

# The Elimination of Self-Reference:

## Generalized Yablo-Series and the Theory of Truth<sup>1,2</sup>

P. Schlenker

(UCLA & Institut Jean-Nicod)

Second Draft (revised and extended, though not quite final<sup>3</sup>; last modified on November 9, 2005)

*Abstract :* Although it was traditionally thought that self-reference is a crucial ingredient of semantic paradoxes, Yablo (1993, 2004) showed that this was not so by displaying an infinite series of sentences none of which is self-referential but which, taken together, are paradoxical. Yablo's paradox consists in a series  $\{ \langle s(i), [\forall k: k > i] \neg Tr(s(k)) \rangle : i \geq 0 \}$ , where  $Tr$  is the truth predicate and where for each number  $i$  the term  $s(i)$  denotes the sentence  $[\forall k: k > i] \neg Tr(s(k))$ . We generalize Yablo's result along two dimensions. **1.** First, we investigate the general behavior of the series  $\{ \langle s(i), [Qk: k > i] f[(s(k))_{k \geq i}] \rangle : i \geq 0 \}$ , where  $Q$  is a generalized quantifier and where  $f$  is some fixed truth function. We show that under broad conditions *all the sentences in the series must have the same value*, and we derive a *characterization of those values of  $Q$  for which the series is paradoxical*. **2.** Second, we show that in the Strong Kleene trivalent logic, Yablo's results are a special case of a much more general phenomenon: under certain conditions, *any semantic phenomenon that involves self-reference can be reproduced without self-reference*. This is shown by way of a translation which associates to each pair  $\langle s, F \rangle$  in a non-quantificational language  $L'$  (where the term  $s$  denotes the sentence  $F$ ) an infinite series of translations  $\{ \langle s(i), [Qk: k > i] [F]_k \rangle : i \geq 0 \}$  in a quantificational language  $L^*$  (where  $[F]_k$  is a modification of  $F$ ). We provide a characterization of those values of  $Q$  for which the translation goes through, and we discuss various extensions of the translation procedure. The paper, which generalizes recent results by Cook (2004), shows that under certain conditions self-reference is not essential to *any* semantic phenomena.

## 1 Simple and Infinite Liars

### 1.1 Semantic Effects of Self-Reference

Natural language includes various means to express self-referential statements. These may be entirely innocuous, as in the case of (1) and (2), which are uncontroversially true and false respectively:

- (1) This very sentence contains six words.
- (2) This very sentence contains ten words.

In other cases, however, self-reference leads to more interesting logical monstrosities such as the Liar and the Truth-Teller, illustrated in (3) and (4) respectively:

- (3) This very sentence is false.
- (4) This very sentence is true.

The Liar is *paradoxical* because it cannot coherently be assigned the value *true* or the value *false*. The Truth-Teller is *pathological* because it can be assigned either the value *true* or the value *false*, but in a way that appears to be utterly arbitrary. As shown by Kripke 1975, some statements may or may not be Liar-like or Truth-Teller-like depending on some empirical facts. Thus (5) as uttered by Smith is paradoxical if it is the only statement made by Smith on

---

<sup>1</sup> I thank the following for discussion of various stages of this work: Denis Bonnay, Serge Bozon, Paul Egré, Marcus Kracht, Tony Martin, Benjamin Spector, Albert Visser, as well as audiences at IHPST, UCLA and U. of Amsterdam. Special thanks to Denis Bonnay, who corrected innumerable errors and to Albert Visser, who commented on this paper at a French-Dutch logic meeting in Amsterdam (PALMYR, June '05).

<sup>2</sup> A shorter and less technical version of the present research can be found in Schlenker 2005.

<sup>3</sup> Section 5 is still rough.

a particular day. But it is simply false if Smith utters on the very same day another sentence which is true, for instance *The earth is round*.

(5) Every sentence I will have uttered today will turn out to be false.

Tarski (e.g. Tarski 1944) observed that a bivalent framework is incapable of providing an adequate theory of truth for a language that includes (i) means of self-reference, and (ii) its own truth predicate. He concluded that the (bivalent) languages he studied in logic could not contain their own truth predicates. As an account of natural language semantics, however, this falls short. Natural language *does* generally contain a truth predicate, and speakers have relatively clear intuitions about its meaning. For this reason, Kripke 1975 suggested that bivalence should be sacrificed, and he showed how an adequate trivalent semantics can be given for a language that contains its own truth predicate and means of self-reference. This move away from bivalence would seem to be linguistically motivated, since speakers often want to classify the Liar as being neither true nor false. We will henceforth say that a sentence is 'indeterminate' or 'has the value #' if it is neither true nor false.

Once a trivalent framework is adopted, Tarski's 'Convention T' must be reformulated to require that the truth predicate *Tr* be (i) true of all the true sentences, and (ii) false of all the false sentences, and (iii) neither true nor false of the sentences that are neither true nor false. If we design the syntax in such a way that that *Tr* can only take as arguments terms that denote sentences, the revised version of Convention T for a language *L* can be expressed as follows (the domain of objects is *D*, and we assume that  $L \subseteq D$ , which means that the sentences of the language belong to the domain of objects):

(6)  $I^+(Tr) = \{d \in D : d \in L \wedge I'(d) = 1\}$   
 $I^-(Tr) = \{d \in D : d \in L \wedge I'(d) = 0\}$

An interpretation that satisfies these conditions is henceforth called a 'fixed point', following Kripke's own terminology<sup>4</sup>. Kripke 1975 gives a procedure (which will be generalized below) to show that some fixed points can indeed be constructed.

## 1.2 Yablo's Paradox

Once this framework is in place, one would like to understand what are the essential ingredients needed to obtain the interesting semantic phenomena we started out with (logical and empirical Liars and Truth-Tellers). Granted, a truth predicate and some devices of self-reference are *sufficient* to generate these phenomena. But are they *necessary*? Some researchers have been concerned with paradoxes that can be generated without a truth predicate (see for instance Egré 2005 for a recent presentation). In this article, we will only be concerned with paradoxes of truth. With this restriction, two results can be established.

a) Yablo (1993, 2004) has shown that logical paradoxes can be obtained without self-reference if quantification over infinite series of sentences is allowed. His conclusion is strengthened by Cook 2004, who studies a very simple system in which every Liar (or Truth-

---

<sup>4</sup> A few words might be in order to justify this terminology. Kripke starts from a classical interpretation *I*, to which the Truth predicate *Tr* (which is itself trivalent) is added. A new (trivalent) interpretation *I'* is determined by the combination of *I* and a specification of a pair  $\langle I^+(Tr), I^-(Tr) \rangle$  of an extension and an anti-extension for *Tr*. If we define  $f(\langle I^+(Tr), I^-(Tr) \rangle) = \langle \text{true sentences according to } I', \text{ false sentences according to } I' \rangle$ , we see that a fixed point in Kripke's sense satisfies  $f(\langle I^+(Tr), I^-(Tr) \rangle) = \langle I^+(Tr), I^-(Tr) \rangle$ .

When the truth predicate can take as arguments terms that denote non-sentences, some decision must be made as to how to classify them. Kripke puts them in the anti-extension of the Truth predicate, which can then be paraphrased as: *is a true sentence* (rather than: *is true*). The revised Convention T must then take the following form:

(i)  $I^+(Tr) = \{d \in D : d \in L \wedge I'(d) = 1\}$ ;  $I^-(Tr) = \{d \in D : d \notin L \vee (d \in L \wedge I'(d) = 0)\}$

Teller) involving self-reference can be 'unwinded' to yield an infinite Liar (or Truth-Teller) without self-reference<sup>5</sup>.

b) On the negative side, it can be shown that a Yabloesque construction does *not* yield paradoxes when effected in a first-order language with a truth predicate and sentence names but no quantifiers<sup>6</sup>.

To see an example of Yablo's construction, let us start from a classical interpretation  $I$  for a language  $L$  that includes some simple arithmetic vocabulary (i.e. variables ranging over natural numbers, a name  $i$  for each natural number  $i$ , and the relation  $>$ ). We will assume throughout that the arithmetic vocabulary is based on the standard interpretation of the integers. We add to  $L$  a list of sentence names, which in the present case will be functional, and which are interpreted by a denotation relation  $N^*$ .  $L$  together with the sentence names and the truth predicate  $Tr$  yields a language  $L^*$ . We will often write  $\{<s(\mathbf{k}), F_k>: k \geq 0\}$  if we study a set of formulas  $F_k$  ( $k \geq 0$ ), and if  $N^*$  specifies that for each natural number  $k$ ,  $s(\mathbf{k})$  denotes  $F_k$  [we will often import sentence names in the meta-language by underlining them, thus using  $\underline{s(\mathbf{k})}$  as our name for  $F_k$ ].

With these assumptions, it can be shown that the following system, which we call the  $\forall$ -Liar, is paradoxical: no fixed point may assign to every sentence a 'classical value' (i.e. 0 or 1).

(7)  $L_\forall := \{<s_i, \forall k (k > i \Rightarrow \neg Tr(s_k))>: i \geq 0\}$

a) Suppose that for each  $i \geq 0$ ,  $s_i$  is false. This leads to a contradiction:  $s_0$  is true iff for each  $k > 0$ ,  $s_k$  is false, which by hypothesis is indeed the case. So  $s_0$  should be true, contrary to hypothesis.

b) Suppose that for some  $i \geq 0$ ,  $s_i$  is true. Then for each  $k > i$ ,  $s_k$  is false. In particular, for each  $k > i+1$ ,  $s_k$  is false. Hence  $s_{i+1}$  is true, which contradicts the fact that for each  $k > i$ ,  $s_k$  is false.

Quite a bit of ink has been spilled to determine whether Yablo's sentences are not subtly self-referential in some indirect fashion. We give in an Appendix an argument to the effect that they aren't, but in the rest of the paper we simply assume that Yablo's observation is correct. The skeptic can take our claims to be conditional: *If Yablo's sentences do not involve self-reference, then \_\_\_\_\_*.

While Yablo's paradox is normally stated in a setting in which the relevant sentences are linearly ordered (as is the case with the denotation relation  $N^*$  we just introduced), much less is necessary to obtain the paradox (as was already observed by Cook 2004 for his special

<sup>5</sup> In addition, Cook's construction departed from Yablo's in relying on infinite conjunction rather than on quantification over sentences. We will not be concerned with this distinction in the present paper.

<sup>6</sup> The argument, which was pointed out to me by Tony Martin, goes as follows.

