



# Non-reflexivity and Revenge

Julien Murzi<sup>1</sup> · Lorenzo Rossi<sup>2</sup>

Received: 12 June 2020 / Accepted: 19 June 2021 / Published online: 12 November 2021  
© The Author(s) 2021

## Abstract

We present a revenge argument for non-reflexive theories of semantic notions – theories which restrict the rule of assumption, or (equivalently) initial sequents of the form  $\varphi \vdash \varphi$ . Our strategy follows the general template articulated in Murzi and Rossi [21]: we proceed via the definition of a notion of paradoxicality for non-reflexive theories which in turn breeds paradoxes that standard non-reflexive theories are unable to block.

**Keywords** Liar paradox · Curry’s Paradox · Validity Curry · Non-reflexive logics · Revenge

Jc Beall and Julien Murzi [2] argue that, if the semantic paradoxes are to be solved via a revision of classical logic, paradoxes of *naïve validity* such as the Validity-Curry require a restriction of one of the classical *structural rules*.<sup>1</sup> Building on [23], [22] consider the possibility of restricting certain instances of *structural reflexivity*

$$\frac{}{\varphi \vdash \varphi} \text{SRef}$$

as a way of treating the paradoxes of naïve validity, and the semantic paradoxes more generally. This paves the way for a consistent and broadly Kripkean account of naïve

---

<sup>1</sup>For more discussion of Beall and Murzi’s argument, see [8] and [22].

✉ Julien Murzi  
j.murzi@gmail.com

Lorenzo Rossi  
lorenzo.rossi@lrz.uni-muenchen.de

<sup>1</sup> University of Salzburg, Salzburg, Austria

<sup>2</sup> Munich Center for Mathematical Philosophy, LMU, München, Germany

validity, *grounded validity*, as Carlo Nicolai and Lorenzo Rossi call it. In this note, we show that a broad family of semantic theories based on *non-reflexive logics*, including Nicolai and Rossi's theory of grounded validity, breed revenge paradoxes they are unable to block.<sup>2</sup> The revenge strategy we follow is but an instance of the general revenge argument discussed in [21]. In a nutshell, the strategy is this. Non-classical theories such as Nicolai and Rossi's restrict the application of classical principles to paradoxical sentences. However, there exist finitely many principles  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  such that a sentence satisfies  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  just in case it satisfies all the principles of classical logic. This allows one to characterise *paradoxical* sentences as the ones that only satisfy  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  on pain of triviality. Unsurprisingly perhaps, such a notion of paradoxicality – which articulates a basic *fact* about non-classical theories – in turn gives rise to new paradoxes that can't be blocked by restricting, in the non-reflexive case, instances of SRef.

## 1 Paradoxes of Truth and Naïve Validity

Let  $\mathcal{L}$  be a first-order language encoding a minimal amount of syntax theory. More specifically, we require that the following two requirements are satisfied:

- (i) There is a function  $\ulcorner \urcorner$  such that for every sentence  $\varphi$ ,  $\ulcorner \varphi \urcorner$  is a closed term – informally, a name of  $\varphi$ .
- (ii) For every open formula  $\varphi(x)$  there is a term  $t_\varphi$  such that  $\ulcorner \varphi(t_\varphi/x) \urcorner$  denotes the same element as  $t_\varphi$ , in a fixed acceptable model of the language, where  $\varphi(t_\varphi/x)$  is the result of replacing every occurrence of  $x$  with  $t_\varphi$  in  $\varphi$ .

Theories that are rich enough to satisfy (i) and (ii) include, for instance, first-order arithmetic in the standard signature  $\{0, S, +, \times\}$  plus primitive symbols and defining formulae for all the primitive recursive functions.<sup>3</sup> We let  $\mathcal{L}_{\text{Tr}}$  be  $\mathcal{L}$  plus a primitive unary predicate  $\text{Tr}$ , for 'is true'.

We now specify the rules of classical propositional logic we will employ (henceforth, CPL).<sup>4</sup> We formulate them in a sequent-style natural deduction calculus where structural rules are explicitly formulated (see [25, Ch. 2]). A *sequent* is an expression

---

<sup>2</sup>Other non-reflexive approaches to truth and paradox have been developed, for instance, in [10, 19], and [9]. The approaches discussed in those papers are somewhat different from the one we discuss (for instance, [9] formulates his theory via many-sided sequents). However, the revenge paradox to be developed in what follows uniformly applies to all the theories (and subtheories) satisfying the (minimal) requirements formulated in Definition 1, Definition 4, and that are encoded in the proof of Proposition 5.

<sup>3</sup>See [12] for more details.

<sup>4</sup>Since our main result only requires propositional logical rules, this suffices for the purposes of this paper. A double line indicates that a rule can be read in both directions. Our axiomatization is redundant, but simplifies the subsequent derivations.

of the form  $\Gamma \vdash \varphi$ , where  $\Gamma$  is a finite multiset of sentences. A rule is an *inference* if its premises are empty, and a *meta-inference* otherwise.

$$\begin{array}{c}
 \frac{}{\varphi \vdash \varphi} \text{SRef} \quad \frac{\Gamma \vdash \chi}{\Gamma, \varphi \vdash \chi} \text{SWeak} \quad \frac{\Gamma, \varphi, \varphi \vdash \chi}{\Gamma, \varphi \vdash \chi} \text{SContr} \\
 \\
 \frac{\Gamma \vdash \varphi \quad \Delta, \varphi \vdash \psi}{\Gamma, \Delta \vdash \psi} \text{Cut} \\
 \\
 \frac{\Gamma \vdash \varphi \quad \Delta \vdash \psi}{\Gamma, \Delta \vdash \varphi \wedge \psi} \wedge\text{-I} \quad \frac{\Gamma \vdash \varphi \wedge \psi}{\Gamma \vdash \varphi} \wedge\text{-E}_1 \quad \frac{\Gamma \vdash \varphi \wedge \psi}{\Gamma \vdash \psi} \wedge\text{-E}_2 \\
 \\
 \frac{\frac{\Gamma \vdash \varphi}{\Gamma \vdash \varphi \vee \psi} \vee\text{-I}_1 \quad \frac{\Gamma \vdash \psi}{\Gamma \vdash \varphi \vee \psi} \vee\text{-I}_2}{\Gamma \vdash \varphi \vee \psi \quad \Delta_0, \varphi \vdash \chi \quad \Delta_1, \psi \vdash \chi} \vee\text{-E} \\
 \frac{}{\Gamma, \Delta_0, \Delta_1 \vdash \chi} \\
 \\
 \frac{\Gamma, \varphi \vdash \psi}{\Gamma \vdash \varphi \rightarrow \psi} \rightarrow\text{-I} \quad \frac{\Gamma \vdash \varphi \quad \Delta \vdash \varphi \rightarrow \psi}{\Gamma, \Delta \vdash \psi} \rightarrow\text{-E} \\
 \\
 \frac{}{\Gamma \vdash \varphi \rightarrow \psi} \text{Contr} \\
 \frac{}{\Gamma \vdash \neg\psi \rightarrow \neg\varphi} \\
 \\
 \frac{\Gamma \vdash \neg\neg\varphi}{\Gamma \vdash \varphi} \neg\neg\text{-I/E} \quad \frac{\Gamma, \varphi \vdash \perp}{\Gamma \vdash \neg\varphi} \neg\text{-I} \quad \frac{\Gamma \vdash \varphi \quad \Delta \vdash \neg\varphi}{\Gamma, \Delta \vdash \perp} \neg\text{-E} \quad \frac{\Gamma \vdash \perp}{\Gamma \vdash \varphi} \perp\text{-E}
 \end{array}$$

