ulcorner \psi_{\mathcal{A}}(\dot{u}, \dot{y}) \urcorner, \lambda w.T_{(w)_1}^{\prec} (\ulcorner \psi((\dot{w})_0, (\dot{w})_1) \urcorner, (w)_1), x, y);$$

note that  $(w)_1 \prec y \wedge T_{(w)_1}^{\prec} \ulcorner \psi^{\neg}((w)_0, (w)_1) \urcorner \equiv T_{(w)_1}^{\prec} (\ulcorner \psi((\dot{w})_0, (\dot{w})_1) \urcorner, (w)_1)$ . By the standard construction of  $\psi_{\mathcal{A}}$ , we can assume that  $\psi_{\mathcal{A}}$  is of the form  $\mathcal{B}(T_y^{\prec}, T_{(w)_1}^{\prec}, x, y)$  for some  $R$ -positive  $\mathcal{B} \in \mathcal{L}_0(R, Q)$  and thus  $\text{PUTB}_{\prec} \vdash \forall x \forall y [T_y^{\prec} \ulcorner \psi_{\mathcal{A}}(\dot{x}, \dot{y}) \urcorner \leftrightarrow \psi_{\mathcal{A}}(x, y)]$ . Then, the embedding  $*$  of  $\widehat{\text{ID}}_{\prec}$  in  $\text{PUTB}_{\prec}$  is obtained by  $(P^{\mathcal{A}}x)^* \mapsto T_y^{\prec} \ulcorner \psi_{\mathcal{A}}((\dot{x})_0, (\dot{x})_1) \urcorner$ , since  $(P_{\prec y}^{\mathcal{A}}w)^* \equiv T_{(w)_1}^{\prec} (\ulcorner \psi((\dot{w})_0, (\dot{w})_1) \urcorner, (w)_1)$ .  $\square$

**Theorem 5.5.5.**  $|\text{PUTB}_{\alpha}| = |\text{PUTB}_{\alpha} + \text{Cons}_{\alpha}| = |\text{PUTB}_{\alpha} + \text{Comp}_{\alpha}| = |\text{KF}_{\alpha}|$ . We also have  $|\text{Aut}(\text{PUTB})| = |\text{Aut}(\text{PUTB} + \text{Cons})| = |\text{Aut}(\text{PUTB} + \text{Comp})| = \varphi 200$ .

## 5.6 Kripke-Feferman Truth with Weak Kleene Logic

In the present section, we consider the iteration and autonomous progression of WKF-truth. I remind the reader that the truth tables for the logical connectives in weak Kleene logic are given as follows:

$\neg$		$\wedge$	T	U	F
T	F	T	T	U	F
U	U	U	U	U	U
F	T	F	F	U	F

The evaluation of the universal quantifier is:  $\forall x\phi(x)$  is true iff  $\phi(x)$  is true for all  $x$ ; and  $\forall x\phi(x)$  is false iff  $\phi(x)$  is false for some  $x$  and  $\phi(x)$  is either true or false for all  $x$ . The other logical primitives are defined in the usual manner in terms of ‘ $\neg$ ’, ‘ $\wedge$ ’ and ‘ $\forall$ ’.

**Definition 5.6.1.** Given a p.r. ordering  $\prec$ , the system  $\text{WKF}_{\prec}$  over  $\mathcal{L}^{\prec}$  consists of  $\text{PA}(U)$ , full induction for  $\mathcal{L}^{\prec}$ ,  $\text{TI}_{\mathcal{L}^{\prec}}(\prec)$ , **K1** $_{\prec}$ –**K4** $_{\prec}$ , and: (instead of **K5**–**K6**)

**W5** $_{\prec}$   $\forall a \in \text{fd}(\prec) \forall x, y \in \text{St}_a^{\prec}$

$$\begin{aligned} & [(T_a^{\prec}(x \wedge y) \leftrightarrow (T_a^{\prec}x \wedge T_a^{\prec}y)) \\ & \wedge [F_a^{\prec}(x \wedge y) \leftrightarrow ((F_a^{\prec}x \wedge F_a^{\prec}y) \vee (T_a^{\prec}x \wedge F_a^{\prec}y) \vee (F_a^{\prec}x \wedge T_a^{\prec}y))]]. \end{aligned}$$

**W6** $_{\prec}$   $\forall a \in \text{fd}(\prec) \forall z \forall x$

$$\begin{aligned} & [\forall z. x \in \text{St}_a^{\prec} \rightarrow [(T_a^{\prec}(\forall z.x) \leftrightarrow \forall y T_a^{\prec}x(y)) \\ & \wedge (F_a^{\prec}(\forall z.x) \leftrightarrow (\exists y F_a^{\prec}x(y) \wedge \forall y (T_a^{\prec}x(y) \vee F_a^{\prec}x(y)))]]. \end{aligned}$$

The system  $\text{Aut}(\text{WKF})$  is  $\text{PA}(U)$  plus full-induction for  $\mathcal{L}$ , and the following rules:

$$\frac{\text{TI}(\prec; U)}{\mathbf{K1}_{\prec} \wedge \cdots \wedge \mathbf{K4}_{\prec} \wedge \mathbf{W5}_{\prec} \wedge \mathbf{W6}_{\prec}} \quad \text{and} \quad \frac{\text{TI}(\prec; U)}{\text{TI}(\prec; \psi)},$$

for each p.r. linear ordering  $\prec$  and arbitrary  $\mathcal{L}$ -formula  $\psi$ .

$\text{WKF}_{\prec}$  is embeddable in  $\widehat{\text{ID}}_{\prec}$  in the same way as in the case of  $\text{KF}_{\prec}$ ; for,  $\text{WKF}_{\prec}$  can be modelled by iterated Kripkean fixed-points with weak Kleene schema (cf. [50]). Hence, we can regard  $\text{WKF}_{\prec} \subset \widehat{\text{ID}}_{\prec}$ ; therefore  $|\text{WKF}_{\prec}| \leq |\text{KF}_{\prec}|$  and  $|\text{Aut}(\text{WKF})| \leq |\text{Aut}(\text{KF})|$ . In the present section, we give a well-ordering proof for  $\text{WKF}_{\alpha}$  which yields  $|\text{WKF}_{\alpha}| \geq |\text{KF}_{\alpha}|$ . For readability, as in §4, we

identify ordinals with their representations in the fixed notation system and drop the superscript  $\triangleleft \upharpoonright_\alpha$  whenever there is no worry of confusion.

The well-ordering proof for  $\text{WKF}_\alpha$  goes almost parallel to the one for  $\text{KF}_\alpha$  ( $\text{SRT}_\alpha$ ) in [43]. The point is that  $\text{KF}_\alpha$  and  $\text{WKF}_\alpha$  behave in the same way on determinate sentences. However, we still need to make slight modifications. In what follows, instead of giving a well-ordering proof specifically for  $\text{WKF}_\alpha$ , we will show a more general theorem.

**Theorem 5.6.2.** Given an ordinal  $\alpha$  such that

$$\alpha = \omega^{1+\alpha_0} + \dots + \omega^{1+\alpha_n} + m,$$

for ordinals  $\alpha_0 \geq \dots \geq \alpha_n$  and  $m < \omega$ , let  $\mathbf{Q}_\alpha$  be a theory over  $\mathcal{L}^{\triangleleft \upharpoonright_\alpha}$  which extends  $\text{PA}(U)$  and contains full-induction for  $\mathcal{L}^{\triangleleft \upharpoonright_\alpha}$  and  $\text{TI}_{\mathcal{L}^{\triangleleft \upharpoonright_\alpha}}(\alpha)$ . We write  $D_\beta^+(x)$  for  $\text{St}_\beta(x) \wedge [F_\beta x \leftrightarrow \neg T_\beta x]$ .

Suppose  $\mathbf{Q}_\alpha$  derives all the following:

- (i) The axioms  $\mathbf{K1}_\alpha$  and  $\mathbf{K2}_\alpha$ ;
- (ii)  $\forall \beta < \alpha \forall x \forall y [(x \in D_\beta^+ \wedge y \in D_\beta^+) \rightarrow (\neg x \in D_\beta^+ \wedge x \wedge y \in D_\beta^+)]$ ;
- (iii)  $\forall \beta < \alpha \forall x \forall z [(\forall z.x \in \text{St}_\beta \wedge \forall y(x(y) \in D_\beta^+)) \rightarrow (\forall z.x \in D_\beta^+)]$ ;
- (iv)  $\forall \beta < \alpha \forall x [\neg x \in D_\beta^+ \rightarrow (T_\beta \neg x \leftrightarrow \neg T_\beta x)]$ ;
- (v)  $\forall \beta < \alpha \forall x \forall y [x \wedge y \in D_\beta^+ \rightarrow (T_\beta(x \wedge y) \leftrightarrow (T_\beta x \wedge T_\beta y))]$ ;
- (vi)  $\forall \beta < \alpha [\forall z.x \in D_\beta^+ \rightarrow (T_\beta \forall z.x \leftrightarrow \forall y T_\beta x(y))]$ ;
- (vii) The axiom  $\mathbf{K3}_\alpha$ .

Then,  $\mathbf{Q}_\alpha \vdash \text{TI}(\xi; U)$  for all  $\xi < \varphi 1 \alpha_0 (\dots (\varphi 1 \alpha_n (\alpha | m)))$ .

Since  $\text{WKF}_\alpha$  meets all these conditions, this theorem suffices for our goal and we will give its proof in the rest of the present section.

Let  $x \in tc_\beta$  denote  $\text{For}_\beta(x) \wedge \forall y (F_\beta x(y) \leftrightarrow \neg T_\beta x(y))$ ,<sup>11</sup> and set:

$$\text{Prog}^\beta(a) := \forall \eta [(\forall \xi < \eta) T_\beta a(\xi) \rightarrow T_\beta a(\eta)]$$

$$\text{TI}^\beta(a, \gamma) := \text{Prog}^\beta(a) \rightarrow \forall \xi < \gamma T_\beta a(\xi)$$

$$I^\beta(\gamma) := (\forall \delta < \beta) (\forall a \in tc_\delta) \text{TI}^\delta(a, \gamma).$$

Correspondingly, given a formula  $\phi$ , we write

$$\text{Prog}(\phi) := \forall \eta (\forall \xi < \eta \phi(\xi) \rightarrow \phi(\eta)) \quad \text{TI}(\phi, \gamma) := \text{Prog}(\phi) \rightarrow \forall \xi < \gamma \phi(\xi).$$

Let  $\text{Lim}(x)$  represent the set of limit ordinals. We can easily check these possess the following desired properties in  $\mathbf{Q}_\alpha$  (in fact, they are provable without using (i)-(vii)).

- $\forall \beta < \alpha \forall a [\forall \gamma (\text{TI}^\beta(a, \gamma) \rightarrow \text{TI}^\beta(a, \gamma+1)) \wedge \forall \gamma \in \text{Lim} (\forall \xi < \gamma \text{TI}^\beta(a, \xi) \rightarrow \text{TI}^\beta(a, \gamma))]$
- $\forall \beta \leq \alpha [\forall \gamma (I^\beta(\gamma) \rightarrow I^\beta(\gamma+1)) \wedge \forall \gamma \in \text{Lim} ((\forall \xi < \gamma) I^\beta(\xi) \rightarrow I^\beta(\gamma))]$
- $\forall \beta \leq \alpha \forall \gamma [\forall \delta \leq \beta (I^\beta(\gamma) \rightarrow I^\delta(\gamma)) \wedge (\beta \in \text{Lim} \wedge (\forall \delta < \beta) I^\delta(\gamma) \rightarrow I^\beta(\gamma))].$

An  $\mathcal{L}^\prec$ -formula  $\phi(x, \vec{z})$  is said to be *typed by*  $x$ , when  $\phi$  is constructed from  $\mathcal{L}_0$ -atomics and  $T_x^\prec(s)$  by logical connectives and quantifiers over variables different from  $x$ . Then, when  $\phi(x, \vec{z}) \in \mathcal{L}^\prec$  is typed by  $x$ ,  $\phi(t, \vec{z})$  is  $\mathcal{L}_t^\prec$  for a closed term  $t$  with  $t \in \text{fd}(\prec)$  and  $\forall a \in \text{fd}(\prec) (\ulcorner \phi(\dot{a}, \vec{z}) \urcorner \in \text{Fml}_a^\prec)$  (provably in  $\text{PA}(U)$ ); in particular,  $\forall \beta < \alpha (\ulcorner \phi(\dot{\beta}, \vec{z}) \urcorner \in \text{Fml}_\beta)$  when  $\phi \in \mathcal{L}^{\prec 1_\alpha}$ . We will write  $\phi_a(\vec{z})$  for a typed formula  $\phi(a, \vec{z})$ . Notice that  $\text{Prog}^\beta$  and  $\text{TI}^\beta$  are typed by  $\beta$ . Now  $\mathbf{Q}_\alpha$  proves the following.

- (a)  $\forall \beta < \alpha \forall \gamma < \beta \forall a [(\text{For}_\gamma(a) \rightarrow a \in tc_\beta) \wedge (\text{St}_\gamma(a) \rightarrow D_\beta^+(a))];$  induction on  $a$  using (i)-(iii).

---

<sup>11</sup>It stands for ‘total and consistent predicate at the  $\beta$ -th level’; cf. [5]. It is also denoted by  $S_\beta$  in [43].

- (b)  $\forall\beta < \alpha \forall\gamma < \beta \forall \vec{z} [T_\beta(\ulcorner \phi_\gamma(\vec{z}) \urcorner) \leftrightarrow \phi_\gamma(\vec{z})]$ , for each formula  $\phi_x(\vec{z})$  typed by  $x$ ; by meta-induction on  $\phi$  using (a), (i) and (iv)-(vi).
- (c)  $\forall\beta < \alpha \forall\gamma < \beta \forall\delta [(\text{Prog}^\beta(\ulcorner \phi_\gamma \urcorner) \leftrightarrow \text{Prog}(\phi_\gamma)) \wedge (\text{TI}^\beta(\ulcorner \phi_\gamma \urcorner, \delta) \leftrightarrow \text{TI}(\phi_\gamma, \delta))]$ , for each formula  $\phi_x(z)$  typed by  $x$ ; use (b).
- (d)  $\forall\beta < \alpha \forall a \in tc_\beta[\ulcorner T_\beta \dot{a}(\cdot) \urcorner \in tc_\beta]$ ; by (vii).
- (e)  $\forall\beta < \alpha \forall a \in tc_\beta[\ulcorner \text{TI}^{\dot{\beta}}(\dot{a}, \cdot) \urcorner \in tc_\beta \wedge \forall\delta (T_\beta(\ulcorner \text{TI}^{\dot{\beta}}(\dot{a}, \delta) \urcorner) \leftrightarrow \text{TI}^\beta(a, \delta))]$ ; all substitution instances of subformulae of  $\ulcorner \text{TI}^{\dot{\beta}}(\dot{a}, \cdot) \urcorner$  is in  $D_\beta^+$  by (d) and (i)-(iii); then the second conjunct follows from (i) and (iv)-(vii).