Consider a series of sentences of the form  $\{<c_k, f_k(c_{k+1}, \dots, c_{k+n_k})>: k \geq 0\}$ , where for each  $k \geq 0$   $f_k$  is a Boolean function. We show that for any such series there exists a bivalent valuation. Let us say that an assignment of truth-values to  $c_0, \dots, c_n$  is *acceptable* just in case for each  $i \leq n$ , (1) or (2) holds:

(1) for some  $k$  such that  $c_k$  is an argument of  $f_i$ ,  $k > n$

(2) (1) fails, and the truth value assigned to  $c_i$  is as required by the value of  $f_i$ .

For each  $n$ , there is an acceptable assignment of bivalent values to  $c_0, \dots, c_n$ . We can simply start with an arbitrary value for  $c_n$  and any other sentence which has at least an argument  $c_m$  for  $m > n$ . We then compute the values of the other sentences as the  $f_i$  dictate.

Thus the binary tree of all acceptable assignments has arbitrary long branches. By Koenig's Lemma, it has an infinite branch, which is the desired valuation.

Albert Visser (p.c.) suggests an alternative argument. Consider the theory  $Th := \{c_k \leftrightarrow f_k(c_{k+1}, \dots, c_{k+n_k}): k \geq 0\}$ , where the  $c_k$ 's are construed as propositional letters rather than as sentence names. Each finite subset of  $Th$  has a model (as in the preceding argument, we assign an arbitrary value to the 'last'  $c_k$ 's, and compute the value of the 'earlier'  $c_k$ 's as the  $f_k$ 's require to make the equivalences true). By the Compactness Theorem for propositional logic,  $Th$  itself has a model, which provides the desired valuation.

language). In fact, all that is needed is that the sentences be ordered according to a relation  $R$  which (i) is non-empty, (ii) is transitive, and (iii) has no end points. Instead of quantifying over natural numbers, we consider sentences that quantify directly over other sentences (in the following example  $s$  is a variable ranging over sentences, and we say that  $s$  sees  $s'$  just in case  $\langle s, s' \rangle$  is in the extension of  $R$ ):

- (8)  $L_R := \{ \langle s_i, \forall s (s_i R s \Rightarrow \neg \text{Tr}(s)) \rangle : i \geq 0 \}$ , where  $R$  is interpreted a set of pairs of sentences of  $L_R$ .
- a) Suppose for each  $i \geq 0$ ,  $s_i$  is false. Then for all  $i \geq 0$  the content of  $s_i$  should be true (for either  $s_i$  'sees' no other sentences according to  $R$ , in which case it is vacuously true; or it does see other sentences, but they by assumption false, as  $s_i$  claims). Contradiction.
  - b) Suppose that for some  $i \geq 0$ ,  $s_i$  is true. Since  $R$  has no endpoints,  $s_i$  sees some sentence  $s_{i^*}$ , which must be false.  $s_{i^*}$  itself sees some sentence or sentences, which by the transitivity of  $R$  must be seen by  $s_{i^*}$  as well, and hence be false. But this should suffice to make  $s_{i^*}$  true, contrary to what we just showed. Contradiction.

The present paper has two primary goals. First, we derive a general result about Yablo-style series in which sentences are linearly ordered, and we characterize those that are paradoxical (the result fails when the sentences are simply ordered by a transitive relation without endpoints). Second, and more importantly, we show that Yablo's and Cook's results are special cases of a more general phenomenon: *under certain conditions, every semantic result that can be obtained with self-reference in Kleene's Strong trivalent logic can be emulated without self-reference*. Yablo's and Cook's constructions had nothing to say about sentences like (5), which are paradoxical only when certain empirical facts hold. Because our result is general, and stated within Kripke's theory of truth, it entails that every *empirical* paradox also has a non-self-referential variant. It also entails that every other semantic phenomenon that one might be interested in can be replicated without self-reference as well.

The rest of this paper is organized as follows. We start by generalizing Yablo's result by studying in greater generality certain infinite series of sentences, which we call *Generalized Yablo Series* (Section 2). We then develop our main construction, which shows how self-reference can be eliminated from a non-quantificational language (Section 3). We discuss some extensions and generalizations (Section 4), and we briefly consider an extension to quantificational languages (Section 5). The paper ends with some perspectives for future research (Section 6), and with an Appendix that provides a sufficient condition of non-self-reference, which suggests that Yablo's sentences are indeed self-reference-free.

## 2 Generalized Yablo-Series

### 2.1 Versions of Yablo's Paradox

Yablo's paradox comes in several varieties. We have already shown that  $L_\forall$ , repeated in (9)a, is paradoxical. Yablo 2004 shows that  $L_\exists$  and  $L_{\exists\forall}$  as defined below are equally paradoxical, and we will show in a second that  $L_{\forall\exists}$  is as well.

- (9) a.  $L_\forall := \{ \langle s_i, \forall k (k > i \Rightarrow \neg \text{Tr}(s_k)) \rangle : i \geq 0 \}$   
 b.  $L_\exists := \{ \langle s_i, \exists k (k > i \wedge \neg \text{Tr}(s_k)) \rangle : i \geq 0 \}$   
 c.  $L_{\exists\forall} := \{ \langle s_i, \exists k (k > i \wedge \forall k' (k' > k \Rightarrow \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$   
 d.  $L_{\forall\exists} := \{ \langle s_i, \forall k (k > i \Rightarrow \exists k' (k' > k \wedge \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$

Before discussing too many special cases, it is worth noting that new paradoxical series can be obtained out of old ones by a kind of duality principle. We will take a series to be paradoxical (given a classical interpretation  $I$  and a denotation function  $N$  for the sentence terms) just in case no fixed point can be found which assigns to all of its members classical truth values. With this definition, the Duality Lemma stated below yields a recipe to obtain new paradoxes out of old ones.

*Definition:* If  $F$  is a formula, let  $F^*$  be the formula obtained by replacing every occurrence of the form  $Tr(.)$  in  $F$  with  $\neg Tr(.)$ . [Note that  $F^{**}$  is equivalent to  $F$ , and that  $\neg(F^*)$  is identical to  $(\neg F)^*$ ].

*Duality Lemma:* If  $\{ \langle s_k, F_k \rangle : k \geq 0 \}$  is paradoxical, so is  $\{ \langle s_k, \neg F_k^* \rangle : k \geq 0 \}$  [ $\{ \langle s_k, F_k \rangle : k \geq 0 \}$  and  $\{ \langle s_k, \neg F_k^* \rangle : k \geq 0 \}$  fully describe the intended denotation functions for sentence terms; in each case  $s_k$  could be any sentence-denoting terms, not just constants].

*Proof:* Suppose that  $\{ \langle s_k, F_k \rangle : k \geq 0 \}$  is paradoxical but that  $\{ \langle s_k, \neg F_k^* \rangle : k \geq 0 \}$  is not. This means that there is some fixed point  $I^*$  which yields a bivalent valuation for  $\{ \langle s_k, \neg F_k^* \rangle : k \geq 0 \}$ . We claim that this suffices to find a fixed point  $I'$  which assigns classical values to the sentences in  $\{ \langle s_k, F_k \rangle : k \geq 0 \}$ . We take  $I'$  to be identical to  $I^*$ , except that:

- (i)  $I'(s) = F$  iff  $I^*(s) = \neg F^*$
- (ii)  $I'^+(Tr) = \{ F \in L' : \neg F^* \in I^*(Tr) \}$  and  $I'^-(Tr) = \{ F \in L' : \neg F^* \in I^{*+}(Tr) \}$ .

We start by observing that:

- (iii) for each formula  $F$ ,  $I'(F) = I^*(F^*)$

This is because it follows from (i) and (ii) that if  $s$  is sentence-denoting and if  $I'(s) = F$ ,

$$\begin{aligned}
 I'(Tr(s)) = 1 & \quad \text{iff} \quad F \in I'^+(Tr) & \quad [\text{by trivalent semantics}] \\
 & \quad \text{iff} \quad \neg F^* \in I^*(Tr) & \quad [\text{by (ii)}] \\
 & \quad \text{iff} \quad I^*(s) \in I^*(Tr) & \quad [\text{by (i)}] \\
 & \quad \text{iff} \quad I^*(Tr(s)) = 0 & \quad [\text{by trivalent semantics}] \\
 & \quad \text{iff} \quad I^*(\neg Tr(s)) = 1
 \end{aligned}$$

(The case  $I'(Tr(s)) = 0$  is symmetric). Since  $F^*$  is identical to  $F$  except that every occurrence of  $Tr(.)$  is replaced by  $\neg Tr(.)$ , we obtain the desired result.

Let us now assume that  $I^*$  is a fixed point, and let us show that  $I'$  as defined is a fixed point as well.

$$\begin{aligned}
 F \in I'^+(Tr) & \quad \text{iff} \quad \neg F^* \in I^*(Tr) & \quad [\text{by (ii)}] \\
 & \quad \text{iff} \quad I^*(\neg F^*) = 0 & \quad [\text{since } I^* \text{ is a fixed point}] \\
 & \quad \text{iff} \quad I'(\neg F) = 0 & \quad [\text{by (iii)}] \\
 & \quad \text{iff} \quad I'(F) = 1
 \end{aligned}$$

The case  $F \in I'^-(Tr)$  is symmetric:  $F \in I'^-(Tr)$  iff  $I'(F) = 0$ . Thus if  $I'$  is a fixed point, so is  $I^*$ , and furthermore for any formula  $F$   $I'(F) = I^*(F^*)$ . It immediately follows that if  $I^*$  is a fixed point that assigns classical values to all the sentences in  $\{ \langle s_k, \neg F_k^* \rangle : k \geq 0 \}$ , this is also the case of  $I'$  with respect to the sentences of  $\{ \langle s_k, F_k \rangle : k \geq 0 \}$ , which shows that the latter system is not paradoxical.

From the Duality Lemma and our earlier observations about the  $\forall$ -Liar, it follows that its dual, namely  $\{ \langle s_i, \neg \forall k (k > i \Rightarrow \neg \neg Tr(s_k)) \rangle : i \geq 0 \}$ , is also paradoxical. But the latter is immediately equivalent to the  $\exists$ -Liar  $\{ \langle s_i, \exists k (k > i \wedge \neg Tr(s_k)) \rangle : i \geq 0 \}$ .