In order to present the Liar and other truth-theoretic paradoxes, we use the following introduction and elimination rules for naïve truth (for convenience, we assume both positive and negative forms):

$$\begin{array}{c}
 \frac{\Gamma \vdash \varphi}{\Gamma \vdash \text{Tr}(\ulcorner \varphi \urcorner)} \text{Tr-I} \quad \frac{\Gamma \vdash \text{Tr}(\ulcorner \varphi \urcorner)}{\Gamma \vdash \varphi} \text{Tr-E} \quad \frac{\Gamma \vdash \neg\varphi}{\Gamma \vdash \neg\text{Tr}(\ulcorner \varphi \urcorner)} \neg\text{-Tr-I} \\
 \\
 \frac{}{\Gamma \vdash \neg\text{Tr}(\ulcorner \varphi \urcorner)} \neg\text{-Tr-E}
 \end{array}$$

Now let  $\neg\text{Tr}(t_\lambda)$  be a sentence – a liar sentence – s.t.  $t_\lambda = \ulcorner \neg\text{Tr}(t_\lambda) \urcorner$  and let  $\lambda$  be a shorthand for  $\neg\text{Tr}(t_\lambda)$ . We may then (paradoxically) reason thus. We first prove  $\text{Tr}(\ulcorner \lambda \urcorner) \vdash \perp$ :

$$\frac{\frac{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \text{Tr}(\ulcorner \lambda \urcorner)}{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \lambda} \text{SRef} \quad \frac{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \text{Tr}(\ulcorner \lambda \urcorner)}{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \lambda} \text{Tr-E} \quad \frac{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \lambda}{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \neg\text{Tr}(\ulcorner \lambda \urcorner)} \text{Definition of } \lambda}{\text{Tr}(\ulcorner \lambda \urcorner), \text{Tr}(\ulcorner \lambda \urcorner) \vdash \perp} \neg\text{-E} \\
 \frac{}{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \perp} \text{SContr}$$

Call the above derivation  $\mathcal{D}_1$ . We then derive  $\text{Tr}(\ulcorner \lambda \urcorner)$  from  $\mathcal{D}_1$ :

$$\begin{array}{c}
 \mathcal{D}_1 \\
 \frac{\text{Tr}(\ulcorner \lambda \urcorner) \vdash \perp}{\vdash \neg\text{Tr}(\ulcorner \lambda \urcorner)} \neg\text{-I} \\
 \frac{}{\vdash \lambda} \text{Definition of } \lambda \\
 \frac{}{\vdash \text{Tr}(\ulcorner \lambda \urcorner)} \text{Tr-I}
 \end{array}$$

Call this derivation  $\mathcal{D}_2$ .  $\mathcal{D}_1$  and  $\mathcal{D}_2$  can now be combined together to yield a proof of absurdity, courtesy of Cut:

$$\frac{\mathcal{D}_2 \quad \mathcal{D}_1}{\frac{\vdash \text{Tr}(\ulcorner \lambda \urcorner) \quad \text{Tr}(\ulcorner \lambda \urcorner) \vdash \perp}{\vdash \perp} \text{Cut}}$$

Given  $\perp$ -E, the foregoing Liar reasoning yields a proof of *any sentence*  $\varphi$ , thus trivialising any theory in which the paradox can be derived.

Triviality can also be established without making use of  $\perp$ -E, via Curry’s Paradox [6]. The paradox involves a sentence  $\text{Tr}(t_\kappa) \rightarrow \psi$  (where  $\psi$  is any sentence) s.t.  $t_\kappa = \ulcorner \text{Tr}(t_\kappa) \rightarrow \psi \urcorner$ . Now let  $\kappa$  etc. be shorthand for  $\text{Tr}(t_\kappa) \rightarrow \psi$ . We can then essentially replicate the Liar reasoning to yield a ‘proof’ of  $\psi$ . We first ‘prove’  $\text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi$  – call this derivation derivation  $\mathcal{D}_3$ :

$$\frac{\frac{\frac{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \text{Tr}(\ulcorner \kappa \urcorner)}{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \text{Tr}(\ulcorner \kappa \urcorner)} \text{SRef} \quad \frac{\frac{\frac{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \text{Tr}(\ulcorner \kappa \urcorner)}{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \kappa} \text{Tr-E} \quad \text{Definition of } \kappa}{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \text{Tr}(\ulcorner \kappa \urcorner) \rightarrow \psi} \rightarrow\text{-E}}{\text{Tr}(\ulcorner \kappa \urcorner), \text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi} \text{SContr}}{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi} \text{SRef}}{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi} \text{SContr}$$

We then use  $\mathcal{D}_3$  to derive  $\text{Tr}(\ulcorner \lambda \urcorner)$ :

$$\frac{\mathcal{D}_3 \quad \frac{\frac{\text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi}{\vdash \text{Tr}(\ulcorner \kappa \urcorner) \rightarrow \psi} \rightarrow\text{-I} \quad \text{Definition of } \kappa}{\vdash \kappa} \text{Tr-I}}{\vdash \text{Tr}(\ulcorner \kappa \urcorner)} \text{Tr-I}$$

Call this derivation  $\mathcal{D}_4$ .  $\mathcal{D}_3$  and  $\mathcal{D}_4$  can now be combined together to yield a proof of absurdity, courtesy of Cut:

$$\frac{\mathcal{D}_4 \quad \mathcal{D}_3}{\frac{\vdash \text{Tr}(\ulcorner \kappa \urcorner) \quad \text{Tr}(\ulcorner \kappa \urcorner) \vdash \psi}{\vdash \psi} \text{Cut}}$$

On a revisionary approach to paradox,<sup>5</sup> naïve principles such as Tr-I and Tr-E are non-negotiable: they’re both required in order for the truth predicate to play its *inferential role*. For instance, the following reasoning is usually taken to motivate the unrestricted rule of Tr-E:

(AGREEMENT) All the theorems of Peano Arithmetic (PA) are true.  $\varphi$  is a theorem of PA. Hence,  $\varphi$  is true. Therefore (by Tr-E),  $\varphi$ .

Similar kinds of reasoning are standardly cited in support of a naïve conception of truth and, therefore, of a revisionary approach to paradox.

<sup>5</sup>See e.g. [7], [1, §§1.1-1.2], [5, 13].

Several revisionary theories of truth modify the logical rules in order to non-trivially retain naïve truth-theoretical principles, such as Tr-I, Tr-E,  $\neg$ -Tr-I,  $\neg$ -Tr-E. In light of the Liar and Curry’s Paradox, one might think that revising *operational* rules is enough to avoid semantic paradox altogether. Indeed, if one’s only naïve semantic notion is truth, then provably non-trivial approaches are available that only restrict logical rules. However, truth is not the only meaningful semantic notion.

In particular, *substructural* approaches to the Liar and the semantic paradoxes, i.e. approaches that restrict one of the classical structural rules, are required if one’s semantic theory is to be rich enough to express a naïve conception of validity, obeying the following principles:

(Validity Proof)

If  $\psi$  follows from  $\varphi$ , then the argument from  $\varphi$  to  $\psi$  is naïvely valid.

(Validity Detachment)

$\psi$  follows from  $\varphi$  and the claim that the argument from  $\varphi$  to  $\psi$  is naïvely valid.

Both principles are intuitively compelling. What is more, they can also be motivated by the need to model reasonings involving expressions of agreement and disagreement, just like the principles for naïve truth. For instance, the following reasoning can be taken to motivate VD (see [22] for more details):

(AGREEMENT\*) All the inferences performed by Mary are naïvely valid. Mary infers  $\chi$  from  $\varphi$  and  $\psi$ . Both  $\varphi$  and  $\psi$  hold. Therefore (by Validity Detachment),  $\chi$  also holds.

