

Iterated reflection over full disquotational truth

MARTIN FISCHER, *MCMP, LMU Muenchen.*

E-mail: M.Fischer@lrz.uni-muenchen.de

CARLO NICOLAI, *Department of Philosophy and Religious Studies, Utrecht University, Janskerkhof 13, 3512 BL Utrecht, The Netherlands.*

E-mail: carlo.nicolai6@gmail.com

LEON HORSTEN, *Cotham House, Bristol BS6 6JL, United Kingdom.*

E-mail: Leon.Horsten@bristol.ac.uk

Abstract

Iterated reflection principles have been employed extensively to unfold epistemic commitments that are incurred by accepting a mathematical theory. Recently this has been applied to theories of truth. The idea is to start with a collection of Tarski-biconditionals and arrive by iterated reflection at strong compositional truth theories. In the context of classical logic, it is incoherent to adopt an initial truth theory in which A and ‘ A is true’ are inter-derivable. In this article, we show how in the context of a weaker logic, which we call Basic De Morgan Logic, we can coherently start with such a fully disquotational truth theory and arrive at a strong compositional truth theory by applying a natural uniform reflection principle a finite number of times.

Keywords: Truth, semantic Paradoxes, Basic de Morgan logic, reflection principles.

1 Introduction

In the article we pursue the strategy of iterating reflection principles on a basic truth theory TS_0 that encapsulates, or so we argue, the fundamental building blocks of truth-theoretic reasoning. Its components are: a theory of the objects of truth (syntax theory in our case), basic truth-theoretic principles that enable us to infer $\top\phi\top$ from a sentence ϕ and vice-versa. To remain faithful to these assertability conditions for truth ascriptions, classical logic cannot be used on account of the liar paradox. In particular, one should be prepared not to have the rule of conditionalization for all sentences. Also, it is desirable to work at the right level of generality by, for instance, reasoning without committing oneself to paracomplete or paraconsistent options. We employ a logic that does this and call it, following [9], *Basic De Morgan logic* (cf. §2 for the definition).

The theory TS_0 , however, may not be all there is to truth. Many authors have discussed further desiderata for truth — such as full compositionality — that are out of reach for our basic theory TS_0 . Recently, Leon Horsten and Graham Leigh have studied iterations of reflection over a basic theory in classical logic [18]. Their classical starting point, however, leads to a loss of the intimate connection between truths and their assertability that is present in TS_0 . In the following we therefore extend Horsten and Leigh’s strategy in the framework of Basic De Morgan logic.

Reflection is rooted in the fundamental intuition that we are committed to the truth of the sentences that are provable in a theory that we accept. This operation can be expressed in different ways: as *global reflection*, where we make explicit use of the truth predicate, or as *uniform reflection*, where

we express this intuition schematically without direct reference to truth. In the classical case, global and uniform reflection are provably different operations: a variety of well-known theories of truth can be closed under uniform reflection but not under global reflection.¹ However, this is not so in the case of Basic De Morgan logics and variations thereof, where global and uniform reflection coincide for a wide class of theories containing the truth principles of TS_0 .

The theories that we study in this work, resulting from iterating reflection over TS_0 , can be characterized as *internal* axiomatizations of Kripke's fixed point construction — specifically, its four-valued version based on Basic De Morgan Logic as it can be found in [14] — because their principles and rules are all sound with respect to this class of models. The paradigmatic case of an internal theory is the theory PKF introduced in [15]. These theories are faithful to the fixed-point models and interact well with the process of reflection. There is an alternative array of theories capturing Kripke's construction *externally*. They are couched in classical logic and they are meant to be faithful to the set-theoretic definition of this class of models. Among the external theories we find the well-known classical axiomatization KF of Kripke's construction [8], and also the iteration of reflection studied by Horsten and Leigh [18]. Although external axiomatizations invoke principles that are not valid in the intended semantics and they usually cannot be closed under global reflection, they are proof-theoretically stronger than the corresponding internal axiomatizations. As a consequence, they also deem true more sentences that belong to the intended extension of the truth predicate than the corresponding internal axiomatizations; for some authors (cf. [14, 16]), this is considered a clear advantage over internal theories.

For instance, KF is proof-theoretically much stronger than its natural internal counterpart PKF. The former proves transfinite induction for all sentences of the language with the truth predicate up to any ordinal smaller than ε_0 , the latter only up to any ordinal smaller than ω^ω . The main result of the present article is that two steps of reflection over TS_0 enables us to recapture all principles of PKF and prove significantly more transfinite induction than what is available in PKF. Moreover, iterated reflection on TS_0 enables us to reach the strength of KF.

2 The core laws of truth

In this section we introduce the main components of the theory TS_0 . We first introduce a two-sided sequent version of Basic De Morgan logic and state some simple properties of this calculus. We then introduce the principles governing the objects to which truth is ascribed, which will amount to the axioms of a very weak arithmetical theory. Finally we state the truth theoretic principles of TS_0 .

2.1 Basic De Morgan logic

We employ a two-sided sequent calculus BDM reminiscent of the one employed in [14]; sequents are expressions of the form $\Gamma \Rightarrow \Delta$ where Γ, Δ are finite sets of formulas. We write $\neg\Gamma$ for $\{\neg A \mid A \in \Gamma\}$. BDM is a subsystem of a suitable two-sided classical calculus; its axioms and rules are listed in Table 1. Intuitively, BDM is obtained from classical logic by replacing the usual clauses for negation with (CP1) and (CP2) below. However, the general negation rules (CP1-2) enable us to derive the sequents $A \Rightarrow \neg\neg A$ and $\neg\neg A \Rightarrow A$ for all formulas A and make the following contraposition rule admissible in BDM:

$$\frac{\Gamma \Rightarrow \Delta}{\neg\Delta \Rightarrow \neg\Gamma} \text{ (Cont).}$$

¹Examples are the system KF from [8] plus the axiom stating that the extension of the truth predicate is consistent and the system FS from [13] originating in [11].

TABLE 1. The system BDM

$A \Rightarrow A$	for $A \in \mathcal{L}_T$
$\frac{\Gamma \Rightarrow \Delta}{A, \Gamma \Rightarrow \Delta}$ (LW)	$\frac{\Gamma \Rightarrow \Delta}{\Gamma \Rightarrow \Delta, A}$ (RW)
$\frac{\Gamma \Rightarrow \Delta, A \quad A, \Gamma \Rightarrow \Delta}{\Gamma \Rightarrow \Delta}$ (Cut)	
$\frac{\neg \Gamma \Rightarrow \Delta}{\neg \Delta \Rightarrow \Gamma}$ (CP ₁)	$\frac{\Gamma \Rightarrow \neg \Delta}{\Delta \Rightarrow \neg \Gamma}$ (CP ₂)
$\frac{A, B, \Gamma \Rightarrow \Delta}{A \wedge B, \Gamma \Rightarrow \Delta}$ (L \wedge)	$\frac{\Gamma \Rightarrow \Delta, A \quad \Gamma \Rightarrow \Delta, B}{\Gamma \Rightarrow \Delta, A \wedge B}$ (R \wedge)
$\frac{A, \Gamma \Rightarrow \Delta \quad B, \Gamma \Rightarrow \Delta}{A \vee B, \Gamma \Rightarrow \Delta}$ (L \vee)	$\frac{\Gamma \Rightarrow \Delta, A, B}{\Gamma \Rightarrow \Delta, A \vee B}$ (R \vee)
$\frac{\Gamma, A(t) \Rightarrow \Delta}{\Gamma, \forall x A \Rightarrow \Delta}$ (L \forall)	$\frac{\Gamma \Rightarrow \Delta, A(t)}{\Gamma \Rightarrow \Delta, \exists x A}$ (R \exists)
$\frac{\Gamma \Rightarrow \Delta, A(x)}{\Gamma \Rightarrow \Delta, \forall x A}$ (R \forall)	$\frac{\Gamma, A(x) \Rightarrow \Delta}{\Gamma, \exists x A \Rightarrow \Delta}$ (L \exists)
x not free in Γ, Δ	x not free in Γ, Δ

The following closely related lemma will be extensively used in what follows:

LEMMA 1

For a language \mathcal{L}_0 with signature $\mathcal{S} = \{P_1, \dots, P_n\}$, if for all atomic formulas A of \mathcal{L}_0 we can prove $\Rightarrow A, \neg A$, then BDM formulated in \mathcal{L}_0 is closed under the following classical rules for negation:

$$\frac{\Gamma \Rightarrow \Delta, A}{\neg A, \Gamma \Rightarrow \Delta} (\neg L) \qquad \frac{A, \Gamma \Rightarrow \Delta}{\Gamma \Rightarrow \Delta, \neg A} (\neg R).$$

BDM enjoys standard properties of Gentzen-type sequent calculi such as substitution, inversion and cut elimination.

A natural semantics for BDM is given in terms of four-valued models, that is we also allow predicates with a partial or a paraconsistent behaviour (gaps and gluts).² The intended satisfaction relation has a double clause: a sequent is satisfied in a model \mathcal{M} just in case if all formulas in the antecedent are true in \mathcal{M} there is a formula in the consequent true in \mathcal{M} , and if all formulas in

²Our logic is close to what is sometimes called FDE. We follow, however, Field's terminology in [9, Ch. 3].

the succedent are false in \mathcal{M} , there is a false in \mathcal{M} formula in the antecedent. **BDM** is sound and complete with respect to the semantics just hinted at (see [3]).