Let us now turn to the  $\exists\forall$ -Liar (discussed in Yablo 2004). We may reason as follows:

(10)  $L_{\exists\forall} := \{ \langle s_i, \exists k (k > i \wedge \forall k' (k' > k \Rightarrow \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$

a) Suppose that for some  $i \geq 0$ ,  $s_i$  is true. Let  $k > i$  be such that  $\forall k' (k' > k \Rightarrow \neg \text{Tr}(s_{k'}))$ . In particular,  $s_{k+1}$  must be false. However the condition  $\forall k' (k' > k \Rightarrow \neg \text{Tr}(s_{k'}))$  suffices to make  $s_{k+1}$  true. Contradiction.

b) Suppose now that for each  $i \geq 0$ ,  $s_i$  is false. Then in particular what  $s_0$  says is true. Contradiction.

It follows from the Duality Lemma that the series  $\{ \langle s_i, \neg \exists k (k > i \wedge \forall k' (k' > k \Rightarrow \neg \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$  is also paradoxical. Since the latter is equivalent  $\{ \langle s_i, \forall k (k > i \Rightarrow \exists k' (k' > k \wedge \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$ , we derive the result that the  $\forall\exists$ -Liar is paradoxical as well, as was announced. (We may also note for completeness that the same results would have held if the sentences had been ordered according to a transitive relation without endpoints rather than a linear ordering.)

It will prove important to observe that there is an interesting conceptual difference between the  $\forall$ - and  $\exists$ -Liars on the one hand and the  $\exists\forall$ - and  $\forall\exists$ -Liars on the other. If  $L_{\forall}$  is evaluated in an interpretation which is not a fixed point, its members need *not* all have the same truth value. For instance, if  $I^+(\text{Tr}) = \{s_1\}$  and  $I^-(\text{Tr}) = \{s_i : i \geq 2\}$ , it will follow that  $I'(s_0) = 0$  and for all  $i \geq 1$ ,  $I'(s_i) = 1$ . A similar argument applies to  $L_{\exists}$ . By contrast, however, all members of  $L_{\exists\forall}$  and  $L_{\forall\exists}$  have the same value in *any* interpretation (not just in fixed points) because they all have the same semantic content. Any sentence  $s(i)$  of  $L_{\exists\forall}$  asserts, in effect, that all but finitely many members of the series beyond rank  $i$  are true. But the modifier *beyond rank  $i$*  turns out to be semantically idle: the claim is utterly insensitive to what happens in any given initial segment of the series, and for this reason all the sentences make the very same claim, namely that all but finitely many members of the series are true (this argument will be fleshed out below). Similarly any sentence  $s(i)$  of  $L_{\forall\exists}$  asserts that infinitely many members of the series *beyond rank  $i$*  are true. But here too the modifier *beyond rank  $i$*  is eliminable, and thus all the members of the series make the same claim. This conceptual difference will have important repercussions when we provide a translation procedure to eliminate self-reference systematically (it will make use of a generalization of the constructions at work in  $\exists\forall$ - and  $\forall\exists$ -Liars, and we will show a similar generalization of the  $\forall$ - and  $\exists$ -Liars fails).

It can shown that the Truth-Teller has a variety of self-reference free versions, which parallel the corresponding constructions for the Liar. Thus the systems defined below have the property that (i) in any fixed point  $I'$ , all the sentences in the series have the same truth value according to  $I'$ ; and (ii) this value can be arbitrarily chosen to be 0, 1, or  $\#$ .<sup>7</sup>

(11) a.  $T_{\forall} := \{ \langle s_i, \forall k (k > i \Rightarrow \text{Tr}(s_k)) \rangle : i \geq 0 \}$

b.  $T_{\exists} := \{ \langle s_i, \exists k (k > i \wedge \text{Tr}(s_k)) \rangle : i \geq 0 \}$

c.  $T_{\exists\forall} := \{ \langle s_i, \exists k (k > i \wedge \forall k' (k' > k \Rightarrow \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$

d.  $L_{\forall\exists} := \{ \langle s_i, \forall k (k > i \Rightarrow \exists k' (k' > k \wedge \neg \text{Tr}(s_{k'}))) \rangle : i \geq 0 \}$

It should be pointed out that even with respect to the various versions of Yablo's paradox, we haven't quite finished the semantic job. At least two questions remain open. (i) First, we have shown that in each case at least one member of the series must have a non-classical value. But do *all* the members of the series have the value undefined? In the case of  $L_{\exists\forall}$  and  $L_{\forall\exists}$  this result of uniformity follows from our earlier observation that all the members of the series make the very same claim. But what about the other cases? We will soon see that a more general result, which we call the Uniformity Property, guarantees that in series of this sort (whether paradoxical or not), all the sentences have the same value in any given fixed point. (ii) Second, we would like to have at least a partial characterization of those series of

<sup>7</sup> To put it more rigorously: for each value  $v \in \{0, 1, \#\}$ , there is a fixed point in which all the sentences in the series have the value  $v$ .

Yablo type which are indeed paradoxical. The Uniformity Property will make the characterization straightforward for an entire class of series.

## 2.2 Yablo Series in a General Setting

We turn to a generalization of the results we have discussed up to this point. The theory of Generalized Quantifiers turns out to offer a versatile tool to study the more general form of the problem.

### 2.2.1 Yablo Series with Generalized Quantifiers

#### □ The Series

We will provide a simple characterization of the behavior of series of sentences of the form  $\{ \langle s(i), [Qk: k > i] \text{ Tr}(s(k)) \rangle : i \geq 0 \}$ , where  $Q$  is a binary generalized quantifier (e.g. *some*, *most*, *no*, *all*, *an odd number of*, etc.) which satisfies Permutation Invariance, Extension and Conservativity, three natural properties that are believed to hold of natural language determiners (Keenan 1996; see the definition below). To see that this is indeed a generalization of the cases we have considered so far, we may observe that for special values of  $Q$  we obtain different version of the Infinite Liar and of the Infinite Truth-Teller:

- (12)  $S_Q = \{ \langle s(i), [Qk: k > i] \text{ Tr}(s_k) \rangle : i \geq 0 \}$ .
- a. For  $Q = \text{All}$ ,  $S_Q$  is the Universal Truth-Teller.
  - b. For  $Q = \text{Some}$ ,  $S_Q$  is the Existential Truth-Teller.
  - c. For  $Q = \text{All but a finite number of}$ ,  $S_Q$  is the Almost Universal Truth-Teller.
  - a'. For  $Q = \text{No}$ ,  $S_Q$  is the Universal Liar.
  - b'. For  $Q = \text{Not all}$ ,  $S_Q$  is the Existential Liar.
  - c'. For  $Q = \text{At most a finite number of}$ ,  $S_Q$  is the Almost Universal Liar.

Permutation Invariance, Extension and Conservativity are defined as follows (Keenan 1996):

- (13) A binary function  $R$  defined for each universe  $D$  and all subsets  $X, Y$  of  $D$  satisfies:
- a. *Permutation Invariance* just in case for each universe  $D$ , for any permutation  $\pi$  of  $D$ , for all  $X, Y \subseteq D$ ,  $R_D(X, Y) = R_D(\pi(X), \pi(Y))$
  - b. *Extension* just in case for any  $X, Y, D, D'$ , if  $X, Y \subseteq D$  and  $X, Y \subseteq D'$ , then  $R_D(X, Y) = R_{D'}(X, Y)$
  - c. *Conservativity* iff for all  $X, Y, D$ :  $R_D(X, Y) = R_D(X, X \cap Y)$

What is important for our purposes is that, taken together, these properties ensure that the truth value of, say, *Most students passed*, only depends on two numbers: the number  $a$  of individuals that are students and did not pass, and the number of individuals that are students and passed (for the determiner *most*, the condition is that  $b > a$ ). More generally, a generalized quantifier  $Q$  that satisfies the conditions in (13) is defined by its 'tree of numbers'  $\underline{Q}$ , which is a set of pairs of numbers (including  $\infty$ ) such that: for any formulas  $F, F'$  with extensions  $\underline{F}$  and  $\underline{F'}$ ,  $Qx F F'$  is true (in a bivalent system) iff  $\langle |\underline{F} - \underline{F'}|, |\underline{F} \cap \underline{F'}| \rangle \in \underline{Q}$  (van Benthem 1986) [Examples:  $\underline{\text{No}} = \{ \langle n, 0 \rangle : n \geq 0 \text{ or } n = \infty \}$ ;  $\underline{\text{Most}} = \{ \langle n, n' \rangle : n' > n \}$ ].

#### □ Generalization of the Tree of Numbers

Since we are interested in systems of sentences that might be paradoxical, we must develop the analysis in a logic that is at least trivalent. The case we consider is quite special, however, because the first argument of  $Q$  is a classical formula, which (given any assignment function) has either the value true or the value false. Let us call an  $n$ -valued logic

'reasonable' if it has the following property, which can be seen as a generalization of the 'tree of numbers' found in the bivalent case:

(14) Reasonableness

An n-valent logic with truth values in E is *reasonable* just in case:

If for any assignment function F has a classical value, then for any generalized quantifier  $Q$ , the value of a formula  $[Qk: F]F'$  under an assignment function s and an interpretation I only depends on  $(\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=e\})_{e \in E}$ .

Examples: (i) In the bivalent case, Reasonableness is just the requirement that the quantifiers should be definable in terms of the tree of numbers. In other words,  $I_s([Qk: F]F')$  only depends on  $\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=1\}$ ,  $\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=0\}$

(ii) In the trivalent case, Reasonableness requires that  $I_s([Qk: F]F')$  only depend on the numbers:  $\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=1\}$ ,  $\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=0\}$ ,  $\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=\#\}$

Throughout our discussion, we will restrict attention to formulas in which every generalized quantifier takes a classical formula as its restrictor, so that Reasonableness will have some 'bite'. We also restrict attention to denumerable domains. With this restriction, the semantics of a formula  $[Qk: F]F'$  in an n-valent logic which is reasonable and has truth values in E is determined by a 'generalized Tree of Numbers' which can be seen as a function from  $(\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=e\})_{e \in E}$  to E. We can assimilate  $(\{d \in D: I_{s[k \rightarrow d]}(F)=1\} \cap \{d \in D: I_{s[k \rightarrow d]}(F')=e\})_{e \in E}$  to a function from E to  $\mathbb{N} \cup \{\infty\}$ : a number is associated to each possible truth value in E. We will thus adopt the following convention:

*Notational Convention 1:* If Q is a quantifier, we call  $\underline{Q}^n$  the function from  $E \rightarrow \mathbb{N} \cup \{\infty\}$  to E which defines its n-valent semantics on denumerable domains.