However, as in the case of truth, *Validity Proof* and *Validity Detachment* give rise to a version of Curry’s Paradox, called v-Curry, that only employs structural and semantic principles ([2]).<sup>6</sup>

To see this, let  $\text{Val}(x, y)$  be understood as ‘the argument from  $x$  to  $y$  is naïvely valid’. Then, the above principles can be respectively formalized as follows:

$$\frac{\varphi \vdash \psi}{\vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner)} \text{VP} \quad \frac{\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner) \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi} \text{VDm.}$$

We can now reason thus. Let  $\pi$  be a sentence identical to  $\text{Val}(\ulcorner \pi \urcorner, \ulcorner \perp \urcorner)$  (so that  $\pi$  says of itself that it entails absurdity), and let  $\mathcal{D}_5$  be the following derivation of  $\text{Val}(\ulcorner \pi \urcorner, \ulcorner \perp \urcorner)$ :

$$\frac{\frac{\frac{\pi \vdash \pi}{\pi \vdash \text{Val}(\ulcorner \pi \urcorner, \ulcorner \perp \urcorner)} \text{SRef} \quad \text{Def. of } \pi}{\pi \vdash \pi} \text{SRef}}{\frac{\pi, \pi \vdash \perp}{\pi \vdash \perp} \text{SContr}} \text{VP}$$

<sup>6</sup>Note that our formalization of ‘naïvely valid’ in AGREEMENT\* as a binary predicate is only one amongst several options. Another possibility is to formalize it as a unary predicate, taking arguments which (in the relevant applications) code arguments, i.e. codes of pairs of sentences. Alternatively, yet another option is to formalize it as a predicate applying to codes of pairs of multisets, for more structured arguments (this is done, for instance, in [23]). Nothing crucial hinges on this choice, however: in each of these cases, one can formulate analogues of VP and VD that give rise to essentially the same (naïve) validity-involving version of Curry’s Paradox.

Then, using  $\mathcal{D}_5$ , we can then ‘prove’  $\perp$ :

$$\frac{\frac{\mathcal{D}_5}{\vdash \text{Val}(\ulcorner \pi \urcorner, \ulcorner \perp \urcorner)} \quad \frac{\mathcal{D}_5 \quad \vdash \text{Val}(\ulcorner \pi \urcorner, \ulcorner \perp \urcorner)}{\vdash \pi} \text{Definition of } \pi}{\vdash \perp} \text{VDm}}$$

Since the argument makes no assumptions about the logic of negation and the conditional, it resists *fully structural* revisionary treatments, i.e. treatments that retain all of SRef, SContr, and Cut. As [2] point out, on a revisionary approach to the semantic paradoxes, the v-Curry Paradox calls for a *substructural* treatment – one on which one between SRef, SContr, and Cut is restricted.<sup>7</sup>

## 2 Fixed Points for Naïve Validity

Substructural approaches to the semantic paradoxes mostly fall into two main camps: SContr-free (see e.g. [16, 34–38]) and Cut-free (see e.g. [4, 5, 26–28]).<sup>8</sup> However, both kinds of approaches have been criticised as being both radical and unnecessary. For, a classically minded theorist might reason, how could a difference in the number of times a given premise is used in the course of an argument affect that argument’s validity? And how could consequence fail to be transitive? Moreover, one might legitimately ask, is there really a coherent notion of naïve validity for which *Validity Proof* and *Validity Detachment* are sound and yet non-trivial? In particular, ([7] pp. 10–11) argues that substructural approaches to paradox hard to get one’s head around and alleges [8] that there can be no coherent notion of validity satisfying both *Validity Proof* and *Validity Detachment*.

However, Field’s assessment is too bleak. Nicolai and Rossi [23] offer a SRef-free theory of naïve truth and naïve validity (called *grounded validity*), based on a family of fixed point constructions that can be seen as a generalisation of Kripke’s original construction for truth. As we’re about to see, the theory validates versions of *Validity Proof* and *Validity Detachment*. Since the new construction is but a natural generalisation of Kripke’s original one, [22] argue, in response to Field, that grounded validity is just as legitimate as a naïve conception of truth based on Kripke’s

<sup>7</sup>Strictly speaking, Cut is not employed in the above derivation. This is because the formulation of VDm essentially already incorporates a form of Cut (or meta-inferential *modus ponens*), and is therefore unavailable to non-transitive theories. More generally, the principle of *Validity Detachment* can be (and has been) formulated in two different ways in the literature: as an inference,

$$\frac{\varphi, \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner) \vdash \psi}{\text{VD}}$$

and as a meta-inference, as in our VDm (the ‘m’ stands for ‘meta-inference’). A derivation of  $\perp$  via a v-Curry reasoning involving both Cut and VD can be easily obtained via a simple modification of the above derivation (see [22], §1.3). We here focus on the meta-inferential VDm, because our target are non-reflexive theories such as the one developed in Section 2, which can be consistently closed under VDm but not VD.

<sup>8</sup>For general background on revisionary approaches to semantic paradox, see [3]. For a comprehensive critical comparison of SContr-free and Cut-free, see [29].

original construction. In what follows, we offer an informal presentation of Nicolai and Rossi’s construction – the KV-construction.<sup>9</sup>

### 2.1 The KV-Construction

Let  $\mathcal{L}_V$  be  $\mathcal{L}$  plus the designated predicate  $\text{Val}(x, y)$ , for ‘is naïvely valid’. We don’t need to add a predicate for truth, since naïve truth can be defined putting  $\text{Tr}(\ulcorner \varphi \urcorner) := \text{Val}(\ulcorner \top \urcorner, \ulcorner \varphi \urcorner)$ . Similarly, negation can be defined via the validity predicate, letting  $\neg\varphi := \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \perp \urcorner)$ , and conditionals can be defined via negation and disjunction, setting  $\varphi \rightarrow \psi := \neg\varphi \vee \psi$ . For this reason, we freely employ the truth predicate, negations, and conditionals in what follows, even though the semantic construction to be developed in this section only contains explicit clauses for the naïve validity predicate.<sup>10</sup>

Let  $\mathcal{M}$  be a countable acceptable model of the non-semantic fragment of  $\mathcal{L}_V$ .<sup>11</sup> We also assume that  $\mathcal{L}_V$  has constants for all the elements in  $\mathcal{M}$ . The KV-construction defines an increasing succession of sets

$$I_\Psi^0 \subsetneq I_\Psi^1 \subsetneq \dots \subsetneq I_\Psi^{\alpha+1} \subsetneq \dots$$

which ultimately determines the extension of the predicate  $\text{Val}$ . However, rather than building a succession of sets of sentences, the KV-construction builds a succession of sets of multiple conclusion *sequents*, i.e. objects of the form  $\Gamma \vdash \Delta$ , which are then internalized via naïve clauses for the predicate  $\text{Val}$ , thus building its extension.<sup>12</sup> The succession of sets of inferences is formalized by an inductive definition, and therefore it reaches a fixed point, i.e. for some ordinal  $\delta$ :

$$I_\Psi^0 \subsetneq I_\Psi^1 \subsetneq \dots \subsetneq I_\Psi^{\alpha+1} \subsetneq \dots \subsetneq I_\Psi^\delta = I_\Psi^{\delta+1} = \dots$$

The starting point of the construction is just the empty set:  $I_\Psi^0 = \emptyset$ . At stage 1, it contains only the inferences that are validated by the base model  $\mathcal{M}$ . Therefore,  $I_\Psi^1$  contains:

- $\Gamma \vdash \varphi, \Delta$  where  $\varphi$  is an atomic  $\mathcal{L}$ -sentence and  $\mathcal{M} \models \varphi$ ,
- $\Gamma, \psi \vdash \Delta$  where  $\psi$  is an atomic  $\mathcal{L}$ -sentence and  $\mathcal{M} \not\models \psi$

<sup>9</sup>Our discussion here presupposes familiarity with Kripke’s theory (strong Kleene version). For a technically informed presentation, see ([17] Chapters 3 and 4); for a less technical presentation, see ([32] Chapters 4 and 5).