**Lemma 5.6.3.** Let  $B^\beta(\gamma)$  be  $\forall\delta[\forall a \in tc_\beta \text{TI}^\beta(a, \delta) \rightarrow \forall a \in tc_\beta \text{TI}^\beta(a, \varphi_\gamma \delta)]$ . Then, we have  $\mathbb{Q}_\alpha \vdash \forall\beta < \alpha \text{Prog}(\lambda\gamma. B^\beta(\gamma))$ ; shown in the standard manner (cf. [17]) using (e) above.<sup>12</sup>

**Corollary 5.6.4.**  $\mathbb{Q}_\alpha \vdash \forall\beta \leq \alpha [\beta \in \text{Lim} \rightarrow \forall\xi (I^\beta(\xi) \rightarrow I^\beta(\varphi_\xi 0))]$ ; by the last lemma and applying (c) to  $B^\delta$  for  $\delta < \beta$  (note that  $B^\delta$  is typed by  $\delta$  and thus a  $tc_{\delta+1}$ -predicate).

**Corollary 5.6.5.**  $\mathbb{Q}_\alpha \vdash \forall\beta \leq \alpha [\beta \in \text{Lim} \rightarrow \text{Prog}(\lambda\delta. I^\beta(\Gamma_\delta))]$ ; use the last corollary.

The arguments so far have been parallel to [43], but we need a slight modification for the next step. Jäger et al. showed in Lemma 4 of [43] that  $\text{KF}_\alpha$  proves

$$\forall\beta \in (\text{Lim} \cap \alpha) \exists a \in tc_\beta \forall x [T_\beta a(x) \leftrightarrow I^\beta(x)], \quad (5.8)$$

which states that  $I^\beta$  can be represented by a  $tc_\beta$ -predicate  $a$ . Although they didn't give an explicit proof, such a representation  $a$  should be the code of the following kind of formula:

$$(\forall\delta < \beta)(\forall b) [T_\beta(\ulcorner tc_\delta(\dot{b}) \urcorner) \rightarrow \text{TI}^\beta(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, x)]. \quad (5.9)$$

<sup>12</sup>The corresponding statement for  $\text{SRT}_\alpha$  is implicit in the proof of Lemma 5 of [43], which then corresponds to our Corollary 5.6.4 below.

For this  $a$ , since (5.9) above is an  $\mathcal{L}_\beta$ -sentence,  $\mathbf{KF}_\alpha$  can derive:

$$T_\beta a(x) \leftrightarrow (\forall \delta < \beta)(\forall b)[T_\beta(\ulcorner T_\beta(\ulcorner tc_\delta(\dot{b}) \urcorner) \urcorner) \rightarrow T_\beta(\ulcorner \text{TI}^\beta(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, \dot{x}) \urcorner)]; \quad (5.10)$$

thus we can focus only on  $\delta < \beta$  and  $b \in tc_\delta$  (note that  $tc_\delta$  is typed by  $\delta$ ) and then obtain (5.8) in  $\mathbf{KF}_\alpha$ . However, (5.10) is no longer valid in  $\mathbf{WKF}_\alpha$ ; in general, even when  $A$  is arithmetical,  $T_\beta(\ulcorner A \rightarrow B \urcorner)$  isn't equivalent to  $A \rightarrow T_\beta(B)$  in  $\mathbf{WKF}_\alpha$ . Since (5.8) is used in the proof of Main Lemma I of [43], we have to make a modification at this point.

**Proposition 5.6.6.** For each term  $t(z, \vec{x})$ , there is an  $\mathcal{L}^\prec$ -formula  $\phi_z(\vec{x})$  typed by  $z$  such that

$$\mathbf{Q}_\alpha \vdash \forall \beta < \alpha \forall \gamma < \beta \forall \vec{x} [T_\beta(\ulcorner \phi_\gamma(\vec{x}) \urcorner) \leftrightarrow I^\gamma(t(\gamma, \vec{x}))].$$

*Proof.* As such  $\phi$ , we can take  $(\forall \delta < z)(\forall b)[T_z \ulcorner tc_\delta(\dot{b}) \urcorner \rightarrow \text{TI}^z(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, t(z, \vec{x}))]$ . This  $\phi$  is typed by  $z$ .

Hence, for  $\gamma < \beta < \alpha$ , it follows from (b) that:

$$\forall \vec{x} [T_\beta \ulcorner \phi_\gamma(\vec{x}) \urcorner \leftrightarrow (\forall \delta < \gamma)(\forall b)(T_\gamma \ulcorner tc_\delta(\dot{b}) \urcorner \rightarrow \text{TI}^\gamma(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, t(\gamma, \vec{x})))];$$
 and then

$$\forall \vec{x} [T_\beta \ulcorner \phi_\gamma(\vec{x}) \urcorner \leftrightarrow (\forall \delta < \gamma)(\forall b \in tc_\delta) \text{TI}^\gamma(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, t(\gamma, \vec{x}))].$$

Since  $\text{TI}^\gamma(\ulcorner T_\delta \dot{b}(\cdot) \urcorner, w) \leftrightarrow \text{TI}^\delta(b, w)$  for  $\delta < \gamma$  and all  $w$  (by **K2** $_\alpha$ ), we obtain the claim.  $\square$

This proposition tells that  $I^\gamma$  can be represented at any higher level  $\beta > \gamma$ . In fact, the next lemma suffices for proving Main Lemma I of [43]. It follows from Proposition 5.6.6 and doesn't require for  $I^\gamma$  to be representable at the *same* level  $\gamma$  (i.e., (5.8)).

**Lemma 5.6.7.**  $\mathbf{Q}_\alpha \vdash \forall \beta \leq \alpha \forall \gamma < \beta \forall \delta \forall \xi [(\text{Lim}(\beta) \wedge I^\beta(\xi) \wedge \text{Prog}(\lambda x. I^\gamma(\varphi 1 \delta x))) \rightarrow I^\gamma(\varphi 1 \delta \xi)].$

*Proof.* Let  $\beta \leq \alpha$  with  $\text{Lim}(\beta)$  and  $\gamma < \beta$ . Suppose  $I^\beta(\xi) \wedge \text{Prog}(\lambda x. I^\gamma(\varphi 1 \delta x))$ . It follows from the last proposition that there is  $a \in \text{For}_\gamma \subset tc_{\gamma+1}$  which represents  $I^\gamma(\varphi 1 \delta x)$ . Then, since  $\gamma + 1 < \beta$ ,

we obtain  $\text{TI}^{\gamma+1}(a, \xi)$  from  $I^\beta(\xi)$ . Finally, since  $\text{Prog}(\lambda x. I^\gamma(\varphi 1 \delta x))$  is equivalent to  $\text{Prog}^{\gamma+1}(a)$ , we obtain  $(\forall \eta < \xi) T_{\gamma+1} a(\varphi 1 \delta \eta)$ ; thus  $I^\gamma(\varphi 1 \delta \xi)$ .  $\square$

Following the notations of [43], let  $\gamma \uparrow \beta$  and  $\text{Main}_\alpha(\delta)$  respectively denote

$$\exists \eta \in \text{Lim} \exists \beta_0 [\beta = \beta_0 + \gamma \cdot \eta] \quad \text{and} \quad \forall \beta \leq \alpha \forall \gamma [\omega^{1+\delta} \uparrow \beta \wedge I^\beta(\gamma) \rightarrow I^\beta(\varphi 1 \delta \gamma)],$$

**Lemma 5.6.8.**  $\mathbb{Q}_\alpha \vdash \text{Prog}(\lambda \delta. \text{Main}_\alpha(\delta))$ .

*Proof.* The proof is divided into three cases: i.e., the base case ( $\delta = 0$ ), successor case and limit case. Our proof goes parallel to Main Lemma I of [43] except that we use the last lemma instead of (5.8). We will briefly illustrate how the last lemma is used.

Assume  $(\forall \eta < \delta) \text{Main}_\alpha(\eta)$  and take  $\beta \leq \alpha$  with  $\omega^{1+\delta} \uparrow \beta$ ; then  $\beta \in \text{Lim}$ . For a suitable fundamental sequence  $\beta[n]$  ( $\lim_{n \rightarrow \infty} \beta[n] = \beta$ ), we can show in the same way as [43] that:

$$(\forall n) \text{Prog}(\lambda x. I^{\beta[n]}(\varphi 1 \delta x)); \tag{5.11}$$

note that we use Corollary 5.6.5 for the base case where  $\delta = 0$  (and IH for the other cases).

We have not yet used the last lemma until this point, and it is the next step that we need it, in which we show  $(\forall n) I^{\beta[n]}(\varphi 1 \delta \gamma)$  from (5.11) under the assumption  $I^\beta(\gamma)$ . This immediately follows from Lemma 5.6.7, since  $\beta \in \text{Lim}$  and  $\beta[n] < \beta$  for all  $n$ .  $\square$

We can show the following in the standard manner:

- $\mathbb{Q}_\alpha \vdash \forall \beta \leq \alpha \forall \delta \forall \beta_0 [\beta = \beta_0 + \omega^{1+\delta} \rightarrow (\forall \eta < \delta) (\forall \gamma) [I^\beta(\gamma) \rightarrow I^\beta(\varphi 1 \eta \gamma)]]$ ; by the last lemma and  $\text{TI}_{\mathcal{L}^{\leq 1} \alpha}(\alpha)$ . This corresponds to Corollary 8 of [43].
- $\mathbb{Q}_\alpha \vdash \forall \beta \leq \alpha \forall \delta \forall \beta_0 [\beta = \beta_0 + \omega^{1+\delta} \rightarrow \text{Prog}(\lambda \gamma. I^\beta(\varphi 1 \delta \gamma))]$ ; by Corollary 5.6.5 (for  $\delta = 0$ ) and the above (for the other cases). This corresponds to Corollary 9 of [43].

- $\mathbf{Q}_\alpha \vdash I^{\omega^{1+\alpha_0}+\dots+\omega^{1+\alpha_n}}(\xi)$  for each  $\xi < (\alpha|m)$ ; by Lemma 5.6.3 and the fact  $\mathbf{Q}_\alpha \vdash \text{TI}_{\mathcal{L}^{\triangleleft \alpha}}(< \varepsilon(\alpha))$  using (a) and (c).

Using Lemma 5.6.7, these finally entail Theorem 5.6.2.

**Theorem 5.6.9.**  $|\text{WKF}_\alpha| = |\text{KF}_\alpha|$  and  $|\text{Aut}(\text{WKF})| = |\text{Aut}(\text{KF})|$ .

## 5.7 Determinate Truth and Feferman Logic

The present section introduces and studies the iteration and autonomous progression of Feferman's determinate truth axiomatized by DT.

As was shown in Chapter 3, DT is identical with the system  $\text{FKF} + \text{Cons}$ . Recall that Feferman logic is a variant of weak Kleene logic in which the conditional ' $\rightarrow$ ' is not definable in terms of ' $\neg$ ' and ' $\vee$ ' and the evaluation of ' $\rightarrow$ ' is determined by the following separate truth table:

$\rightarrow$	T	U	F
T	T	U	F
U	U	U	U
F	T	T	T

The evaluations of all the other connectives and quantifiers in Feferman logic are the same as weak Kleene logic. Hence, in order to deal with the iteration of DT-truth and FKF-truth in our current framework, we assume that the conditional ' $\rightarrow$ ' is assigned its own Gödel number independently of the other connectives (recall that we have excluded ' $\rightarrow$ ' from the primitive logical connectives in the present chapter); thus  $x \rightarrow y \neq (\neg x) \vee y$ . Then, the systems of iterated DT-truth and FKF-truth are defined as follows.