For perspicuity, we will sometimes make use a further convention:

*Notational Convention 2:* If f is a function in  $E \rightarrow \mathbb{N} \cup \{\infty\}$ , we will sometimes write  $\underline{Q}^n(e: f(e))_{e \in E}$  instead of  $\underline{Q}^n(f)$ . We also write  $\underline{Q}^n(e: a, -e: b)$  for  $\underline{Q}^n(f)$  where  $f(e)=a$  and  $f(e')=b$  for all  $e' \neq e$ .

### 2.2.2 The Uniformity Property

It might be helpful to start by considering a special case of the Uniformity Property. We will thus restrict attention to the series  $S_Q := \{ \langle s(i), [Qk: k > i] \text{Tr}(s(k)) \rangle : i \geq 0 \}$ , and we will prove the following result:

(15) Uniformity Property (Special Case): For any fixed point I for a language that includes  $S_Q$ , for each m,  $n \geq 0$ ,  $I(s(\underline{\mathbf{m}})) = I(s(\underline{\mathbf{n}}))$ .

Given a valuation for  $S_Q$ , the truth value of each sentence  $s(\underline{\mathbf{i}})$  in  $S_Q$  is determined by three numbers  $-_i$ ,  $+_i$  and  $\#_i$ , which are the numbers of integers  $k > i$  that make the nuclear scope  $\text{Tr}(s(k))$  false, true and indeterminate respectively. We further define  $/i/ = \langle -_i, +_i, \#_i \rangle$ ; the value of  $s_i$  is thus equal to  $\underline{Q}^3(/i/)$ , which we also write as  $\underline{Q}^3(-_i, +_i, \#_i)$ .



With these conventions, we can give an easy proof of the Property (this argument was greatly simplified by Denis Bonnay<sup>8</sup>):

Proof: Consider any sentence  $\underline{s(i)}$  of  $S_Q$ . The restrictor of  $\underline{s(i)}$  holds true of an infinite number of natural numbers (because there are infinitely many numbers that are greater than  $i$ ). Therefore  $\neg_i$ ,  $+_i$ , or  $\#_i$  (where 'or' is inclusive) are infinite. Suppose for instance that  $\neg_i = \infty$ . Consider any  $i' \geq 0$  for which  $I(\underline{s(i'+1)}) = 0$ . Since  $I(\underline{s(i'+1)}) = 0$ ,  $Q^3(\neg_{i'+1}, +_{i'+1}, \#_{i'+1}) = 0$ . Since  $\neg_{i'+1} = \infty$ , we have (with slight abuses of notation):

$/i'/ = \langle \infty + 1, +_{i'+1}, \#_{i'+1} \rangle$  (because  $\underline{s(i')}$  is followed by the same sentences as  $\underline{s(i'+1)}$ , plus  $\underline{s(i')}$  itself, which is false). Thus  $/i'/ = \langle \infty, +_{i'+1}, \#_{i'+1} \rangle = \langle \neg_{i'+1}, +_{i'+1}, \#_{i'+1} \rangle = /i'+1/$ , and  $Q^3(/i'/) = Q^3(/i'+1/) = 0$ . By iterating this reasoning, we see that for all  $i'' \leq i'$   $I(\underline{s(i'')}) = 0$ . Since there are infinitely many false sentences in the series, we can show that all sentences in the series are false. By similar reasoning, we can show that if  $+_i = \infty$  all sentences are true, and that if  $\#_i = \infty$  all sentences are indeterminate.

It may be observed that the above argument did not depend on the precise semantics of the quantifiers (except for the fact that they satisfied the generalized tree of numbers), nor even on the number of truth values in the logic. This suggests that a more general result can be proven, which is indeed the case.

(16) Uniformity Property (General Case): Let  $E$  be a finite set of truth values and let  $f$  a propositional formula that denotes a Boolean function  $\underline{f}$  from  $E^\infty$  to  $E$ , i.e.

$$\underline{f}: E^\infty \rightarrow E$$

Let  $Q$  be a binary generalized quantifier of a reasonable logic with values in  $E$ . Let  $S$  be a naming relation defined by:

$$S = \{ \langle s(i), [Qk: k > i] f[(s(k'))_{k \geq k}] \rangle : i \geq 0 \}$$

[note that  $f[(s(k'))_{k \geq k}]$  may be an infinitely long formula; note also that  $f$  does not depend on  $i$ ].

Then if  $I$  is a fixed point, for all  $i, i' \geq 0$ ,  $I(\underline{s(i)}) = I(\underline{s(i')})$

Proof:

For each  $i \geq 0$  and for each  $e \in E$ , let us call:

$$a) f_i := I(f[(s(k'))_{k \geq i}])$$

$$b) e_i := |\{k > i: f_k = e\}|$$

To illustrate, if  $E = \{0, 1\}$ , we have:

$0_i = |\{k > i: f_k = 0\}|$  (this is the number of objects that satisfy the restrictor but not the nuclear scope of the statement  $s_i$ )

$1_i = |\{k > i: f_k = 1\}|$  (this is the number of objects that satisfy the restrictor and the nuclear scope of the statement  $s_i$ )

(i) Because the logic is *reasonable*,  $I(\underline{s(i)})$  only depends on  $(e_i)_{e \in E}$

(ii) Since  $E$  is a finite set of truth values, there is a non-negative integer  $i^*$  such that for each  $i > i^*$ , if  $f_i = e$ , then there are infinitely many natural numbers  $i'$  such that  $f_{i'} = e$ .

[Note that this condition would not hold if  $E$  were not finite].

<sup>8</sup> A prior version considered many more cases than was necessary.

Specifically, we define  $E^f := \{e \in E: \text{only finitely many } k\text{'s are such that } f_k = e\}$ . Define  $i^*$  as:  $i^* := \text{Max } \{k \geq 0: \exists e \in E^f f_k = e\}$ .  $E^f$  is finite and only finitely many sentences have values in  $E^f$  hence  $i^*$  is well-defined.

iii) For all  $i > i^*$ , if  $e \in E^f$ ,  $e_i = 0$ ; if  $e \notin E^f$ ,  $e_i = \infty$

iv) By i) and iii), for all  $i, i' > i^*$ ,  $I(\underline{s(i)}) = I(\underline{s(i')})$ . So  $E - E^f$  is a singleton. Let us call  $e^*$  its only member. The situation can be pictured as follows, where the right-hand column represents the value of  $f_i$  for various values of  $i$ .

Value of $i$	Value of $f_i$
...	...
$i^*+3$	$e^*$
$i^*+2$	$e^*$
$i^*+1$	$e^*$
$i^*$	
...	...
0	

Let us define  $s^* := \underline{Q^n}[(n(e))_{e \in E}]$  where for each  $e \neq e^*$   $n(e) = 0$  and  $n(e^*) = \infty$ . We can now complete the above picture:

Value of $i$	Value of $f_i$	Value of $I(\underline{s(i)})$
...	...	...
$i^*+3$	$e^*$	$s^*$
$i^*+2$	$e^*$	$s^*$
$i^*+1$	$e^*$	$s^*$
$i^*$		
...	...	...
0		

v) Let us now compute the truth value of  $\underline{s(i^*)}$ .

For all  $e \in E$ ,  $e_{i^*} = |\{k > i^*: f_k = e\}|$ . Given iv), if  $e = e^*$ ,  $e_{i^*} = \infty$ ; and if  $e \neq e^*$ ,  $e_{i^*} = 0$ . Thus  $(e_{i^*})_{e \in E} = (e_{i^*+1})_{e \in E}$ , and by i)  $I(\underline{s(i^*)}) = I(\underline{s(i^*+1)}) = s^*$   
 $f_{i^*} = I(f[(s(\mathbf{k}'))_{k \geq i^*}]) = I(f[(s(\mathbf{k}'))_{k \geq i^*+1}]) = f_{i^*+1} = e^*$

By iterating the reasoning, we see that for each  $i \geq 0$ ,  $I(\underline{s(i)}) = s^*$ . This reasoning shows that in any fixed point  $I$ , all the sentences in the series have the same truth value  $e_i$ . Furthermore, the argument shows that  $e_i$  must satisfy the following equation:

$$(17) \quad \underline{Q^n}(f[(e_i)^\infty]: \infty, -f[(e_i)^\infty]: 0) = e_i$$

Conversely, any truth value  $e_i$  which satisfies this equation yields a coherent valuation for the entire series (it follows from a result to be discussed shortly that this valuation can then be extended to a fixed point for the entire language)<sup>9</sup>.

<sup>9</sup> To see an application of the Uniformity Property, let us consider the following series:

(i)  $V = \{ \langle s_i, [\forall k: k > i](\neg \text{Tr}(s(\mathbf{k}+2)) \vee \text{Tr}(s(\mathbf{k}+1))) \rangle \}$

### 2.2.3 Behavior of infinite series in a trivalent system

Let us now consider a trivalent system which is an extension of a classical logic with generalized quantifiers, in the sense that for each Quantifier  $Q$ , for any numbers (including  $\infty$ )  $a, b$ ,

$$(18) \quad \underline{Q}^3(0: a, 1: b, \#: 0) = \underline{Q}^2(0: a, 1: b)$$

With these assumptions, we study the case of Yablo-series of the form:

$$S_Q = \{ \langle s(i), [Qk: k > i] \text{ Tr}(s(k)) \rangle : i \geq 0 \}.$$

Applied to the present case, the equation in 0 gives a necessary and sufficient condition for  $e$  to be a possible value of (all) the sentences in  $S_Q$ :

$$(19) \quad \underline{Q}^3(e: \infty, -e: 0) = e$$

If  $e$  is classical,  $\underline{Q}^3(e: \infty, -e: 0) = \underline{Q}^2(e: \infty, -e: 0)$  [by (18)]. In other words:

-The sentences in  $S_Q$  can be coherently assigned the value 1 iff  $\underline{Q}^2(0: 0, 1: \infty) = 1$

-The sentences in  $S_Q$  can be coherently assigned the value 0 iff  $\underline{Q}^2(0: \infty, 1: 0) = 0$

The results we have obtained so far can be summarized as follows:

(20) Let  $Q$  be a binary generalized quantifiers satisfying Permutation Invariance, Extension and Conservativity. Then:

- a. A binary valuation can be found in which  $S_Q$  has the value *true* iff  $\underline{Q}^2(0: 0, 1: \infty) = 1$
- b. A binary valuation can be found in which  $S_Q$  has the value *false* iff  $\underline{Q}^2(0: \infty, 1: 0) = 0$
- c.  $S_Q$  is paradoxical iff no binary valuation can be found in which  $S_Q$  has the value *true* and no binary valuation can be found in which  $S_Q$  has the value *false*  
iff  $\underline{Q}^2(0: 0, 1: \infty) = 0$  and  $\underline{Q}^2(0: \infty, 1: 0) = 1$

This result allows us to determine without further reasoning whether a Yablo series is paradoxical or not.