2.2 The theory \mathbf{TS}_0

The basic theory \mathbf{TS}_0 and all its extensions will be formulated in the language of arithmetic $\mathcal{L} = \{0, \mathbf{S}, +, \times, \mathbf{exp}\}$ with exponentiation. With x not occurring in the \mathcal{L} -term t , expressions of the form $(\forall x \leq t) \varphi(x)$ and $(\exists x \leq t) \varphi(x)$ are said to be obtained from $\varphi(v)$ by *bounded quantification*: Δ_0 -formulas (or *elementary*) of \mathcal{L} are formulas that contain only bounded quantifiers. The classes of Σ_n and Π_n formulas of \mathcal{L} are then defined in the usual way. In practice, we work in the expansion of \mathcal{L} with finitely many function symbols corresponding to suitable elementary operations and the truth predicate \mathbf{T} . We call the resulting language $\mathcal{L}_{\mathbf{T}}$. We call a formula of $\mathcal{L}_{\mathbf{T}}$ *arithmetical* if it does not contain occurrences of \mathbf{T} .

To formulate \mathbf{TS}_0 , we first consider identity, which is governed by usual principles:

$$\Rightarrow t = t \quad (\text{Id1})$$

$$s = t, A(s) \Rightarrow A(t). \quad (\text{Id2})$$

\mathbf{TS}_0 will also contain initial sequents $\Rightarrow A$ for all basic axioms A of a suitable system of arithmetic, in our case Kalmar's elementary arithmetic **EA** formulated in $\mathcal{L}_{\mathbf{T}}$ (cf. [2, 12]).³ In addition, our basic theory features an induction rule

$$\frac{\Gamma, A(x) \Rightarrow A(x+1), \Delta}{\Gamma, A(0) \Rightarrow A(t), \Delta} \quad (\Delta_0\text{-IND})$$

where x is not free in $A(0)$, Δ , Γ , t is arbitrary, and A is a Δ_0 -formula of the language \mathcal{L} of arithmetic without the truth predicate. We call the resulting system **Basic**.

The core principles of truth capture the fundamental idea that one is justified in asserting a sentence A precisely when she is justified in asserting that A is true.

DEFINITION 1 (The system \mathbf{TS}_0)

\mathbf{TS}_0 is obtained by extending **Basic** with the initial sequents

$$\mathbf{T}(\ulcorner A \urcorner) \Rightarrow A \quad (\mathbf{T1})$$

$$A \Rightarrow \mathbf{T}(\ulcorner A \urcorner) \quad (\mathbf{T2})$$

for all $\mathcal{L}_{\mathbf{T}}$ -sentences A .

\mathbf{TS}_0 stands for 'truth sequents'. The subscript 0 indicates a restriction of induction to Δ_0 -formulas; its absence indicates full induction. The semantic conservativeness — and therefore the

³The class of *elementary functions* \mathcal{E} is obtained by closing the initial functions $\mathbf{zero}(\cdot)$, $\mathbf{suc}(\cdot)$, $+$, \times , 2^x , $\mathbf{P}_i^n(x_1, \dots, x_n) = x_i$ with $(1 \leq i \leq n)$, truncated subtraction $x \dot{-} y$ under the operations of composition and bounded minimalization:

$$H(\vec{x}) = F(G_1(\vec{x}), \dots, G_n(\vec{x})); \quad (\mu t \leq y) P(\vec{x}, t) = \begin{cases} \text{the least } t \leq y \text{ s.t. } P(\vec{x}, t) \\ 0, \text{ if there is no such } t \end{cases}$$

where F, G are elementary functions and P an elementary predicate. **EA** has sufficient resources to naturally introduce new relations corresponding to the elementary functions by proving their defining equations. In particular, the functions in \mathcal{E} are exactly the functions that can be Σ_1 -defined in **EA** (see [25, §3.1]).

consistency — of TS_0 over **Basic** can be obtained by expanding any model of the latter with an interpretation of the truth predicate resulting from a positive inductive definition along the lines of Kripke’s fixed-point construction (cf. [4, §5]). The following observation can be found in [15, Lem. 16].

LEMMA 2

$\Rightarrow A, \neg A$ is derivable in TS_0 for A arithmetical.

The principles of TS_0 are just right to capture the desired assertability conditions for truth ascriptions: its basic truth-theoretic principles (T1)–(T2) are in fact not as strong as the classical Tarski-biconditionals: otherwise they would lead to inconsistency. But they are also stronger than mere inference rules, as the latter do not allow for conditionalization for arithmetical sentences.

We can think of TS_0 as a minimal internal axiomatization of a fixed-point construction along the lines of Kripke’s [20]. The fixed-points we are interested in are in fact fixed-points of a monotone operator \mathcal{J} associated with the Basic De Morgan evaluation scheme.⁴ The crucial property of Kripke style fixed points S , i.e. sets of sentences S such that $\mathcal{J}(S)=S$, is that every sentence A is in S iff $\top(\ulcorner A \urcorner)$ is in S . By combining this fact with the notion of satisfaction introduced on p. 2633 we can easily see that for a fixed point S , the model (\mathbb{N}, S) satisfies TS_0 when S is taken to be the extension of the truth predicate. Moreover, for (\mathbb{N}, S) to satisfy TS_0 , S has to contain the same sentences as $\mathcal{J}(S)$. This means that S is a fixed-point of \mathcal{J} iff (\mathbb{N}, S) satisfies TS_0 , and therefore it is an internal axiomatization of the fixed-point construction in the sense of §1. Moreover any \mathcal{L}_\top -theory in **BDM** satisfying the adequacy condition just considered will contain the principles of TS_0 .

2.3 Intermezzo on arithmetization

In what follows, we assume a canonical Gödel numbering for \mathcal{L}_\top -expressions. For a fixed expression e of \mathcal{L}_\top , we will use the usual Gödel corners for the closed term of \mathcal{L}_\top representing Gödel number $\#e$ of e . Therefore, for formulas A of \mathcal{L}_\top , we will have $\ulcorner A \urcorner = \#A$. Similarly, for sequents $\Gamma \Rightarrow \Delta$, $\ulcorner \Gamma \Rightarrow \Delta \urcorner = \#\Gamma \Rightarrow \Delta$, where the Gödel code of $\Gamma \Rightarrow \Delta$ is taken to be an ordered pair whose components are the codes of the finite sets Γ and Δ .⁵ Closed terms standing for specific Gödel codes of \mathcal{L}_\top -expressions contrast with open terms standing for templates to generate such closed terms: a well-known example of such a template is the open \mathcal{L}_\top -term $\text{sub}(\ulcorner A(v) \urcorner, \ulcorner v \urcorner, \text{num}(x))$, standing for the result of formally substituting, in the formula $A(v)$, the free variable v with the numeral for x . To distinguish these open terms from specific codes, we use square brackets instead of Gödel codes, so that, for instance, $[A(x)]$ stands for $\text{sub}(\ulcorner A(v) \urcorner, \ulcorner v \urcorner, \text{num}(x))$.⁶ This distinction clearly generalizes to sequents and formulas with more than one free variable: $[\Gamma \vec{x} \Rightarrow \Delta \vec{x}]$ refers to the simultaneous substitution in $\ulcorner \Gamma \Rightarrow \Delta \urcorner$ of the variables in the strings \vec{x} with their corresponding numerals, where of course $[\Gamma x \Rightarrow \Delta x]$ is short for $\text{sub}(\ulcorner \Gamma \urcorner, \ulcorner \Delta \urcorner, \ulcorner x \urcorner, \text{num}(x))$. When it is clear from the context which free variable we are formally substituting, we will omit it and treat **sub** as a binary function.

In **EA** we can easily carry out an elementary arithmetization of the standard syntactic notions and operations such as the notion of being a closed term of \mathcal{L}_\top (formally, $\text{ct}(x)$), the notion of being a sentence of the language \mathcal{L}_\top (formally, $\text{Sent}_{\mathcal{L}_\top}(x)$), the operation of prepending a truth predicate to x (formally, $\text{tr}(x)$), and so on (see e.g. [2]). Theories will be elementary presented sets of axioms

⁴For a definition of the operator and the evaluation scheme we refer the reader to Halbach [14, section 15.1].

⁵We assume that the code of the finite set Γ is the code of the sequence of codes of formulas in Γ in ascending order.