**Definition 5.7.1.** Let  $\prec$  be a p.r. ordering.  $D_a^\prec(x)$  denotes an  $\mathcal{L}^\prec$ -formula  $T_a^\prec x \vee F_a^\prec x$ ; it expresses ' $x$  is determinate at the  $a$ -th level'. The system  $\text{DT}_\prec$  over  $\mathcal{L}^\prec$  consists of  $\text{PA}(U)$ -axioms, full induction for  $\mathcal{L}^\prec$ ,  $\text{TI}_{\mathcal{L}^\prec}(\prec)$ , and:

$$\mathbf{D1}_{\prec} \quad \forall a \in \text{fd}(\prec) [\forall x (\text{AtSt}_0(x) \rightarrow D_a^{\prec}(x)) \wedge \forall b \prec a \forall x (\text{AtSt}_b^{\prec}(x) \rightarrow D_a^{\prec}(x))]$$

$$\mathbf{D2}_{\prec} \quad \forall a, x \in \text{CT} [a^\circ \in \text{fd}(\prec) \rightarrow (D_{a^\circ}^{\prec}(T_a x) \leftrightarrow D_a^{\prec}(x^\circ))]$$

$$\mathbf{D3}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, y \in \text{St}_a^{\prec} [(D_a^{\prec}(x \rightarrow y) \leftrightarrow (D_a^{\prec}x \wedge (T_a^{\prec}x \rightarrow D_a^{\prec}y))) \\ \wedge (D_a^{\prec}(\neg x) \leftrightarrow D_a^{\prec}x) \wedge (D_a^{\prec}(x \wedge y) \leftrightarrow (D_a^{\prec}x \wedge D_a^{\prec}y))]$$

$$\mathbf{D4}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall z \forall x [\forall z. x \in \text{St}_a^{\prec} \rightarrow (D_a^{\prec}(\forall z. x) \leftrightarrow \forall y D_a^{\prec}x(y))]$$

$$\mathbf{D5}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall \vec{x} \in \text{CT} [T_a^{\prec}(R\vec{x}) \leftrightarrow R\vec{x}^\circ], \text{ for each predicate symbol } R \text{ of } \mathcal{L}_0$$

$$\mathbf{D6}_{\prec} \quad \forall a, b, x \in \text{CT} [b^\circ \preceq a^\circ \rightarrow ((T_a^{\prec}(T_b^{\prec}x) \leftrightarrow T_b^{\prec}x^\circ) \wedge (T_a^{\prec}(T_a^{\prec}x) \leftrightarrow T_a^{\prec}x^\circ))]$$

$$\mathbf{D7}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, y \in D_a^{\prec} \cap \text{St}_a^{\prec} [(T_a^{\prec}(x \rightarrow y) \leftrightarrow (T_a^{\prec}x \rightarrow T_a^{\prec}y)) \\ \wedge (T_a^{\prec}(\neg x) \leftrightarrow \neg T_a^{\prec}x) \wedge (T_a^{\prec}(x \wedge y) \leftrightarrow (T_a^{\prec}x \wedge T_a^{\prec}y))]$$

$$\mathbf{D8}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall z \forall x [\forall z. x \in D_a^{\prec} \cap \text{St}_a^{\prec} \rightarrow (T_a^{\prec}(\forall z. x) \leftrightarrow \forall y T_a^{\prec}x(y))],$$

Then, the system  $\text{Aut}(\text{DT})$  is defined as  $\text{PA}(U)$  plus full-induction for  $\mathcal{L}$ , and the following two rules: for each p.r. linear ordering  $\prec$  and arbitrary  $\mathcal{L}$ -formula  $\psi$ ,

$$\frac{\text{TI}(\prec; U)}{\mathbf{D1}_{\prec} \wedge \cdots \wedge \mathbf{D8}_{\prec}} \quad \text{and} \quad \frac{\text{TI}(\prec; U)}{\text{TI}(\prec; \psi)}.$$

**Definition 5.7.2.**  $\text{FKF}_{\prec}$  over  $\mathcal{L}^{\prec}$  is defined as  $\text{WKF}_{\prec}$  plus the compositional axiom for ‘ $\rightarrow$ ’:

$$\mathbf{F7}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, y \in \text{St}_a^{\prec}$$

$$[(T_a^{\prec}(x \rightarrow y) \leftrightarrow ((T_a^{\prec}x \wedge T_a^{\prec}y) \vee F_a^{\prec}x)) \wedge (F_a^{\prec}(x \rightarrow y) \leftrightarrow (T_a^{\prec}x \wedge F_a^{\prec}y))].$$

Then,  $\text{Aut}(\text{FKF})$ ,  $\text{Aut}(\text{FKF} + \text{Cons})$  and  $\text{Aut}(\text{FKF} + \text{Comp})$  are defined analogously to  $\text{Aut}(\text{KF})$ ,  $\text{Aut}(\text{KF} + \text{Cons})$  and  $\text{Aut}(\text{KF} + \text{Comp})$  respectively.

**Proposition 5.7.3.**  $\text{FKF}_{\prec} + \text{Cons}_{\prec}$  and  $\text{DT}_{\prec}$  are identical (as theories). This is shown in the same manner as Theorem 21 of [24] with obvious modifications. Thus,  $\text{Aut}(\text{DT})$  is also identical with  $\text{Aut}(\text{FKF} + \text{Cons})$ .

Consequently, we can apply the cut-elimination technique developed in §4 to  $\text{DT}_\alpha$  with straightforward modifications and this gives the upper bound of  $|\text{DT}_\alpha|$ . However, we can even show that  $\text{DT}_\prec$  is syntactically embeddable in  $\text{KF}_\prec$  in general.

**Lemma 5.7.4.**  $\text{FKF}_\prec + \text{Cons}_\prec (= \text{DT}_\prec)$  is syntactically embeddable in  $\text{KF}_\prec$ .

*Proof.* Let us drop the superscript ‘ $\prec$ ’ from ‘ $T^\prec$ ’ for readability. First, we define a binary function  $I$  by the recursion theorem (or diagonalization) such that:  $I(b, a) :=$

$$\left\{ \begin{array}{ll} b & \text{if } b \in \text{Fml}_0 \\ \ulcorner T_c I(d, c) \wedge \neg T_c I(\neg d, c) \wedge F_c I(\neg d, c) \wedge \neg F_c I(d, c) \urcorner & \text{if } \exists c, d \in \text{Tm}(b = T_c d \wedge c^\circ \prec a) \\ \ulcorner T_c I(\neg d, c) \vee \neg T_c I(d, c) \vee F_c I(d, c) \vee \neg F_c I(\neg d, c) \urcorner & \text{if } \exists c, d \in \text{Tm}(b = \neg T_c d \wedge c^\circ \prec a) \\ \ulcorner T_c I(d, c) \urcorner & \text{if } \exists c, d \in \text{Tm}(b = T_c d \wedge c^\circ = a) \\ \ulcorner T_c I(\neg d, c) \urcorner & \text{if } \exists c, d \in \text{Tm}(b = \neg T_c d \wedge c^\circ = a) \\ I(c, a) & \text{if } b = \neg \neg c \\ (I(c, a) \wedge I(d, a) \wedge \neg I(\neg c, a) \wedge \neg I(\neg d, a)) & \text{if } b = c \wedge d \\ I(\neg c, a) \wedge I(\neg d, a) \wedge \neg I(c, a) \wedge \neg I(d, a) & \\ \vee (I(\neg c, a) \wedge I(d, a) \wedge \neg I(c, a) \wedge \neg I(\neg d, a)) & \\ \vee (I(c, a) \wedge I(\neg d, a) \wedge \neg I(\neg c, a) \wedge \neg I(d, a)) & \text{if } b = \neg(c \wedge d) \\ (I(c, a) \wedge I(d, a) \wedge \neg I(\neg c, a) \wedge \neg I(\neg d, a)) & \\ \vee (I(\neg c, a) \wedge \neg I(c, a)) & \text{if } b = c \rightarrow d \\ I(c, a) \wedge I(\neg d, a) \wedge \neg I(\neg c, a) \wedge \neg I(d, a) & \text{if } b = \neg(c \rightarrow d) \\ \forall c. (I(d, a) \wedge \neg I(\neg d, a)) & \text{if } b = \forall c. d \\ \forall c. [(I(d, a) \wedge \neg I(\neg d, a)) \vee (I(\neg d, a) \wedge \neg I(d, a))] & \\ \wedge \exists c. I(\neg d, a) & \text{if } b = \neg \forall c. d; \end{array} \right.$$

otherwise, we set  $I(b, a) = 0$ . The embedding  $\mathcal{T}$  of  $\text{DT}_\prec$  in  $\text{KF}_\prec$  is given by:

$$T_z x \mapsto T_z I(x, z) \wedge \neg T_z I(\neg x, z) \wedge F_z I(\neg x, z) \wedge \neg F_z I(x, z).$$

$\mathcal{T}$  is of essentially the same form as the syntactical embedding of  $\text{DT}$  in  $\text{KF}$  given in [24] and the proof is its straightforward generalization except the case for  $\mathbf{K2}_\prec$ . Thus, we only illustrate the

case for  $\mathbf{K2}_{\prec}$ : i.e.,  $\mathbf{KF}_{\prec} \vdash (\mathbf{K2}_{\prec})^{\mathcal{T}}$ . For  $b, x \in \text{CT}$  with  $b^{\circ} \prec a$ ,

$$\begin{aligned}
(F_a(T_b x))^{\mathcal{T}} &\leftrightarrow T_a I(\neg T_b x, a) \wedge \neg T_a I(T_b x, a) \wedge F_a I(T_b x, a) \wedge \neg F_a I(\neg T_b x, a) \\
&\leftrightarrow [T_a \ulcorner T_b I(\neg x, b) \urcorner \vee T_a \ulcorner \neg T_b I(x, b) \urcorner \vee T_a \ulcorner F_b I(x, b) \urcorner \vee T_a \ulcorner \neg F_b I(\neg x, b) \urcorner] \\
&\quad \wedge [\neg T_a \ulcorner T_b I(x, b) \urcorner \vee \neg T_a \ulcorner \neg T_b I(\neg x, b) \urcorner \vee \neg T_a \ulcorner F_b I(\neg x, b) \urcorner \vee \neg T_a \ulcorner \neg F_b I(x, b) \urcorner] \\
&\quad \wedge [F_a \ulcorner T_b I(x, b) \urcorner \vee F_a \ulcorner \neg T_b I(\neg x, b) \urcorner \vee F_a \ulcorner F_b I(\neg x, b) \urcorner \vee F_a \ulcorner \neg F_b I(x, b) \urcorner] \\
&\quad \wedge [\neg F_a \ulcorner T_b I(\neg x, b) \urcorner \vee \neg F_a \ulcorner \neg T_b I(x, b) \urcorner \vee \neg F_a \ulcorner F_b I(x, b) \urcorner \vee \neg F_a \ulcorner \neg F_b I(\neg x, b) \urcorner] \\
&\leftrightarrow T_{b^{\circ}} I(\neg x^{\circ}, b^{\circ}) \vee \neg T_{b^{\circ}} I(x^{\circ}, b^{\circ}) \vee F_{b^{\circ}} I(x^{\circ}, b^{\circ}) \vee \neg F_{b^{\circ}} I(\neg x^{\circ}, b^{\circ}) \equiv (\neg T_{b^{\circ}}(x^{\circ}))^{\mathcal{T}};
\end{aligned}$$

note that each of the four conjuncts of the third formula is equivalent to the fourth. We can dually show  $(T_a \ulcorner T_b x \urcorner)^{\mathcal{T}} \leftrightarrow (T_b x^{\circ})^{\mathcal{T}}$  in the parallel manner.  $\square$

We can easily show that  $\mathbf{WKF}_{\prec} + \mathbf{Comp}_{\prec}$  and  $\mathbf{KF}_{\prec} + \mathbf{Comp}_{\prec}$  are identical; cf. Lemma 51 of [24]. It is also observed that, for all  $a \in \text{fd}(\prec)$  and  $x, y \in \text{St}_a$ ,  $\mathbf{FKF}_{\prec} + \mathbf{Comp}_{\prec}$  derives  $T_a(x \rightarrow y) \leftrightarrow T_a(\neg x \vee y)$  and  $F_a(x \rightarrow y) \leftrightarrow F_a(\neg x \vee y)$ . Namely, in the presence of  $\mathbf{Comp}_{\prec}$ , the inner logics of  $\mathbf{WKF}$  and  $\mathbf{FKF}$  collapse into strong Kleene logic. Now, by the recursion theorem, we can take a function  $h$  such that  $h(\ulcorner T_a t \urcorner) = \ulcorner T_a h(t) \urcorner$ ,  $h(\ulcorner \phi \rightarrow \psi \urcorner) = \neg h(\ulcorner \phi \urcorner) \vee h(\ulcorner \psi \urcorner)$  and commutes with the other connectives. Then, the embedding of  $\mathbf{FKF} + \mathbf{Comp}_{\prec}$  in  $\mathbf{KF} + \mathbf{Comp}_{\prec}$  is given by  $T_z x \mapsto T_z h(x)$ . The main theorem of this section is thus obtained:

**Theorem 5.7.5.** (1)  $|\text{DT}_{\alpha}| = |\mathbf{FKF}_{\alpha}| = |\mathbf{WKF}_{\alpha} + \mathbf{Cons}_{\alpha} \text{ (or } \mathbf{Comp}_{\alpha})| = |\mathbf{KF}_{\alpha}|$ .

(2)  $|\text{Aut}(\text{DT})| = |\text{Aut}(\mathbf{FKF})| = |\text{Aut}(\mathbf{WKF} + \mathbf{Cons})| = |\text{Aut}(\mathbf{WKF} + \mathbf{Comp})| = \varphi_{200}$ .

## 5.8 Cantini's VF-Truth

In the present section, we consider the iteration of this VF-truth and its autonomous progression.

### 5.8.1 The Systems $\mathbf{VF}_{\prec}$ and $\mathbf{Aut}(\mathbf{VF})$

Given a system  $\mathbf{S}$ ,  $\mathbf{Ax}_{\mathbf{S}}$  represents the set of the axioms of  $\mathbf{S}$  and  $\mathbf{Prv}_{\mathbf{S}}$  is a canonical provability predicate of  $\mathbf{S}$ . Let  $\mathbf{PL}_{\mathcal{L}^{\prec}}$  be pure classical predicate logic for the language  $\mathcal{L}^{\prec}$  (thus with equality); then  $\mathbf{Ax}_{\mathbf{PL}_{\mathcal{L}^{\prec}}}$  represents the set of logical axioms (for  $\mathcal{L}^{\prec}$ ). Our following formulation of  $\mathbf{VF}_{\prec}$  is based on the simpler version of  $\mathbf{VF}$  in [7] (or [45]) rather than the original one in [6].