Examples: (i)  $\underline{No}^2(0: 0, 1: \infty) = 0$  and  $\underline{No}^2(0: \infty, 1: 0) = 1$  and therefore  $S_{No}$ , i.e. the Universal Liar, is indeed paradoxical.

(ii)  $\underline{Some}^2(0: 0, 1: \infty) = 1$  and  $\underline{Some}^2(0: \infty, 1: 0) = 0$  and therefore  $S_{\exists}$  is an infinite Truth-Teller.

(iii)  $\underline{All}^2(0: 0, 1: \infty) = 1$  and  $\underline{All}^2(0: \infty, 1: 0) = 0$  and therefore  $S_{\forall}$  is an infinite Truth-Teller.

(iv) All but a finite number of<sup>2</sup>(0: 0, 1:  $\infty$ ) = 1 and All but a finite number of<sup>2</sup>(0:  $\infty$ , 1: 0) = 0 and therefore  $S_{\forall}$  is an infinite Truth-Teller.

By the Uniformity Property (General Case), if  $I$  is a fixed point compatible with  $V$ , all the sentences in  $V$  have the same value according to  $I$ .

*Case 1.* All sentences have the value 0. This is immediately absurd: for each  $i > 0$ ,  $\underline{s(i)}$  asserts something true, since for each  $k > 0$ ,  $\underline{s(k+2)}$  is false, and hence  $\neg \text{Tr}(\underline{s(k+2)})$  is true, as is  $\neg \text{Tr}(\underline{s(k+2)}) \vee \text{Tr}(\underline{s(k+1)})$ .

*Case 2.* All sentences have the value 1. No contradiction follows: for each  $i > 0$ , for each  $k > i$ ,  $\underline{s(k+1)}$  is true, and hence so is  $\neg \text{Tr}(\underline{s(k+2)}) \vee \text{Tr}(\underline{s(k+1)})$ .

Thus we see that all the sentences in  $V$  have the value *true* in any fixed point in which they have a classical value. Furthermore none of these sentences involves self-reference.

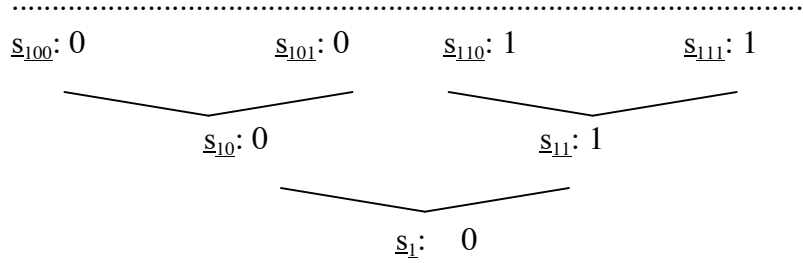
### 2.2.4 Failure of the Uniformity Property when the sentences are not linearly ordered

We already observed (following Cook 2004) that Yablo's paradox can be produced with much less than a linear ordering of sentences: a transitive relation without endpoints was shown to be sufficient to derive the paradox. Interestingly, the Uniformity Condition fails in some cases of this sort. Let us consider the following series, which is a modified  $\forall$ -Truth-Teller ( $y$  is a variable ranging over sentences, and for each  $i \geq 0$   $s_i$  is a sentence-denoting constant):

$$(21) \quad T_{R, \forall} := \{ \langle s_i, \forall s (s_i R s \Rightarrow \text{Tr}(s)) \rangle : i \geq 0 \}$$

It is possible to interpret  $R$  as a transitive relation  $\underline{R}$  without endpoints, and yet to find a valuation for  $T_{R, \forall}$  which does not assign the same value to all the sentences. Here is an example:

- (22) For convenience we index sentences with natural numbers written in binary notation, i.e. as strings of 0's and 1's. We stipulate that  $\underline{s}_i \underline{R} \underline{s}_k$  iff  $k$  written in binary notation is  $i$  concatenated with 0 or 1. We consider the valuation that assigns the value 0 to sentences whose index is 1 or starts with 10. It is a coherent valuation for the series  $\{ \langle s_{ij}, \forall s (s_{ij} R s \Rightarrow \text{Tr}(s)) \rangle : i, j \geq 0 \}$ .

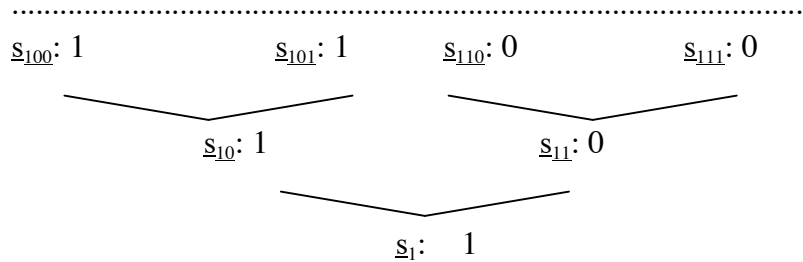


In fact, the valuation also yields a counter-example to the Uniformity Property when applied to the modified  $\forall\exists$ -Truth-Teller:

$$(23) \quad \{ \langle s_i, \forall s (s_i R s \Rightarrow (\exists s' (s R s' \wedge \text{Tr}(s')))) \rangle : i \geq 0 \}$$

By permuting the semantic values '0' and '1' in the valuation, we also obtain a counter-example to the modified  $\exists$ - and  $\exists\forall$ -Truth-Tellers:

- (24) a.  $\{ \langle s_i, \exists s (s_i R s \wedge \text{Tr}(s)) \rangle : i \geq 0 \}$   
 b.  $\{ \langle s_i, \exists s (s_i R s \wedge \forall s' (s R s' \Rightarrow \text{Tr}(s'))) \rangle : i \geq 0 \}$



### 2.3 Extending Local Fixed Points to Global Fixed Points

We have been proceeding throughout as if one could interchangeably talk about the existence of a *valuation* for a Yablo-style series and the existence of a *fixed point* for it. The

first perspective is local: it only requires that the truth predicate  $Tr$  be coherently interpreted *when we restrict attention to the sentences of the series*. The second perspective is global, and seeks to find a fixed point for the entire language. There is indeed an equivalence between the two notions, but only because the series we have considered are 'semantically autonomous', in the sense that all the sentences whose semantic status could potentially affect their truth value are themselves included in the series. Let us now prove this result.

It is clear that a global fixed point will yield a valuation for the series. To show the converse, we will proceed in two steps.

-First, we start from a trivalent valuation  $V$  for the Yablo series  $S$ . We then define a trivalent interpretation  $I'$  by:  $I'^+(Tr) = \{s \in S: V(s) = 1\}$ ;  $I'^-(Tr) = \{s \in S: V(s) = 0\}$ . This will yield a 'local fixed point' for the series itself, in the sense that when one restricts attention to sentences in  $S$ ,  $Tr$  is true of those sentences that are in fact true, and is false of those sentences that are false. This is defined formally in (25)

(25) Local fixed point

Let  $I'$  be a trivalent interpretation.  $I'$  is a local fixed point for a set  $S$  of sentences just in case:

$$\forall s \in S [(I'(s) = 1 \leftrightarrow s \in I'^+(Tr)) \wedge (I'(s) = 0 \leftrightarrow s \in I'^-(Tr))]$$

-Second, an Extension Lemma will show how this 'local' fixed point can be extended to a global fixed point for the entire language. The argument is a modification of the proof that Kripke 1975 gave to show that some fixed point always exists for certain evaluation schemes. As in Kripke's construction, the crucial assumption is that the trivalent logic is monotonic, in the sense that when we start from an interpretation  $I'$  based on an extension  $E$  and anti-extension  $A$  for the truth predicate, if we 'add' sentences to  $E$  or to  $A$  we will only have to revise the value of sentences that heretofore had the value  $\#$  (that is, the 'classical' values won't have to be revised when we resolve some of the indeterminates as being true or false). This condition is stated more precisely in (26).

(26) Monotonicity

An evaluation scheme is monotonic if each language  $L$  and for all interpretations  $I_1$  and  $I_2$  of  $L$  the following holds:

if: for each predicate  $P$  of  $L$ ,  $I_1^+(P) \subseteq I_2^+(P)$  and  $I_1^-(P) \subseteq I_2^-(P)$

then:  $\{s \in L: I_1(s) = 1\} \subseteq \{s \in L: I_2(s) = 1\}$  and  $\{s \in L: I_1(s) = 0\} \subseteq \{s \in L: I_2(s) = 0\}$

The key observation is that the series we study only 'talk about' other sentences in the series. These series are thus 'semantically autonomous', in the following sense:

(27) Semantic Autonomy

Let  $I$  be a classical interpretation.  $S$  is *semantically autonomous* in  $I$  just in case:

for all trivalent interpretations<sup>10</sup>  $I'$ ,  $I''$  that extend  $I$  (by assigning an extension and anti-extension to the truth predicate), if  $\forall s \in S [(s \in I'^+(Tr)) \leftrightarrow s \in I''^+(Tr)] \wedge (s \in I'^-(Tr) \leftrightarrow s \in I''^-(Tr))$ , then  $\forall s \in S I'(s) = I''(s)$ .

We could prove in greater detail that Yablo-Series are semantically autonomous in any interpretation compatible with the naming relation that they define, but for brevity we will accept that this is indeed so<sup>11</sup>. We turn directly to the main Lemma:

<sup>10</sup> Note that the requirement holds of *all* interpretations, not just fixed points.

<sup>11</sup> In a nutshell, the argument is that all the sentences we have considered can be expressed in terms of a special form of restricted quantification, with each quantifier being of the form  $[\forall x: F]_{\_}$  or  $[\exists x: F]_{\_}$  for some *classical* formulas  $F$ ,  $F'$  (henceforth called *restrictors*). Conservativity implies that only the objects satisfying the restrictors  $F$ ,  $F'$  need to be taken into account when evaluating the truth of these formulas. If the restrictors

(28) *Extension Lemma*

Consider a monotonic evaluation scheme. Let  $I'$  be a local fixed point for a set  $S$  of sentences of  $L$  which is semantically autonomous in  $I'$ . Then there is a global fixed point for  $L$  which agrees with  $I'$  on  $S$ .