⁶The square brackets notation is often replaced by the so-called Feferman dot notation, in which, for instance, $\text{sub}(\ulcorner A(v) \urcorner, \ulcorner v \urcorner, \text{num}(x))$ is abbreviated with $\ulcorner A(\dot{x}) \urcorner$.

and rules, and we write $\text{Ax}_T(x)$ for ‘ x is an axiom of T ’. Variations of canonical provability will be particularly relevant in what follows, especially in our discussion of reflection. We list them here accompanied by their intuitive meaning:

$\text{Prf}_T(x, y)$	‘ x is a proof in T of the sequent y ’
$\text{Prv}_T(x, y)$	‘the sequent y is provable in T with a proof of length no greater than x ’
$\text{Pr}_T(x)$	‘ x is a provable sequent in T ’
$\text{Thm}_T(x)$	‘ x is a theorem of T , that is the ‘sequent’ $(\ulcorner \varnothing \urcorner, x)$ is provable in T ’
$\text{Pr}_T^2(x, y)$	‘the inference from the sequent x to the sequent y is provably admissible in T ’

We conclude by commenting on the predicate $\text{Pr}_T^2(x, y)$. In addition to being admissible, a *provably admissible rule* in a theory T requires the existence of a T -provable proof transformation of the proof of the premise into a proof of the conclusion of the rule.⁷ Pr_T^2 enjoys generalized versions of some of the properties usually ascribed to provability predicates:

If $\frac{\Gamma(x) \Rightarrow \Delta(x)}{\Theta(x) \Rightarrow \Lambda(x)}$ is admissible in T , provably in **Basic**, then (Pr1)

Basic $\vdash \Rightarrow \text{Pr}_T^2([\Gamma(x) \Rightarrow \Delta(x)], [\Theta(x) \Rightarrow \Lambda(x)])$

If the sequents (Pr2)

$\Rightarrow \text{Pr}_T^2([\Gamma(x) \Rightarrow \Delta(x)], [\Theta(x) \Rightarrow \Lambda(x)]);$

and

$\Rightarrow \text{Pr}_T([\Gamma(x) \Rightarrow \Delta(x)])$

are derivable in **Basic**, then also

$\Rightarrow \text{Pr}_T([\Theta(x) \Rightarrow \Lambda(x)])$

is derivable in **Basic**.

2.4 *The weakness of TS_0 and the advantages of reflection in **BDM***

To evaluate the theories of truth considered in this work we list a number of desiderata. Many of them have been already proposed and discussed by truth-theorists — see for instance [21] and [17]. As already emphasized, we consider the intersubstitutivity of φ and $\ulcorner \varphi \urcorner$ as guiding principle and therefore we do not argue for it but take it as primitive.⁸ In addition, we require our truth predicate to be *compositional*, therefore enabling us to explain how we can understand complex sentences only on the basis of an understanding of its compounds and its logical structure. In particular, we

⁷For instance, although the rule of cut applied to ‘geometric’ formulations of Robinson’s arithmetic **Q** is admissible in it, cut is not provably admissible in **Q** as this procedure is of hyperexponential growth rate. (For a geometric presentation of Robinson’s arithmetic and for the cut elimination for it, see [22].)

⁸This immediately yields the identity of inner and outer logic as defined in [14, 21].

prefer theories that, given a set of logical constants, allow the truth predicate to commute with these constants for *all* sentences of the language, that is we require compositionality in quantified form and not only compositionality in schematic form. Finally, we aim at *proof-theoretically strong* theories of truth. The desideratum of strength can be motivated in several ways: following the broadly abductive picture sketched in [27], one might simply argue that a theory that proves more, all other things being equal, has a higher scientific status than a weaker one. Moreover, truth has served a foundational role in several philosophical programmes, such as Feferman’s predicativism or Aczel’s attempt to recover Frege’s foundations [1, 8], and therefore a stronger theory might provide alternative, and perhaps more appealing, formulations of theories of sets or of other mathematical objects.⁹

In the light of these criteria TS_0 , although natural, simple, and fully disquotational, is clearly not completely adequate. For one thing, it is not compositional in the required sense. For a universally quantified sentence such as $\forall x A(x)$, for instance, TS_0 cannot show how its truth value depends on the truth values of its compounds $A(t)$ because the sequent $\forall x \text{T}([A(x)]) \Rightarrow \text{T}(\ulcorner \forall x A(x) \urcorner)$ is not derivable for all A .

Closely related to the criterion of strength is the capability of a theory of proving the soundness of its base theory. This claim, in our setting, takes the form of the global reflection principle for **Basic**, i.e.

$$\text{Thm}_{\text{Basic}}(x) \Rightarrow \text{T}(x), \tag{GRF}_{\text{Basic}}$$

which is not derivable in TS_0 . The underderivability of the global reflection principle directly follows from the fact that TS_0 is a conservative extension of **Basic**. Therefore, in terms of strength measured with respect to **Basic**, TS_0 is as bad as it can get. Conservativeness, however, is only a boundary that one can use to partition truth theories into conservative and non-conservative ones. To distinguish between non-conservative extensions of **Basic**, finer-grained measures are required: one option is to consider how many arithmetical sentences a theory of truth proves (or equivalently, how much arithmetical transfinite induction the theory proves). This, however, does not directly take into account general claims involving the truth predicate that we would like to consider. Therefore, to measure the strength of our theories of truth we will consider the amount of transfinite induction for the language \mathcal{L}_{\top} that is provable in them. This criterion has the advantage of giving us direct information about how many truth iterations are provable in a theory (see §4).

Following a strategy already proposed and defended in [18] for theories formulated in classical logic, one may think of TS_0 as *implicitly* containing stronger principles, including compositional ones and principles of transfinite induction. This relation of implicit containment can be unfolded via postulating a hierarchy of reflection principles over TS_0 . Traditionally, reflection principles for a theory T are explicit soundness assertions (*‘whatever is provable in T , is true’*). The soundness of T is naturally expressed via GRF_T .¹⁰ However, by Tarski’s undefinability theorem, GRF_T can only be formulated if the expressive resources of T are increased with a fresh truth predicate. Therefore, if one wants to express soundness in an arithmetical language, one must resort to schemata. A well-known candidate is what is widely known as the *uniform reflection principle* for T :

$$\forall x(\text{Thm}_T([A(x)]) \rightarrow A(x)). \tag{RFN}_T$$

RFN_T states that, for every number x , if A is satisfied by the numeral for x , provably in T , then A is satisfied by x . However, we are mainly concerned with languages that *do* contain a truth predicate.

⁹Our criteria more or less correspond to the ones listed in [21], except of course the one requiring classical logic.

¹⁰This reading of reflection is ubiquitous in the literature. See for instance the classical handbook entry [26] and [14]. Kreisel and Lévy in [19] clearly states that global reflection is the intended soundness claim for a theory T .

In this context, therefore, the most natural way to express the soundness of a theory is by means of global reflection. In fact, if the truth predicate satisfies minimal conditions, the global reflection principles implies all instances of uniform reflection for a theory T .

There is therefore an intuitive connection between uniform and global reflection: both are intended to express the soundness of the base theory. It turns out, however, that this connection is lost in the classical axiomatizations of Kripke's fixed point construction considered by Leigh and Horsten in [18]. For T an axiomatization of Kripke's fixed point construction in classical logic, in fact, the result of adding GRF_T to it determines a severe restriction of the class of acceptable models: all *consistent* fixed points are excluded, i.e., if (\mathbb{N}, S) models $T + \text{GRF}_T$ with S a fixed point, then S is inconsistent.¹¹ In contrast, $T + \text{RFN}_T$ can have models of the form (\mathbb{N}, S) for S a consistent fixed point (in fact all consistent fixed points).

There is a natural explanation for the internal inconsistency of $T + \text{GRF}_T$: classical theories T of the sort just mentioned are in fact *unsound* with respect to the notion of truth captured by T , and GRF_T makes this explicit. In fact, many theorems involving the truth predicate in a classical axiomatization T of the fixed-point construction based on Basic De Morgan Logic are outside the extension of the truth predicate given by consistent fixed points. The classical tautology $\lambda \vee \neg\lambda$ involving a liar sentence is one such example. Uniform reflection alone, in theories such as T , does not suffice to uncover their unsoundness:¹² this is the sense in which the intimate connection between global and uniform reflection is lost in the classical setting considered here.¹³

The close connections between the two forms of reflection just considered, however, can be restored by moving to internal axiomatizations of (Basic De Morgan) Kripke fixed points such as extensions of TS_0 . To see this, we first reformulate both principles in rule form and adapt them to the sequent-style formulation of TS_0 that we have chosen.

$$\frac{\Rightarrow \text{Thm}_T([Ax])}{\Rightarrow A(x)} (\text{RFN}_T^R) \qquad \frac{\Rightarrow \text{Sent}_{\mathcal{L}_T}(x) \wedge \text{Thm}_T(x)}{\Rightarrow \text{T}x} (\text{GRF}_T^R).$$

Finally, we introduce an extension of TS_0 obtained by replacing the axioms (T1) and (T2) with

- (i) $A(x) \Rightarrow \text{T}[Ax]$;
- (ii) $\text{T}[Ax] \Rightarrow A(x)$.

We call the resulting system UTS_0 ('*uniform* TS_0 '). We can now establish that not only uniform and global reflection are connected in Basic De Morgan logic, but that they actually *coincide*.

PROPOSITION 1

Let T contain UTS_0 . Then $T + \text{RFN}_T^R$ and $T + \text{GRF}_T^R$ are identical theories.

PROOF. We start by showing that global reflection entails uniform reflection. Reasoning in $T + \text{GRF}_T^R$, we assume that the sequent $\Rightarrow \text{Thm}_T([Ax])$ is derivable in it. Then, by GRF_T^R , we have $\text{T}[Ax]$ and therefore Ax by (ii) above.

