# ON LANGUAGES THAT CONTAIN THEIR OWN
## UNGROUNDEDNESS PREDICATE[*]

### STEFAN WINTEIN

*Abstract*

Kripke's fixed point construction deals with paradoxical sentences such as the Liar or elements of a Yablo sequence by declaring them neither true nor false. We say that these sentences receive a third truth value, which we call *ungrounded*. We specify circumstances under which an interpreted language can — besides its truth and falsity predicate — contain its own *ungroundedness predicate*; in such a language the assertion '$\sigma$ is ungrounded' is true just in case $\sigma$ is in fact ungrounded. Then, our result is applied to shed light on a dubious claim that has recently been advanced in the literature with respect the so called *Hardest Logic Puzzle Ever*.

## 1. *Introduction*

The notorious Liar paradox causes serious problems for our intuitive understanding of truth. The literature's reactions to the Liar constitute a densely populated inhomogeneous area of theories. Formally, a boundary can be drawn between *axiomatic* and *model-theoretic* theory constructions while philosophically such a boundary can be drawn between theories that study the Liar phenomenon in an *ordinary language environment* or in an *environment of mathematical language*. In the philosophical area, theories of (self-referential) truth 'range from attempts to explicate our intuitive notion of truth to assigning the truth predicate a rôle in the foundations of mathematics'[1]. The contribution of this paper is a model-theoretic result that is obtained in the attempt to make sense of a particular phenomenon of self-referentiality occurring in ordinary language and thought.

[1] (Halbach & Horsten 2006), slight change of context.

*The model-theoretic result.* A very important technique for the construction of theories of self-referential truth is *Kripke's fixed point construction.* Starting with a classical structure (called a *ground structure*) for the truth-free fragment[2] of a language $L_T$ which contains a truth predicate '$T$' the construction generates, upon specification of a monotonic valuation scheme $V$, a partial structure for $L_T$. The associated interpreted language $\mathcal{L}_T$[3] has the so called *fixed point property*; the truth value — true (t), false (f) or ungrounded (u) — of a sentence $\sigma$ equals the truth value of the sentence which expresses that $\sigma$ is true; $\sigma$ and $T(\ulcorner\sigma\urcorner)$ are *semantically intersubstitutable.* When (Kripke 1975) writes that 'Being a fixed point $\mathcal{L}_T$ contains its own truth predicate' it is arguable that his reason for calling $\mathcal{L}_T$ a language that contains its own truth predicate is precisely that it has the mentioned semantic intersubstitutability property.

In fact, Kripke's construction may also be applied to a ground structure to obtain an language $\mathcal{L}_{TF}$ which, in the sense in which $\mathcal{L}_T$ is a language that contains its own truth predicate, is a language that contains its own truth and falsity predicate. However, the reason for calling $\mathcal{L}_{TF}$ a language that contains its own falsity predicate is obviously not the semantic intersubstitutability of $\sigma$ with $F(\ulcorner\sigma\urcorner)$. As '$T$' and '$F$' are both *truth value predicates*, used to express that a sentence has truth value t respectively f, one may ask for a specification of general conditions that have to be fulfilled for an interpreted language to contain its own truth value predicate(s) in the sense alluded to by Kripke. I take it that a *necessary* condition for a language to contain its own truth value predicate with respect to a certain truth value is that the language is *Truth Value Correct* ($TVC$) with respect to that value. A language is $TVC$ with respect to truth value v just in case whenever a sentence $\sigma$ has truth value v, the sentence which expresses that $\sigma$ has truth value v is true. As $\mathcal{L}_{TF}$ is $TVC$ with respect to t and f, the question arises whether a language can also be truth value correct with respect to the truth value *ungrounded*.

This paper's model-theoretic result, called the *paradoxical $TVC$ theorem* (partially) answers this question. The theorem states that any $\Delta$-*neutral* ground structure can be expanded to a structure for a language $L_{TFU}$ — containing a truth, falsity and Ungroundedness predicate — such that the associated interpreted language $\mathcal{L}_{TFU}$ is truth value correct with respect to t, f and u. In a $\Delta$-neutral structure, we have the ability to form Liar sentences, Yablo

---

[2] In this paper, we will only consider Kripkean fixed point constructions that are carried out starting with an empty extension and anti-extension of the truth predicate; we only consider *minimal* fixed points.

[3] $\mathcal{L}_T = \langle L_T, M, V \rangle$, with $V$ a monotonic valuation scheme and $M$ the partial structure for $L_T$ that is obtained by carrying out the minimal fixed point construction with $V$.

sequences ((Yablo 1993)) and, in fact, to form any self-referential construction[4] whatsoever using the predicates '$T$' and '$F$'. However, $\Delta$-neutrality excludes the formation of self-referential constructions which use the ungroundedness predicate '$U$'. Sentences like 'this sentence is ungrounded' or 'this sentence is ungrounded or it is false' have no formal counterpart in a $\Delta$-neutral structure.

The paper is organized as follows. Section 2 sets up notation and Section 3 is used to review two theorems, one due to Gupta and one due to Kripke, in light of our notion of truth value correctness. Section 4 is devoted to the proof of the paradoxical $TVC$ theorem, which involves a combination of Kripkean fixed point techniques with *revisionist techniques* and which is inspired by Gupta's proof of a theorem which can be found in Section 3 (Theorem 1) of this paper. In Section 5 we sketch an application of our theorem; we show how the paradoxical $TVC$ theorem sheds light on the status of an interesting — though obscure — argument involving self-referentiality that has recently been advanced in the literature ((Rabern & Rabern 2008)) with respect to the *Hardest Logic Puzzle Ever* ($HLPE$).