<sup>10</sup>As it turns out, the so-defined material conditional has the same truth-conditions of the corresponding  $\text{Val}$ -statement in the KV-construction. In other words,  $\varphi \rightarrow \psi$  is validated by the KV-construction whenever  $\text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner)$  is. This means that, on the semantics on offer, the material conditional can simply be taken to be the ‘connective version’ of the naïve validity predicate. To be sure, *qua* predicate, the validity predicate can be used by itself to formulate paradoxical sentences (unlike the material conditional).

<sup>11</sup>An acceptable model has an in-built coding scheme and a sub-structure which is isomorphic to  $\mathbb{N}$  (the standard model of arithmetic) and  $\leq_{\mathbb{N}}$  (its standard well-ordering). See ([20] p. 22) for more details.

<sup>12</sup>Using multiple conclusion sequents simplifies some of the proofs of the properties of the KV-construction.

At stage  $\alpha + 1$ , the construction only contains the sequents that are obtained from sequents in stage  $\alpha$ , by closing it under all the classically valid meta-inferences. More precisely,  $\Gamma \vdash \Delta \in I_{\Psi}^{\alpha+1}$  if:

- $\Gamma \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma \vdash \varphi \wedge \psi, \Delta_0$  and  $\Gamma \vdash \varphi, \Delta_0$  is in  $I_{\Psi}^{\alpha}$  and  $\Gamma \vdash \psi, \Delta_0$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma_0, \varphi \wedge \psi \vdash \Delta$  and  $\Gamma_0, \varphi, \psi \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma \vdash \varphi \vee \psi, \Delta_0$  and  $\Gamma \vdash \varphi, \psi, \Delta_0$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma_0, \varphi \vee \psi \vdash \Delta$  and  $\Gamma_0, \varphi \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$  and  $\Gamma_0, \psi \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma \vdash \forall x \varphi(x), \Delta_0$  and for every closed  $\mathcal{L}_V$ -term  $s$ :  $\Gamma \vdash \varphi(s), \Delta_0$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma_0, \forall x \varphi(x) \vdash \Delta$  and for some closed  $\mathcal{L}_V$ -term  $s$ :  $\Gamma_0, \varphi(s) \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner), \Delta_0$  and  $\Gamma, \varphi \vdash \psi, \Delta_0$  is in  $I_{\Psi}^{\alpha}$ , or
- $\Gamma \vdash \Delta$  is  $\Gamma_0, \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner) \vdash \Delta$  and  $\Gamma_0 \vdash \varphi, \Delta$  is in  $I_{\Psi}^{\alpha}$  and  $\Gamma_0, \psi \vdash \Delta$  is in  $I_{\Psi}^{\alpha}$

Finally, at limit stages we take unions. For  $\beta$  a limit ordinal:

$$I_{\Psi}^{\beta} = \bigcup_{\alpha < \beta} I_{\Psi}^{\alpha}$$

We indicate with  $I_{\Psi}$  the (smallest) fixed point of the above construction.

Just like Kripke’s original construction,  $I_{\Psi}$  validates the intersubstitutivity of  $\varphi$  and  $\text{Tr}(\ulcorner \varphi \urcorner)$  in all non-opaque contexts. More precisely, recalling that  $\text{Tr}(\ulcorner \varphi \urcorner)$  is defined as  $\text{Val}(\ulcorner \top \urcorner, \ulcorner \varphi \urcorner)$ , the construction clearly validates the following equivalences:<sup>13</sup>

$$\begin{aligned} \Gamma \vdash \text{Tr}(\ulcorner \varphi \urcorner), \Delta \in I_{\Psi} \text{ iff } \Gamma \vdash \text{Val}(\ulcorner \top \urcorner, \ulcorner \varphi \urcorner), \Delta \in I_{\Psi} \text{ iff } \Gamma \vdash \varphi, \Delta \in I_{\Psi} \\ \Gamma, \text{Tr}(\ulcorner \varphi \urcorner) \vdash \Delta \in I_{\Psi} \text{ iff } \Gamma, \text{Val}(\ulcorner \top \urcorner, \ulcorner \varphi \urcorner) \vdash \Delta \in I_{\Psi} \text{ iff } \Gamma, \varphi \vdash \Delta \in I_{\Psi} \end{aligned}$$

In addition, it also validates the following versions of *Validity Proof* and *Validity Detachment*:

- (VP) if  $\Gamma, \varphi \vdash \psi$  is in  $I_{\Psi}$ , then  $\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner)$  is in  $I_{\Psi}$ .
- (VDM) if  $\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner)$  is in  $I_{\Psi}$  and  $\Delta \vdash \varphi$  is in  $I_{\Psi}$ , then  $\Gamma, \Delta \vdash \psi$  is in  $I_{\Psi}$ .

### 2.2 A Non-reflexive Theory of Naïve Validity

The set of sequents and the extension of  $\text{Val}$  provided by  $I_{\Psi}$  yields a *non-reflexive* theory of naïve validity (and truth). Structural reflexivity is the only classically valid structural rule that in general fails in  $I_{\Psi}$ , since  $\pi \vdash \pi$  (or equivalently  $\emptyset \vdash \text{Val}(\ulcorner \pi \urcorner, \ulcorner \pi \urcorner)$ ) is not in  $I_{\Psi}$ . More generally,  $I_{\Psi}$  does not validate any inference: it is *not* the case that, if a schematic inference  $\Gamma \vdash \Delta$  is classically valid, then all its instances are in  $I_{\Psi}$ . However,  $I_{\Psi}$  is closed under all the classically valid *meta-inferences*: if the meta-inference

$$\frac{\Gamma_0 \vdash \Delta_0, \dots, \Gamma_n \vdash \Delta_n}{\Gamma \vdash \Delta}$$

<sup>13</sup>Proving this requires establishing some closure properties of  $I_{\Psi}$ . See [23] for details.



is classically valid, then whenever  $\Gamma_0 \vdash \Delta_0, \dots, \Gamma_n \vdash \Delta_n$  are in  $I_\Psi$ , so is  $\Gamma \vdash \Delta$ . Since structural reflexivity is the only inference among the classically valid structural rules (i.e. the other structural rules are all meta-inferences), it is the only structural rule that's not unrestrictedly validated.

The theory of naïve validity provided by  $I_\Psi$  implicitly adopts a *tolerant-strict* (or TS) notion of consequence [4, 5]. Simplifying a bit, a TS-consequence relation interprets the logical constants as in a strong Kleene evaluation, and therefore makes use of three semantic values (0, 1/2, and 1). A set of sentences  $\Delta$  is said to *TS-follow* from a set of sentences  $\Gamma$  if, for any given strong Kleene evaluation, whenever all the sentences in  $\Gamma$  have value 1 or 1/2, then at least one of the sentences in  $\Delta$  has value 1. Once validity is so defined, the failure of structural reflexivity is easy to see: since  $\pi$  receives value 1/2 in the strong Kleene evaluation associated with  $I_\Psi$ , the sequent  $\pi \vdash \pi$  is not TS-valid.