¹¹By the diagonal lemma, the arithmetical part of T already proves $(\lambda \wedge \neg \text{T}^\Gamma \lambda^\neg) \vee (\neg \lambda \wedge \text{T}^\Gamma \lambda^\neg)$ for λ a liar sentence. Therefore $T + \text{GRF}_T$ proves $\text{T}(\text{T}^\Gamma (\lambda \wedge \neg \text{T}^\Gamma \lambda^\neg) \vee (\neg \lambda \wedge \text{T}^\Gamma \lambda^\neg)^\neg)$. Since T is an axiomatization of the class of Kripke fixed points, we can use compositional and truth-iteration principles to obtain, still in $T + \text{GRF}_T$, $\text{T}(\text{T}^\Gamma \lambda \wedge \neg \lambda^\neg)$. A well-known example of such a theory T is a the theory KF from [8] — see also [14].

¹²This is also the reason why Horsten and Leigh could consider iterations of uniform reflection without restrictions on the fixed-points models.

¹³This phenomenon is not confined to axiomatizations of Kripke fixed points based on De Morgan Logic. Also in supervaluational fixed points uniform and global reflection come apart.

For the other direction, we reason in $T + \text{RFN}_T^R$ and assume that the sequent $\Rightarrow \text{Sent}_{\mathcal{L}_T}(x) \wedge \text{Thm}_T(x)$ is derivable in it. Also, we know that (i) and (ii) are derivable sequents of UTS_0 — and then also of T —, and therefore $\Rightarrow \text{Pr}_T(x, [\top x])$ is a derivable sequent of T . By combining this latter fact with our assumption, we obtain $\Rightarrow \text{Thm}_T([\top x])$. By RFN_T^R , therefore, we can conclude $\Rightarrow \top x$, as desired. ■

Proposition 1 suggests that Leigh and Horsten’s project can be more coherently carried out in the context of *non-classical* theories of truth. In the next section, we will in fact employ strengthenings of the reflection principles considered in Proposition 1 to unfold the truth-theoretic and mathematical content implicit in the acceptance of TS_0 .

3 Reflecting on TS_0

This section introduces the main results of the article: in Section 3.1 we discuss several alternative reflection rules and motivate the choice of a particular form of reflection on admissible rules that turns out to be stronger than simple reflection on derivable sequents. In Section 3.2 we show that the closure under two applications of our rule of reflection suffices to recover the strong internal axiomatization of Kripke’s fixed point PKF. Finally, in Section 3.3 we show that the result of reflecting twice on TS_0 proves more transfinite induction for the language with the truth predicate than PKF itself. We conclude the section by investigating further iterations of reflection.

3.1 Reflection on sequents and rules

As we have seen, in the classical setting the uniform reflection schema and rule take the form:

$$\text{Thm}_T([A(x)]) \Rightarrow A(x) \tag{RFN}_T^R$$

$$\frac{\Rightarrow \text{Thm}_T([A(x)])}{\Rightarrow A(x)}. \tag{URFN}_T^R$$

Over EA, URFN_T^R and RFN_T^R are equivalent, as shown by Feferman in [5]. In the non-trivial direction, i.e. going from the rule to the initial sequent, one shows that **Basic** suffices to formalize the fact that the sentence $\text{Prf}_T(\bar{n}, \ulcorner \Rightarrow A(\bar{m}) \urcorner) \rightarrow A(\bar{m})$ is provable in T for any $m, n \in \omega$. Therefore one application of (URFN_T^R) yields (RFN_T^R) .

In the non-classical setting the situation is different. Whereas with classical logic on the background we can formulate (URFN_T^R) and (RFN_T^R) in a one-sided sequent calculus, there are good reasons to stick with a two-sided calculus for Basic De Morgan logic. In a one-sided classical system, in fact, sequents $A, \neg A$ play the role that initial sequents $A \Rightarrow A$ play in a two-sided setting. In our system this correspondence breaks down: first of all, $A, \neg A$ is not generally valid in our intended semantics — if A is a liar sentence, for instance —, whereas $A \Rightarrow A$ are initial sequents of our system. Moreover, there is no conditional naturally corresponding to the sequent arrow since $\Rightarrow A \rightarrow A$ is just a notational variant of $\Rightarrow A \vee \neg A$.

We therefore opt for a formulation of our first reflection principle as applying to provable (two-sided) sequents. As a consequence, basic syntactic considerations force us to formulate reflection in rule-form. The *uniform reflection principle for sequents* of T takes the following form:

$$\frac{\Rightarrow \text{Pr}_T([\Gamma \vec{x} \Rightarrow \Delta \vec{y}])}{\Gamma \Rightarrow \Delta}. \tag{r}_T$$

But the simple rule of reflection (r_T) is not the only form of reflection that will be relevant for what follows. A suitable conditional — such as the classical or the intuitionistic conditional — enables one to compress in one sequent chains of reasoning featuring embedded implications. In our setting, the highly *meta-theoretic* nature of the sequent arrow forces us to capture these chains of reasoning explicitly via suitable extensions of the simple reflection rule (r_T). One way to achieve this is to focus not only on provable sequents, but also to take into account rules admissible in T via the provability predicate $\text{Pr}_T^2(x, y)$ introduced above.

We define the *uniform reflection principle for provably admissible rules* in T :

$$\frac{\Rightarrow \text{Pr}_T^2([\Gamma(x) \Rightarrow \Delta(x)], [\Theta(x) \Rightarrow \Lambda(x)]) \quad \Gamma(x) \Rightarrow \Delta(x)}{\Theta(x) \Rightarrow \Lambda(x)} \quad (\mathbf{R}_T)$$

Obviously, in the context of any reasonable theory T , the theory $T + r_T$ is a subtheory of $T + \mathbf{R}_T$. The rule \mathbf{R}_T states that if there is a uniform T -provable proof transformation of a proof of $\Gamma(\bar{n}) \Rightarrow \Delta(\bar{n})$ into a proof of $\Theta(\bar{n}) \Rightarrow \Lambda(\bar{n})$ for each $n \in \omega$, and moreover $\Gamma(x) \Rightarrow \Delta(x)$ is provable, then also $\Theta(x) \Rightarrow \Lambda(x)$ is provable. The rule \mathbf{R}_T is in a sense supplementing **BDM** with the possibility of dealing with chains of reasoning involving sequents. It is based on a formalization of admissible rules, which in the context of classical logic are easily compressed into sentence form and combined, informally, via conditionalization and transitivity. Although it increases the expressive power of the theories in **BDM**; however, it is clear that \mathbf{R}_T does not amount to restoring classical logic, as these chains of reasoning can only be managed in the safe and controlled environment guaranteed by the formalization in the support theory: the rule \mathbf{R}_T relies essentially on the syntactic capabilities of the background non-logical principles and is clearly not a logical principle. More importantly, the acceptability of our reflection principles, just like the acceptability of any other reflection principle, is based on preservation of soundness: as we shall see later on, the addition of \mathbf{R}_T to our theories of truth preserves not only soundness but all semantic properties of the theories of truth in **BDM**, intersubstitutivity *in primis*, and therefore it can be safely assumed regardless of how it interacts with the underlying logic.

If T is an axiomatizable theory, then the *reflection on T* is the closure of **Basic** under the reflection rules $r(T)$ and $\mathbf{R}(T)$:

$$\begin{aligned} r(T) &:= \text{Basic} + (r_T) \\ \mathbf{R}(T) &:= \text{Basic} + (\mathbf{R}_T). \end{aligned}$$

Theories obtained by iterating our reflection rules are then defined in a standard manner: for instance, $\mathbf{R}(\mathbf{R}(T))$ is the result of closing $\mathbf{R}(T)$ under $\mathbf{R}_{\mathbf{R}(T)}$. We abbreviate $\mathbf{R}(\mathbf{R}(T))$ as $\mathbf{R}^2(T)$, and similarly for more iterations.

We have introduced (\mathbf{R}_T) as a generalization of (r_T) . A natural question is whether (\mathbf{R}_T) is actually stronger than the simpler rule. We will not answer to this question in this article but we will prove some facts that may be relevant for a future answer. For instance, we now provide an upper bound for the strength of $r(\text{UTS}_0)$; later — cf. Proposition 3 — we will show that the resulting theory is a proper subtheory of $\mathbf{R}^2(\text{TS}_0)$.

The upper bound for $r(\text{UTS}_0)$ that we now provide is given in terms of the theory **PKF** that was mentioned in the introduction. **PKF** is also formulated in the language \mathcal{L}_T , and its axioms and rules are displayed in Table 2.

In $\mathbb{T}_{=1-2}$, the function symbol \equiv represents the elementary syntactic operation of forming an identity statement out of two terms. A similar notation will be applied for other syntactic operations.