## 2. *Preliminaries*

We identify a first order language $L$ with its set of non-logical constants and we assume that '$=$' is a logical constant, expressing the identity relation. With $n \geq 1$, $Con(L)$, $Pred^n(L)$ and $Fun^n(L)$ are used to denote the set of all constant symbols, $n$-ary predicate symbols and $n$-ary function symbols of $L$ respectively. $Pred(L)$ and $Fun(L)$ denote the set of all predicate respectively function symbols so that $L = Con(L) \cup Pred(L) \cup Fun(L)$. The set of sentences of $L$ (constructed in the usual manner) will be denoted as $Sen(L)$, its set of closed terms as $Cterm(L)$. A *structure* for $L$ is a pair $M = \langle D, I \rangle$ consisting of a domain $D$ and a function $I$ that interprets $L$. With $c \in Con(L)$, $f \in Fun^n(L)$ we have $I(c) \in D$ and $I(f) \in D^{D^n}$. With $R \in Pred^n(L)$, $I(R) = (R^+, R^-) \in \mathcal{P}(D^n) \times \mathcal{P}(D^n)$ *such that* $R^+ \cap R^- = \emptyset$.[5] Whenever $R^- = D^n - R^+$ for each $n \geq 1$ and $R \in Pred^n(L)$, we say that $M$ is *classical*, otherwise $M$ is *non-classical*. A *valuation scheme* $V$ assigns a function $V_M : Sen(L) \to \{t, f, u\}$ to each structure $M$ for $L$. Here $\{t, f, u\}$ is the set of *truth values*; $\sigma$ can be true (t),

---

[4] Whether or not a Yablo sequence is a genuine manifestation of self-reference is a controversial issue.

[5] $R^+$ is called the *extension*, $R^-$ the *anti-extension* of $R$. Indeed, the definition of a structure in this paper is such that the extension and anti-extension are always disjoint.

false (f) or *ungrounded* (u). The classical (Tarskian) valuation scheme ($\mathcal{C}$), the Strong Kleene scheme ($SK$), the Weak Kleene scheme ($WK$) and the Supervaluation scheme ($SV$) — and only these schemes — we call *appropriate (valuation) schemes*. Note that $\mathcal{C}$ is only defined for classical structures, whereas the other appropriate schemes are defined for all structures. Any appropriate scheme $V$ is *normal*[6] meaning that whenever $M$ is a classical structure for some language $L$, $V_M(\sigma) = \mathcal{C}(\sigma)$ for all $\sigma \in Sen(L)$. We will use $den_M \subseteq Cterm(L) \times D$ for the denotation relation in structure $M = \langle D, I \rangle$; $\langle t, d \rangle \in den_M$ just in case $t$ denotes $d$ in $M$. Whenever we write '$\langle t, \sigma \rangle \in den_M$', let it be understood that $t$ denotes a *sentence $\sigma$* of the language under consideration.

*Definition 1*: Quotational closure
*Let $L$ be an arbitrary first order language. We set $L^0 = L$ and define:*

- $L^{n+1} = L^n \cup \{[\sigma] \mid \sigma \in Sen(L^n)\}$, $n \geq 0$
- $\bar{L} = \bigcup_{i=0}^{\infty} L^i$

*When $\sigma$ is a sentence of $L^n$, $[\sigma]$ is a* constant symbol *of $L^{n+1}$. $\bar{L}$ is the* quotational closure *obtained from $L$ and $\{L^n\}_{n \in \mathbb{N}}$ is the* quotational hierarchy *of $\bar{L}$. Note that $m \leq n \Rightarrow L^m \subseteq L^n$. Any language $\bar{L}$, obtained as the quotational closure of some language $L$, is called a* quotational language.

A quotational language $\bar{L}$ will be interpreted by a *sentence structure*, which is a structure that has the sentences of $\bar{L}$ as objects in its domain and which has *at least* one closed term (the quotational constant $[\sigma]$) referring to each sentence $\sigma$ of $\bar{L}$.

*Definition 2*: Sentence structures
*A* sentence structure *$M = \langle D, I \rangle$ is a structure for a quotational language $\bar{L}$ such that:*

(1) $Sen(\bar{L}) \subseteq D$
(2) $I([\sigma]) = \sigma$ for all $\sigma \in Sen(\bar{L})$.

*Thus the domain of a sentence structure $M = \langle D, I \rangle$ for $\bar{L}$ consists of the sentences of $\bar{L}$ and $\mathcal{O}$ther objects. We use $\mathcal{O}_M = D - Sen(\bar{L})$ for the set of non-sentential objects in $M$'s domain.*[7]

---

[6] Throughout the paper, 'appropriate scheme' is interchangeable with 'monotonic normal scheme'.

[7] So the sentences themselves — rather than their (Gödel) codes — populate our domain. This approach is not uncommon in the literature. For instance, see (Gupta 1982) or (Gupta & Belnap 1993).

When $\bar{L}$ is some quotational language, we use $\mathcal{L}$ to range over all triples $\langle \bar{L}, M, V \rangle$, where $M$ is a sentence structure for $\bar{L}$ and where $V$ is an appropriate scheme *that is defined for* $M$.[8] When $V = \mathcal{C}$, we say that $\mathcal{L}$ is classical, otherwise, we say that $\mathcal{L}$ is non-classical.[9]

*Definition 3*: Ground structures and their expansions
*Let $\bar{L}$ be a quotational language and let $P \subseteq Pred(\bar{L})$. We say that $\hat{M} = \langle D, I \rangle$ is a* ground structure *for $\bar{L} - P$ just in case $\hat{M}$ is a* classical *structure for $\bar{L} - P$ such that:*
  (1) *$Sen(\bar{L}) \subseteq D$*
  (2) *$I([\sigma]) = \sigma$ for all $\sigma \in Sen(\bar{L})$.*
*When $\hat{M}$ is a ground structure for $\bar{L} - P$, $M$ is an $\bar{L}$-expansion of $\hat{M}$ when $M$ is a structure for $\bar{L}$ such that the domains of $\hat{M}$ and $M$, as well as their respective interpretations of $\bar{L} - P$, are identical.*

As an example, let $L_T$ be a language containing the truth predicate symbol '$T$' and, amongst others, the constant symbol '$\lambda$'. A ground structure for $\bar{L}_T - \{T\}$ has all the sentences of $\bar{L}_T$ in its domain and so in particular the sentence '$\neg T(\lambda)$'. When the ground structure is such that $I(\lambda) = \neg T(\lambda)$, the sentence '$\neg T(\lambda)$' intuitively says of itself that it is not true; '$\neg T(\lambda)$' is a (strengthened) Liar sentence. In the next section, we review two theorems, one due to Gupta and one due to Kripke, that can be interpreted as specifying the conditions under which a ground structure $\hat{M}$ for $\bar{L}_T - \{T\}$ can be $\bar{L}_T$-expanded to $M$ such that $\mathcal{L}_T = \langle \bar{L}_T, M, V \rangle$ contains its own truth predicate. In Section 4, we extend the results of Gupta and Kripke by specifying the conditions under which a ground structure $\hat{M}$ for $\bar{L}_{TFU} - \{T, F, U\}$ can be $\bar{L}_{TFU}$-expanded to $M$ such that $\mathcal{L}_{TFU} = \langle \bar{L}_{TFU}, M, V \rangle$ contains its own truth, falsity and ungroundedness predicate; we prove our *paradoxical $TVC$ theorem*.