**Definition 5.8.1.** Given a p.r. ordering  $\prec$ , the system  $\mathbf{VF}_{\prec}$  over  $\mathcal{L}^{\prec}$  comprises  $\mathbf{PA}(U)$ , full-induction for  $\mathcal{L}^{\prec}$ ,  $\mathbf{TI}_{\mathcal{L}^{\prec}}(\prec)$  and the following truth axioms:

$$\mathbf{V1}_{\prec} \quad \forall \vec{x} \forall a \in \text{fd}(\prec) [T_a^{\prec}(\ulcorner \phi(\vec{x}) \urcorner) \rightarrow \phi(\vec{x})], \text{ for each } \mathcal{L}^{\prec}\text{-formula } \phi(\vec{x}).$$

$$\mathbf{V2}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall \vec{x} \in \text{CT} [(T_a^{\prec}(R\vec{x}) \leftrightarrow R\vec{x}^{\circ}) \wedge (F_a^{\prec}(R\vec{x}) \leftrightarrow \neg R\vec{x}^{\circ})], \text{ for each } R \in \mathcal{L}_0$$

$$\mathbf{V3}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall b, x \in \text{CT} [b^{\circ} \prec a \rightarrow (T_a^{\prec}(T_b^{\prec}x) \leftrightarrow T_b^{\prec}x^{\circ}) \wedge (F_a^{\prec}(T_b^{\prec}x) \leftrightarrow \neg T_b^{\prec}x^{\circ})]$$

$$\mathbf{V4}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x \in \text{St}^{\prec} [\mathbf{Ax}_{\mathbf{PL}_{\mathcal{L}^{\prec}}}(x) \rightarrow T_a^{\prec}(x)]$$

$$\mathbf{V5}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall z \forall x [\text{St}^{\prec}(\forall z.x) \rightarrow (\forall y T_a^{\prec}x(y) \rightarrow T_a^{\prec}\forall z.x)]$$

$$\mathbf{V6}_{\prec} \quad \forall a \in \text{fd}(\prec) \forall x, y \in \text{St}^{\prec} [T_a^{\prec}(x \rightarrow y) \rightarrow (T_a^{\prec}x \rightarrow T_a^{\prec}y)]$$

The system  $\mathbf{Aut}(\mathbf{VF})$  is defined as  $\mathbf{PA}(U)$  plus full-induction for  $\mathcal{L}$ , and the following rules:

$$\frac{\mathbf{TI}(\prec; U)}{\phi}, \quad \frac{\mathbf{TI}(\prec; U)}{\mathbf{V2}_{\prec} \wedge \dots \wedge \mathbf{V6}_{\prec}} \quad \text{and} \quad \frac{\mathbf{TI}(\prec; U)}{\mathbf{TI}(\prec; \psi)},$$

for each p.r. linear ordering  $\prec$ , instance  $\phi$  of  $\mathbf{V1}_{\prec}$  and  $\mathcal{L}$ -formula  $\psi$ .

**Proposition 5.8.2.** The following are all provable in  $\mathbf{VF}_{\prec}$ .

$$(i) \quad \forall a \in \text{fd}(\prec) \forall x \in \text{St}^{\prec} [\mathbf{Prv}_{\mathbf{PL}_{\mathcal{L}^{\prec}}}(x) \rightarrow T_a^{\prec}(x)].$$

$$(ii) \quad \forall a \in \text{fd}(\prec) \forall x, y \in \text{St}^{\prec}$$

$$[(T_a^{\prec}(x \wedge y) \leftrightarrow (T_a^{\prec}x \wedge T_a^{\prec}y)) \wedge ((T_a^{\prec}x \vee T_a^{\prec}y) \rightarrow T_a^{\prec}(x \vee y))].$$

(iii)  $\forall a \in \text{fd}(\prec) \forall x \forall z$

$$[\forall z. x \in \text{St}^\prec \rightarrow ((\forall y T_a^\prec x(y) \leftrightarrow T_a^\prec \forall z. x) \wedge (\exists y T_a^\prec x(y) \rightarrow T_a^\prec \exists z. x))].$$

(iv)  $\forall a \in \text{fd}(\prec) \forall x \in \text{St}^\prec [T_a^\prec x \leftrightarrow T_a^\prec (\neg \neg x)]$ .

(v)  $\forall x \in \text{For}^\prec \forall z \in \text{CT} [T_a^\prec \text{sb}(x; z) \leftrightarrow T_a^\prec \text{sb}(x; \text{nm} \circ \text{val}(z))]$ .

(vi) For each  $\mathcal{L}^\prec$ -formula  $\phi_x(\vec{z})$  typed by  $x$  (see §6 for its definition),

$$\forall a \in \text{fd}(\prec) \forall b, \vec{z} \in \text{CT} [b^\circ \prec a \rightarrow (T_a^\prec(\ulcorner \phi_b(\vec{z})^\urcorner) \leftrightarrow \phi_{b^\circ}(\vec{z}^\circ)) \wedge (F_a^\prec(\ulcorner \phi_b(\vec{z})^\urcorner) \leftrightarrow \neg \phi_{b^\circ}(\vec{z}^\circ))].$$

We notice that (i) says ‘every logically valid sentence is true at any level’ and (v) expresses ‘the truth of a formula  $A(t)$  depends only on the value of the closed term  $t$  (at any level)’.

*Proof.* We can show (i)-(v) in a parallel manner to Proposition 2.1 of [6]. For (vi), we show the both conjuncts simultaneously by meta-induction on  $\phi$ ; one direction immediately follows from **V1** $_\prec$  using (v); the converse is shown by **V2** $_\prec$  and **V3** $_\prec$  for the base step and by (ii) and (iii) for the induction steps.  $\square$

In what follows, we first embed the system  $\text{ID}_\prec$  of iterated inductive definitions along  $\prec$  in  $\text{VF}_\prec$  and then embed  $\text{VF}_\prec$  in the system  $\Pi_1^1\text{-CA}_\prec + (\text{Bi})$  of iterated  $\Pi_1^1$ -comprehensions along  $\prec$  plus the Bar induction (Bi) (with the proviso that  $\prec$  has no end point); thus we have  $\text{ID}_\prec \subset \text{VF}_\prec \subset \Pi_1^1\text{-CA}_\prec + (\text{Bi})$ . Fix a natural notation system up to a sufficiently large ordinal suitable for ordinal analyses of impredicative systems. Then, together with already known results on  $\text{ID}_\alpha$  and  $\Pi_1^1\text{-CA}_\alpha$  (and also  $\text{KPI}_\alpha$ ) [3, 12, 65], these syntactical embeddings entail that  $\text{VF}_\alpha$  is proof-theoretically equivalent to  $\text{ID}_\alpha$  and  $\Pi_1^1\text{-CA}_\alpha + (\text{Bi})$  for a limit  $\alpha$ ; in the same ordinal notation as in [3] and [13], we also have

$$|\text{VF}_\alpha| = \Theta \varepsilon_{\Omega_\alpha + 1} 0 \quad (\text{for a limit } \alpha) \quad \text{and} \quad |\text{Aut}(\text{VF})| = \Theta \Omega_{\Omega_1} 0.$$

As in the preceding sections, we drop the superscript  $\prec$  in what follows for readability whenever there is no worry of confusion.

### 5.8.2 Lower Bound

Kahle [45] gave a syntactical embedding of  $\text{ID}_1$  in  $\text{VF}$ . In the present subsection, we generalize it to an embedding of  $\text{ID}_\prec$  in  $\text{VF}_\prec$ .

Given a p.r. ordering  $\prec$ , the language  $\mathcal{L}_{\text{ID}}$  of  $\text{ID}_\prec$  is  $\mathcal{L}_0$  plus new unary predicates  $J^{\mathcal{A}}$  associated to each inductive operator form  $\mathcal{A}$ . We write  $J_a^{\mathcal{A}}(b)$  and  $J_{\prec a}^{\mathcal{A}}(b)$  for  $J^{\mathcal{A}}(\langle b, a \rangle)$  and  $(b)_1 \prec a \wedge J^{\mathcal{A}}(b)$  respectively. The theory  $\text{ID}_\prec$  consists of  $\text{PA}(U)$ , full induction for  $\mathcal{L}_{\text{ID}}$ ,  $\text{TI}_{\mathcal{L}_{\text{ID}}}(\prec)$  and the axiom schemata of inductive definitions:

$$(\text{ID}_\prec^1) \quad \forall a \in \text{fd}(\prec) \forall x [\mathcal{A}(J_a^{\mathcal{A}}, J_{\prec a}^{\mathcal{A}}, x, a) \rightarrow J_a^{\mathcal{A}}(x)].$$

$$(\text{ID}_\prec^2) \quad \forall a \in \text{fd}(\prec) [\forall x (\mathcal{A}(\psi, J_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \psi(x)) \rightarrow \forall x (x \in J_a^{\mathcal{A}} \rightarrow \psi(x))],$$

for arbitrary  $\mathcal{L}_{\text{ID}}$ -formula  $\psi$ .

Let  $\mathcal{A}(R, Q, x, y)$  be an inductive operator form, and let  $\psi(x)$ ,  $\theta(x)$ ,  $\chi(x)$  be  $\mathcal{L}^\prec$ -formulae possibly with other parameters  $\vec{v}$ . By straightforward meta-induction on  $\mathcal{A}$ , we can show:

$$\text{PL}_{\mathcal{L}^\prec} \vdash \forall \vec{v} \forall y \forall x [(\mathcal{A}(\psi, \theta, x, y) \wedge \forall z (\psi(z) \rightarrow \chi(z))) \rightarrow \mathcal{A}(\chi, \theta, x, y)]. \quad (5.12)$$

**Lemma 5.8.3.** Let  $\theta_x(z, \vec{v})$  be a  $\mathcal{L}^\prec$ -formula typed by  $x$  with a distinguished variable  $z$  and parameters  $\vec{v}$  and let  $\theta_{\prec x}(z, \vec{v})$  denote  $(z)_1 \prec x \wedge \theta_{(z)_1}(z, \vec{v})$ . Then,  $\text{VF}_\prec$  derives:

$$\forall \vec{v} \forall a \in \text{fd}(\prec) \forall z [(T_a \ulcorner \theta_{\prec a}(\dot{z}, \vec{v}) \urcorner \leftrightarrow \theta_{\prec a}(z, \vec{v})) \wedge (F_a \ulcorner \theta_{\prec a}(\dot{z}, \vec{v}) \urcorner \leftrightarrow \neg \theta_{\prec a}(z, \vec{v}))]$$

*Proof.* Fix  $a \in \text{fd}(\prec)$  and  $\vec{v}$ . For both conjuncts, the one direction immediately follows from  $\mathbf{V1}_\prec$  and the converse follows from  $\mathbf{V2}_\prec$  and Proposition 5.8.2 (ii) and (vi).  $\square$

**Lemma 5.8.4.** Let  $\mathcal{A}(R, Q, x, y)$  be an inductive operator form and let  $\psi(x, z)$  and  $\theta(x, z)$  be  $\mathcal{L}^<$ -formulae possibly with parameters  $\vec{v}$ . Suppose  $\mathbf{VF}_{<}$  derives:

$$\forall \vec{v} \forall a \in \text{fd}(<) \forall z [(T_a \ulcorner \theta(\dot{a}, \dot{z}, \vec{v}) \urcorner \leftrightarrow \theta(a, z, \vec{v})) \wedge (F_a \ulcorner \theta(\dot{a}, \dot{z}, \vec{v}) \urcorner \leftrightarrow \neg \theta(a, z, \vec{v}))]. \quad (5.13)$$

We write  $\psi_{a, \vec{v}}(z)$  and  $\theta_{a, \vec{v}}(z)$  for  $\psi(a, z, \vec{v})$  and  $\theta(a, z, \vec{v})$  respectively. Then,

$$\mathbf{VF}_{<} \vdash \forall \vec{v} \forall a \in \text{fd}(<) \forall z [\mathcal{A}(\lambda u. T_a \ulcorner \psi_{\dot{a}, \vec{v}}(\dot{u}) \urcorner, \theta_{a, \vec{v}}, z, a) \rightarrow T_a \ulcorner \mathcal{A}(\psi_{\dot{a}, \vec{v}}, \theta_{\dot{a}, \vec{v}}, \dot{z}, \dot{a}) \urcorner].$$

*Proof.* Meta-induction on  $\mathcal{A}$ ; use  $\mathbf{V2}_{<}$ , (5.13) and Proposition 5.8.2. □

We introduce two notations which generalize  $\text{Clos}_{\mathcal{A}}(\psi)$  and  $\psi'$  in [45]:

$$\text{Clos}_{\mathcal{A}}[\psi; \theta] := \forall x [\mathcal{A}(\psi, \theta, x, y) \rightarrow \psi(x)]$$

$$\psi'_{\theta}(x) := \text{Clos}_{\mathcal{A}}[\psi; \theta] \rightarrow \psi(x)$$

**Lemma 5.8.5.** Let  $\mathcal{A}(R, Q, x, y)$  be an inductive operator form and  $\psi(x)$  and  $\theta(x)$  be arbitrary  $\mathcal{L}^<$ -formulae possibly with parameters  $y$  and  $\vec{v}$ . Then,  $\text{PL}_{\mathcal{L}^<} \vdash (\forall y)(\forall \vec{v}) \text{Clos}_{\mathcal{A}}[\psi'_{\theta}; \theta]$ .

*Proof.* A straightforward generalization of Lemma 33 of [45]. We have to show that

$$\forall y \forall \vec{v} \forall x [\mathcal{A}(\psi'_{\theta}, \theta, x, y) \rightarrow (\text{Clos}_{\mathcal{A}}[\psi; \theta] \rightarrow \psi(x))]$$

Fix the parameters. Take any  $x$  and suppose  $\mathcal{A}(\psi'_{\theta}, \theta, x, y)$  and  $\text{Clos}_{\mathcal{A}}[\psi; \theta]$ . By  $\text{Clos}_{\mathcal{A}}[\psi; \theta]$ , we have  $\forall z [\psi'_{\theta}(z) \rightarrow \psi(z)]$ . Then, it follows from (5.12) that  $\mathcal{A}(\psi'_{\theta}, \theta, x, a) \rightarrow \mathcal{A}(\psi, \theta, x, y)$ ; and thus  $\mathcal{A}(\psi, \theta, x, y)$ . Again by  $\text{Clos}_{\mathcal{A}}[\psi; \theta]$ , we finally obtain  $\psi(x)$ . □

In what follows, we suppress all irrelevant parameters for simplicity.

**Lemma 5.8.6.** Let  $\mathcal{A}(R, Q, x, y)$  be an inductive operator form and  $\theta_z(x)$  be the same as Lemma 5.8.4. Then,  $\mathbf{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) \text{Clos}_{\mathcal{A}}[T_a \ulcorner \psi'_{\theta_a} \urcorner; \theta_a]$ , for any  $\mathcal{L}^{\prec}$ -formula  $\psi$ .