$I_\Psi$  guarantees the consistency and non-triviality of several non-reflexive axiomatic theories of naïve truth and validity. Here, we (informally) present one very simple such theory:<sup>14</sup>

**Definition 1 (TSTV)** TSTV is provided by:

- Primitive name-forming;
- CPL (defined as above), *minus* SRef, *plus* the following principle:

$$\frac{\Gamma \vdash \varphi \rightarrow \varphi}{\Gamma, \varphi \vdash \varphi} \text{WSRef}$$

- The naïve truth rules Tr-I, Tr-E,  $\neg$ Tr-I,  $\neg$ Tr-E, and the following two rules for naïve validity:

$$\frac{\Gamma, \varphi \vdash \psi}{\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner)} \text{VP} \qquad \frac{\Gamma \vdash \text{Val}(\ulcorner \varphi \urcorner, \ulcorner \psi \urcorner) \quad \Delta \vdash \varphi}{\Gamma, \Delta \vdash \psi} \text{VDm}$$

TSTV is an extremely weak axiomatic theory: to wit, it does not have compositional rules for truth or validity, nor does it have strong base-theoretic axioms, such as the instances of the induction schema. For all its weakness, however, TSTV is prone to revenge – and so are, *a fortiori*, all of its stronger extensions.

### 2.3 Towards Revenge

As far as we can see, *pace* [8],  $I_\Psi$  and TSTV characterise a coherent notion of validity – one that satisfies forms of *Validity Proof* and *Validity Detachment* and that also justifies, for this reason, the adoption of a substructural logic of paradox.<sup>15</sup> To be sure, restrictions of SRef might seem pretty drastic. After all, SRef simply codifies, in a sequent setting, the natural deduction *rule of assumption*, which says that we may assume,

<sup>14</sup>A much richer, axiomatic non-reflexive theory of naïve truth and validity is developed in Nicolai and Rossi, [23]. Systems for non-reflexive consequence (under review).

<sup>15</sup>We defend this claim at length in [22]. See also footnote 7, where we explicitly consider different formalisations of *Validity Detachment*.

for the sake of argument, any sentence we like. And, it might be objected, if we can't make assumptions for the sake of argument for any sentence, how are we to reason? In response to this (understandable) worry, we'd like to make two observations.

To begin with, consider non-classical theories of naïve truth whose underlying logic is the three-valued strong Kleene logic  $K3$ , or some extension thereof [7, 13]. In these theories, validity is defined as preservation of value 1, and meta-inferences such as  $\rightarrow$ -I and  $\neg$ -I fail. Although in  $K3$ -based theories all instances of SRef hold, the assumptions that are thereby licensed cannot really be used for  $\rightarrow$ -I and  $\neg$ -I, since these rules do not hold in  $K3$ -based theories. Thus, it seems to us, our capacity to make assumptions is *already compromised* in well-known and widely accepted non-classical theories, such as the axiomatic theory PKF [11, 13, 14].

The second point is that, although SRef isn't valid in general, many of its instances are perfectly harmless. More precisely, in the foregoing setting, SRef is correctly applied to any sentence with a classical value. To be sure, we don't always know if a sentence actually has a classical value. For instance, natural languages allow the construction of *contingent liars* – sentences such as 'What Jane said is not true', whose paradoxical status depends on contingent facts, such as whether 'What Jane said is not true' is the semantic value of 'what Jane said' (see [15] p. 692 and ff.). If it is, then, in the foregoing framework, 'What Jane said is not true' has value  $1/2$ , and one may not apply SRef to it. However, one might reason, what's the worst that could happen if we applied SRef to a paradoxical instance of 'What Jane said is not true'? Arguably, just this: *on such an assumption*, we'd derive (via familiar Liar-like reasoning) that what Jane said is both true and not true, i.e. (courtesy of  $\neg$ -E) we'd be in a position to derive  $\perp$ . But rather than trivialising our theory, this would simply give us a reason to discharge the assumption that 'What Jane said is not true' satisfies SRef and assert that 'What Jane said is not true' is paradoxical, and shouldn't be reasoned with classically. That is, a principle to the effect that if  $\perp$  can be derived from the assumption that  $\varphi$  satisfies  $\varphi \vdash \varphi$  (or some equivalent principle), then  $\varphi$  is paradoxical, would help us deciding *which sentences* are paradoxical and hence cannot be reasoned with classically. For instance, if the assumption that  $\lambda$  satisfies  $\lambda \vdash \lambda$  (or some equivalent principle) leads to absurdity, then we know that  $\lambda$  is paradoxical and hence may not be reasoned with classically (we return to this point in Section 4).

However, this is precisely where TSTV and its extension reveal their limits. Murzi and Rossi [21] show that four families of revisionary theories (including theories that restrict SContr and theories that restrict Cut) are committed to validating principles about paradoxicality and unparadoxicality under which these theories can only be closed on pain of triviality. In a nutshell, the basic idea is this. Non-classical theories restrict the application of classical principles to unparadoxical sentences. However, for each such non-classical theory, there exist finitely many principles  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  such that a sentence  $\varphi$  satisfies  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  if and only if it satisfies all the principles of classical logic. This allows one to characterise *paradoxical* sentences as the ones that only satisfy each of  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  on pain of triviality and *unparadoxical* sentences as the ones that satisfy each of  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$ . Such a characterisation cannot be explicitly articulated in the theories in question, though, on pain of paradox-induced triviality. In the next section, we apply the recipe to TSTV and its extensions.

*Remark* We should stress that our distinction between paradoxical and unparadoxical sentences (in a given theory  $T$ ) is not assumed to be exhaustive or exclusive. As Lucas Rosenblatt [30] points out, paracomplete Kripkean theories of truth, for instance, can feature sentences – such as the truth-teller  $\tau$ , identical to its own truth predication  $\text{Tr}(\ulcorner \tau \urcorner)$  – which do *not* entail  $\perp$  if reasoned with classically, and which yet fail to satisfy all the classically valid principles. In our framework, then, such sentences are neither paradoxical nor unparadoxical. Rosenblatt [30] uses arguments in the style of [18] to show that *every* classically valid principle  $\mathfrak{P}$  has sets of instances which are mutually inconsistent – essentially, because any consistent set of instances of  $\mathfrak{P}$  can be expanded in incompatible ways, e.g. by adding to it sentences such as  $\tau$ , or their negations. These arguments, however, shouldn't be seen as showing that our notions of paradoxicality and unparadoxicality for a theory  $T$  are unacceptable. Let  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  be  $T$ 's classical recapturing principles and suppose  $\tau$  does not satisfy, in  $T$ , all the principles of classical logic. Now let's further assume that  $T$  is closed under principles to the effect that  $\varphi$  is paradoxical if and only if it satisfies  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$  only on pain of triviality and *unparadoxical* if and only if it satisfies each of  $\mathfrak{P}_1, \dots, \mathfrak{P}_n$ . Then,  $T$  doesn't prove that  $\tau$  is unparadoxical and, since  $\tau$  doesn't entail  $\perp$  in  $T$ , it also doesn't prove that it is paradoxical. That is,  $T$  is *silent* about  $\tau$ . But, by our lights, this need not be a problem: we shouldn't expect sentences such as  $\tau$  and 'Salzburg is East of Munich' to be assigned the same semantic status, if one behaves non-classically and the other doesn't.<sup>16</sup>

### 3 Revenge

Non-classical theories such as TSTV all share a common feature: despite their non-classicality, they have *fully classical* fragments. That is, they limit their restrictions to classical logic to *some* sentences. This is not only a basic fact about those theories; it also allows one to apply those theories to mathematics and science more generally. As it is sometimes said, non-classical theories can *recapture* classical reasoning when needed.<sup>17</sup>

To see how non-classical theories recapture classical theories, our starting point is a particularly simple way of characterising the classical fragment of non-classical theories such as TSTV and their extensions. Such theories all enjoy the following informal property:

(CLASSICALITY PRINCIPLES) There are finitely many classically valid principles such that a sentence satisfies such principles only if it satisfies all classical principles.

We can then say that a theory recaptures classical logic if it is closed under weaker versions of classical rules which, whenever some extra conditions are satisfied,

<sup>16</sup>We thank an anonymous referee for helping us appreciate this point.