## 3. *Truth Value Correctness*

A *truth value predicate* is a predicate that is used to express the assertion that a sentence has a certain truth value; the unary predicate symbols $T$, $F$ and $U$ will be used to express that a sentence is true, false or ungrounded respectively.

---

[8] Thus, when $M$ is non-classical, $V \neq \mathcal{C}$.

[9] Note that a non-classical $\mathcal{L}$ may have classical $M$.

*Definition 4*:  Truth value correctness
*Let L be a language and let* $L_{TFU} = L \cup \{T, F, U\}$. $\mathcal{L}_{TFU} = \langle \bar{L}_{TFU}, M, V \rangle$
*is said to be* Truth Value Correct *with respect to* $\mathsf{v} \in \{\mathsf{t}, \mathsf{f}, \mathsf{u}\}$ *just in case*
$TVC_{\mathsf{v}}$ *holds:*

- $TVC_{\mathsf{t}}$ :     $V_M(\sigma) = \mathsf{t} \Leftrightarrow V_M(T(t)) = \mathsf{t}$,    *for all* $\langle t, \sigma \rangle \in den_M$
- $TVC_{\mathsf{f}}$ :     $V_M(\sigma) = \mathsf{f} \Leftrightarrow V_M(F(t)) = \mathsf{t}$,    *for all* $\langle t, \sigma \rangle \in den_M$
- $TVC_{\mathsf{u}}$ :     $V_M(\sigma) = \mathsf{u} \Leftrightarrow V_M(U(t)) = \mathsf{t}$,    *for all* $\langle t, \sigma \rangle \in den_M$

$\mathcal{L}_{TFU}$ *is* truth value correct *(TVC) just in case it is truth value correct with
respect to each* $\mathsf{v} \in \{\mathsf{t}, \mathsf{f}, \mathsf{u}\}$. *When* $TVC_{\mathsf{t}}$ *($TVC_{\mathsf{f}}$, $TVC_{\mathsf{u}}$) holds, we say
that* $\mathcal{L}_{TFU}$ *contains its own truth (falsity, ungroundedness) predicate.*[10]

Note that, for classical $\mathcal{L}_{TFU}$, $TVC_{\mathsf{t}}$ is equivalent to (1):

$$\mathcal{C}_M(T(t) \leftrightarrow \sigma) = \mathsf{t} \qquad \text{for all } \langle t, \sigma \rangle \in den_M \tag{1}$$

So, for classical $\mathcal{L}_{TFU}$, $TVC_{\mathsf{t}}$ is equivalent to the truth of all instances of
the notorious *T-scheme*. Not every classical $\mathcal{L}_{TFU}$ can contain its own truth
predicate. For instance, let $\mathcal{L}_{TFU}$ be such that $\langle \lambda, \neg T(\lambda) \rangle \in den_M$. The
sentence '$\neg T(\lambda)$', which intuitively says of itself that it is not true is a
(strengthened) *Liar sentence*. Instantiating (1) with the Liar sentence gives
'$\mathcal{C}_M(T(\lambda) \leftrightarrow \neg T(\lambda)) = \mathsf{t}$', which is impossible. In order to define an in-
teresting class of classical $\mathcal{L}_{TFU}$ that do contain their own truth predicate,
we need to define the notion of *X-neutrality*. Let $M = \langle D, I \rangle$ be a sentence
structure — for some $\bar{L}$ — and let $X \subseteq D$. With $\vec{d} = \langle d_1, \ldots, d_n \rangle \in D^n$,
we say that $\vec{d'} \in D^n$ is an *X-swap* of $\vec{d}$ just in case for every $1 \leq i \leq n$
we have that $d_i \notin X \Rightarrow d_i' = d_i$ and that $d_i \in X \Rightarrow d_i' \in X$. We use
$X(\vec{d}) \subseteq D^n$ to denote the set of all *X*-swaps of $\vec{d}$. We are now ready to
define the notion of *X*-neutrality.

*Definition 5*:  $X$-neutrality
*Let* $\bar{L}$ *be a quotational language, let* $M = \langle D, I \rangle$ *be a sentence structure for*
$\bar{L}$ *and let* $X \subseteq Sen(\bar{L})$. *We say that* $M$ *is* $X$-neutral *just in case, for every*
$f \in Fun(\bar{L})$ *and every* $R \in Pred(\bar{L})$ *we have that:*

(1) $\sigma \in X$ *and* $\langle t, \sigma \rangle \in den_M \Rightarrow t = [\sigma]$
(2) $\vec{d} \in R^+ \Leftrightarrow X(\vec{d}) \subseteq R^+$,    $\vec{d} \in R^- \Leftrightarrow X(\vec{d}) \subseteq R^-$
(3) $I(f)(\vec{d}) = I(f)(\vec{d'})$          *for all* $\vec{d'} \in X(\vec{d})$

---

[10] Although it can be argued that $TVC_{\mathsf{v}}$ is merely a *necessary* condition for a language to
contain the corresponding truth value predicate, we will not discuss this issue in this paper
and use the phrase 'containing its own truth (falsity, ungroundedness) predicate' as inter-
changeable with $TVC_{\mathsf{t}}$ ($TVC_{\mathsf{f}}$, $TVC_{\mathsf{u}}$).

Thus, in an $X$-neutral structure, 1) the *only* closed term that refers to a sentence $\sigma \in X$ is its quotational name $[\sigma]$, 2) interchanging $X$ members for (other) $X$ members in a tuple $\langle d_1, \ldots, d_n \rangle$ does not change the membership relation of the tuple with respect to extensions and anti-extensions of $n$-ary predicates, 3) nor does it change the output of an $n$-ary function on the tuple.

(Gupta 1982) showed that every ground structure $\hat{M}$ for $\bar{L}_T - \{T\}$[11] that is $Sen(\bar{L}_T)$-*neutral* can be $\bar{L}_T$-expanded, using a *revision process*, to a classical structure $M$ such that the classical $\mathcal{L}_T$ associated with $M$ contains its own truth predicate. However, his results are easily generalizable, delivering the following theorem.