*Proof.* It follows from the last lemma and Proposition 5.8.2 (i) that

$$\mathbf{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) T_a \ulcorner \text{Clos}_{\mathcal{A}}[\psi'_{\theta_a}; \theta_a] \urcorner; \text{i.e.,}$$

$$\mathbf{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) T_a \ulcorner \forall x (\mathcal{A}(\psi'_{\theta_a}, \theta_a, x, \dot{a}) \rightarrow \psi'_{\theta_a}(x)) \urcorner; \text{by Prop 5.8.2(iii) and } \mathbf{V6}_{\prec},$$

$$\mathbf{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) \forall x [T_a \ulcorner \mathcal{A}(\psi'_{\theta_a}, \theta_a, \dot{x}, \dot{a}) \urcorner \rightarrow T_a \ulcorner \psi'_{\theta_a}(\dot{x}) \urcorner].$$

We also have  $\forall a \in \text{fd}(\prec) \forall x [\mathcal{A}(T_a \ulcorner \psi'_{\theta_a} \urcorner, \theta_a, x, a) \rightarrow T_a \ulcorner \mathcal{A}(\psi'_{\theta_a}, \theta_a, \dot{x}, \dot{a}) \urcorner]$  by Lemma 5.8.4. Combining these, we obtain  $\forall a \in \text{fd}(\prec) \forall x [\mathcal{A}(T_a \ulcorner \psi'_{\theta_a} \urcorner, \theta_a, x, a) \rightarrow T_a \ulcorner \psi'_{\theta_a}(\dot{x}) \urcorner]$ .  $\square$

Given an inductive operator form  $\mathcal{A}(R, Q, x, y)$ , let  $\Phi^{\mathcal{A}}(x, y, w)$  denote:

$$\forall z \in \text{For}^{\prec} \left( \text{Clos}_{\mathcal{A}}[T_y z(\cdot); \lambda u. ((u)_1 \prec y \wedge T_{(u)_1} w((u)_0, (u)_1))] \rightarrow T_y z(x) \right).$$

By diagonalization (on  $w$ ), there is a closed term  $t$  such that  $t = \ulcorner \Phi^{\mathcal{A}}(x, y, t) \urcorner$  and thus  $\Phi^{\mathcal{A}}(x, y, t) \leftrightarrow \Phi^{\mathcal{A}}(x, y, \ulcorner \Phi^{\mathcal{A}}(x, y, t) \urcorner)$ . Set  $S^{\mathcal{A}}(x, y) \equiv \Phi^{\mathcal{A}}(x, y, t)$  and  $P^{\mathcal{A}}(x, y) \equiv T_y(\ulcorner S^{\mathcal{A}}(\dot{x}, \dot{y}) \urcorner)$ . Let  $\Psi^{\mathcal{A}}(u, y)$  be  $T_y t((u)_0, (u)_1)$ . Then  $\Psi^{\mathcal{A}}(u, y)$  is typed by  $y$  and  $\Psi^{\mathcal{A}}_{\prec y}(u) \equiv (u)_1 \prec y \wedge T_{(u)_1} t((u)_0, (u)_1)$  in the notation of Lemma 5.8.3. By definition, we have

$$S^{\mathcal{A}}(x, y) \equiv \forall z \in \text{For}^{\prec} (\text{Clos}_{\mathcal{A}}[T_y z(\cdot); \lambda u. \Psi^{\mathcal{A}}_{\prec y}(u)] \rightarrow T_y z(x)) \quad (5.14)$$

$$P^{\mathcal{A}}(x, y) \equiv T_y \left( \ulcorner \forall z \in \text{For}^{\prec} (\text{Clos}_{\mathcal{A}}[T_{\dot{y}} z(\cdot); \lambda u. \Psi^{\mathcal{A}}_{\prec \dot{y}}(u)] \rightarrow T_{\dot{y}} z(\dot{x})) \urcorner \right). \quad (5.15)$$

We define an embedding  $*$  of  $\text{ID}_{\prec}$  in  $\text{VF}_{\prec}$  by  $J^{\mathcal{A}}(x) \mapsto P^{\mathcal{A}}((x)_0, (x)_1)$ . Then,

$$(J_{\prec y}^{\mathcal{A}}x)^* \equiv (x)_1 \prec y \wedge P^{\mathcal{A}}((x)_0, (x)_1) \quad (5.16)$$

$$\equiv (x)_1 \prec y \wedge T_{(x)_1} \ulcorner S^{\mathcal{A}}((x)_0, (x)_1) \urcorner \quad (5.17)$$

$$\leftrightarrow (x)_1 \prec y \wedge T_{(x)_1} t((x)_0, (x)_1) \urcorner \equiv \Psi_{\prec y}^{\mathcal{A}}(x). \quad (5.18)$$

Let us write  $P_{\prec y}^{\mathcal{A}}x$  for  $(x)_1 \prec y \wedge P^{\mathcal{A}}((x)_0, (x)_1)$ , i.e.,  $(J_{\prec y}^{\mathcal{A}}x)^*$ . By (5.16), we have

$$\mathcal{A}^*(J_y^{\mathcal{A}}, J_{\prec y}^{\mathcal{A}}, x, y) \equiv \mathcal{A}(P_y^{\mathcal{A}}, P_{\prec y}^{\mathcal{A}}, x, y) \leftrightarrow \mathcal{A}(P_y^{\mathcal{A}}, \Psi_{\prec y}^{\mathcal{A}}, x, y) \quad (5.19)$$

The following two lemmata will establish that the translation  $*$  is indeed an embedding.

**Lemma 5.8.7.**  $\text{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) \forall x [\mathcal{A}(P_a^{\mathcal{A}}, P_{\prec a}^{\mathcal{A}}, x, a) \rightarrow P_a^{\mathcal{A}}(x)]$ , for each inductive operator form  $\mathcal{A}(R, Q, x, y)$ .

*Proof.* Fix arbitrary  $a \in \text{fd}(\prec)$ . By logic, it follows from (5.14) that

$$\text{PL}_{\mathcal{L}^{\prec}} \vdash \forall x \forall z \in \text{For}^{\prec} [S_a^{\mathcal{A}}(x) \rightarrow (\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}] \rightarrow T_a z(x))]; \text{ and thus}$$

$$\text{PL}_{\mathcal{L}^{\prec}} \vdash \forall x \forall z \in \text{For}^{\prec} [\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}] \rightarrow (S_a^{\mathcal{A}}(x) \rightarrow T_a z(x))].$$

Hence, it follows from (5.12) that, in  $\text{PL}_{\mathcal{L}^{\prec}}$ ,

$$\forall x \forall z \in \text{For}^{\prec} [\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}] \rightarrow (\mathcal{A}(S_a^{\mathcal{A}}, \Psi_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \mathcal{A}(T_a z(\cdot), \Psi_{\prec a}^{\mathcal{A}}, x, a))].$$

Since  $z$  is not free in  $\mathcal{A}(S_a^{\mathcal{A}}, \Psi_{\prec a}^{\mathcal{A}}, x, a)$ , we equivalently have: in  $\text{PL}_{\mathcal{L}^{\prec}}$ ,

$$\forall x [\mathcal{A}(S_a^{\mathcal{A}}, \Psi_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \forall z \in \text{For}^{\prec} (\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}] \rightarrow \mathcal{A}(T_a z(\cdot), \Psi_{\prec a}^{\mathcal{A}}, x, a))].$$

Recall that  $\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}]$  denotes  $\forall x[\mathcal{A}(T_a z(\cdot), \Psi_{\prec a}^{\mathcal{A}}, x, a) \rightarrow T_a z(x)]$ ; thus,

$$\text{PL}_{\mathcal{L}^{\prec}} \vdash \forall x[\mathcal{A}(S_a^{\mathcal{A}}, \Psi_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \forall z \in \text{For}^{\prec}(\text{Clos}_{\mathcal{A}}[T_a z(\cdot); \Psi_{\prec a}^{\mathcal{A}}] \rightarrow T_a z(x))].$$

Since all these have been proved in  $\text{PL}_{\mathcal{L}^{\prec}}$ ,  $\text{VF}_{\prec}$  proves that, for all  $x$ ,

$$T_a \ulcorner \mathcal{A}(S_{\dot{a}}^{\mathcal{A}}, \Psi_{\prec \dot{a}}^{\mathcal{A}}, \dot{x}, \dot{a}) \urcorner \rightarrow T_a \ulcorner \forall z \in \text{For}^{\prec}(\text{Clos}_{\mathcal{A}}[T_{\dot{a}} z(\cdot); \Psi_{\prec \dot{a}}^{\mathcal{A}}] \rightarrow T_{\dot{a}} z(\dot{x})) \urcorner. \quad (5.20)$$

On the other hand, since  $\Psi_a^{\mathcal{A}}$  is typed, it follows from Lemmata 5.8.3 and 5.8.4,

$$\mathcal{A}(P_a^{\mathcal{A}}, \Psi_{\prec a}^{\mathcal{A}}, x, a) \equiv \mathcal{A}(T_a \ulcorner S_{\dot{a}}^{\mathcal{A}} \urcorner, \Psi_{\prec a}^{\mathcal{A}}, x, a) \rightarrow T_a \ulcorner \mathcal{A}(S_{\dot{a}}^{\mathcal{A}}, \Psi_{\prec \dot{a}}^{\mathcal{A}}, \dot{x}, \dot{a}) \urcorner. \quad (5.21)$$

Thus, our claim finally follows from (5.19), (5.21), (5.20) and (5.14).  $\square$

**Lemma 5.8.8.** For each inductive operator  $\mathcal{A}(R, Q, x, y)$  and an arbitrary  $\mathcal{L}^{\prec}$ -formula  $\psi(x)$ ,

$$\text{VF}_{\prec} \vdash \forall a \in \text{fd}(\prec) [\forall x(\mathcal{A}(\psi, P_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \psi(x)) \rightarrow \forall y(P_a^{\mathcal{A}}(y) \rightarrow \psi(y))].$$

*Proof.* Fix  $a \in \text{fd}(\prec)$  and take any  $y$ . Suppose  $P_a^{\mathcal{A}}(y)$ . By  $\mathbf{V1}_{\prec}$ , we obtain  $S_a^{\mathcal{A}}(y)$  and thus  $\forall z \in \text{For}^{\prec}(\text{Clos}_{\mathcal{A}}[T_a z(\cdot); P_{\prec a}^{\mathcal{A}}] \rightarrow T_a z(y))$  by (5.16). Let  $\phi_a$  be  $\psi'_{P_{\prec a}^{\mathcal{A}}}$ ; then we have

$$\text{Clos}_{\mathcal{A}}[T_a(\ulcorner \phi_{\dot{a}}(\cdot) \urcorner); P_{\prec a}^{\mathcal{A}}] \rightarrow T_a(\ulcorner \phi_{\dot{a}}(\dot{y}) \urcorner). \quad (5.22)$$

Since  $P_a^{\mathcal{A}}$  is typed by  $a$ , we already have  $\text{Clos}_{\mathcal{A}}[T_a(\ulcorner \phi_{\dot{a}}(\cdot) \urcorner); P_{\prec a}^{\mathcal{A}}]$  by Lemmata 5.8.3 and 5.8.6. Then, (5.22) entails  $T_a(\ulcorner \phi_{\dot{a}}(\dot{y}) \urcorner)$  and thus  $\phi_a(y)$  again by  $\mathbf{V1}_{\prec}$ . Finally,  $\psi(y)$  follows under the condition  $\forall x(\mathcal{A}(\psi, P_{\prec a}^{\mathcal{A}}, x, a) \rightarrow \psi(x))$  by the definition of  $\phi_a$  ( $\equiv \psi'_{P_{\prec a}^{\mathcal{A}}}$ ).  $\square$

**Theorem 5.8.9.** We can regard  $\text{ID}_{\prec} \subset \text{VF}_{\prec}$ . In particular, the system  $\overline{\text{ID}}$  of autonomously iterated

inductive definitions (see [3]) is embeddable in  $\text{Aut}(\text{VF})$ .

### 5.8.3 Upper Bound

In the present subsection, we embed  $\text{VF}_{\prec}$  in  $\Pi_1^1\text{-CA}_{\prec} + (\text{Bi})$ . Let  $\mathcal{L}^2$  be the second-order language of arithmetic obtained from  $\mathcal{L}_0$ ; we assume that the classes  $\Pi_n^1$  ( $\Sigma_n^1$ ) of  $\mathcal{L}^2$ -formulae are suitably defined so that  $\Pi_0^1 = \Sigma_0^1$  (i.e., the class of ‘arithmetical’ formulae) includes  $\mathcal{L}_0$ . Given a p.r. ordering  $\prec$ , the axiom schema  $\Pi_1^1\text{-CA}_{\prec}$  is defined by:

$$\forall z \in \text{fd}(\prec) \exists X \forall y \prec z \forall x [x \in X^y \leftrightarrow \phi(X^{\prec y}, x, y)],$$

for each  $\Pi_1^1$ -formula  $\phi(Z, x, y)$  only with the displayed parameters, where  $x \in X^y$  and  $x \in X^{\prec y}$  respectively denote  $\langle x, y \rangle \in X$  and  $(x)_1 \prec y \wedge x \in X$ . Then, the system  $\Pi_1^1\text{-CA}_{\prec}$  over  $\mathcal{L}^2$  is a subsystem of second-order arithmetic defined as  $\Pi_1^1\text{-CA}$  plus  $\text{Wf}(\prec) := \forall X [\text{Prog}(\prec; X) \rightarrow \forall x \in \text{fd}(\prec)(x \in X)]$  (‘ $\prec$  is well-founded’) and  $\Pi_1^1\text{-CA}_{\prec}$ . We notice that this definition of  $\Pi_1^1\text{-CA}_{\prec}$  is different from the one in [64].<sup>13</sup> The Bar induction (Bi) is defined in the standard manner; see, [64]. Feferman [13] showed that (Bi) is equivalent, modulo ACA, to the so-called quantifier schema:

**(QS)**  $\forall X \phi(X) \rightarrow \phi[\psi/X]$ , for any arithmetical  $\phi$  and arbitrary  $\psi$ .

A  $\Pi_1^1$ -formula  $\mathcal{A}(X, Y, x, y)$  possibly with other parameters is called a *monotone*  $\Pi_1^1$  operator (with respect to  $X$ ), when it holds that

$$\forall X \forall X' [X \subset X' \rightarrow (\mathcal{A}(X, Y, x, y) \rightarrow \mathcal{A}(X', Y, x, y))].$$

In particular, when  $\mathcal{A}(X, Y, x, y)$  is  $X$ -positive arithmetical, it is monotone  $\Pi_1^1$ .