TABLE 2. The theory PKF

LOGIC	logical initial sequents and rules of BDM
IDENTITY	Id1, Id2
ARITHMETIC	$\Rightarrow A$ for A a basic axiom of EA plus the <i>full</i> induction rule for \mathcal{L}_T : $\frac{\Gamma, A(x) \Rightarrow A(x + \bar{1}), \Delta}{\Gamma, A(0) \Rightarrow A(t), \Delta} \text{ (IND)}$
ATOMIC TRUTH	(T=1) $\text{ct}(x), \text{ct}(y), \text{val}(x) = \text{val}(y) \Rightarrow T(x \doteq y)$ (T=2) $\text{ct}(x), \text{ct}(y), T(x \doteq y) \Rightarrow \text{val}(x) = \text{val}(y)$ (TT1) $T[Tx] \Rightarrow Tx$ (TT2) $Tx \Rightarrow T[Tx]$
TRUTH PRINCIPLES FOR CONNECTIVES	(T∧1) $\text{Sent}_{\mathcal{L}_T}(x \wedge y), T(x) \wedge T(y) \Rightarrow T(x \wedge y)$ (T∧2) $\text{Sent}_{\mathcal{L}_T}(x \wedge y), T(x \wedge y) \Rightarrow T(x) \wedge T(y)$ (T∨1) $\text{Sent}_{\mathcal{L}_T}(x \vee y), T(x) \vee T(y) \Rightarrow T(x \vee y)$ (T∨2) $\text{Sent}_{\mathcal{L}_T}(x \vee y), T(x \vee y) \Rightarrow T(x) \vee T(y)$ (T¬1) $\text{Sent}_{\mathcal{L}_T}(x), T(\neg x) \Rightarrow \neg T(x)$ (T¬2) $\text{Sent}_{\mathcal{L}_T}(x), \neg T(x) \Rightarrow T(\neg x)$
TRUTH PRINCIPLES FOR QUANTIFIERS	(T∀1) $\text{Sent}_{\mathcal{L}_T}(\forall yx), \forall y T x(y/v) \Rightarrow T(\forall yx)$ (T∀2) $\text{Sent}_{\mathcal{L}_T}(\forall yx), T(\forall yx) \Rightarrow \forall y T x(y/v)$ (T∃1) $\text{Sent}_{\mathcal{L}_T}(\exists yx), \exists y T x(y/v) \Rightarrow T(\exists yx)$ (T∃2) $\text{Sent}_{\mathcal{L}_T}(\exists yx), T(\exists yx) \Rightarrow \exists y T x(y/v)$

As mentioned earlier, PKF is an internal axiomatization of Kripke’s theory of truth. Crucially, PKF is fully compositional as also negation commutes with the truth predicate. Halbach and Horsten in [15] have measured the proof-theoretic strength of PKF by showing that PKF proves arithmetical transfinite induction up to the ordinal $\varphi_\omega 0$. Therefore we can use PKF as means of comparison for our theories of iterated reflection.

PROPOSITION 2

$r(\text{UTS}_0)$ is a subtheory of PKF.

PROOF. To prove Proposition 2 we only need to check that PKF can handle reflection. We first establish that PKF is strong enough to prove the soundness of UTS_0 . To this end, for ordinal codes α , we define a hierarchy of predicates $\text{Tr}_\alpha(\cdot)$ as $T(\cdot) \wedge \text{Sent}_{\mathcal{L}^{<\alpha}}(\cdot)$, where $\mathcal{L}^{<\alpha}$ is defined as \mathcal{L} plus all

predicates Tr_β for $\beta < \alpha$ (cf. [14, Ch. 9]).¹⁴ Halbach and Horsten establish in [15] that PKF proves the predicates Tr_β to behave like Tarskian truth predicates for $\beta < \omega^\omega$: that is, for formulas of \mathcal{L}_\top that are in $\mathcal{L}_\top^{<\omega}$, the classical commutation conditions for typed truth predicates hold, while for \mathcal{L}_\top -formulas A not in \mathcal{L}_\top^ω , we can prove $\neg \text{Tr}_\omega \ulcorner A \urcorner$.

We can extend this definition to sequents via the predicate $\text{TR}_\alpha(\cdot)$ in the following way:

$$\text{TR}_\alpha(\ulcorner \Gamma \Rightarrow \Delta \urcorner) := ((\text{Tr}_\alpha(\ulcorner \Gamma \wedge \Gamma \urcorner) \rightarrow \text{Tr}_\alpha(\ulcorner \Gamma \vee \Delta \urcorner)) \wedge (\text{Tr}_\alpha(\ulcorner \Gamma \rightarrow \Delta \urcorner) \rightarrow \text{Tr}_\alpha(\ulcorner \neg \wedge \Gamma \urcorner))).$$

From the proof-theoretic analysis of PKF in [15] we know that the truth predicates up to ω^ω behave classically in it, and that we can employ the material conditional to carry out the inductive proof of the following:

$$\text{PKF} \vdash \Rightarrow \text{Pr}_{\text{UTS}_0}([\Gamma x \Rightarrow \Delta x]) \rightarrow \text{TR}_\omega([\Gamma x \Rightarrow \Delta x]). \quad (1)$$

The proof employs the induction rule of PKF. It suffices, therefore, to establish

$$\Rightarrow \text{Prv}_{\text{UTS}_0}(0, [\Gamma x \Rightarrow \Delta x]) \rightarrow \text{TR}_\omega([\Gamma x \Rightarrow \Delta x]) \quad (2)$$

$$\text{Prv}_{\text{UTS}_0}(u, [\Gamma x \Rightarrow \Delta x]) \rightarrow \text{TR}_\omega([\Gamma x \Rightarrow \Delta x]) \Rightarrow \quad (3)$$

$$\text{Prv}_{\text{UTS}_0}(u+1, [\Gamma x \Rightarrow \Delta x]) \rightarrow \text{TR}_\omega([\Gamma x \Rightarrow \Delta x]),$$

where $\text{Prv}_T(x, y)$ expresses that y is provable in T with a proof of length less or equal to x (see page 2636).

We consider the crucial case of the characterizing principles of UTS_0 , (i)–(ii) on page 2638. Reasoning classically in PKF, we assume

$$\Rightarrow \text{Prv}_{\text{UTS}_0}(0, [\top[Ax] \Rightarrow Ax]). \quad (4)$$

We need to show

$$\Rightarrow \text{Tr}_\omega[\top[Ax]] \rightarrow \text{Tr}_\omega[Ax] \quad (5)$$

$$\Rightarrow \text{Tr}_\omega[\neg Ax] \rightarrow \text{Tr}_\omega[\neg \top[Ax]]. \quad (6)$$

We start with (5) and reason informally: if $\text{Tr}_\omega[\top[Ax]]$, then for some $n \in \omega$ and $m < n$, $\text{Tr}_n[\top[Ax]]$. Therefore, since Tr_n and Tr_m are Tarskian truth predicates, also $\text{Tr}_n[Ax]$.

Similarly for (6), if $\text{Tr}_\omega[\neg Ax]$, then $\text{Tr}_n[\neg Ax]$ for some $n \in \omega$, and, since $\text{Tr}_n[Ax]$ is in \mathcal{L}_\top^{n+1} , also $\text{Tr}_{n+1} \neg[\top[Ax]]$ and therefore $\text{Tr}_\omega[\neg \top[Ax]]$. ■

3.2 *Recovering compositionality by reflection*

One of the goals of this section is to show that by reflecting on our core laws of truth we can recover desirable compositional principles. More specifically, reflecting on TS_0 is sufficient to recover the initial sequents and the full induction rule of PKF.

In a first step we show that adding the reflection principle for TS_0 to **Basic** allows us to derive the initial sequents of UTS_0 .

¹⁴The truth predicates Tr_α can be defined for as many ordinals as we can code in our theory. In Section 3.3, in particular, we will employ a coding for ordinals smaller than Γ_0 .

LEMMA 3

$\text{UTS}_0 \subseteq r(\text{TS}_0)$.

PROOF. For all $n \in \omega$ we have $\text{TS}_0 \vdash \text{T}(\ulcorner A(n) \urcorner) \Rightarrow A(n)$. Therefore, **Basic** proves:

$$\Rightarrow \forall y (\text{Sent}_{\mathcal{L}_T}(y) \rightarrow \text{Ax}_{\text{TS}_0}(\text{sub}(\ulcorner \text{T}x \urcorner, \text{num}(y)), y)) \quad (7)$$

and

$$\Rightarrow \forall x (\text{Sent}_{\mathcal{L}}([Bx]), \quad (8)$$

where, we recall, $[Bx] := \text{sub}(\ulcorner Bv \urcorner, \text{num}(x))$ for all \mathcal{L}_T -formulas B with one free variable. Therefore, by combining (7) and (8), we also have in **Basic**

$$\Rightarrow \text{Ax}_{\text{TS}_0}(\text{sub}(\ulcorner \text{T}x \urcorner, \text{num}([Ax]), [Ax])). \quad (9)$$

Therefore, by definition of the canonical provability predicate Pr_{TS_0} ,

$$\Rightarrow \text{Pr}_{\text{TS}_0}(\text{sub}(\ulcorner \text{T}x \urcorner, \text{num}([Ax]), [Ax])). \quad (10)$$

Recall that $\text{tr}(x)$ is the elementary function that formally prefixes $\ulcorner \text{T} \urcorner$ to the numeral for x . Then **Basic** also proves the equation:

$$\Rightarrow \text{sub}(\ulcorner \text{T}x \urcorner, \text{num}([Ax])) = \text{sub}(\text{tr}([Av]), \text{num}(x)). \quad (11)$$

By performing the appropriate substitution in (10), we have

$$\Rightarrow \text{Pr}_{\text{TS}_0}(\text{sub}(\text{tr}([Av]), \text{num}(x)), [Ax]). \quad (12)$$

The other direction is analogous.