*Theorem 1*:  Non-paradoxical $TVC$ theorem (Gupta)
*Let $V$ be an appropriate scheme and let $\hat{M} = \langle D, I \rangle$ be a ground structure for $\bar{L}_{TFU} - \{T, F, U\}$ that is $Sen(\bar{L}_{TFU})$-neutral. There exists a classical $\bar{L}_{TFU}$-expansion $M$ of $\hat{M}$ such that $\mathcal{L}_{TFU} = \langle \bar{L}_{TFU}, M, V \rangle$ is truth value correct.*

*Proof.* See (Gupta 1982) or (Gupta & Belnap 1993) for a proof in terms of $\mathcal{L}_T$ and carry out the necessary modifications, interpreting $U$ with $(\emptyset, D)$ to obtain a proof for classical $\mathcal{L}_{TFU}$. As the expansion of $\hat{M}$ is classical and as any appropriate valuation scheme is normal, it follows that the theorem in fact holds for any $\mathcal{L}_{TFU}$.                                □

The reason that I baptized this result of Gupta the *non-paradoxical $TVC$* theorem is that the conditions for the theorem (i.e. $Sen(\bar{L}_{TFU})$-neutrality) explicitly forbid the formation of well-known "paradoxical" sentences such as the Liar or elements of a *Yablo sequence* ((Yablo 1993)). In a $Sen(\bar{L}_{TFU})$-neutral structure we can only refer to a sentence $\sigma$ of $\bar{L}_{TFU}$ via its quote name $[\sigma]$ and so the formation of a Liar sentence using an interpretation function $I$ and constant $\lambda$ which are such that $I(\lambda) = I([\neg T(\lambda)]) = \neg T(\lambda)$, is excluded in a $Sen(\bar{L}_{TFU})$-neutral structure. Conditions 2 and 3 of Definition 5 in terms of $Sen(\bar{L}_{TFU})$-neutrality guarantee that a Liar sentence can neither be constructed in a manner that is more common in the literature, which is via a *substitution function*.

(Kripke 1975) showed that, in the presence of paradoxical sentences as the Liar, a language can still be truth value correct with respect to t and f.

---

[11] $L_T = L \cup \{T\}$ for some language $L$.

*Theorem 2*: Paradoxical $TVC$ theorem for $\{\mathsf{t},\mathsf{f}\}$ (Kripke)
*Let $L$ be a language, $L_{TF} = L \cup \{T, F\}$ and let $\hat{M}$ be any ground structure for $\bar{L}_{TF} - \{T, F\}$. When $V$ is any non-classical appropriate scheme, there exists a $\bar{L}_{TF}$-expansion $M$ of $\hat{M}$ such that $\mathcal{L}_{TF} = \langle \bar{L}_{TF}, M, V \rangle$ is truth value correct with respect to $\mathsf{t}$ and $\mathsf{f}$.*

*Proof.* See (Kripke 1975). $\qquad\qquad\qquad\qquad\qquad\qquad\square$

A language $\mathcal{L}_{TF}$ that is obtained via Theorem 2 declares the Liar sentence to be ungrounded, just as the sentences that ascribe truth or falsity to the Liar sentence. However, the languages $\mathcal{L}_{TF}$ considered by Kripke do not contain an ungroundedness *predicate*, so that an assertion like 'the Liar sentence is ungrounded' has no formal representation in $\mathcal{L}_{TF}$. Thus, the question arises whether a language $\mathcal{L}_{TFU}$ can, in the presence of paradoxical sentences, be truth value correct *tout court*. In the next section, we will prove a theorem that specifies conditions under which the answer to this question is 'yes'. The proof of this *paradoxical $TVC$ theorem* is inspired by Gupta's proof of Theorem 1.

4. *The paradoxical $TVC$ theorem*

Let $L$ be a first order language and let $\bar{L}_T$ and $\bar{L}_{TU}$ denote the quotational closure of $L \cup \{T\}$ and $L \cup \{T, U\}$ respectively. We let $\Delta = Sen(\bar{L}_{TU}) - Sen(\bar{L}_T)$ and for each $n \in \mathbb{N}$, we let $\Delta_n = \Delta \cap Sen(L_{TU}^n)$. The paradoxical $TVC$ theorem will be immediate, once we have proven the following lemma.

*Lemma 1*: The paradoxical $TVC$ lemma for $\{\mathsf{t}, \mathsf{u}\}$
*Let $V$ be a non-classical appropriate valuation scheme and $\hat{M}$ a $\Delta$-neutral ground structure for $\bar{L}_{TU} - \{T, U\}$. Then, $\hat{M}$ can be $\bar{L}_{TU}$-expanded to $M$ such that $\mathcal{L}_{TU} = \langle \bar{L}_{TU}, M, V \rangle$ is truth value correct with respect to $\mathsf{t}$ and $\mathsf{u}$.*

The $\bar{L}_{TU}$-expansion referred to in Lemma 1 will be constructed from the ground structure $\hat{M}$ via a two stage process. The first stage, called $FP$, uses a fixed point construction, the second stage, called $RP$, uses revisionist techniques.

$FP$. Let $\hat{M}$ be a $\Delta$-neutral structure and let $V$ be a non-classical appropriate valuation scheme. For any ordinal $\alpha$, let $M_\alpha = \hat{M}(T_\alpha^+, T_\alpha^-)$ denote the $\bar{L}_{TU}$-expansion of $\hat{M}$ that interprets $T$ as $(T_\alpha^+, T_\alpha^-)$ and that (classically)

interprets $U$ as $(\emptyset, D)$. Let $M_0 = \hat{M}(\emptyset, \mathcal{O}_{\hat{M}})$ and, for $\alpha > 0$, define $M_\alpha$ as follows.

$$SUC: \quad \alpha = \beta + 1 : \left\{ \begin{array}{l} T_\alpha^+ = \{\sigma \in Sen(\bar{L}_T) \mid V_{M_\beta}(\sigma) = \mathsf{t}\} \\ T_\alpha^- = \{\sigma \in Sen(\bar{L}_T) \mid V_{M_\beta}(\sigma) = \mathsf{f}\} \cup \mathcal{O}_{\hat{M}} \end{array} \right.$$

$$LIM: \quad \alpha \text{ is limit} : \left\{ \begin{array}{l} T_\alpha^+ = \bigcup_{\beta < \alpha} T_\beta^+ \\ T_\alpha^- = \bigcup_{\beta < \alpha} T_\beta^- \end{array} \right.$$

By well-known arguments, the sequence of structures $\{M_\alpha\}_{\alpha \in On}$ has a fixed point, i.e. there exists an ordinal after which further applications of $SUC$ and $LIM$ do not change the resulting structures anymore. We call this fixed point structure $M^*$.