---

<sup>13</sup>It is observed from Theorem 3.2.4.2 and the proof of Lemma 3.3.5.2 of [64] that  $\text{KPI}_{\alpha} \vdash \Pi_1^1\text{-CA}_{\alpha}$ . It is also known that  $\text{KPI} \vdash \Pi_1^1\text{-CA} + (\text{Bi})$ ; their proofs are found in [64]. Hence,  $\Pi_1^1\text{-CA}_{\alpha} + (\text{Bi})$  in our definition can be embedded in  $\text{KPI}_{\alpha}$ . Then, the proof-theoretic equivalence of  $\text{ID}_{\alpha}$  and  $\text{KPI}_{\alpha}$  (Theorem 3.4.6.18 of [64]) finally yields our claim.

Given a monotone  $\Pi_1^1$  operator  $\mathcal{A}(X, Y, x, y)$ , define

$$\text{Clos}_{\mathcal{A}}(X, Y, y) \equiv (\forall x)[\mathcal{A}(X, Y, x, y) \rightarrow x \in X].$$

$$H_{\mathcal{A}}(Y, x, y) \equiv (\forall Z)[\text{Clos}_{\mathcal{A}}(Z, Y, y) \rightarrow x \in Z]$$

Informally,  $H_{\mathcal{A}}(Y, x, y)$  expresses  $x \in \bigcap \{Z \mid \text{Clos}_{\mathcal{A}}(Z, Y, y)\}$ .  $H_{\mathcal{A}}$  doesn't appear to be  $\Pi_1^1$  at a first glance. However, we can show that  $\text{Clos}_{\mathcal{A}}$  is  $\Sigma_1^1$  in  $\Pi_1^1\text{-CA}$  using  $\Sigma_1^1\text{-AC}$  ( $\neg \Pi_1^1\text{-CA}$ ); therefore  $H_{\mathcal{A}}$  is  $\Pi_1^1$  in  $\Pi_1^1\text{-CA}$ . Let  $S$  be the set defined by  $H_{\mathcal{A}}$ . Then, we can show in the standard manner that  $S$  is the least fixed-point of  $\mathcal{A}$ ; see [13] or [64].

When  $\mathcal{A}(X, Y, x, y)$  contains no other parameters, we can apply  $\Pi_1^1\text{-CA}_{\prec}$  to  $H_{\mathcal{A}}$  and obtain:

$$\Pi_1^1\text{-CA}_{\prec} \vdash \forall z \in \text{fd}(\prec)(\exists X)(\forall y \prec z)(\forall x)[x \in X^y \leftrightarrow H_{\mathcal{A}}(X^{\prec y}, x, y)].$$

Given  $z \in \text{fd}(\prec)$ , let  $S$  be a set such that  $(\forall x)(\forall y \prec z)[x \in S^y \leftrightarrow H_{\mathcal{A}}(S^{\prec y}, x, y)]$ . Then, for each  $y \prec z$ , we have

$$(\forall X)[\text{Clos}_{\mathcal{A}}(X, S^{\prec y}, y) \rightarrow S^y \subset X] \quad \text{and} \quad \forall x[\mathcal{A}(S^y, S^{\prec y}, x, y) \rightarrow x \in S^y].$$

These imply  $\forall x[\mathcal{A}(S^y, S^{\prec y}, x, y) \leftrightarrow x \in S^y]$  ( $S^y$  is a fixed-point of  $\mathcal{A}(\cdot, S^{\prec y}, x, y)$ ).

Let us fix a p.r. ordering  $\prec$ . We will construct an embedding of  $\text{VF}_{\prec}$  in  $\Pi_1^1\text{-CA}_{\prec} + (\text{Bi})$ . As was shown by Cantini [6],  $\text{VF}$  can be modelled by a Kripkean fixed-point with van Fraassen's supervaluation schema. Indeed,  $\text{VF}_{\prec}$  can be similarly modelled by iterating such structures. In what follows, we will formalize this argument within  $\Pi_1^1\text{-CA}_{\prec} + (\text{Bi})$ .

Given  $V \subset \mathbb{N}$ , let  $(\mathbb{N}, V)$  be a  $\mathcal{L}^{\prec}$ -structure in which the truth predicate  $T^{\prec}$  is interpreted by  $(V)_0 \times (V)_1$  where  $(V)_i = \{(x)_i \mid x \in V\}$  for  $i = 0, 1$ ; we simply write  $V \models \phi$  for  $(\mathbb{N}, V) \models \phi$  for each  $\mathcal{L}^{\prec}$ -sentence  $\phi$ . It is known that the set  $\{\sigma \in \mathcal{L}^{\prec} \mid (\mathbb{N}, V) \models \sigma\}$  is  $\Pi_1^1$  (indeed  $\Delta_1^1$ ) in  $V$ ; see [6].

This fact can be formalized in  $\Pi_1^1$ -CA. We define an  $Z$ -positive arithmetical formula  $\mathcal{B}(Z, X, Y, x, a)$  as follows:

$$\begin{aligned}
& \bigwedge_{R \in \mathcal{L}_0} \exists \vec{y} \in \text{CT} [(x = R\vec{y} \wedge R(\vec{y}^\circ)) \vee (x = \neg R\vec{y} \wedge \neg R(\vec{y}^\circ))] \\
& \vee \exists b, y \in \text{CT} [b^\circ \prec a \wedge [(x = T_b y \wedge y^\circ \in Y^{b^\circ}) \vee (x = \neg T_b y \wedge y^\circ \notin Y^{b^\circ})]] \\
& \vee \exists b, y \in \text{CT} [b^\circ \not\prec a \wedge [(x = T_b y \wedge y^\circ \in X^{b^\circ}) \vee (x = \neg T_b y \wedge y^\circ \notin X^{b^\circ})]] \\
& \vee \exists y, z \in \text{St}^\prec [(x = y \wedge z \wedge y \in Z \wedge z \in Z) \vee (x = \neg(y \wedge z) \wedge (\neg y \in Z \vee \neg z \in Z))] \\
& \vee \exists y \exists z [(x \in \text{St}^\prec \wedge x = \forall y.z \wedge \forall v(z(v) \in Z)) \\
& \qquad \qquad \qquad \vee (x \in \text{St}^\prec \wedge x = \neg \forall y.z \wedge \exists v(\neg z(v) \in Z))] \\
& \vee \exists y \in \text{St}^\prec (x = \neg \neg y \wedge y \in Z).
\end{aligned}$$

Let  $\text{Tr}(X, Y, x, a)$  be the  $\Pi_1^1$ -formula defining the least fixed-point of  $\mathcal{B}$  (i.e.,  $H_{\mathcal{B}}(X, Y, x, a)$ ). We sometimes write  $x \in \text{Tr}(X, Y, a)$  for  $\text{Tr}(X, Y, x, a)$ . As we shall see below,  $\text{Tr}(X, Y, x, a)$  represents the set  $\{\ulcorner \phi \urcorner \mid (X^{\not\prec a} \cup Y^{\prec a}) \models \phi\}$  as is expected, where  $X^{\not\prec a}$  denotes the set  $\{\langle x, y \rangle \in X \mid y \not\prec a\}$  ( $\supset X^a$ ).

The next is shown by using the fixed-point property of  $\text{Tr}(X, Y, a)$ .

**Proposition 5.8.10.** The following hold in  $\Pi_1^1$ -CA: for all  $X, Y$  and  $a \in \text{fd}(\prec)$ ,

- (1)  $\forall \vec{x} \in \text{CT} [(R\vec{x} \in \text{Tr}(X, Y, a) \leftrightarrow R\vec{x}^\circ) \wedge (\neg R\vec{x} \in \text{Tr}(X, Y, a) \leftrightarrow \neg R\vec{x}^\circ)]$ , for  $R \in \mathcal{L}_0$
- (2)  $\forall x, b \in \text{CT} [b^\circ \prec a \rightarrow [(T_b x \in \text{Tr}(X, Y, a) \leftrightarrow x^\circ \in Y^{b^\circ}) \wedge (\neg T_b x \in \text{Tr}(X, Y, a) \leftrightarrow x^\circ \notin Y^{b^\circ})]]$
- (3)  $\forall x, b \in \text{CT} [b^\circ \not\prec a \rightarrow [(T_b x \in \text{Tr}(X, Y, a) \leftrightarrow x^\circ \in X^{b^\circ}) \wedge (\neg T_b x \in \text{Tr}(X, Y, a) \leftrightarrow x^\circ \notin X^{b^\circ})]]$
- (4)  $\forall x, y \in \text{St}^\prec [[x \wedge y \in \text{Tr}(X, Y, a) \leftrightarrow (x \in \text{Tr}(X, Y, a) \wedge y \in \text{Tr}(X, Y, a))]$

$$\wedge[\neg(x \wedge y) \in \text{Tr}(X, Y, a) \leftrightarrow (\neg x \in \text{Tr}(X, Y, a) \vee \neg y \in \text{Tr}(X, Y, a))]$$

$$(5) \forall x \forall z [\forall z. x \in \text{St}^{\prec} \rightarrow [\forall z. x \in \text{Tr}(X, Y, a) \leftrightarrow \forall y (x(y) \in \text{Tr}(X, Y, a))]]$$

$$\wedge[\neg \forall z. x \in \text{Tr}(X, Y, a) \leftrightarrow \exists y (\neg x(y) \in \text{Tr}(X, Y, a))]$$

$$(6) \forall x \in \text{St}^{\prec} (\neg \neg x \in \text{Tr}(X, Y, a) \leftrightarrow x \in \text{Tr}(X, Y, a))$$

$$(7) \forall x \in \text{St}^{\prec} (\neg x \in \text{Tr}(X, Y, a) \leftrightarrow x \notin \text{Tr}(X, Y, a)); \text{ by induction on } x \text{ using (1)-(6)}.$$

For an  $\mathcal{L}^{\prec}$ -formula  $\phi$  and a set  $M$ ,  $\phi^M$  will denote the  $\mathcal{L}^2$ -formula obtained from  $\phi$  by replacing each occurrence of  $T_z^{\prec}(x)$  by  $x \in M^z$ . Then the next proposition follows from the last one by meta-induction on  $\phi$ .

**Proposition 5.8.11.** For each  $\mathcal{L}^{\prec}$ -formula  $\phi(\vec{x})$ , we have

$$\Pi_1^1\text{-CA} \vdash \forall X \forall Y \forall a \forall \vec{x} [\ulcorner \phi(\vec{x}) \urcorner \in \text{Tr}(X, Y, a) \leftrightarrow \phi^{(X^{\neq a} \cup Y^{\prec a})}(\vec{x})].$$

We can formalize the soundness of  $\text{PL}_{\mathcal{L}^{\prec}}$  (i.e., soundness of classical logic):

**Proposition 5.8.12.**  $\Pi_1^1\text{-CA} \vdash \forall X \forall Y \forall a \forall x \in \text{St}^{\prec} (\text{Ax}_{\text{PL}_{\mathcal{L}^{\prec}}}(x) \rightarrow x \in \text{Tr}(X, Y, a))$ .

Define an operator  $\text{FV}(X, Y, x, a)$  (on  $x$  with the parameters  $X, Y$  and  $a$ ) by

$$\text{FV}(X, Y, x, a) \text{ :} \equiv \forall Z [X \subset Z^a \rightarrow x \in \text{Tr}(Z, Y, a)].$$

$\text{FV}(X, Y, x, a)$  informally expresses the set

$$\{\ulcorner \phi \urcorner \mid \phi \in \mathcal{L}^{\prec} \text{ and } (Z^{\neq a} \cup Y^{\prec a}) \models \phi \text{ for all } Z \text{ s.t. } Z^a \supset X \}.$$

Since  $\text{Tr}$  is  $\Pi_1^1$  and  $\text{FV}(X, Y, x, a)$  is  $X$ -positive,  $\text{FV}(X, Y, x, a)$  is a monotone  $\Pi_1^1$  operator (w.r.t.  $X$ ). Hence, we can define the least fixed-point  $H_{\text{FV}}(Y, x, a)$  of  $\text{FV}$ :

$$\begin{aligned}\text{Clos}_{\text{FV}}(X, Y, a) &::= (\forall x)[\text{FV}(X, Y, x, a) \rightarrow x \in X] \\ H_{\text{FV}}(Y, x, a) &::= (\forall Z)[\text{Clos}_{\text{FV}}(Z, Y, a) \rightarrow x \in Z].\end{aligned}$$

Recall that  $\text{Clos}_{\text{FV}}(X, Y, a)$  is  $\Sigma_1^1$  and  $H_{\text{FV}}(Y, x, a)$  is  $\Pi_1^1$ . In what follows, we assume that  $\prec$  has no end point, i.e.,  $\forall a \in \text{fd}(\prec)\exists b \in \text{fd}(\prec)(a \prec b)$ ; therefore, if  $\prec$  is well-ordered then  $\text{otyp}(\prec)$  is a limit ordinal. Then, we have

$$\Pi_1^1\text{-CA}_{\prec} \vdash \forall a \in \text{fd}(\prec)\exists M\forall b \preceq a\forall x[x \in M^b \leftrightarrow H_{\text{FV}}(M^{\prec b}, x, b)]. \quad (5.23)$$

Now, we define a translation  $*$  from  $\mathcal{L}^{\prec}$  to  $\mathcal{L}^2$  as follows:

$$T_z x \mapsto \exists M[\forall w \preceq z\forall y(y \in M^w \leftrightarrow H_{\text{FV}}(M^{\prec w}, y, w)) \wedge x \in M^z].$$

Our goal is to show that  $*$  is indeed an embedding of  $\text{VF}_{\prec}$  in  $\Pi_1^1\text{-CA}_{\prec} + (\text{Bi})$ .