In $r(\text{TS}_0)$ — and a fortiori in $\text{R}(\text{TS}_0)$ — therefore, we obtain

$$\text{T}[Ax] \Rightarrow A(x)$$

$$A(x) \Rightarrow \text{T}[Ax]$$

as desired. ■

As a consequence of the previous lemma, in $r(\text{TS}_0)$ we can already prove the full truth sequents for atomic arithmetical formulas and for truth ascriptions containing free variables ($\text{T}=1$), ($\text{T}=2$), (TT_1), and (TT_2). That (TT_1) and (TT_2) are direct instances of the initial truth sequents of UTS_0 is immediate. For the identity initial sequents, a slightly more general version of the Tarski sequents would be required, namely one in which at least two free variables appear. However, since we are working over **Basic**, we can always assume that the free variable in the truth sequents of UTS_0 stands for (the code of) a string of free variables of finite length.

Also the initial, compositional sequents of PKF for the propositional connectives \wedge, \vee, \neg can be proved in $r(\text{TS}_0)$:

LEMMA 4

In $r(\text{TS}_0)$ we can derive $(\text{T}\wedge_{1-2})$, $(\text{T}\vee_{1-2})$, $(\text{T}\neg_{1-2})$.

PROOF. In TS_0 we can directly prove the schematic form of the compositional clauses, for example $\text{TS}_0 \vdash \text{T}(\ulcorner A \urcorner) \wedge \text{T}(\ulcorner B \urcorner) \Rightarrow (\ulcorner A \wedge B \urcorner)$ for all \mathcal{L}_\top -sentences A, B . By formalizing this fact in **Basic**, we obtain

$$\Rightarrow \text{Pr}_{\text{TS}_0}([\text{Sent}_{\mathcal{L}_\top}(x \wedge y), \text{T}x \wedge \text{T}y \Rightarrow \text{T}(x \wedge y)]). \quad (13)$$

In $\text{r}(\text{TS}_0)$, therefore, we can then move from the formalization to the full quantifiable statement

$$\text{Sent}_{\mathcal{L}_\top}(x \wedge y), \text{T}x \wedge \text{T}y \Rightarrow \text{T}(x \wedge y) \quad (14)$$

as desired. The cases of the other connectives are analogous. \blacksquare

However, by looking at Table 2 one realizes that compositional initial sequents for the propositional connectives by themselves are not enough to capture all truth principles of PKF: we also need initial sequents for quantifiers and full induction for \mathcal{L}_\top (IND). We will first show how to recover full induction.

It is a well-known result that (uniformly) reflecting on EA suffices to obtain the full induction schema for \mathcal{L} . Kreisel and Lévy, in [19], proved the equivalence of uniform reflection and full induction over EA — that is the equivalence of EA plus uniform reflection for EA and Peano Arithmetic (PA).¹⁵ We will apply Kreisel and Lévy's strategy to our setting. To do so, however, their original argument has to be modified in several respects. First of all, we allow formulas of \mathcal{L}_\top and not just of \mathcal{L} to appear in instances of the induction schema. In addition, we have to consider an induction rule because the induction axiom involving the material conditional fails to be sound in the setting of Basic De Morgan Logic. Finally, already in this step, we shall employ our generalized reflection rule R_T instead of the basic reflection rule r_T . In what follows we denote as PA_\top the version of PA formulated in \mathcal{L}_\top whose logic is BDM and in which the truth predicate can appear in instances of induction.

LEMMA 5

$\text{PA}_\top \subseteq \text{R}(\text{Basic})$.

PROOF. Let $A(x)$ be a formula in \mathcal{L}_\top with one free variable. We want to show that in $\text{R}(\text{Basic})$ the full induction rule

$$\frac{\Gamma, A(x) \Rightarrow A(x + \bar{1}), \Delta}{\Gamma, A(0) \Rightarrow A(t), \Delta} \quad (15)$$

for formulas of \mathcal{L}_\top is admissible. The following inference is admissible in **Basic** – and in fact in predicate logic in \mathcal{L}_\top only — for any $n \in \omega$:

$$\frac{\Gamma, A(x) \Rightarrow A(x + \bar{1}), \Delta}{\Gamma, A(0) \Rightarrow A(\bar{n}), \Delta} \quad (16)$$

By (Pr1), since the proof transformation in (16) is elementary, **Basic** proves

$$\Rightarrow \text{Pr}_{\text{Basic}}^2(\ulcorner \Gamma, A(x) \Rightarrow A(x + \bar{1}), \Delta \urcorner, \ulcorner \Gamma, A(0) \Rightarrow A(\bar{y}), \Delta \urcorner). \quad (17)$$

Now by assumption, (17), and R_{Basic} we conclude

$$\Gamma, A(0) \Rightarrow A(y), \Delta. \quad \blacksquare$$

¹⁵See Beklemishev [2, Cor. 4.3] for a proof of this fact.

The full set of compositional sequents of PKF is obtained by complementing the clauses for the connectives by the ones for quantifiers. This can be achieved by closing the theory $R(TS_0)$ under $R_{R(TS_0)}$, that is, by performing one iteration of the general reflection rule.

LEMMA 6

$R^2(TS_0)$ proves $(T\forall_{1-2})$ and $(T\exists_{1-2})$.

PROOF. We prove $T\forall_1$; the other cases are treated similarly. For all \mathcal{L}_T -formulas $A(v)$ with only v free, $R(TS_0)$ proves

$$\begin{array}{ll} T[Ay] \Rightarrow A(y) & \text{by Lemma 3} \\ \forall y T[Ay] \Rightarrow \forall y A(y) & \text{by logic} \\ \forall y T[Ay] \Rightarrow T \ulcorner \forall y A(y) \urcorner & \text{by (T2).} \end{array}$$

The argument just carried out in $R(TS_0)$ can uniformly be formalized in Basic, i.e., Basic proves:

$$\Rightarrow \text{Pr}_{R(TS_0)}(\ulcorner \text{Sent}_{\mathcal{L}_T}(\forall y \dot{x}), \forall y T \dot{x}(y/v) \Rightarrow T(\forall y \dot{x}(y/v)) \urcorner).$$

Therefore $R^2(TS_0)$ suffices to conclude

$$\text{Sent}_{\mathcal{L}_T}(\forall y x), \forall y T x(y/v) \Rightarrow T(\forall y x(y/v)),$$

as desired. ■

COROLLARY 1

$\text{PKF} \subseteq R^2(TS_0)$.

Corollary 1 shows that two iterations of the generalized rule R_{TS_0} over our basic theory TS_0 suffice to recover all compositional truth laws that were not immediately provable in the original theory as well as the full induction rule for the language \mathcal{L}_T . If reflection is considered to be a procedure already implicit in the acceptance of TS_0 , then the laws of PKF follow naturally from a few applications of this process. However, it is natural to ask whether the inclusion established in Cor. 1 is proper.

These questions translate, on the conceptual side, into the task of approximating the set of sentences that are valid in the intended models of our theories, which are the Kripke fixed-point models. In doing so, we gather information on how many truth iterations and general claims involving truth we are permitted to assert upon accepting TS_0 (after reflection) and how many mathematical patterns of reasoning we regain in the form of transfinite induction.

3.3 Recovering transfinite induction by reflection

In this section, we investigate the question of how much transfinite induction for \mathcal{L}_T can be recovered in iterations of the generalized reflection rule over TS_0 . One of the upshots of our analysis will be that $R^2(TS_0)$ properly extends PKF.

To carry out our proofs, we need to assume a notation system $(OT, <)$ for ordinals up to the Feferman–Schütte ordinal Γ_0 as it can be found, for instance, in [24]. OT is a primitive recursive set of ordinal codes and $<$ a primitive recursive relation on OT that is isomorphic to the usual ordering of ordinals up to Γ_0 . We distinguish between fixed ordinal codes, which we denote with $\alpha, \beta, \gamma, \dots$, and $\zeta, \eta, \theta, \dots$ as abbreviations for variables ranging over elements of OT . From the results in [15] it

follows that PKF proves transfinite induction for \mathcal{L}_T only up to any ordinal smaller than ω^ω . If we focus only on \mathcal{L} -formulas, however, PKF proves transfinite induction for much higher ordinals. In particular, PKF proves the same arithmetical sentences as PA plus transfinite induction for \mathcal{L} up to any ordinal smaller than $\varphi_\omega 0$.

Before analysing how much transfinite induction can be proved in $R^2(\text{TS}_0)$, we introduce some notation. The schema of transfinite induction up to α for the formula $A(v)$ of a language \mathcal{L}_1 containing \mathcal{L} is the rule

$$\frac{\forall \xi < \eta A(\xi) \Rightarrow A(\eta)}{\Rightarrow \forall \xi < \alpha A(\xi)} \text{TI}_{\mathcal{L}_1}(A, \alpha).$$

We then denote transfinite induction up to some ordinal α with $\text{TI}_{\mathcal{L}_1}(<\alpha)$, standing for the closure under all rules $\text{TI}_{\mathcal{L}_1}(A, \beta)$ for $A \in \mathcal{L}_1$ and $\beta < \alpha$. Analogously, we write $\text{TI}_{\mathcal{L}_1}(\alpha)$ for the closure under all rules $\text{TI}_{\mathcal{L}_1}(A, \alpha)$ for $A \in \mathcal{L}_1$. In what follows, we will only deal with the cases in which \mathcal{L}_1 is either \mathcal{L} itself or \mathcal{L}_T .