$RP$.  For any ordinal $\alpha$, let $M_\alpha^* = \hat{M}(T_\alpha^+, T_\alpha^-, U_\alpha^+)$ denote the $\bar{L}_{TU}$-expansion of $\hat{M}$ that interprets $T$ as $(T_\alpha^+, T_\alpha^-)$ and that interprets $U$ (classically) as $(U_\alpha^+, D - U_\alpha^+)$. We set $M_0^* = M^*$ and define for each $\alpha > 0$, $M_\alpha^*$ as follows.

$$SUC': \quad \alpha = \beta+1 : \left\{ \begin{array}{l} T_\alpha^+ = \{\sigma \in Sen(\bar{L}_{TU}) \mid V_{M_\beta^*}(\sigma) = \mathsf{t}\} \\ T_\alpha^- = \{\sigma \in Sen(\bar{L}_{TU}) \mid V_{M_\beta^*}(\sigma) = \mathsf{f}\} \cup \mathcal{O}_{\hat{M}} \\ U_\alpha^+ = \{\sigma \in Sen(\bar{L}_{TU}) \mid V_{M_\beta^*}(\sigma) = \mathsf{u}\} \end{array} \right.$$

$$LIM': \quad \alpha \text{ is limit:} \left\{ \begin{array}{l} T_\alpha^+ = \{\sigma \in Sen(\bar{L}_{TU}) \mid \exists \beta : \sigma \in \bigcap_{\beta < \gamma < \alpha} T_{M_\gamma^*}^+)\} \\[2mm] T_\alpha^- = \{\sigma \in Sen(\bar{L}_{TU}) \mid \exists \beta : \sigma \in \bigcap_{\beta < \gamma < \alpha} T_{M_\gamma^*}^-\} \\[2mm] U_\alpha^+ = \{\sigma \in Sen(\bar{L}_{TU}) \mid \exists \beta : \sigma \in \bigcap_{\beta < \gamma < \alpha} U_{M_\gamma^*}^+\} \end{array} \right.$$

In order to prove Lemma 1, we will need the following lemma.

*Lemma 2*: Stabilization lemma
*Let $V$ be a non-classical appropriate valuation scheme. Let $\hat{M}$ be a $\Delta$-neutral structure for $\bar{L}_{TU} - \{T, U\}$ and let $\{M_\alpha^*\}_{\alpha \in On}$ be the series of structures generated from $\hat{M}$ via $FP$ and $RP$. Then, for all $n \in \mathbb{N}$ and $\alpha \in On$ such that $\alpha \geq n + 1$*

$$\sigma \in Sen(L_{TU}^n) \Rightarrow V_{M_{n+1}^*}(\sigma) = V_{M_\alpha^*}(\sigma) \tag{2}$$

*Proof.*  Suppose that the lemma is false. Then there has to be a least natural number, say $n'$ for which it fails and, given $n'$ there has to be a least ordinal $\geq n' + 1$, say $\alpha'$ such that $M_{n'+1}^*$ and $M_{\alpha'}^*$ disagree about the truth value of

$\sigma \in Sen(L_{TU}^{n'})$. Thus from the hypothesis that the lemma is false and the minimality of $n'$ and $\alpha'$ we get:

$\mathsf{C}_1$ : For all $n < n', \alpha \geq n + 1, \sigma \in Sen(L_{TU}^{n}) : V_{M_{n+1}^*}(\sigma) = V_{M_\alpha^*}(\sigma)$

$\mathsf{C}_2$ : For all $\alpha$ s.t. $n' + 1 \leq \alpha < \alpha', \sigma \in Sen(L_{TU}^{n'}) : V_{M_{n'+1}^*}(\sigma) = V_{M_\alpha^*}(\sigma)$

$\mathsf{C}_3$ : $\exists \sigma \in Sen(L_{TU}^{n'}) : V_{M_{n'+1}^*}(\sigma) \neq V_{M_{\alpha'}^*}(\sigma)$

We will show that these 3 conditions can not (jointly) hold, contradicting the hypothesis of the falsity of the lemma. From the definition of $RP$ it follows that $\alpha'$ has to be a successor ordinal, say $\alpha' = \beta + 1$. The structures $M_{n'+1}^*$ and $M_{\beta+1}^*$ only differ with respect to the interpretation of the predicate symbols $T$ and $U$. By definition of $RP$, these interpretations are fully determined by the functions $V_{M_{n'}^*}(\cdot)$ and $V_{M_\beta^*}(\cdot)$ respectively. As these functions valuate $\aleph_0$ different sentences of $\Delta - \Delta_{n'-1}$ to be true (false),[12] there exists a bijection $\chi : \Delta - \Delta_{n'-1} \to \Delta - \Delta_{n'-1}$ such that — with $X = T^+, T^-$ or $U^+$:

$$\forall \sigma \in (\Delta - \Delta_{n'-1}) : \sigma \in X_{n'+1} \Leftrightarrow \chi(\sigma) \in X_{\alpha'} \tag{3}$$

We extend $\chi$ to a bijection[13] from $D$ to $D$, by specifying that $\chi$ acts as the identity function on objects in $D - (\Delta - \Delta_{n'-1})$. We will show that $\chi$ is an isomorphism between the structures $M_{n'+1}^*$ and $M_{n'+\alpha'}^*$ in the language $L_{TU}^{n'}$. From the fact that isomorphic structures in a language are elementary equivalent w.r.t. the sentences of that language, it then follows that there cannot be a $\sigma \in Sen(L_{TU}^{n'})$ such that $M_{n'+1}^*$ and $M_{n'+\alpha'}^*$ disagree about the truth value of $\sigma$. Hence, we establish a contradiction with $\mathsf{C}_3$ and, consequently, with the hypothesis that the lemma is false. By definition of an isomorphism between structures, in order to show that $\chi$ is an isomorphism between $M_{n'+1}^*$ and $M_{\alpha'}^*$ in the language $L_{TU}^{n'}$, we need to establish that, for every $n \in \mathbb{N}$ and $\langle d_1, \ldots, d_n \rangle \in D^n$:

---

[12] Let $\Delta_{-1} = \emptyset$.