Let  $\text{VF}_{\prec}(z, M)$  denote a formula  $\forall w \preceq z\forall x[x \in M^w \leftrightarrow H_{\text{FV}}(M^{\prec w}, x, w)]$ ; thus (5.23) is rewritten as  $\Pi_1^1\text{-CA}_{\prec} \vdash \forall a \in \text{fd}(\prec)\exists M\text{VF}_{\prec}(a, M)$ . For simplicity, we make the following abuses of notation: for each  $\mathcal{L}^{\prec}$ -formula  $\phi$ ,

$$\phi^{\prec a}(x) ::= (x)_1 \prec a \wedge \phi((x)_0) \quad \& \quad \phi^{\not\prec a}(x) ::= (x)_1 \not\prec a \wedge \phi((x)_0).$$

**Proposition 5.8.13.** One can straightforwardly verify the following; cf. [13, 64].

$$\Pi_1^1\text{-CA} + \text{Wf}(\prec) \vdash \forall a \in \text{fd}(\prec) \forall M \forall N$$

$$[\text{VF}_{\prec}(a, M) \wedge \text{VF}_{\prec}(a, N) \rightarrow (\forall b \preceq a) \forall x [x \in M^b \leftrightarrow x \in N^b].$$

**Lemma 5.8.14.**  $\Pi_1^1\text{-CA}_{\prec} + (\text{Bi}) \vdash (\mathbf{V1}_{\prec})^*$ .

*Proof.* Suppose  $(T_a \ulcorner \phi(\vec{x}) \urcorner)^*$  for  $a \in \text{fd}(\prec)$  and  $\phi \in \mathcal{L}^{\prec}$ . Take  $N$  with  $\text{VF}_{\prec}(a, N)$ . Then, by Proposition 5.8.13, we have  $\ulcorner \phi(\vec{x}) \urcorner \in N^a$ . Since  $N^a$  is a fixed-point of  $\text{FV}(\cdot, N^{\prec a}, a)$ ,

$$(\forall Z)[N^a \subset Z^a \rightarrow \ulcorner \phi(\vec{x}) \urcorner \in \text{Tr}(Z, N^{\prec a}, a)].$$

By Proposition 5.8.11, we obtain

$$(\forall Z)[N^a \subset Z^a \rightarrow \phi^{(Z^{\prec a} \cup N^{\prec a})}(\vec{x})].$$

Since the part ' $N^a \subset Z^a \rightarrow \phi^{(Z^{\prec a} \cup N^{\prec a})}(\vec{x})$ ' is arithmetical, by (QS), we obtain:

$$N^a \subset ((\lambda y.Ty)^*)^a \rightarrow \phi^{(((\lambda x.Tx)^*)^{\prec a} \cup N^{\prec a})}(\vec{x});$$

here, we are treating  $(\lambda y.Ty)^*$  as if it were a set, but it is in fact a mere abbreviation. Then, since

$N^a = (((\lambda y.Ty)^*)^a)^a$  by Proposition 5.8.13, we obtain

$$\phi^{(((\lambda x.Tx)^*)^{\prec a} \cup N^{\prec a})}(\vec{x}).$$

Since  $((\lambda y.Ty)^*)^{\prec a} \cup N^{\prec a} = (\lambda y.Ty)^*$  by Proposition 5.8.13, we obtain  $\phi^*(\vec{x})$ . □

**Proposition 5.8.15.**  $\Pi_1^1\text{-CA}_{\prec} \vdash (\mathbf{V3}_{\prec})^* \wedge (\mathbf{V2}_{\prec})^*$ .

*Proof.* We only show  $(\mathbf{V3}_{\prec})^*$ .  $(\mathbf{V2}_{\prec})^*$  is shown similarly. Given  $a \in \text{fd}(\prec)$ , let  $N$  be such that  $\text{VF}(a, N)$  and take  $b, x \in \text{CT}$  with  $b^\circ \prec a$ . Suppose  $(T_a T_b x)^*$ . Then, by Proposition 5.8.13, we have  $T_b x \in N^a$ . Since  $N^a$  is a fixed-point of  $\text{FV}(\cdot, N^{\prec a}, a)$ , we have

$$(\forall Z)[N^a \subset Z^a \rightarrow T_b x \in \text{Tr}(Z, N^{\prec a}, a)].$$

In particular, we have  $T_b x \in \text{Tr}(N^a, N^{\prec a}, a)$ . By Proposition 5.8.10-(2), we have  $x^\circ \in N^{b^\circ}$  and thus  $(T_b x^\circ)^*$ . For the converse, suppose  $(T_b x^\circ)^*$ . Then  $x^\circ \in N^{b^\circ}$  by Proposition 5.8.13 and  $\forall Z[T_b x \in \text{Tr}(Z, N^{\prec a}, a)]$  by the definition of  $\text{Tr}$ ; in particular, we have  $T_b x \in \text{FV}(N^a, N^{\prec a}, a)$ . Since  $N^a$  is a fixed-point, we have  $T_b x \in N^a$  and thus  $(T_a T_b x)^*$ . The other conjunct is shown similarly.  $\square$

**Proposition 5.8.16.**  $\Pi_1^1\text{-CA}_{\prec} \vdash (\mathbf{V4}_{\prec})^*$ .

*Proof.* Take  $a \in \text{fd}(\prec)$  and suppose  $\text{Ax}_{\text{PL}_{\prec}}(x)$ . Let  $N$  be such that  $\text{VF}(a, N)$ . It suffices to show that  $x \in N^a$ ; this follows from Proposition 5.8.12.  $\square$

**Proposition 5.8.17.**  $\Pi_1^1\text{-CA}_{\prec} \vdash (\mathbf{V5}_{\prec})^* \wedge (\mathbf{V6}_{\prec})^*$ .

*Proof.* We only show  $(\mathbf{V6}_{\prec})^*$ ;  $(\mathbf{V5}_{\prec})^*$  is similarly shown. Let  $a \in \text{fd}(\prec)$  and  $x \rightarrow y \in \text{St}^{\prec}$ . Suppose  $(T_a x \rightarrow y)^*$  and  $(T_a x)^*$ . Let  $N$  be such that  $\text{VF}(a, N)$ . By Proposition 5.8.13, we have  $x, x \rightarrow y \in N^a$ . Then, since  $N^a$  is a fixed-point of  $\text{FV}$ ,  $(\forall Z)[N^a \subset Z^a \rightarrow x, x \rightarrow y \in \text{Tr}(Z, N^{\prec a}, a)]$ . By Proposition 5.8.10, for all  $Z$  with  $Z^a \supset N^a$ ,

$$x \in \text{Tr}(Z, N^{\prec a}, a) \rightarrow y \in \text{Tr}(Z, N^{\prec a}, a) \quad \text{and} \quad x \in \text{Tr}(Z, N^{\prec a}, a);$$

therefor  $y \in \text{Tr}(Z, N^{\prec a}, a)$ . Hence we have  $\text{FV}(N^a, N^{\prec a}, y, a)$ . Again since  $N^a$  is fixed-point, we obtain  $y \in N^a$ . This implies  $(T_a y)^*$ .  $\square$

**Theorem 5.8.18.** Let  $\prec$  be a p.r. ordering with no end point. Then,  $\text{VF}_\prec$  is syntactically embeddable in  $\Pi_1^1\text{-CA}_\prec + (\text{Bi})$ .

## 5.9 Summary and Concluding Remarks of Chapter 5

Let us summarize the results obtained in the present chapter.

$$\begin{aligned}
|\text{KF}_\alpha| &= |\text{KF}_\alpha + \text{Cons}_\alpha \text{ (or Comp}_\alpha)| = |\text{WKF}_\alpha| = |\text{WKF}_\alpha + \text{Cons}_\alpha \text{ (or Comp}_\alpha)| \\
&= |\text{PUTB}_\alpha| = |\text{PUTB}_\alpha + \text{Cons}_\alpha \text{ (or Comp}_\alpha)| \\
&= |\text{FKF}_\alpha| = |\text{DT}_\alpha (= \text{FKF}_\alpha + \text{Cons}_\alpha)| = |\text{FKF}_\alpha + \text{Comp}_\alpha| = |\widehat{\text{ID}}_\alpha|;
\end{aligned}$$

in fact, they are all proof-theoretically equivalent. We also have

$$\begin{aligned}
|\text{Aut}(\text{WKF})| &= |\text{Aut}(\text{WKF} + \text{Cons})| = |\text{Aut}(\text{WKF} + \text{Comp})| \\
&= |\text{Aut}(\text{PUTB})| = |\text{Aut}(\text{PUTB} + \text{Cons})| = |\text{Aut}(\text{PUTB} + \text{Comp})| \\
&= |\text{Aut}(\text{FKF})| = |\text{Aut}(\text{DT})| = |\text{Aut}(\text{FKF} + \text{Comp})| \\
&= |\text{Aut}(\text{KF} + \text{Cons})| = |\text{Aut}(\text{KF} + \text{Comp})| = |\text{Aut}(\text{KF})| = \varphi_{200};
\end{aligned}$$

they are all proof-theoretically equivalent as well. Also, for a limit  $\alpha$ ,

$$|\text{VF}_\alpha| = |\text{ID}_\alpha| = |\Pi_1^1\text{-CA}_\alpha + \text{Bi}| = \Theta_{\varepsilon_{\Omega_\alpha+1}0} \quad \text{and} \quad |\text{Aut}(\text{VF})| = \Theta_{\Omega_1}0.$$

There are at least two directions for further studies from what we have obtained. First, we may further investigate the iterations and autonomous progressions of other self-applicable truths. Second, as far as I know, the system  $\text{Aut}(\text{VF})$  is the strongest system ever among the so far presented truth systems which have natural and philosophical motivations behind them. However, its proof-

theoretic strength is far below that of, e.g.,  $\Pi_2^1$ -CA. Seeking for a well-motivated and stronger system of self-applicable truth might be of some significance from the foundational point of view.

## Chapter 6

# Further Prospects

I conclude my thesis by making a very brief comment on the results and discussions so far developed.

I have given proof-theoretic analyses to a variety of axiomatic systems of truth and compared them by some different means. I dare not deduce from them a hasty conclusion about which system is the best or most adequate system of truth. The research I have so far conducted is not of the kind by which one can directly and immediately deduce such a conclusion. Rather, as I stated at the end of Chapter 1, I would like to suggest that we now need to conduct tedious but down-to-earth ‘field work’ of examining various systems of truth from various points of view. As long as we have decided to renounce our very naïve conception of truth, we have lost general guiding principles for constructing systems of truth and it seems to me unlikely that we can obtain the answer to ‘what is *the* system of truth’ merely by the traditional kind of philosophical speculations. Nonetheless, I still hope and would like to claim that my results have provided philosophers with more information on the respective notions of truth and on how they are related and compared, by which one might be able to select some systems as good and sieve out others as bad for their purpose.

The axiomatic study of truth is relatively a young area of research, which has been attracting more and more interest from philosophers and logicians in the last two decades. I expect that the

axiomatic study of truth will open a new rich and fruitful scientific field of study and I hope that my research in the present thesis will contribute to it. Instead of drawing a definite conclusion, I raise some further prospects and perspectives for future studies of axiomatic truth in the rest of this short chapter.

Since Feferman [17] and Friedman and Sheard [23], many subsequent researchers have been following the same arithmetical setting that they adopted and the study of axiomatic truths has so far been concentrated mainly on those over arithmetic; the present thesis of mine also follows this tradition and setting.<sup>1</sup> Accordingly, axiomatic systems of truth have been a target of traditional proof theory and ordinal analysis of arithmetic. However, we can of course consider another setting and adopt a different kind of base system over which truth is formulated. The most natural alternative is probably set theory. It might even be argued that set theory is a more suitable setting for investigating the notion of mathematical truth, since set theory is often considered to be the foundation of mathematics. However, not much attention has so far been paid to this direction of study and, as far as I know, no research paper has been published on this topic. There are some plausible reasons for having put set theory aside from the axiomatic study of truth. First, from a foundational point of view, while set theory is often considered to be a sufficiently rich framework within which we can develop contemporary mathematics, it is sometimes suspected to be too strong and doubted to be necessary as the overall framework for mathematics. For these reasons, some people do not accept strong set theories such as ZF. However, there are still many who accept ZF or its equivalents and take it as the foundation of mathematics, and thus it may well be of a good significance to study axiomatic truths over set theory and to investigate what would be brought about by them. Second, from a practical and technical point of view, the study of axiomatic truth systems over ZF is subsumed in that of the realm below inaccessibles (but above ZF), to which

---

<sup>1</sup>In fact, Feferman [17] suggests possible formulations of truth systems over set theory as well, although his study there is principally devoted to those over arithmetic.

not much attention has been paid and for which not enough proof-theoretic techniques have been developed. As far as the I see, any so far presented (natural) system of truth for arithmetic, when suitably reformulated for ZF, cannot reach the strength of the postulate of the existence of inaccessible cardinals; for, they have a ‘standard’ model over  $\mathbb{N}$  (with some exceptions like Friedman and Sheard’s FS [23, 31]) and thus, when they are formulated over ZF, we can construct their ‘standard’ models over the set  $V_\kappa$  for an inaccessible  $\kappa$  in completely parallel manners. However, some new perspectives and research interests about this realm below inaccessibles have recently been proposed and versatile proof-theoretic techniques have consequently been developed by Jäger and his colleagues [40, 41, 42, 44] which turn out to enable us to analyze axiomatic truth systems for set theory.<sup>2</sup> With these backgrounds, it may well be a good time to start the axiomatic study of truth for set theory.

Second, of course, we may well push forward the axiomatic study of truth over arithmetic under various directions. I suggest two directions among them: I propose to search for and study metapredicative and impredicative systems of truth. As for the former, the notion of metapredicativity was recently introduced by Jäger, Strahm and their colleagues. Roughly speaking, metapredicative systems are those which go beyond the predicative realm but can be analysed by predicative methods without resorting to the typical techniques in impredicative proof theory such as the technique of collapsing functions. Some metapredicative systems of truth have been already presented by in Chapter 5. However, these systems are of iterated truth and we haven’t yet known any natural metapredicative system of non-iterated single truth. I also expect that the study of metapredicative truth will contribute to better understanding of metapredicativity and its connection with impred-

---

<sup>2</sup>Jäger’s initial motivation seems to have come from his analysis of Feferman’s operational set theory OST. He showed in [41] that a certain extension of OST is equivalent to the subsystem  $\text{NBG}_{<E_0}$  of MK (cf. §3.3). The study of class theory along this line is kept vigorously developing by Jäger and other logicians. Besides these recent approaches to class theory by Feferman, Jäger and their colleagues, Sato’s recent works on his system SS [73] bring another perspective on class theory. Sato proposed the system SS as a versatile framework covering a broad range of mathematical logic from bounded arithmetic to set theory even beyond ZF. He showed that SS plus the existence of an inaccessible cardinal (Inacc) is equivalent to ZF (by private communication). In SS, its axiom schemata are restricted to fixed bounded complexities; then if we strengthen them by admitting more complex formulae, the resulting systems together with Inacc become stronger than ZF and correspond to some subsystems of MK.

icativity. As for impredicative truth, we already have Cantini's system VF [6] of truth. We have also seen that the impredicative system  $\text{Aut}(\text{VF})$  is fairly strong and goes beyond  $\Delta_2^1\text{-CA}$ ; to repeat,  $\text{Aut}(\text{VF})$  is the strongest system of truth so far to my knowledge, but it still falls far short of, say,  $\Pi_2^1\text{-CA}$ . According to my view explained in Ch.5-§1, finding a strong but natural system of truth contributes to the justification of richer mathematics. I thus propose to pursue stronger but natural impredicative systems of truth.