*Proof:* The proof extends Kripke's technique to the case at hand. The desired interpretation is defined as the union of a (transfinite) series of interpretations which need not themselves be fixed points.

We define a series of interpretations  $I'_i$  (for ordinal  $i$ ) each of which is fully determined by the pair  $\langle E_i, A_i \rangle$  of the extension and anti-extension of the truth predicate. If  $I'_i$  is determined by  $\langle E_i, A_i \rangle$ , we further define

$$f(\langle E_i, A_i \rangle) := \langle \{d \in L: I'_i(d)=1\}, \{d \in L: I'_i(d)=0\} \rangle$$

(note that we still assume that the truth predicate  $Tr$  only takes sentence-denoting terms as arguments; without this assumption we would have to give a slightly more complicated definition for  $f$ ).

To define the interpretations  $I'_i$ , it suffices to define the pairs  $\langle E_i, A_i \rangle$ . We do so by the following induction:

$$1) \langle E_0, A_0 \rangle := \langle I'^+(Tr) \cap S, I'^-(Tr) \cap S \rangle$$

2) If  $i$  is a successor ordinal  $k+1$ , we set:

$$\langle E_i, A_i \rangle := f(\langle E_k, A_k \rangle)$$

3) If  $i$  is a limit ordinal, we set:

$$\langle E_i, A_i \rangle := \langle \bigcup_{k < i} E_k, \bigcup_{k < i} A_k \rangle$$

We prove by induction that the property  $\pi(i)$  holds of all ordinals  $i$ :

$\pi(i) : \text{for all } i', i'' \text{ for which } i'' \leq i' \leq i, (a) \langle E_{i''}, A_{i''} \rangle \subseteq \langle E_{i'}, A_{i'} \rangle \text{ and } (b) \langle E_{i''} \cap S, A_{i''} \cap S \rangle = \langle E_{i'} \cap S, A_{i'} \cap S \rangle$

(1)  $\pi(0)$  is trivially true.

(2) Suppose that  $i$  is a successor ordinal  $k+1$ . Then  $\langle E_i, A_i \rangle = f(\langle E_k, A_k \rangle)$ . By the Induction Hypothesis, for each  $k' \leq k$ ,  $\langle E_{k'}, A_{k'} \rangle \subseteq \langle E_k, A_k \rangle$ . By the monotonicity of  $f$ , it follows that for each  $k' \leq k$ ,  $f(\langle E_{k'}, A_{k'} \rangle) \subseteq f(\langle E_k, A_k \rangle)$ , i.e. that  $\langle E_{k'+1}, A_{k'+1} \rangle \subseteq \langle E_{k+1}, A_{k+1} \rangle$ . If  $k$  is a successor ordinal  $k'+1$ ,  $k' \leq k$  and  $\langle E_k, A_k \rangle \subseteq \langle E_{k+1}, A_{k+1} \rangle$ . If  $k$  is a limit ordinal,  $\langle E_k, A_k \rangle = \bigcup_{k' < k} \langle E_{k'}, A_{k'} \rangle \subseteq \bigcup_{k' < k} \langle E_{k'+1}, A_{k'+1} \rangle \subseteq \langle E_{k+1}, A_{k+1} \rangle$ . In all cases  $\langle E_k, A_k \rangle \subseteq \langle E_{k+1}, A_{k+1} \rangle$ , which together with the Induction Hypothesis yields part (a) of  $\pi(i)$ .

To prove part (b), we observe that by assumption (' $S$  is semantically autonomous'), the semantic value of members of  $S$  is fixed by the restriction of the interpretation of  $Tr$  to  $S$ . From the Induction Hypothesis it follows that for each  $k' \leq k$ ,  $\langle E_{k'} \cap S, A_{k'} \cap S \rangle = \langle E_k \cap S, A_k \cap S \rangle$ , whence  $\langle E_{k'+1} \cap S, A_{k'+1} \cap S \rangle = \langle E_{k+1} \cap S, A_{k+1} \cap S \rangle$ . If  $k$  is a successor ordinal  $k'+1$ ,  $k' \leq k$  and  $\langle E_{k'+1} \cap S, A_{k'+1} \cap S \rangle = \langle E_{k+1} \cap S, A_{k+1} \cap S \rangle$ , i.e.  $\langle E_k \cap S, A_k \cap S \rangle = \langle E_{k+1} \cap S, A_{k+1} \cap S \rangle$ . If  $k$  is a limit ordinal,  $\langle E_k \cap S, A_k \cap S \rangle = \bigcup_{k' < k} (\langle E_{k'}, A_{k'} \rangle \cap \langle S, S \rangle) = \langle E_0, A_0 \rangle$  (again thanks to the Induction Hypothesis), and thus  $\langle E_{k+1} \cap S, A_{k+1} \cap S \rangle = \langle E_k \cap S, A_k \cap S \rangle$ .

(3) Suppose that  $i$  is a limit ordinal  $k$ . Then  $\langle E_i, A_i \rangle = \bigcup_{k' < i} \langle E_{k'}, A_{k'} \rangle$ , which given the Induction Hypothesis immediately yields  $\pi(i)$ .

The series  $\langle E_i, A_i \rangle$  is increasing on the ordinals and thus it must have a fixed point<sup>12</sup>  $\langle E_{i^*}, A_{i^*} \rangle$ , which is the interpretation we were looking for.

---

only hold true of other sentences in the series, Semantic Autonomy is guaranteed. A development of this idea is used in the Appendix to find a sufficient condition of non-self-reference.

<sup>12</sup> Otherwise, at each successor ordinal there should be a sentence that is first declared true at this level; we could thus define a function  $f$  that associates to each sentence the ordinal at which it is first declared true, if this ordinal exists; and 0 otherwise. The domain of  $f$  is the set of sentences of the language. Its range is a proper class of ordinals (since it includes the class of all successor ordinals). But this contradicts the Axiom of Replacement.

### 3 Elimination of Self-Reference for a Language without Quantifiers I: Basic Case

As was mentioned earlier, there is an important conceptual difference between the  $\exists\forall$ - and the  $\forall\exists$ -Liars on the one hand and the  $\forall$ - and  $\exists$ -Liars on the other. In this section we will show that the mechanism at work in the former can be generalized to give a procedure by which a language  $L'$  can be translated into a self-reference-free fragment of a language  $L^*$  which includes quantifiers over natural numbers and functional names of sentences. We will show that our construction *fails* if we try to generalize the procedure at work in the  $\forall$ - and  $\exists$ -Liars, and we will provide a characterization of those Generalized Quantifiers that can be used in the translation. For simplicity we assume throughout this section that the initial language has no quantifiers.

#### 3.1 Preliminaries

##### 3.1.1 Languages

Let us start by presenting in greater detail the languages we will be working with. We start from a classical language  $L$ , to which we add sentence-denoting constants and a truth predicate  $Tr$  to obtain an enriched language  $L'$ . We stipulate in the syntax of  $L'$  that the  $Tr$  can only take sentence-denoting constants as arguments, which will simplify a bit the definition of what a 'fixed point' should be<sup>13</sup>.

(29) Syntax of the base language  $L$

-Object-denoting Terms:	$o := c_i$	[the $c_i$ 's are constants]
-Object-denoting Predicates:	$P := P_i^n$	
-Formulas:	$F := P_i^n(o_1, \dots, o_n) \mid \neg F \mid (F \wedge F) \mid (F \vee F)$	

(30) Syntax of the enriched language  $L'$

-Object-denoting Terms:	$o := c_i$	[the $c_i$ 's are constants]
-Sentence-denoting terms:	$s := s_i$	[the $s_i$ 's are constants]
-Object-denoting Predicates:	$P := P_i^n$	
-Sentence-denoting Predicate:	$Tr$	[ $Tr$ is the Truth predicate]
-Formulas:	$F := P_i^n(o_1, \dots, o_n) \mid Tr(s) \mid \neg F \mid (F \wedge F) \mid (F \vee F)$	

The translation language  $L^*$  is richer. Besides the truth predicate  $Tr$ , it extends  $L$  with (i) some simple arithmetic vocabulary, and (ii) function symbols, which take as arguments number-denoting terms to yield sentence-denoting terms. Concretely, whenever  $s_i$  is a sentence-denoting constant in  $L'$ ,  $s_i(\mathbf{n})$  is a sentence-denoting functional term in  $L^*$ . (We also include in the definition of  $L^*$  sentence-denoting variables, but these won't be used except in the Appendix).

(31) Syntax of the target language  $L^*$

-Object-denoting terms:	$o := c_i$	[the $c_i$ 's are constants]
-Number-denoting terms:	$n := k_i \mid 0 \mid S n$	[the $k_i$ 's are variable]
	We abbreviate $S^n(0)$ as $\mathbf{n}$	

<sup>13</sup> If  $Tr$  could take other terms as arguments, we would have to stipulate whether, say, a chair or a table fall under the extension or the anti-extension of the truth predicate; the problem does not arise given our definitions.

-Sentence-denoting functions:	$S := s_i$
-Sentence-denoting terms:	$s := S(n) \mid x_i$ [the $x_i$ 's are variables] We abbreviates $s_i(k)$ as $s_{i,k}$ or even $s_{ik}$ .
-Object-denoting Predicates:	$P := P_i^n$
-Truth Predicate:	$Tr$
-Number-denoting predicate:	$<$
-Formulas:	$F := P_i^n(o_1, \dots, o_n) \mid Tr(s) \mid \neg F \mid (F \wedge F) \mid$ $(F \vee F) \mid \exists k_i F \mid \forall k_i F \mid \exists x_i F \mid \forall x_i F$

The semantics is given in the usual way given the specification of (i) an evaluation scheme (we will generally restrict attention to Kleene's Strong Kleene Logic), (ii) a base interpretation  $I$  for  $L$ , (iii) a denotation function for the sentence-denoting terms, and (iv) the specification of an extension and an anti-extension for the truth predicate.

### 3.1.2 Goal

Unless otherwise noted, we work with the Strong Kleene evaluation scheme, whose main clauses are as follows ( $s$  is an assignment function):

- (32) a.  $I_s(\neg F)=1$  iff  $I_s(F)=0$ ;  $=1$  iff  $I(F)=1$   
 b.  $I_s(F \wedge G)=1$  iff  $I(F)=I(G)=1$ ;  $=0$  iff  $I(F)=0$  or  $I(G)=0$   
 c.  $I_s(\forall x F)=1$  iff for each  $d$  in the domain,  $I_{s[x \rightarrow d]}(F)=1$ ;  $=0$  iff for some  $d$  in the domain,  $I_{s[x \rightarrow d]}(F)=0$ .