As a measure of strength of the theories obtained via iteration of reflection, we will mainly focus on how much transfinite induction for \mathcal{L}_T is derivable in such theories. However, there is often a direct connection between the amount of transfinite induction for \mathcal{L}_T and \mathcal{L} derivable in a truth theory. Both in the case of KF and PKF, for instance, the amount of transfinite induction for \mathcal{L}_T available in the systems — that is $\text{TI}_{\mathcal{L}_T}(<\varphi_1 0)$ and $\text{TI}_{\mathcal{L}_T}(<\varphi_0 \omega)$ respectively — can be used to define classical, Tarskian truth predicates indexed by these ordinals with the crucial contribution of the compositional truth principles of the two theories (see Feferman's [8] for KF and Halbach and Horsten's [15] for PKF). This gives a lower bound for the systems in terms of ramified truth hierarchies up to $\varphi_1 0$ (or ε_0) and $\varphi_0 \omega$ (or ω^ω) respectively, which — by a classical result by Feferman (cf. [7]) — yields that KF and PKF are proof-theoretically as strong as at least $\text{PA} + \text{TI}_{\mathcal{L}_T}(<\varphi_{\varepsilon_0} 0)$ and $\text{PA} + \text{TI}_{\mathcal{L}_T}(<\varphi_\omega 0)$ respectively.

The following proposition shows that iterating the generalized reflection rule twice over TS_0 enables us to go beyond PKF. This also gives us more information about the question that was posed on page 2640 about the comparison between the rules (r_T) and (R_T) . By Proposition 2, the theory $r(\text{UTS}_0)$ is a subtheory of PKF. The next will entail that $R^2(\text{TS}_0)$ is indeed stronger than PKF.

PROPOSITION 3

$R^2(\text{Basic}) \vdash \text{TI}_{\mathcal{L}_T}(\omega^\omega)$.

PROOF. We first prove in $R(\text{Basic})$ that, for all $n \in \omega$,

$$\frac{\Gamma, \forall \zeta < \eta A(\zeta) \Rightarrow A(\eta), \Delta}{\Gamma \Rightarrow \forall \zeta < \omega^n A(\zeta), \Delta}. \quad (18)$$

To prove (18), we first prove in $R(\text{Basic})$, for all $n \in \omega$:

$$\frac{\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)}{\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^n A(\zeta)}. \quad (19)$$

We reason as follows in $\mathbf{R}(\mathbf{Basic})$:

$$\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta) \tag{20}$$

$$\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^0 A(\zeta) \quad \text{by (20)} \tag{21}$$

$$\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^n A(\zeta) \quad \text{external ind. hyp.} \tag{22}$$

$$\forall \zeta < \eta + (\omega^n \times x) A(\zeta) \Rightarrow \forall \zeta < \eta + (\omega^n \times x) + \omega^n A(\zeta) \quad \text{from (22)} \tag{23}$$

$$\forall \zeta < \eta + (\omega^n \times 0) A(\zeta) \Rightarrow \forall x \forall \zeta < \eta + (\omega^n \times x) A(\zeta) \quad \text{by (IND)} \tag{24}$$

$$\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^{n+1} A(\zeta) \tag{25}$$

The last two lines give us the induction step and therefore (19) by, possibly, a series of cuts.

Now in \mathbf{Basic} ,

$$\Rightarrow \text{Pr}_{\mathbf{R}(\mathbf{Basic})}^2([\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)], [\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^x A(\zeta)]). \tag{26}$$

Therefore, in $\mathbf{R}^2(\mathbf{Basic})$,

$$\frac{\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)}{\forall \zeta < \eta A(\zeta) \Rightarrow \forall x \forall \zeta < \eta + \omega^x A(\zeta)}. \tag{27}$$

That is

$$\frac{\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)}{\forall \zeta < \eta A(\zeta) \Rightarrow \forall \zeta < \eta + \omega^\omega A(\zeta)}. \tag{28}$$

From (28), by letting η to be 0, we get

$$\frac{\Gamma, \forall \zeta < \eta A(\zeta) \Rightarrow A(\eta), \Delta}{\Gamma \Rightarrow \forall \zeta < \omega^\omega A(\zeta), \Delta}. \tag{29}$$

■

By the proof theoretic analysis of PKF we know that it can only prove transfinite induction for \mathcal{L}_\top for ordinals smaller than ω^ω . But this fact is not dependent in any way on the truth theoretic principles of PKF: already PA_\top , in fact, proves $\text{TI}_{\mathcal{L}_\top}(< \omega^\omega)$. This is also reflected by the fact that Proposition 3 does not rely on the truth principles of TS_0 . However, by Corollary 1, we have:

COROLLARY 2

PKF is a proper subtheory of $\mathbf{R}^2(\text{TS}_0)$.

Transfinite induction up to ω^ω , however, is clearly not the limit of what we can achieve in $\mathbf{R}^2(\mathbf{Basic})$. By using similar methods to the ones employed in Proposition 3, and starting from (28), we can verify that the following rule is admissible in $\mathbf{R}^2(\mathbf{Basic})$:

$$\frac{\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)}{\forall \zeta < \theta A(\zeta) \Rightarrow \forall \zeta < \theta + \omega^{\omega+k} A(\zeta)}.$$

Generalizing this strategy it is possible to show the following:

LEMMA 7

In $\mathbf{R}^{(n+1)}(\mathbf{Basic})$ the following rule is admissible:

$$\frac{\forall \zeta < \eta A(\zeta) \Rightarrow A(\eta)}{\forall \zeta < \theta A(\zeta) \Rightarrow \forall \zeta < \theta + \omega^{\omega \times n} A(\zeta)}.$$

PROOF. By external induction on n . We have established the claim for $n = 1$. Assume that it holds for n . Then we can argue in $\mathbf{R}^{(n+1)}(\mathbf{Basic})$: Assume

$$\forall \zeta \prec \eta A(\zeta) \Rightarrow A(\eta).$$

Then by the induction hypothesis we have

$$\forall \zeta \prec \theta A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times n} A(\zeta)$$

and

$$\forall \zeta \prec \theta + \omega^{\omega \times n} \times x A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times n} \times x + \omega^{\omega \times n} A(\zeta).$$

By the induction principle (IND) we obtain

$$\forall \zeta \prec \theta + \omega^{\omega \times n} \times 0 A(\zeta) \Rightarrow \forall x \forall \zeta \prec \theta + \omega^{\omega \times n} \times x A(\zeta)$$

giving us

$$\forall \zeta \prec \theta A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times n} \times \omega A(\zeta)$$

which is

$$\forall \zeta \prec \theta A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times n + 1} A(\zeta).$$

Therefore, by iterating this argument m -times, we can obtain, for each m :

$$\forall \zeta \prec \theta A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times n + m} A(\zeta).$$

In $\mathbf{R}^{(n+2)}(\mathbf{Basic})$, therefore, we can conclude

$$\forall \zeta \prec \theta A(\zeta) \Rightarrow \forall \zeta \prec \theta + \omega^{\omega \times (n+1)} A(\zeta).$$

■

Lemma 7 immediately entails that $\mathbf{R}^{n+1}(\mathbf{Basic})$ proves $\text{TI}_{\mathcal{L}_\top}(< \omega^{\omega \times n})$. Therefore, if we reflect on TS_0 instead of \mathbf{Basic} , we are able to define in $\mathbf{R}^n(\text{TS}_0)$ ramified truth predicates for any ordinal smaller than $\omega^{\omega \times n}$ by following the strategy employed by Halbach and Horsten in [15] and hinted at on page 2646.

The strategy employed in Lemma 7 can be iterated even further. Ideally, we would like to reach, by as little reflection iterations as possible, the amount of transfinite induction for \mathcal{L}_\top — and therefore of ramified truth predicates — that are available in \mathbf{KF} , the classical counterpart of \mathbf{PKF} . However, we conclude this section by providing only a first, and presumably rather inefficient, approximation to this task.

By letting $\mathbf{R}^\omega(\mathbf{Basic}) := \bigcup_{n \in \omega} \mathbf{R}^n(\mathbf{Basic})$, a direct consequence of Lemma 7 is that

COROLLARY 3

In $\mathbf{R}^\omega(\mathbf{Basic})$ we have $\text{TI}_{\mathcal{L}_\top}(< \omega^{(\omega^2)})$

Therefore the theory $\mathbf{R}^\omega(\text{TS}_0)$ can define ramified truth predicates indexed by all ordinals $\omega^{\omega \times n}$ for all natural numbers n .

Although ω may seem to be a natural stopping point, the procedure can be iterated even further into the transfinite. Following a well-known tradition initiated by Feferman in [6], the theories $\mathbf{R}^n(\mathbf{Basic})$ can all be shown to be recursively enumerable. Moreover, the notion of being a proof in $\mathbf{R}^n(\mathbf{Basic})$

is recursive. We can then find a primitive recursive function enumerating all those proof predicates. By employing the recursion theorem, therefore, we can find an index for this enumeration that can be used to formalize, via a recursive predicate, the notion of being a proof employing rules proper of one of the theories $\mathbf{R}^n(\mathbf{Basic})$. This, however, suffices to formulate the notion of being a proof in $\mathbf{R}^\omega(\mathbf{Basic})$: clearly, similar procedure can be extended at least to ordinals smaller than ε_0 .