<sup>17</sup>See e.g. ([24] p. 221), [7], ([1] pp. 111-2), [34].

reduce to their classical counterparts. In the case of TSTV, the weaker rule is WSRef and the extra condition is given by the following axiom:

$$(Id) \vdash \varphi \rightarrow \varphi$$

The following definition formally captures these ideas.

**Definition 2** (Classical recapture) Let  $S$  be a non-trivial theory. Then,  $S$  enjoys a classical recapturing property if it is  $\mathfrak{P}$ -classical recapturing for some principle  $\mathfrak{P}$  invalid in  $S$ . In particular,  $S$  is *ld-classical recapturing* if it is closed under the rules of CPL, where SRef is replaced by WSRef.

It's immediate to see why TSTV is ld-classical recapturing: if  $\varphi$  satisfies Id, then, courtesy of WSRef,  $\varphi$  also satisfies  $\varphi \vdash \varphi$ , whence by definition of CPL it also satisfies all the principles of classical logic.

Now, since  $\vdash \varphi \rightarrow \varphi$  is classical recapturing in TSTV, a sentence  $\varphi$  is paradoxical if the assumption that it satisfies  $\vdash \varphi \rightarrow \varphi$  leads to absurdity. Conversely, the claims that  $\varphi$  is paradoxical and that it satisfies  $\vdash \varphi \rightarrow \varphi$  must be jointly inconsistent. We can articulate these intuitive thoughts by means of the following higher-order rules [31]:<sup>18</sup>

$$\begin{array}{c} \square \frac{}{\Gamma \vdash \varphi \rightarrow \varphi} \text{ } n \\ \vdots \\ \frac{\Gamma \vdash \perp}{\Gamma \vdash \text{Par}(\ulcorner \varphi \urcorner)} \text{ld-Par-I, } n \end{array} \qquad \frac{\Gamma \vdash \text{Par}(\ulcorner \varphi \urcorner) \quad \Delta \vdash \varphi \rightarrow \varphi}{\Gamma, \Delta \vdash \perp} \text{ld-Par-E}$$

We employ a ‘higher-order’ rule *à la* Schroeder-Heister [see 31], since some non-reflexive logics that have actually been employed in the literature have no inferences, but only meta-inferences. For this reason, we need to be able to assume, and discharge, *sequents* in order to determine whether a sentence is paradoxical in our sense. Let’s now explicitly include these rules in our theory of naïve validity and truth.

**Definition 3** (TSTVP) TSTVP is the result of formulating TSTV in a language featuring a primitive unary predicate Par (call it  $\mathcal{L}_{V,P}$ ), and closing it under Id-Par-I and Id-Par-E.

We can now show that TSTVP is trivial.<sup>19</sup>

<sup>18</sup>Following [33], the box here indicates that the assumption  $\Gamma \vdash \varphi \rightarrow \varphi$  may not be discharged vacuously.

<sup>19</sup>As an anonymous referee points out, the following result could also be proven via a context-free version of our rule Id-Par-I, thus still weakening our assumptions about revenge-prone theories.

**Proposition 4** TSTVP is trivial, and so is the closure under Id-Par-I and Id-Par-E of any theory extending TSTV.

To prove this, we first prove the following Lemma.

**Lemma 5** Let  $\rho$  be a sentence identical to  $\text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)$ . Then, for any theory as strong as TSTVP, the following holds: if  $\vdash \rho \rightarrow \rho$ , then  $\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Tr}(\ulcorner \rho \urcorner)$ .

*Proof* Let  $\rho$  be a sentence identical to  $\text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)$ . We now reason in TSTVP, as follows:

$$\begin{array}{c}
 \frac{\frac{\vdash \rho \rightarrow \rho}{\rho \vdash \rho} \text{WSRef}}{\rho \vdash \text{Tr}(\ulcorner \rho \urcorner)} \text{Tr-I} \quad \frac{\frac{\vdash \rho \rightarrow \rho}{\rho \vdash \rho} \text{WSRef}}{\rho \vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)} \text{Def. of } \rho \\
 \hline
 \frac{\rho, \rho \vdash \text{Par}(\ulcorner \rho \urcorner)}{\rho \vdash \text{Par}(\ulcorner \rho \urcorner)} \text{SContr} \quad \vdash \rho \rightarrow \rho \quad \text{Id-Par-E} \\
 \hline
 \frac{\frac{\frac{\rho \vdash \perp}{\vdash \neg \rho} \neg\text{-I}}{\vdash \neg \text{Tr}(\ulcorner \rho \urcorner)} \neg\text{Tr-I}}{\neg \text{Tr}(\ulcorner \rho \urcorner) \vdash \neg \text{Tr}(\ulcorner \rho \urcorner)} \text{SWeak} \\
 \frac{\vdash \neg \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \neg \text{Tr}(\ulcorner \rho \urcorner)}{\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Tr}(\ulcorner \rho \urcorner)} \text{Contr}
 \end{array}$$

□

Making use of the Lemma, can now prove Proposition 5.

*Proof* Again, we reason in TSTVP. We begin by deriving  $\text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)$  from  $\vdash \rho \rightarrow \rho$ :

$$\begin{array}{c}
 \frac{\frac{\vdash \rho \rightarrow \rho}{\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Tr}(\ulcorner \rho \urcorner)} \text{Lemma 2}}{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Tr}(\ulcorner \rho \urcorner)} \text{WSRef} \\
 \frac{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Tr}(\ulcorner \rho \urcorner)}{\text{Tr}(\ulcorner \rho \urcorner) \vdash \rho} \text{Tr-E} \\
 \frac{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)}{\text{Tr}(\ulcorner \rho \urcorner), \text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Par}(\ulcorner \rho \urcorner)} \text{Def. of } \rho \\
 \frac{\text{Tr}(\ulcorner \rho \urcorner), \text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Par}(\ulcorner \rho \urcorner)}{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Par}(\ulcorner \rho \urcorner)} \text{SContr} \\
 \frac{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Par}(\ulcorner \rho \urcorner)}{\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)} \text{Contr}
 \end{array}$$

Call this derivation from the open assumption  $\vdash \rho \rightarrow \rho$  (taken twice),  $\mathcal{D}_6$ . We now use two copies of  $\mathcal{D}_6$  to prove that  $\rho$  is paradoxical:

$$\begin{array}{c}
 \frac{}{\vdash \rho \rightarrow \rho} 1 \\
 \mathcal{D}_6 \\
 \frac{\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)}{\vdash \rho} \text{Tr-I} \quad \text{Def. of } \rho \quad \frac{}{\vdash \rho \rightarrow \rho} 1 \\
 \frac{}{\vdash \text{Par}(\ulcorner \rho \urcorner)} \rightarrow\text{-E} \quad \frac{}{\vdash \rho \rightarrow \rho} 1 \\
 \frac{}{\vdash \perp} \text{Id-Par-I, 1}
 \end{array}$$

Call this second derivation  $\mathcal{D}_7$ . We now use two copies of  $\mathcal{D}_7$  to derive  $\perp$  from the empty set of assumptions:

$$\begin{array}{c}
 \mathcal{D}_7 \\
 \frac{}{\vdash \text{Par}(\ulcorner \rho \urcorner)} \text{SWeak} \\
 \frac{\text{Tr}(\ulcorner \rho \urcorner) \vdash \text{Par}(\ulcorner \rho \urcorner)}{\vdash \text{Tr}(\ulcorner \rho \urcorner) \rightarrow \text{Par}(\ulcorner \rho \urcorner)} \rightarrow\text{-I} \\
 \frac{}{\vdash \rho} \text{Def. of } \rho \\
 \frac{}{\rho \vdash \rho} \text{SWeak} \\
 \frac{}{\vdash \rho \rightarrow \rho} \rightarrow\text{-I} \\
 \frac{}{\vdash \perp} \text{Id-Par-E}
 \end{array}$$

□

### 4 Discussion

It may be objected that  $\mathcal{D}_7$  is unacceptable for someone accepting TSTV or some of its extensions. After all, one might argue, one cannot in general start a derivation from the assumption  $\vdash \rho \rightarrow \rho$  in any theory whose underlying logic is TS, for the simple reason that this would be tantamount to relying on the validity of SRef, which is however TS-invalid. To be sure, if we know that a sentence is unparadoxical, then we can use the corresponding instance of SRef. But that’s not the case with  $\rho$ : here *we don’t yet know*  $\rho$ ’s semantic status; whence we cannot reason classically about it – in particular, disallowing SRef means that we can’t start a derivation from  $\rho \vdash \rho$ , or, from that matters, from its TSUWSRef-equivalent,  $\vdash \rho \rightarrow \rho$ . Or so one might argue.