[13] Note that the existence of this bijection depends *only* on the fact that $M_0^*$ valuates $\aleph_0$ sentences as t and $\aleph_0$ sentences as f, which make $V_{M_{n'}^*}(\cdot)$ and $V_{M_\beta^*}(\cdot)$ do so too. $M_0^*$ either valuates $\aleph_0$ or $0$ sentences as u and so $V_{M_{n'}^*}(\cdot)$ and $V_{M_\beta^*}(\cdot)$ respectively valuate either $\aleph_0$ or $0$ sentences as u. In both cases, cardinality considerations show that the function $\chi$ satisfying (3) can be found and so the sought for bijection exists.

(1) For every $R \in Pred^n(L_{TU}^{n'})$:
  $\langle d_1, \ldots, d_n \rangle \in R_{n'+1}^+$ iff $\langle \chi(d_1), \ldots, \chi(d_n) \rangle \in R_{\alpha'}^+$
  $\langle d_1, \ldots, d_n \rangle \in R_{n'+1}^-$ iff $\langle \chi(d_1), \ldots, \chi(d_n) \rangle \in R_{\alpha'}^-$

(2) For every $f \in Fun^n(L_{TU}^{n'})$:
  $\chi(I(f)(d_1, \ldots, d_n)) = I(f)(\chi(d_1), \ldots, \chi(d_n))$

(3) For every $c \in Con(L_{TU}^{n'})$: $\chi(I(c)) = I(c)$

*ad 1.* When $R \notin \{T, U\}$, the claim readily follows from the fact that $\chi$ acts as the identity function on $D - \Delta$ and that $\hat{M}$ is $\Delta$-neutral structure. So let $R \in \{T, U\}$. Observe that $d \in D$ implies that $d$ is either an element of $\mathcal{O}_{\hat{M}}$, $\Delta_{n-1}$, $\Delta - \Delta_{n'-1}$ or $Sen(\bar{L}_T)$. When $d \in \mathcal{O}_{\hat{M}}$, the claim follows from the fact that $d \in T_{n'+1}^- \cap T_{\alpha'}^-$ and that $\chi$ acts as the identity function on $\mathcal{O}_{\hat{M}}$. When $d \in \Delta_{n'-1}$ the claim follows from $C_1$ and the definition of $RP$. When $d \in \Delta - \Delta_{n'-1}$ the claim follows from (3). Finally, let $d \in Sen(\bar{L}_T)$. Observe that, as $M^*$ results from $FP$, $M^*$ is a fixed point structure "with respect to the sentences of $\bar{L}_T$", i.e.:

$$V_{M^*}(\sigma) = V_{M_\alpha^*}(\sigma) \text{ for all } \alpha \in On, \sigma \in Sen(\bar{L}_T) \tag{4}$$

From (4) and the definition of $RP$ it follows, — with $X = T^+, T^-$ or $U^+$ — that:

$$\sigma \in X_{M_1^*} \Leftrightarrow \sigma \in X_{M_{1+\alpha}^*} \text{ for all } \alpha \in On, \sigma \in Sen(\bar{L}_T) \tag{5}$$

Now the claim follows from (5) and the fact that $\chi$ acts as the identity function on $Sen(\bar{L}_T)$.

*ad 2.* The claim follows from the fact that $M$ is an $\Delta$-neutral structure and that $\chi$ only permutes elements of $\Delta$.

*ad 3.* When $c$ denotes an element $d \notin \Delta$, the claim follows from the fact that $\chi(d) = d$. Whenever $c \in Con(L_{TU}^{n'})$ denotes an element of $\Delta$, $\Delta$-neutrality of $\hat{M}$ guarantees that it denotes an element of $\Delta_{n'-1}$, on which $\chi$ also acts as the identity function. $\qquad\qquad \square$

Lemma 1 and, in fact, the paradoxical $TVC$ theorem now follow easily.

*Theorem 3*: The paradoxical $TVC$ theorem
*Let $L$ be a first order language, $L_{TFU} = L \cup \{T, F, U\}$ and $L_{TF} = L \cup \{T, F\}$. Let $\bar{\Delta} = Sen(\bar{L}_{TFU}) - Sen(\bar{L}_{TF})$. Let $V$ be a non-classical*

*appropriate valuation scheme and let $\hat{M}$ be a $\bar{\Delta}$-neutral ground structure for $\bar{L}_{TFU} - \{T, F, U\}$. Then, $\hat{M}$ can be $\bar{L}_{TFU}$-expanded to a structure $M$ such that $\mathcal{L}_{TFU} = \langle \bar{L}_{TFU}, M, V \rangle$ is truth value correct.*

*Proof.* Apply $FP$ and $RP$, modified in the obvious way, to expand $\hat{M}$ by filling the extensions and anti-extensions of $T, F$ and $U$. From the (modified) Stabilization Lemma it follows that the generated series of structures $\{M^*_\alpha\}_{\alpha \in On}$ has a fixed point at $\omega$, i.e. $M^*_\omega = M^*_{\omega+1}$. From the definition of (modified) $RP$, it now immediately follows that $\mathcal{L}_{TFU} = \langle \bar{L}_{TFU}, M^*_\omega, V \rangle$ is truth value correct.                                                                □

## 5. *HLPE and the paradoxical $TVC$ theorem*

HLPE. In this section, the paradoxical $TVC$ theorem will be applied to shed light on the status of a proof which appeared in (Rabern & Rabern 2008). The background of their proof is the so called *Hardest Logic Puzzle Ever* (*HLPE*). Originally devised by Raymond Smullyan, $HLPE$ was first discussed in an academic journal by (Boolos 1996). Neglecting a detail which is irrelevant for our purposes[14], $HLPE$ may be presented as follows.

> The Puzzle: Three gods A, B and C are called, in some order, True, False, and Random. True always speaks truly, False always speaks falsely, but whether Random speaks truly or falsely is a completely *random* matter. Your task is to determine the identities of A, B, and C by asking three yes-no questions; each question must be put to exactly one god.                                  ((Boolos 1996), p. 62)

Boolos gives the following solution. *First*, ask 'you are true iff $A$ is Random' to $B$. If $B$'s answer is 'yes', we can conclude that $C$ is not Random, while if $B$'s answer is 'no', we can conclude that $A$ is not Random. *Second*, go the god, $A$ or $C$, that you now know to be not Random. Ask him a tautology to find out whether he is True or False. *Third*, ask the non Random god whether $B$ is Random to find out the identity of all three gods. (Roberts 2001) criticized Boolos' solution for being 'unnecessarily complicated'[15] and came up with an alternative solution, which also involves *three* questions.