But once a recursive formalization of transfinite iterations of our reflection rules is available, it becomes clear that enough iterations of reflection over \mathbf{Basic} will lead us to the amount of transfinite induction for \mathcal{L}_\top available in \mathbf{KF} . By letting $\omega_0 := 1$, and $\omega_{n+1} := \omega^{\omega_n}$, we have, rather unsurprisingly,

OBSERVATION 1
 $\mathbf{R}^{\omega_n+1}(\mathbf{Basic}) \vdash \text{TI}_{\mathcal{L}_\top}(\omega_n)$.

4 Conclusion

Starting with principles that are minimally constitutive of the notion of truth, such as the initial sequents of the theory \mathbf{TS}_0 , we have investigated the result of iterating reflection rules over them. A similar project, in the context of classical logic and therefore without the basic principles of \mathbf{TS}_0 , has been recently pursued by Horsten and Leigh [18]. We claim that our non-classical setting provides a more coherent framework for such a project for two main reasons. First, in a *classical* setting the intersubstitutivity of A and $\top A^\top$ (which is the defining characteristic of \mathbf{TS}_0) cannot be consistently maintained. Second, following a theme by Kreisel, the global reflection \mathbf{GRF}_T for a theory T is the *intended* soundness extension of T . Other proof theoretic reflection principles, including the uniform reflection principle \mathbf{RFN}_T , are only justified by an appeal to global reflection. However, as shown in Section 2.3, in classical axiomatizations of Kripke’s fixed point constructions, the use of the global reflection principle is at odds with the overall strategy of iterating reflection rules.

One way to understand the results of this article is by asking which statements \mathbf{TS}_0 and the result of iterating reflection rules over it can prove to be true, i.e., by considering their provable sequents of the form $\Rightarrow \top A^\top$ for A in \mathcal{L}_\top or, in short, their *truth theorems*. The logic \mathbf{BDM} in itself — i.e. without identity — has no theorems at all. When initial sequents for identity are added to it as well as arithmetical initial sequents, even if the truth predicate is in the signature of the theory, one only obtains arithmetical theorems but no truth theorems. \mathbf{TS}_0 , in contrast, does prove truth theorems, but only truth theorems of the form

$$\underbrace{\top \dots \top}_{n\text{-times}} A^\top,$$

where A is an arithmetical theorem of \mathbf{Basic} . This shortcoming of \mathbf{TS}_0 is accompanied by the lack of other desirable properties of the theory, such as full compositionality (see again Section 2.3). By adding a uniform or global *reflection rule* to \mathbf{TS}_0 , we restore our full capability of reasoning inductively with the truth predicate, and several compositional truth sequents. Full compositionality, together with the possibility of establishing theorems of the form

$$\underbrace{\top \dots \top}_{\omega^{\omega+n}\text{-times}} A^\top$$

for A again an arithmetical theorem of \mathbf{Basic} , is reached when we consider the theory $\mathbf{R}^2(\mathbf{TS}_0)$, i.e., via a further iteration of the generalized reflection rule \mathbf{R}_T over \mathbf{TS}_0 . At this stage, we already recapture and surpass all truth theorems of the full compositional theory \mathbf{PKF} . A natural goal for the

process of iteration may be to reach the truth theorems of the classical theory KF (or equivalently, $\text{PKF} + \text{TI}_{\mathcal{L}_T}(< \varepsilon_0)$, as shown in [23]). This can be achieved via suitable transfinite progressions of theories obtained by reflection over TS_0 .

From a semantic perspective, there is a tight match between the truth theorems of our theories and the levels of the construction of the minimal fixed point of Kripke's construction from [20]. By extending TS_0 with an ω -rule, this connection can be made explicit: the theorems of TS_0 plus the ω -rule are exactly the \mathcal{L}_T -sentences that are in the extension of the truth predicate in the minimal fixed point of Kripke's theory (see [10] for a recent proof). Uniform reflection principles are recursive approximations of the ω -rule. Therefore iterations of reflection, and the corresponding truth theorems of the resulting theories, can be seen as approximations to the full ω -rule added to TS_0 as they represent initial stages of the construction of the minimal fixed point. It is also clear that all the theories that we have considered are *internal* axiomatizations of Kripke fixed points. Therefore the hierarchy that we have studied can also be seen as an attempt to capture, via recursively axiomatized theories, the set of grounded sentences first isolated by Kripke.

Nonetheless our work leaves many open questions and possibilities for improvement. From a technical point of view, a sharper proof-theoretic analysis of the theories obtained by iterated reflection would be desirable to see clearly, for instance, how much one can obtain with finite iterations of reflection. Moreover, it would be interesting to see whether the reflection rules can be strengthened via 'higher-order' reflection rules in such a way that only finitely many iterations of them can suffice to reach the truth theorems of KF. Finally, there remains the question whether the gap between $\text{TI}_{\mathcal{L}_T}(< \omega^\omega)$ and $\text{TI}_{\mathcal{L}_T}(< \varepsilon_0)$ — which is determined by whether PA in the signature of \mathcal{L}_T is formulated in BDM or classical logic respectively — can be closed by supplementing BDM with a suitable conditional in such a way that the conceptual advantages of the treatment of truth in TS_0 are preserved.

Acknowledgements

This research was supported by the DFG project "Syntactical Treatments of Interacting Modalities" and by the European Commission (grant 658285 - FOREMOTIONS). The authors would like to thank Volker Halbach, Graham Leigh, Johannes Stern, and the two anonymous referees for their comments and suggestions.

References

- [1] P. Aczel. Frege structures and the notions of proposition, truth and set. In *Journal of Symbolic Logic*, J. Barwise, H. J. Keisler and K. Kunen, eds, pp. 244–246. North-Holland, 1980.
- [2] L. Beklemishev. Reflection principles and provability algebras in formal arithmetic. *Russian Mathematical Surveys*, **60**, 197–268, 2005.
- [3] S. Blamey. Partial logic. In *Handbook of Philosophical Logic*, Vol. 5, D. Gabbay and F. Guenther, eds, pp. 261–353. 2 edn, 2002.
- [4] A. Cantini. Notes on formal theories of truth. *Zeitschrift für mathematische Logik und Grundlagen der Mathematik*, **35**, 97–130, 1989.
- [5] S. Feferman. Arithmetization of metamathematics in a general setting. *Fundamenta Mathematicae*, **XLIX**, 35–92, 1960.
- [6] S. Feferman. Transfinite recursive progression of axiomatic theories. *The Journal of Symbolic Logic*, **27**, 259–316, 1962.
- [7] S. Feferman. Systems of predicative analysis. *Journal of Symbolic Logic*, **29**, 1–30, 1964.

- [8] S. Feferman. Reflecting on incompleteness. *The Journal of Symbolic Logic*, **56**, 1–47, 1991.
- [9] H. Field. *Saving Truth from Paradox*. Oxford University Press, 2008.
- [10] M. Fischer and N. Gratzl. Infinitary proof systems and partial truth. Submitted.
- [11] H. Friedman and M. Sheard. An axiomatic approach to self-referential truth. *Annals of Pure and Applied Logic*, **33**, 1–21, 1987.
- [12] P. Hájek and P. Pudlák. *Metamathematics of First-Order Arithmetic*. Springer, 1993.
- [13] V. Halbach. A system of complete and consistent truth. *Notre Dame Journal of Formal Logic*, **35**, 311–27, 1994.
- [14] V. Halbach. *Axiomatic Theories of Truth*. Cambridge University Press, revised edition, 2014.
- [15] V. Halbach and L. Horsten. Axiomatizing Kripke’s theory of truth. *The Journal of Symbolic Logic*, **71**, 677–712, 2006.
- [16] V. Halbach and C. Nicolai. On the costs of nonclassical logic. *Journal of Philosophical Logic*, 2017, doi:10.1007/s10992-017-9424-3.
- [17] L. Horsten and V. Halbach. Norms for theories of reflexive truth. In *Unifying the Philosophy of Truth*, K. Fujimoto, J. M. Fernández, H. Galinon and T. Achourioti, eds. Springer, 2015.
- [18] L. Horsten and G. E. Leigh. Truth is simple. *Mind*, **126**, 195–232, 2017.
- [19] G. Kreisel and A. Lévy. Reflection principles and their use for establishing the complexity of axiomatic systems. *Zeitschrift für mathematische Logik und Grundlagen der Mathematik*, **14**, 97–142, 1968.
- [20] S. Kripke. Outline of a theory of truth. *The Journal of Philosophy*, **72**, 690–716, 1975.
- [21] H. Leitgeb. What theories of truth should be like (but cannot be). *Philosophy Compass*, **2**, 276–290, 2007.
- [22] S. Negri and J. Von Plato. Cut elimination in the presence of axioms. *Bulletin of Symbolic Logic*, **4**, 418–435, 1998.
- [23] C. Nicolai. Provably true sentences across axiomatizations of Kripke’s theory of truth. *Studia Logica*, pp. 1–30, 2017. doi:10.1007/s11225-017-9727-y.
- [24] W. Pohlers. *Proof Theory, The First Step into Impredicativity*. Springer, 2009.
- [25] H. Schwichtenberg. *Proofs and Computations*. Cambridge University Press, 2012.
- [26] C. Smoryński. The incompleteness theorems. In *Handbook of Mathematical Logic*, J. Barwise, ed., pp. 821–865. Dordrecht, 1977.
- [27] T. Williamson. Semantic paradoxes and abductive methodology. In *Reflections on the Liar*, B. Armour-Garb, ed. Oxford University Press, 2017.

Received 6 March 2017