Let’s take this thought seriously for a moment. As we observed in Section 2, whether a sentence is paradoxical may depend on features of the context rather than on the sentence itself. This means, though, that in general we don’t know if a sentence is unparadoxical. We might perhaps know that the sentences of the language of arithmetic, or set theory, are not paradoxical. But we certainly lack such a knowledge for sentences of languages that are rich enough to express contingent Liars – including potentially harmless sentences such as ‘What Jane said is not true’. Thus, if we block derivations like  $\mathcal{D}_7$  on the grounds that they may involve contingent Liars, then we simply cannot in general assume, and hence reason about, sentences such as ‘What Jane said is not true’. However, disallowing the assumptions of such

potentially harmless sentences would not only cripple – as Field would put it – ordinary reasoning involving the truth predicate. It would also make little sense from a revisionary perspective.

Even if the logic is non-classical, one should be able to assume that a given sentence satisfies all the principles of classical logic, and see what follows from this. For instance, let our background logic be  $TS \cup WSRef$  and suppose we assume that the Liar sentence  $\lambda$  satisfies  $\vdash \lambda \rightarrow \lambda$ , and hence all the principles of classical logic. We can now proceed to derive  $\vdash \perp$  via familiar Liar-like reasoning – call his derivation  $\mathcal{D}_8$ :

$$\begin{array}{c}
 \frac{\frac{\vdash \lambda \rightarrow \lambda}{\lambda \vdash \lambda} \text{ WSRef} \quad \frac{\vdash \lambda \rightarrow \lambda}{\lambda \vdash \lambda} \text{ WSRef}}{\lambda \vdash \text{Tr}(\ulcorner \lambda \urcorner)} \text{ Tr-I} \quad \frac{\vdash \lambda \rightarrow \lambda}{\lambda \vdash \text{Tr}(\ulcorner \lambda \urcorner)} \text{ Def. of } \lambda \quad \frac{\frac{\vdash \lambda \rightarrow \lambda}{\lambda \vdash \lambda} \text{ WSRef} \quad \frac{\vdash \lambda \rightarrow \lambda}{\lambda \vdash \text{Tr}(\ulcorner \lambda \urcorner)} \text{ WSRef}}{\lambda \vdash \neg \text{Tr}(\ulcorner \lambda \urcorner)} \text{ Def. of } \lambda \\
 \frac{\lambda, \lambda \vdash \perp}{\lambda \vdash \perp} \text{ SContr} \quad \frac{\vdash \neg \lambda}{\vdash \neg \text{Tr}(\ulcorner \lambda \urcorner)} \neg\text{-I} \quad \frac{\lambda, \lambda \vdash \perp}{\lambda \vdash \perp} \text{ SContr} \quad \frac{\vdash \neg \lambda}{\vdash \neg \text{Tr}(\ulcorner \lambda \urcorner)} \neg\text{-I} \\
 \frac{\vdash \neg \lambda}{\vdash \neg \text{Tr}(\ulcorner \lambda \urcorner)} \neg\text{-Tr-I} \quad \frac{\vdash \lambda}{\vdash \text{Tr}(\ulcorner \lambda \urcorner)} \text{ Tr-I} \\
 \hline
 \vdash \perp
 \end{array}$$

Of course, the natural conclusion of such a reasoning is that  $\lambda$  may not satisfy all the principles of classical logic after all. That is,  $\lambda$  is paradoxical. We can capture this intuitive, and sound, diagnosis by means of our rule  $Id\text{-Par-I}$ :

$$\frac{\overline{\vdash \lambda \rightarrow \lambda}^1 \quad \overline{\vdash \lambda \rightarrow \lambda}^1 \quad \overline{\vdash \lambda \rightarrow \lambda}^1 \quad \overline{\vdash \lambda \rightarrow \lambda}^1}{\vdash \perp} \mathcal{D}_8 \quad \frac{\vdash \perp}{\vdash \text{Par}(\ulcorner \lambda \urcorner)} \text{ Id-Par-I, 1}$$

Here, assuming  $\vdash \lambda \rightarrow \lambda$  has caused us no harm. On the contrary, it has allowed us, courtesy of  $Id\text{-Par-I}$ , to establish that  $\lambda$  is paradoxical, and hence that it cannot be reasoned with classically. We now *know* that we shouldn't trust derivations which rely on undischarged copies of the sequent  $\vdash \lambda \rightarrow \lambda$ . However, this doesn't mean that such a sequent can never be used. It can be *assumed* and later discharged, as in the above derivation, in a proof of  $\vdash \text{Par}(\ulcorner \lambda \urcorner)$ .

This diagnosis can be extended to models of naïve truth and validity. For example, one can capture the extension of the intended notion of paradoxicality (according to which  $\varphi$  is paradoxical just in case it satisfies the classical recapturing principles, and hence full classical logic, only on pain of triviality) within  $\mathcal{I}_\psi$  as follows:

$$\text{Paradoxical in } \mathcal{L}_V := \{\varphi \in \mathcal{L}_V \mid \text{if } \vdash \varphi \rightarrow \varphi \in \mathcal{I}_\psi, \text{ then } \vdash \perp \in \mathcal{I}_\psi\}$$

Intuitively paradoxical sentences, such as  $\lambda$  and  $\pi$ , are clearly paradoxical in  $\mathcal{L}_V$  in the above sense. To be sure, the above notion of paradoxicality is unproblematic because it is limited to  $\mathcal{L}_V$  (and does not apply to the full language  $\mathcal{L}_{V,P}$ ). That is, it doesn't apply to sentences which themselves contain the paradoxicality predicate.

However, by our lights, this nonetheless establishes that the motivation behind the paradoxicality predicate is intuitively coherent.

But why shouldn't such a predicate be self-applicable, i.e. why shouldn't Par apply to sentences in which it itself figures? Just like typed notions of truth are usually disregarded as excessively restricted, especially by revisionary theorists, the same goes with paradoxicality. Consider again  $\mathcal{D}_7$ . Here, too, just like in  $\mathcal{D}_8$ , we assume  $\vdash \rho \rightarrow \rho$ , derive  $\vdash \perp$ , and conclude that  $\rho$  must be paradoxical after all, i.e.  $\vdash \text{Par}(\ulcorner \rho \urcorner)$ . Thus, we contend, the uses of SRef in derivations such as  $\mathcal{D}_7$  are motivated if the uses of SRef in  $\mathcal{D}_8$  are. To be sure, our revenge paradox, i.e. our proof of Proposition 5, shows that some of these seemingly natural uses of Par actually lead to triviality. Thus, the paradox points to some unavoidable limits of non-reflexive approaches. Indeed, since similar paradoxes can be produced for virtually any non-classical theory satisfying CLASSICALITY PRINCIPLES, the paradoxes of paradoxicality and unparadoxicality show the limits of revisionary approaches in general.