# Bibliography

- [1] Jon Barwise and Lawrence Moss. *Vicious Circle*. CSLI Publications, Stanford, 1996.
- [2] Lev Beklemishev. Proof-theoretic analysis by iterated reflection. *Archive for Mathematical Logic*, 42:515–552, 2003.
- [3] Wilfried Buchholz and Wolfram Pohlers. Provable wellorderings of formal theories for transfinitely iterated inductive definitions. *The Journal of Symbolic Logic*, 43:118–125, 1978.
- [4] Andrea Cantini. A note on predicatively reducible theory of iterated elementary induction. *Bollettino Unione Matematica Italiana*, 4-B 6:413–430, 1985.
- [5] Andrea Cantini. Notes on formal theories of truth. *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik*, 35:97–130, 1989.
- [6] Andrea Cantini. A theory of formal truth arithmetically equivalent to ID1. *The Journal of Symbolic Logic*, 55:244–259, 1990.
- [7] Andrea Cantini. *Logical Frameworks for Truth and Abstraction*. Elsevier, Amsterdam, 1996.
- [8] Cezary Cieśliński. Deflationism, conservativeness and maximality. *Journal of Philosophical Logic*, 36:695–705, 2007.
- [9] Solomon Feferman. Reflecting on incompleteness. Handwritten Notes, March 1987, 39 pages.
- [10] Solomon Feferman. Systems of predicative analysis. *The Journal of Symbolic Logic*, 29:1–30, 1964.
- [11] Solomon Feferman. Systems of predicative analysis II: Representations of ordinals. *The Journal of Symbolic Logic*, 33:193–220, 1968.
- [12] Solomon Feferman. Formal theories for transfinite iteration of generalized inductive definitions and some subsystems of analysis. In Akiko Kino, John Myhill, and Richard Vesley, editors, *Intuitionism and Proof Theory*, pages 303–326, Amsterdam, 1970. North-Holland.
- [13] Solomon Feferman. Formal theories for transfinite iteration of generalized inductive definitions and some subsystems of analysis. In Akiko Kino, John Myhill, and Richard Vesley, editors, *Intuitionism and Proof Theory*, pages 303–326, Amsterdam, 1970. North-Holland.
- [14] Solomon Feferman. Iterated inductive fixed-point theories: Application to hancock’s conjecture. In George Metakides, editor, *Patras Logic Symposium*, pages 171–196. North-Holland, Amsterdam, 1982.
- [15] Solomon Feferman. Toward useful type-free theories. *Journal of Symbolic Logic*, 49:75–111, 1984.

- [16] Solomon Feferman. Hilbert’s program relativized: Proof-theoretical and foundational reductions. *The Journal of Symbolic Logic*, 53:364–384, 1988.
- [17] Solomon Feferman. Reflecting on incompleteness. *The Journal of Symbolic Logic*, 56:1–49, 1991.
- [18] Solomon Feferman. Does reductive proof theory have a viable rationale? *Erkenntnis*, 53:63–96, 2000.
- [19] Solomon Feferman. Axioms for determinateness and truth. *The Review of Symbolic Logic*, 1:204–217, 2008.
- [20] Solomon Feferman. Operational set theory and small large cardinals. *Annals of Pure and Applied Logic*, 207:971–979, 2009.
- [21] Hartly Field. Deflating the conservativeness argument. *The Journal of Philosophy*, 96:533–540, 1999.
- [22] Hartly Field. A revenge-immune solution to the semantic paradoxes. *Journal of Philosophical Logic*, 32:139–177, 2003.
- [23] Harvey Friedman and Michael Sheard. An axiomatic approach to self-referential truth. *Annals of Pure and Applied Logic*, 33:1–21, 1987.
- [24] Kentaro Fujimoto. Relative truth definability of axiomatic theories of truth. *The Bulletin of Symbolic Logic*, 16:305–344, 2010.
- [25] Kurt Gödel. Russell’s mathematical logic. In Paul Benacerraf and Hilary Putnam, editors, *Philosophy of Mathematics*, pages 447–469. Cambridge University Press, Cambridge, 1983.
- [26] Anil Gupta and Nuel Belnap. *The Revision Theory of Truth*. MIT Press, Cambridge, Massachusetts, 1993.
- [27] Petr Hájek. *Metamathematics of Fuzzy Logic*. Kluwer, Dordrecht, 1998.
- [28] Petr Hájek, Jeffrey Paris, and John Shepherdson. The liar paradox and fuzzy logic. *The Journal of Symbolic Logic*, 65:339–346, 2000.
- [29] Petr Hájek and Pavel Pudlak. *Metamathematics of First-Order Arithmetic*. Springer, Berlin, 1993.
- [30] Volker Halbach. Axiomatic theories of truth. In Edward Zalta, editor, *Stanford Encyclopedia of Philosophy*. URL: <http://plato.stanford.edu/entries/truth-axiomatic>.
- [31] Volker Halbach. A system of complete and consistent truth. *Notre Dame Journal of Formal Logic*, 35:311–327, 1994.
- [32] Volker Halbach. *Axiomatische Wahrheitstheorien*. Akademie Verlag, Berlin, 1996.
- [33] Volker Halbach. Conservative theories of classical truth. *Studia Logica*, 62:353–370, 1999.
- [34] Volker Halbach. Disquotationalism and infinite conjunctions. *Mind*, 108:1–22, 1999.
- [35] Volker Halbach. How innocent is deflationism? *Synthese*, 126:167–194, 2001.
- [36] Volker Halbach. Reducing compositional to disquotational truth. *The Review of Symbolic Logic*, 2:786–798, 2009.

- [37] Volker Halbach. *Axiomatic Theories of Truth*. Cambridge University Press, 2010.
- [38] Volker Halbach and Leon Horsten. Axiomatizing Kripke’s theory of truth. *The Journal of Symbolic Logic*, 71:677–712, 2006.
- [39] Paul Horwich. The minimalist conception of truth. In Simon Blackburn and Keith Simmons, editors, *Truth*, pages 239–263. Oxford University Press, Oxford, 1999.
- [40] Gerhard Jäger. On Feferman’s operational set theory OST. *Annals of Pure and Applied Logic*, 150:19–39, 2007.
- [41] Gerhard Jäger. Full operational set theory with unbounded existential quantification and power set. *Annals of Pure and Applied Logic*, 160, 2009.
- [42] Gerhard Jäger. Operations, sets and classes. In Clark Glymour, Wei Wang, and Dag Westerstahl, editors, *Logic, Methodology, and Philosophy of Science: Proceedings of the Thirteenth International Congress*. College Publications, 2009.
- [43] Gerhard Jäger, Reinhard Kahle, Anton Setzer, and Thomas Strahm. The proof-theoretic analysis of transfinitely iterated fixed point theories. *The Journal of Symbolic Logic*, 64:53–67, 1999.
- [44] Gerhard Jäger and Jürg Krähenbühl.  $\Sigma_1^1$  choice in a theory of sets and classes. In Ralf Schindler, editor, *Ways of Proof Theory*, pages 283–314. Ontos Verlag, Frankfurt, 2010.
- [45] Reinhard Kahle. Truth in applicative theories. *Studia Logica*, 68:103–128, 2001.
- [46] Reinhard Kahle. Universes over frege structures. *Annals of Pure and Applied Logic*, 119:191–223, 2003.
- [47] Richard Kaye. *Models of Peano Arithmetic*. Clarendon Press, Oxford, 1991.
- [48] Jeffrey Ketland. Deflationism and tarski’s paradise. *Mind*, 108:69–94, 1999.
- [49] Henryk Kotlarski, Stanislaw Krajewski, and Alistair Lachlan. Construction of satisfaction classes for nonstandard models. *Canadian Mathematical Bulletin*, 24:283–293, 1981.
- [50] Saul Kripke. Outline of a theory of truth. *Journal of Philosophy*, 72:690–716, 1975.
- [51] Alistair Lachlan. Full satisfaction classes and recursive saturation. *Canadian Mathematical Bulletin*, 24:295–297, 1981.
- [52] Graham Leigh and Michael Rathjen. An ordinal analysis for theories of self-referential truth. *Archive for Mathematical Logic*, 49:213–247, 2010.
- [53] Hannes Leitgeb. What theories of truth should be like (but cannot be). *Philosophy Compass*, 2:276–290, 2007.
- [54] Per Lindström. *Aspects of Incompleteness*. Lecture Notes in Logic 10. A K Peters, Ltd, Massachusetts, second edition, 2003.
- [55] Grzegorz Malinowski. *Many-Valued Logics*. Oxford University Press, Oxford, 1993.
- [56] Robert L. Martin and Peter W. Woodruff. On representing ‘true-in-L’ in L. *Philosophia*, 5:213–217, 1975.

- [57] Vann McGee. How truthlike can a predicate be? *Journal of Philosophical Logic*, 14:399–410, 1985.
- [58] Vann McGee. *Truth, Vagueness, and Paradox*. Hackett, Indianapolis, 1991.
- [59] Vann McGee. Maximal consistent sets of instances of tarski’s schema (t). *Journal of Philosophical Logic*, 21:235–241, 1992.
- [60] Vann McGee. In praise of the free lunch: Why disquotationalists should embrace compositional semantics. In Thomad Bolander, Vincent F. Hendricks, and Stig A Pedersen, editors, *Self-Reference*, pages 95–120. CSLI Publications, Stanford, 2006.
- [61] Sheard Michael. Truth, provability and naive criteria. In Volker Halbach and Leon Horsten, editors, *Principles of truth*, pages 169–181. Ontos Verlag, Frankfurt, 2002.
- [62] Andrzej Mostowski. Some impredicative definitions in the axiomatic set-theory. *Fundamenta Mathematicae*, 37:87–110, 1950.
- [63] Karl-Georg Niebergall. On the logic of reducibility: Axioms and examples. *Erkenntnis*, 53:27–61, 2000.
- [64] Wolfram Pohlers. Subsystems of set theory and second order number theory. In Samuel Buss, editor, *Handbook of Proof Theory*, pages 209–336. Elsevier, Amsterdam, 1998.
- [65] Wolfram Pohlers. Subsystems of set theory and second order number theory. In Samuel Buss, editor, *Handbook of Proof Theory*, pages 209–336. Elsevier, Amsterdam, 1998.
- [66] Wolfram Pohlers. *Proof Theory*. Springer, Berlin, 2009.
- [67] Graham Priest. Paraconsistent logic. In Dov Gabbay and Franz Guenther, editors, *Handbook of Philosophical Logic*, volume 6, pages 287–393. Kluwer, Dordrecht, 2002.
- [68] Michael Rathjen. The role of parameters in bar rule and bar induction. *The Journal of Symbolic Logic*, 56:715–730, 1991.
- [69] Michael Rathjen. The realm of ordinal analysis. In S. Barry Cooper and John K. Truss, editors, *Sets and Proofs*, pages 219–279. Cambridge University Press, Cambridge, 1999.
- [70] William N Reinhardt. Remarks on significance and meaningful applicability. In Luiz Paulo de Alcantara, editor, *Mathematical Logic and Formal Systems*, pages 227–242. Marcel Dekker, New York, 1985.
- [71] William N Reinhardt. Some remarks on extending and interpreting theories with a partial predicate for truth. *Journal of Philosophical Logic*, 15:219–251, 1986.
- [72] Greg Restall. *An Introduction to Substructural Logic*. Routledge, London, 2000.
- [73] Kentaro Sato. The strength of extensionality II –weak weak set theories without infinity–. To appear in *Annals of Pure and Applied Logic*, 2010.
- [74] Stewart Shapiro. Proof and truth: Through thick and thin. *The Journal of Philosophy*, 95:493–521, 1998.
- [75] Moh Shaw-Kwei. Logical paradoxes for many-valued systems. *The Journal of Symbolic Logic*, 19:37–40, 1954.

- [76] Michael Sheard. A guide to truth predicates in the modern era. *The Journal of Symbolic Logic*, 59:1032–1054, 1994.
- [77] Stephen G. Simpson. *Subsystems of Second Order Arithmetic*. Cambridge University Press, Cambridge, 2009.
- [78] Thomas Strahm. Autonomous fixed point progressions and fixed point transfinite recursion. In Samuel Buss, editor, *Logic Colloquium '98*, volume 13 of *ASL Lecture Notes in Logic*, pages 449–464. A K Peters, 2000.
- [79] Gaisi Takeuti. *Proof Theory*. North-Holland, Amsterdam, second edition, 1987.
- [80] Alfred Tarski. The semantic conception of truth and the foundation of semantics. *Philosophy and Phenomenological Research*, 4:341–376, 1944. Reprinted in: Simon Blackburn and Keith Simmons eds. *Truth*. Oxford University Press, Oxford, 1999: 115–143.
- [81] Alfred Tarski. Truth and proof. *Scientific American*, 220:63–77, 1969.
- [82] Alfred Tarski. The concept of truth in formalized languages. In John Corcoran, editor, *Logic, Semantics, Meta-mathematics*, pages 152–278. Hackett, Indianapolis, second edition, 1983.
- [83] Alfred Tarski and Robert Vaught. Arithmetical extensions of relational systems. *Compositio Mathematica*, 13:81–102, 1956.
- [84] Bas van Fraassen. Presupposition, implication, and self-reference. *The Journal of Philosophy*, 65:136–152, 1968.
- [85] Jouko Väänänen. Second-order logic and foundation of mathematics. *The Bulletin of Symbolic Logic*, 7:504–520, 2001.
- [86] Stephen Yablo. Paradox without self-reference. *Analysis*, 53:251–252, 1993.