---

[14] In Boolos' original formulation, the gods understand English but answer with 'da' and 'ja', which mean 'yes' and 'no' but not necessarily in that order. Although you do not know the meaning of 'da' and 'ja', the da-ja version and yes-no version of $HLPE$ can be solved in the same number of questions, and therefore the da-ja details are irrelevant for our purposes.

[15] Roberts says that the unnecessary complications of Boolos' solution are due to his use of the "iff" construct, which "While well-known to logicians, this is the sort of thing that makes

Rabern and Rabern (R&R) observe that $HLPE$'s instructions do not forbid one to ask self-referential questions to the gods and reflect on the consequences of doing so. On Liar like questions such as 'is it the case that your answer to this question is 'no'?', True cannot answer with either 'yes' or 'no' without lying and he has to show a different reaction accordingly, say that True *explodes*. Interestingly, R&R claim that by asking self-referential questions to the gods they can solve $HLPE$ in just *two* questions. The crucial step in this two question solution is what R&R call the 'Tempered Liar Lemma' ($TLL$), the content of which can be illustrated by means of the following example.

TLL. Suppose that there is an object, $o$, that is either black all over, yellow all over or red all over. You do not know $o$'s color but there is a god, True, who knows $o$'s color and who answers all and only yes-no questions *truthfully*. What is the minimum number $n$ of yes-no questions that you have to ask to True in order to be sure that, no matter how True answers them, you can determine $o$'s color? One may reason as follows. First asking whether $o$ is black and then whether $o$ is yellow shows that $n \leq 2$ and as *obviously*, $n \neq 1$ we have $n = 2$. However, R&R give a proof, in natural language, that claims to show that this appeal to our "$n \neq 1$-intuitions" is unjustified; they claim to prove that $n = 1$. The statement that $n = 1$ will be called the *Tempered Liar Lemma* ($TLL$) and the question by which R&R claim to establish $TLL$ is $Q$, in which 'this' refers to the question as a whole.[16]

$Q$: Is it the case that (your answer to this question is 'no' and $o$ is black) or $o$ is yellow?

R&R, argue that if $Q$ is answered with $a$) 'yes' then $o$ is yellow, $b$) 'no' then $o$ is red while $c$) an explosion indicates that $o$ is black.

The proof. The three material implications $a$), $b$) and $c$), are established via *reductio ad absurdum* as follows.
$a$) Assume that True answers 'yes' and that $o$ is not yellow. Then True says 'yes' to the left disjunct of $Q$ and so in particular to 'your answer to this question is 'no''. This is impossible as True tells the truth.
$b$) Assume that True answers 'no' and that $o$ is not red. Then, as True answered 'no' to $Q$, he denies the left and the right disjunct of $\theta$, from which

most laymen despair of logicians, and wonder why they ever tried to solve such puzzles in the first place." ((Roberts 2001), p. 610)

[16] I will be sloppy in not distinguishing between a yes-no question and its associated declarative sentence; for instance I will speak of the truth-value of a yes-no question.

STEFAN WINTEIN

it respectively follows that $o$ is not black and that $o$ is not yellow and so $o$ is red. Contradiction.
$c$) Assume that True explodes and $o$ is not black. Then $o$ is not yellow either, for otherwise True would answer 'yes'. Hence, as $o$ is neither black nor yellow, True denies both disjuncts of $Q$ and hence answers $Q$ with 'no'. Contradiction.

The paradox. This argument of R&R is — though interesting — obscure, for nowhere in (Rabern & Rabern 2008) are the principles by which True reasons specified. At first sight — at least to me — the proof looks fine. But consider the following argument to the conclusion that True does not explode on $Q$ which is, so it seems, obtained by the same principles as those implicit in R&R's proof. Suppose that $o$ is black and that True explodes on $Q$. Then, the left disjunct of $Q$ is false (as True does not answer 'no'), and so $Q$ is false (as $o$ is black the second disjunct is also false) and hence True should answer $Q$ with 'no'! Also, what would happen if we asked True: 'is the case that your answer to this question is 'no' or that you explode on this question?' Such *strengthened Liar objections* show that R&R's proof is suspect, to say the least, and that an explanation of the assumptions involved is needed.

The truth value-answer link. A possible defense against such strengthened Liar objections is that they are based on a wrong conception of "how True works". For instance, True may answer with 'yes'or 'no' iff answering 'yes' or 'no' is truthful *and* if True can do so *without contradicting himself* and otherwise, True explodes. Such an 'inferential conception' of True may be combined with the thought that if $o$ is black, True explodes on $Q$ and that this explosion renders $Q$ false; the inferential conception of True then gives up the link between the truth value of a sentence and True's answer to it.[17]
In contrast, the paradoxical $TVC$ theorem can be seen as a specification of the conditions under which one can make sense of the argument of R&R when the answer of True to $\sigma$ is understood as a reaction to the truth value of $\sigma$. We sketch two distinct ways to do so, called the meta-language approach and the object-language approach respectively.

[17] In personal communication Brain Rabern explained, as a reaction to my strengthened Liar objections, that his conception of True gives up the link between the truth value of $\sigma$ and True's answer to it; True may very well explode on false sentences. In (Wintein 2009) I develop a formal approach to capture what I call an *inferential conception* of True. There an answer of True to a sentence $\sigma$ is determined by the outcome of an *inferential process* of True which takes $\sigma$ as input. In this formalization of an 'inferential True' one can, due to the assumption that True reacts differently to Truthtellers than to Liar sentences, also determine the color of an object which has 1 out of 4 possible colors by asking a single question to True.