The Strong Kleene scheme is not usually defined for generalized quantifiers. But there is a natural way to extend it to the cases we will be concerned with. To see what the basic idea is, let us remember that the Strong Kleene scheme can be taken to capture the following intuition: the value # represents 'lack of knowledge', which could be resolved as 0 or 1 upon consideration of further information. The rules are designed in such a way that a formula is true in the trivalent system if it can be guaranteed to be true in classical logic no matter how the #'s are resolved. If we now consider a generalized quantifier  $Q$  whose restrictor is classical, we have already posited (under 'Reasonableness' in (14)) that its semantics should be determined by three numbers: the number  $a$  of elements that satisfy the restrictor and of which the nuclear scope is false; the number  $b$  of elements that satisfy the restrictor and of which the nuclear scope is false; and the number  $c$  of elements that satisfy the restrictor and of which the nuclear scope is neither true nor false. To extend the Strong Kleene scheme, it is natural to ask whether no matter how the members of the last group (corresponding to  $c$ ) are redistributed among the first two, we do or do not obtain the same classical truth value. If we do, the formula must have the value in question (0 or 1); if we don't, the formula should have the value #. This recipe is captured by the following rule:

- (33)  $\underline{Q}^3(0: a, 1: b, #: c)=\#$  iff for some  $a', a'', b', b''$  satisfying  $a'+b'=a''+b''=c$ ,  $\underline{Q}^2(0: a+a', 1: b+b') \neq \underline{Q}^2(0: a+a'', 1: b+b'')$ .  
 If  $\neq \#$ ,  $\underline{Q}^3(0: a, 1: b, #: c)=\underline{Q}^2(0: a+a', 1: b+b')$ , where  $a'+b'=c$

The translation procedure will keep constant the interpretation  $I$  for the classical language  $L$ . But it will simultaneously (i) assign to each sentence of the extended language  $L'$  a translation in the quantificational language  $L^*$ , and (ii) replace the (old) denotation function  $N'$  for sentence names of  $L'$  with a new denotation function  $N^*$  for the sentence-denoting terms of  $L^*$ . For simplicity, we will call interpretation of  $L'$  and  $L^*$  *admissible* if they extend



$I$  and are compatible with  $N'$  and  $N^*$  respectively. We will assign to each sentence  $F$  of  $L'$  an infinite series of translations  $h_k(F)$  ( $k \geq 0$ ), and we will call  $h(L')$  the set of all translations of all sentences of  $L'$ . The procedure will be shown to have two main properties:

**Property 1** will guarantee that in any admissible fixed point  $I^*$  of  $L^*$ , all the translations of a given sentence  $F$  of  $L'$  have the same value according to  $I^*$  [i.e. for every sentence  $F$  of  $L'$ , for all  $k, k' \geq 0$ ,  $I^*(h_k(F)) = I^*(h_{k'}(F))$ ]. Property 1 is important to ensure that any translation of a given sentence  $F$ , or alternatively the equivalence class of all its translations, can be taken as 'the' translation of  $F$ . We will sometimes refer to Property 1 as the Uniformity Condition, because it can be seen as a generalization of the Uniformity Property which was shown earlier to hold of any Yablo-series. Specifically, the Yablo-series we discussed earlier will turn out to be the set of translations of the simple Truth-Teller under our scheme. The Uniformity Property only requires that the translations of the Truth-Teller have a uniform value; the Uniformity Condition requires that the translations of *all* sentences receive a uniform truth value in any fixed point.

**Property 2**, which is a bit more cumbersome to state, will guarantee that there is a kind of isomorphism  $j$  between the fixed points of  $L'$  and those of  $L^*$ , and that for any sentence  $F$  of  $L'$ ,  $F$  has the same value in a fixed point  $I'$  as  $h_i(F)$  does in  $j(I')$ . This property is important to guarantee that the semantic behavior of  $L'$  is indeed reflected in the translation. However because only part of  $L^*$  serves to translate  $L'$ , we will have to treat as equivalent fixed points of  $L^*$  that agree on the translation of  $L'$  (i.e. on  $h(L')$ ). As a result, the isomorphism must be defined between the set of fixed points of  $L'$  and the set of *equivalence classes* of fixed points of  $L^*$ . The relevant notions are defined with greater precision below.

### 3.2 Translation and Examples

We start by defining the translation procedure and by illustrating it with some examples.

#### 3.2.1 Translation

As before, we call  $I'$  the interpretation of the initial language  $L'$ , which comprises: (i) a classical interpretation  $I$ , (ii) a denotation function  $N'$  for the sentence-denoting names of  $L'$ , and (iii) a specification of the extension and anti-extension of the truth predicate  $Tr$ . An interpretation  $I^*$  for the target language  $L^*$  will be defined by (i') the same classical interpretation  $I$ , (ii'a) a denotation relation function  $N^*$  for the (functional) sentence-denoting terms of  $L^*$ , (ii'b) an interpretation of the arithmetic vocabulary of  $L^*$  [in the standard model of the natural numbers], and (iii') a specification of the extension and anti-extension of  $Tr$ . We will build simultaneously a translation and a specification of  $N^*$  from  $N'$ , and we will show that this suffices to force the translations to mirror perfectly the behavior of the originals. For brevity we use  $[Qk': k' > k]F$  to abbreviate  $\exists k'' (k'' > k \wedge \forall k' (k' \geq k'' \rightarrow F))$  (later in the paper we will study other conceivable choices of  $Q$  when  $Q$  is a Generalized Quantifier).

- (34) a. *Translation*: For each positive integer  $i$ ,  $h_i(F) = [Qk': k' > i] [F]_k$ , where  $k$  and  $k'$  are 'fresh' number-denoting variables, and where  $[F]_k$  is obtained from  $F$  by replacing every atomic formula of the form  $Tr(c)$  with  $Tr(c(k'))$ .  
 b. *Denotation*:  $s$  denotes  $F$  according to  $N'$  iff  $s(i)$  denotes  $h_i(F)$  according to  $N^*$ .

As before, we write  $\langle s, F \rangle$  for a pair of a formula  $F$  denoted by a sentence-denoting term  $s$ , and we write the set of translations-cum-denotation relation as  $h(\langle s, F \rangle) = \{\langle s(i), h_i(F) \rangle : i \geq 0\}$ .

As was observed earlier, our choice of  $Q$  guarantees that all of the sentences of the form  $[Qk': k' > i] [F]_{k'}$  for various values of  $i$  have exactly the same semantic content and hence the same truth value in any fixed point. This will be essential to guarantee that the Uniformity Condition is satisfied.

### 3.2.2 Examples

Before we study the general properties of the translation procedure, let us illustrate some of its effects.

1. First, we check that the translation is adequate for sentences that do not contain the truth predicate, say *It is raining*, symbolized as  $R$ , and named by a constant  $r$  (we henceforth call a sentence *Tr-free* if it does not contain the truth predicate). Since  $R$  contains no occurrence of the truth predicate, the translation procedure yields a sentence with vacuous quantification, as follows:

$$(35) \quad h(\langle r, R \rangle) = \{\langle r(i), [Qk': k' > i] R \rangle : i \geq 0\}$$

The quantification is vacuous, and given the semantics of  $Q$  ('all but finitely many'), it is immediate that in any interpretation all the translations are equivalent to  $R$ , as is desired.

2. Second, we observe that a sentence that talks about the truth of a *Tr-free* sentence is correctly translated. Let us consider a sentence (named by a constant  $r'$ ) which says that  $r$  is true, yielding a pair  $\langle r', \text{Tr}(r) \rangle$ . The translation procedure yields:

$$(36) \quad h(\langle r', \text{Tr}(r) \rangle) = \{\langle r'(i), [Qk': k' > i] \text{Tr}(r(k')) \rangle : i \geq 0\}$$

We have already established that all the sentences  $\underline{r}(i)$  (for  $i \geq 0$ ) are equivalent to  $\underline{r}$ . It follows that in any fixed point all the sentences  $\underline{r'}(i)$  are also equivalent to  $R$ , and hence to  $\text{Tr}(r)$ , as is desired.

3. Third, let us consider the Liar. We have no new work to do, since we already discussed its translation when we introduced the modified version of Yablo's paradox. As is desired, the Liar  $\langle s, \neg \text{Tr}(s) \rangle$  gets translated as a Yablo-like series which is itself paradoxical, namely  $\{\langle s(i), [Qk': k' > i] \neg \text{Tr}(s(k')) \rangle : i \geq 0\}$ . Since this series has a uniform value, we immediately obtain the result that in any fixed point each sentence in the series should be neither true nor false.

4. Fourth, we should consider the Truth-Teller  $\langle t, \text{Tr}(t) \rangle$ . It is translated as  $\{\langle t(i), \forall k (k > i \rightarrow [Qk': k' > i] \text{Tr}(t(k')) \rangle : i \geq 0\}$ . As before, the form of the translations guarantees that in any interpretation they must all share the same value. It is then easy to see that there are fixed points in which these sentences are true, others in which they are false, and yet others in which they are undefined.

5. Fifth, let us reconsider our empirical versions of the Liar and of Truth-Teller, which we gave respectively as  $\langle e, R \wedge \neg \text{Tr}(e) \rangle$  and  $\langle e', R \wedge \text{Tr}(e') \rangle$ . They are translated as  $\{\langle e(i), [Qk': k' > i] (R \wedge \neg \text{Tr}(e(k')) \rangle : i \geq 0\}$  and as  $\{\langle e'(i), [Qk': k' > i] (R \wedge \text{Tr}(e'(k')) \rangle : i \geq 0\}$ . Reasoning by cases, we see that if  $R$  is false we simply obtain two series of false sentences; and if  $R$  is true, we obtain an infinite Liar and infinite Truth-Teller, as we wished.

### 3.3 Properties of the Construction

We now study the main properties of this translation scheme. We start by observing that  $h(L')$  is semantically autonomous in interpretation of  $L^*$  which is compatible with  $N^*$ . The proof, which is fastidious, would rely on the observation that all the translations can be seen as instances of restricted quantifications of the form  $[Qs: R(s)] \text{ ---}$ , where  $s$  ranges over sentences and where  $R(s)$  is a classical formula that holds true solely of other members of  $h(L')$ . As a result, all interpretations that agree on the restriction of the extension and anti-extension of  $Tr$  to  $h(L')$  must also agree on the values they assign to the members of  $h(L')$ . This is just to say that  $h(L')$  is semantically autonomous (related ideas are implemented in greater detail in the Appendix).