It might again be objected that to assume the sequent  $\vdash \rho \rightarrow \rho$  is effectively to make use of SRef in the metatheory. That is, where  $\models$  is the consequence relation of our metatheory, it might be argued that to assume the sequent  $\vdash \rho \rightarrow \rho$  is to assume the metatheoretic sequent

$$\vdash \rho \rightarrow \rho \models \vdash \rho \rightarrow \rho,$$

which is however just *another* instance of SRef. And, one might reason, if SRef fails in general, it must also fail in the metatheory. There's two problems with this rejoinder, however. To begin with, we haven't yet been given a reason to think that the metatheory  $M$  of a theory such as TSTV must be non-classical. Second, if it *is* non-classical, and one cannot in general make assumptions in  $M$ , there will be very little, if anything, any such metatheory will be able to prove.

Could either Id-Par-I or Id-Par-E be faulted instead? From a revisionary perspective, we see two problems with this. The first is that these principles are seemingly justified by reasonings such as  $\mathcal{D}_8$ , just like naïve principles for naïve truth and validity are seemingly justified by reasonings such as AGREEMENT and AGREEMENT\*. The second problem problem is that there simply are, in theories whose logic is TS, sentences which satisfy  $\varphi \rightarrow \varphi$ , and hence all the principles of classical logic, and sentences which only satisfy  $\varphi \rightarrow \varphi$  (and hence the principles of classical logic) on pain of triviality. Now, it seems to us, it's entirely legitimate to *assign labels* to these two kinds of sentences. For want of a better word, we've called sentences of the first kind unparadoxical and sentences of the second kind paradoxical. Given this terminological choice, for a sentence  $\varphi$  to be paradoxical in a theory such as TSTV *is* to entail  $\perp$  if reasoned with classically. And since we know that  $\varphi$  satisfies  $\varphi \rightarrow \varphi$  in TSTV if and only if it satisfies all the principles of classical logic, we also know that  $\varphi$  is paradoxical if and only if the assumption that it satisfies  $\varphi \rightarrow \varphi$  entails  $\perp$ . As far as we can see, Id-Par-I or Id-Par-E simply articulate this fact about TSTV. Thus, a theory's failure to justify these principles should be interpreted as a serious expressive limitation, of the same kind as the ones revealed by the Liar Paradox, Curry's Paradox, and the semantic paradoxes more generally.



**Acknowledgements** We wish to thank an anonymous reviewer for helpful comments on a previous draft as well as the FWF (project number 29716-G24) for generous financial support.

**Funding** Open access funding provided by Paris Lodron University of Salzburg.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Beall, J. (2009). *Spandrels of truth*. Oxford: Oxford University Press.
2. Beall, J., & Murzi, J. (2013). Two flavors of Curry's Paradox. *The Journal of Philosophy*, *CX*(3), 143–65.
3. Beall, J., & Ripley, D. (2014). Non-classical theories of truth. In M. Glanzberg (Ed.) *The Oxford Handbook of Truth*. Oxford: Oxford University Press.
4. Cobreros, P., Egré, P., Ripley, D., & van Rooij, R. (2012). Tolerant, classical, strict. *Journal of Philosophical Logic*, *41*(2), 347–85.
5. Cobreros, P., Egré, P., Ripley, D., & van Rooij, R. (2013). Reaching transparent truth. *Mind*, *122*, 841–866.
6. Curry, H. (1942). The inconsistency of certain formal logics. *Journal of Symbolic Logic*, *7*, 115–7.
7. Field, H. (2008). *Saving truth from paradox*. Oxford: Oxford University Press.
8. Field, H. (2017). Disarming a paradox of validity. *Notre Dame Journal of Formal Logic*, *58*(1), 1–19.
9. Fjellstad, A. (2017). Non-classical elegance for sequent calculus enthusiasts. *Studia Logica*, *105*(1), 93–119.
10. French, R. (2016). Structural reflexivity and the paradoxes of self-reference. *Ergo, an Open Access Journal of Philosophy*, *3*.
11. Halbach, V., & Horsten, L. (2006). Axiomatizing Kripke's theory of truth. *Journal of Symbolic Logic*, *71*, 677–712.
12. Heck, Jr. R. G. (2007). Self-reference and the languages of arithmetic. *Philosophia Mathematica*, *15*, 1–29.
13. Horsten, L. (2009). Levity. *Mind*, *118*(471), 555–581.
14. Horsten, L. (2012). *The Tarskian Turn. Deflationism and axiomatic truth*. Cambridge: MIT Press.
15. Kripke, S. (1975). Outline of a theory of truth. *Journal of Philosophy*, *72*, 690–716.
16. Mares, E., & Paoli, F. (2014). Logical consequence and the paradoxes. *Journal of Philosophical Logic*, *43*, 439–469.
17. McGee, V. (1991). *Truth, vagueness and paradox*. Indianapolis: Hackett Publishing Company.
18. McGee, V. (1992). Maximal consistent sets of instances of Tarski's schema (T). *Journal of Philosophical Logic*, *21*(3), 235–241.
19. Meadows, T. (2014). Fixed points for consequence relations. *Logique et Analyse*, 333–357.
20. Moschovakis, Y. (1974). Elementary induction on abstract structures. North-Holland.
21. Murzi, J., & Rossi, L. (2021). Generalised revenge. *Australasian Journal of Philosophy*, *98*(1), 153–177.
22. Murzi, J., & Rossi, L. (2018). Naïve validity, Synthese. Online first: <https://doi.org/10.1007/s11229-017-1541-6>.
23. Nicolai, C., & Rossi, L. (2018). Principles for object-linguistic validity: from logical to irreflexive. *Journal of Philosophical Logic*, *47*, 549–577.
24. Priest, G. (2006). *Doubt Truth to be a Liar*. Oxford: Oxford University Press.
25. Restall, G. (2000). *An introduction to substructural logics*. New York: Routledge.

26. Ripley, D. (2012). Conservatively extending classical logic with transparent truth. *Review of Symbolic Logic*, 354–78.
27. Ripley, D. (2013). Paradoxes and failures of cut. *Australasian Journal of Philosophy*, 91(1), 139–64.
28. Ripley, D. (2013). Revising up. *Philosophers' Imprint*, 13(5).
29. Ripley, D. (2015). Comparing substructural theories of truth. *Ergo*, 2(13), 299–328.
30. Rosenblatt, L. (2020). Maximal non-trivial sets of instances of your least favorite logical principle. *The Journal of Philosophy*, 117(1), 30–54.
31. Schroeder-Heister, P. (1984). A natural extension of natural deduction. *Journal of Symbolic Logic*, 49, 1284–1299.
32. Soames, S. (1999). *Understanding truth*. Oxford: Oxford University Press.
33. Tennant, N. (2012). Cut for core logic. *Review of Symbolic Logic*, 5(3), 450–479.
34. Zardini, E. (2011). Truth without contra(di)ction. *Review of Symbolic Logic*, 4, 498–535.
35. Zardini, E. (2013). It is not the case that [p and 'it is not the case that p' is true] nor is it the case that [p and 'p' is not true]. *Thought*, 1(4), 309–19.
36. Zardini, E. (2013). Naïve logical properties and structural properties. *The Journal of Philosophy*, 110(11), 633–44.
37. Zardini, E. (2014). Naïve truth and naïve logical properties. *Review of Symbolic Logic*, 7(2), 351–384.
38. Zardini, E. (2015). Getting one for two, or the contractors' bad deal. Towards a unified solution to the semantic paradoxes. In T. Achourioti, K. Fujimoto, H. Galinon, & J. Martinez-Fernandez (Eds.) *Unifying the Philosophy of Truth* (pp. 461–93). Springer.

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