The meta-language approach. Assuming a link between the truth value of $\sigma$ and True's answer to it, the interpretation of '$F(x)$' as '$x$ is false' is extensionally equivalent to its interpretation as 'True's answer to $x$ is 'no''. Modulo this shift of interpretation, question $Q$ can be represented via a constant $\theta$ and an interpretation function $I$ as follows.

$$I(\theta) = (F(\theta) \wedge B(o)) \vee Y(o) \tag{6}$$

Let $L_{TFU}$ be a language that contains, besides the three truth value predicates (equivalently, "answering predicates") the color predicates $B$, $Y$ and $R$ and the constants $\theta$ and $o$. Let your ignorance about the color of the object be represented by $K \in Sen(L_{TFU})$:

$$\begin{aligned} K := {}& (B(o) \wedge \neg Y(o) \wedge \neg R(o)) \\ & \vee (\neg B(o) \wedge Y(o) \wedge \neg R(o)) \\ & \vee (\neg B(o) \wedge \neg Y(o) \wedge R(o)) \end{aligned}$$

Any $\Delta$-neutral ground structure for $L_{TFU} - \{T, F, U\}$ which interprets $\theta$ as (6), which interprets $o$ with a non-sentential object and which interprets the color predicates such that $K$ is valuated as t we call a *K-ground structure* and the $L_{TFU}$-expansion of a $K$-ground structure via the *Strong Kleene version* of the paradoxical $TVC$ theorem construction, we call a *possible world*. Note that, corresponding to the three possible colors of the object, the class of possible worlds $\mathcal{M}$ allows for a tripartition. A possible way to give a valid reconstruction of R&R's argument is to understand them as reasoning in a classical meta-language about $\mathcal{M}$. We define $\models_{\mathcal{M}} \subseteq \mathcal{P}(Sen(L_{TFU})) \times Sen(L_{TFU})$ by stipulating that $\Delta \models_{\mathcal{M}} \sigma$ just in case in every $M \in \mathcal{M}$ in which all members of $\Delta$ are valuated as t, $\sigma$ is also valuated as t. Neglecting parenthesis for singleton sets, the three claims of R&R may be translated as follows:

$$a') \; T(\theta) \models_{\mathcal{M}} Y(o) \qquad b') \; F(\theta) \models_{\mathcal{M}} R(o) \qquad c') \; U(\theta) \models_{\mathcal{M}} B(o)$$

As the reader may verify, $a')$, $b')$ and $c')$ are true, while the associated object-language counterparts of, $a')$ and $b')$ in terms of material implication do not hold. For instance, we do not have that $\models_{\mathcal{M}} T(\theta) \to Y(o)$; in a world in which the object is black, '$T(\theta)$' is valuated as u, '$Y(o)$' as f and hence the material implication as u.

The object-language approach. If we slightly alter R&R's natural language claims, we can have a correct object language representation of those claims.

STEFAN WINTEIN

For observe that the fact that the ungroundedness predicate is a classical predicate gives us:

$a''$) $\models_{\mathcal{M}} (\neg U(\theta) \wedge T(\theta)) \rightarrow Y(o)$
$b''$) $\models_{\mathcal{M}} (\neg U(\theta) \wedge F(\theta)) \rightarrow R(o)$
$c''$) $\models_{\mathcal{M}} U(\theta) \wedge \rightarrow B(o)$

Conclusion. We used the paradoxical $TVC$ theorem to give a rough sketch of two possible reconstructions of the reasoning of R&R. Although a lot more can be said about the details of both reconstructions, I do not think that either of them can be fruitfully converted into a genuine proof of $TLL$, the reason being that the condition of $\Delta$-neutrality is too restrictive. We would like to know the principles by which True answers questions as 'do you explode on this question?' and the like, which are excluded by $\Delta$-neutrality. It is my conjecture that, in order to get a systematic account of True's answers to such questions, the truth value-answer link has to be traded in for an inferential conception of True. Be that as it may, the paradoxical $TVC$ theorem itself is a nice little result which can be added to our ever growing stock of truths about truth.

Tilburg University, The Netherlands
E-mail: `s.wintein@uvt.nl`

## REFERENCES

Arieli, O. & Avron, A. 1996. Reasoning with Logical Bilattices. *Journal of Logic Language and Information*, 5:25–63.

Belnap, N. 1977. A useful Four-Valued logic. (*In* Dunn, J.M. & Epstein, G. (eds.), *Modern Uses of Multiple-Valued Logic*.)

Boolos, G. 1996. The hardest logic puzzle ever. *The Harvard Review of Philosophy*, 6:62–65.

Dutilh Novaes, C. 2008. A comparative taxonomy to medieval and modern approaches to Liar sentences. *History and Philosophy of Logic*, 29:227–261.

Goldstein, L. & Blum, A. 2008. When is a statement not a statement? When it's a liar. *The Reasoner*, 2:4–6.

Gupta, A. 1982. Truth and Paradox. *Journal of Philosophical Logic*, 11:1–60.

Gupta, A. & Belnap, N. 1993. *The Revision Theory of Truth*. Cambridge: MIT Press.

Halbach, V. & Horsten, L. 2006. Axiomatizing Kripke's theory of truth. *Journal of Symbolic Logic*, 71:677–712.

Kripke, S. 1975. Outline of a Theory of Truth. *Journal of Philosophical Logic*, 72:690–716.

Maudlin, T. 2006. *Truth and Paradox*. Oxford University Press.

Muskens, R. 1999. On partial and Paraconsistent Logics. *Notre Dame Journal of Formal Logic*, 40:352–373.

Rabern, B. & Rabern, L. 2008. A simple solution to the hardest logic puzzle ever. *Analysis*, 68:105–112.

Roberts, T. 2001. Some thoughts about the hardest logic puzzle ever. *Journal of Philosophical Logic*, 30:609–612.

Smullyan, R. 1995. *First-order Logic*. New York: Dover.

Soames, S. 1999. *Understanding Truth*. Oxford University Press.

Thalos, M. 2005. From Paradox to Judgement: towards a metaphysics of expression. *Australian Journal of Logic*, 3:76–107.

Visser, A. 1984. Four-valued semantics and the Liar. *The Journal of Philosophical Logic*, 13:181–212.

Wintein, S. 201x. Assertoric Semantics and the Computational Power of Self-Referential Truth. *Journal of Philosophical Logic*. forthcoming.

Wittgenstein, L. 1939. *Lectures on the Foundations of Mathematics, Cambridge 1939*. Harvestor, ed. Cora Diamond: Hassocks.

Woodward, J. Scientific Explanation. (*In* Zalta, E.N. (ed.), *The Stanford Encyclopedia of Philosophy*, <http://plato.stanford.edu/archives/spr2009/entries/scientific-explanation/>.)

Yablo, S. 1993. Paradox without self-reference. *Analysis*, 53:251–252.