*Definition:*  $I'$  is an *admissible fixed point for  $L'$*  just in case  $I'$  is a fixed point for  $L'$  based on (i) the base interpretation  $I$  and (ii) the denotation function  $N'$  for sentence-names.  $I^*$  is an *admissible fixed point for  $L^*$*  just in case  $I^*$  is a fixed point for  $L^*$  based on (i') the base interpretation  $I$ , (ii'a) the denotation function  $N^*$  for functional sentence names, and (ii'b) the standard interpretation of the arithmetic vocabulary.

**Property 1 (=Uniformity Condition).** For every sentence  $F$  of  $L'$ , for all  $k, k' \geq 0$ ,  $I^*(h_k(F)) = I^*(h_{k'}(F))$ .

*Proof (Sketch):* In the Strong Kleene Scheme, it can be checked that for all  $k \geq 0$ ,

$I^*(h_k(F)) = 1$  iff  $[F]_{k'}$  is true of all but a finite number of the non-negative integers.

$I^*(h_k(F)) = 0$  iff  $[F]_{k'}$  is false of an infinite number of non-negative integers.

$I^*(h_k(F)) = \#$  iff  $I^*(h_k(F)) \neq 1$  and  $I^*(h_k(s)) \neq 0$

The right-hand sides make no reference to  $k$ , and therefore for all  $k, k' \geq 0$ ,  $I^*(h_k(F)) = I^*(h_{k'}(F))^{14}$ .

In order to state Property 2, we must define an equivalence relation over admissible fixed points of  $L^*$ , and an ordering over them:

*Definitions:* (i) If  $I^*_1$  and  $I^*_2$  are admissible fixed points of  $L^*$ ,  $I^*_1 \approx I^*_2$  iff  $I^*_1$  and  $I^*_2$  agree on  $h(L')$ . We write  $[I^*_1]$  for the equivalence class of  $I^*_1$ .

(ii) For any set of sentences  $S$ , we define a partial order on interpretations by stipulating that  $i \leq_s j$  just in case every sentence of  $S$  that has a classical truth value in  $i$  has the same value in  $j$ .

**Property 2 (=Isomorphism Condition)** There is an isomorphism  $j$  between the set of admissible fixed points of  $L'$  ordered by  $\leq_{L'}$  and the set of equivalence classes of admissible

---

<sup>14</sup> Here is a more complete argument. Suppose that  $k$  denotes an integer  $K$ . Then for any formula  $F$  with one free variable  $k'$ ,

1)  $I^*(\exists k'' (k'' > k \wedge \forall k' (k' \geq k'' \rightarrow F[k']))) = 1$  iff for some  $K'' > K$ ,  $I^*(\forall k' (k' \geq K'' \rightarrow F[k'])) = 1$ , iff for some  $K'' > K$ , for each  $K' \geq K''$ ,  $I^*(F[K']) = 1$ . If this condition is interpreted in a standard model of the integers, it is equivalent to:  $F$  is true of all but a finite number of the integers.

2)  $I^*(\exists k'' (k'' > k \wedge \forall k' (k' \geq k'' \rightarrow F[k']))) = 0$  iff for all  $K''$ ,  $I^*(K'' > k \wedge \forall k' (k' \geq k'' \rightarrow F[k'])) = 0$ , iff for all  $K'' > K$ ,  $I^*(\forall k' (k' \geq K'' \rightarrow F[k'])) = 0$ , iff for all  $K'' > K$ , for some  $K' \geq K''$ ,  $I^*(F[K']) = 0$ . If this condition is interpreted in a standard model of the integers, it is equivalent to:  $F$  is false of an infinite number of integers.

fixed points of  $L^*$  ordered by  $\leq_{h(L')}$  and  $j$  guarantees that for every sentence  $F$  of  $L'$ , for every fixed point  $I'$  of  $L'$ ,  $I'(F)=j(I')(h_i(F))$ .

We note for future reference that the only properties of the quantifier  $Q$  which matter in the proof are that (Q1)  $Q$  satisfies the Uniformity Condition, and that (Q2) when  $F$  contains no bound variables,  $[Qk': k'>i] F$  has the same value as  $F$ .

Proof (Sketch): We write that  $J(I', [I^*]) = \text{just in case } I' \text{ and } I^* \text{ are admissible fixed points for } L' \text{ and } L^* \text{ respectively, and } I^{*+}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^+(\text{Tr})\}, I^{*-}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^-(\text{Tr})\}.$

1) Let  $I'$  be an admissible fixed point for  $L'$ . We show that there is exactly one equivalence class of admissible fixed points  $[I^*]$  for  $L^*$  satisfying  $J(I', [I^*])$ .

- 'At most one': given  $N^*$  and  $I$ , the truth value of any member of  $h(L')$  is fixed by the restriction of the interpretation of  $Tr$  to  $h(L')$  [because  $h(L')$  is semantically autonomous]. As a result, once  $I^{*+}(\text{Tr}) \cap h(L')$  and  $I^{*-}(\text{Tr}) \cap h(L')$  are fixed, so is the value of each of the members of  $h(L')$ .

- 'At least one': we show how to construct an admissible fixed point  $I^*$  for  $L^*$  which satisfies  $J(I', [I^*])$ . We start by defining an interpretation  $I^*_0$  which is a local fixed point for  $h(L')$ , and we extend to a global fixed point for  $L^*$  thanks to the Extension Lemma.

(i)  $I^*_0$  is defined by:

$$I^{*+}_0(\text{Tr}) = \{h_k(F): k \geq 0 \text{ and } F \in I'^+(\text{Tr})\}$$

$$I^{*-}_0(\text{Tr}) = \{h_k(F): k \geq 0 \text{ and } F \in I'^-(\text{Tr})\}$$

$I^*_0$  can be shown to be a fixed point of  $h(L')$  because for each sentence  $F$  of  $L'$ ,

$$h_i(F) \in I^{*+}_0(\text{Tr}) \text{ (resp. } I^{*-}_0(\text{Tr})) \text{ iff } F \in I'^+(\text{Tr}) \text{ (resp. } I'^-(\text{Tr}))$$

$$\text{iff } I'(F) = 1 \text{ (resp. } = 0) \text{ [because } I' \text{ is a fixed point]}$$

iff for every  $k \geq 0$ ,  $I^*_0([F]_k) = 1$  (resp.  $= 0$ ), where  $[F]_k$  is obtained from  $F$  by replacing each occurrence of the form  $Tr(c)$  with  $Tr(c(k))$  [this follows because given the definition of  $I^*_0$ ,  $I'(\text{Tr}(c)) = I^*_0(\text{Tr}(c(k)))$  for arbitrary  $k \geq 0$ ]

iff  $I^*_0([Qk': k'>i] [F]_k) = 1$  (resp.  $= 0$ ) [this follows because (i)  $Q$  satisfies property (Q2): when  $F$  contains no bound variables,  $[Qk': k'>i] F$  has the same value as  $F$ , and (ii) for all  $k', k'' \geq 0$ ,  $I^*_0([F]_{k'}) = I^*_0([F]_{k''})$ ]

$$\text{iff } I^*_0(h_i(F)) = 1 \text{ (resp. } = 0)$$

(ii) By the Extension Lemma, this local fixed point can be extended to a global fixed point.

2) Let  $I^*$  be an admissible fixed point for  $L^*$ . We show that there is exactly one admissible fixed point  $I'$  for  $L'$  satisfying  $J(I', [I^*])$ .

Given the Uniformity Condition, for all  $k, k' \geq 0$ ,  $I^*(h_k(s)) = I^*(h_{k'}(s))$ . Given  $I$  and  $N'$ , we can thus define an interpretation  $I'$  by  $I'^+(\text{Tr}) = \{s: \text{for some } k \geq 0, h_k(s) \in I^{*+}\}$  and  $I'^-(\text{Tr}) = \{s: \text{for some } k \geq 0, h_k(s) \in I^{*-}\}$ . It is then immediate that  $I^{*+}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^+(\text{Tr})\}$ ,  $I^{*-}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^-(\text{Tr})\}$ . All that remains to be shown is that  $I'$  is a fixed point.



#	#	#	...
1	#	#	...
1	1	#	...
1	1	1	...
...	...	...	...

Each column represents the values of the translations  $s_i(\mathbf{0})$ ,  $s_i(\mathbf{1})$ ,  $s_i(\mathbf{2})$ , ... of a given sentence  $s_i$  of the original language (for instance, the left-most column indicates that  $s_0(\mathbf{0})$  has the value #, that  $s_0(\mathbf{1})$  has the value 1, that  $s_0(\mathbf{2})$  has the value 1, etc). It can be checked that this valuation is indeed coherent (i.e. that it defines a local fixed point). For instance,  $s_0(\mathbf{0})$  must indeed have the value # because  $s_1(\mathbf{1})$  has the value #. Specifically,  $s_0(\mathbf{0})$  is the formula  $[Qk': k' > 0] Tr(s_i(k'))$ . But there is an element (namely 1) that satisfies the restrictor, and for which the nuclear scope has the value #. Hence by (37)  $s_0(\mathbf{0})$  must have the value #. Although it is coherent, this valuation clearly violates the Uniformity Condition, since in each column we find both the value # and the value 1<sup>16</sup>.

## 4.2 Other Generalized Quantifiers

Let us now restrict attention to the Strong Kleene evaluation scheme, extended to Generalized Quantifiers in accordance with the following rule, already discussed above:

- (38)  $\underline{Q}^3(0: a, 1: b, #: c) = \#$  iff for some  $a', a'', b', b''$  satisfying  $a' + b' = a'' + b'' = c$ ,  $\underline{Q}^2(0: a + a', 1: b + b') \neq \underline{Q}^2(0: a + a'', 1: b + b'')$ . If  $\neq \#$ ,  $\underline{Q}^3(0: a, 1: b, #: c) = \underline{Q}^2(0: a + a', 1: b + b')$ , where  $a' + b' = c$

In our construction  $Q$  ended up having the meaning of 'all but finitely many' (because of our assumption that the sentences were denumerable and linearly ordered, this turned to be expressible in first-order logic). Could we have used other Generalized Quantifiers? We will now give a characterization of those quantifiers that can be used in the construction, by proceeding in two steps:

1. We revisit the translation scheme given above, treating  $Q$  as a parameter rather than as the abbreviation of a particular string of first-order quantifiers. We find necessary and sufficient conditions that  $Q$  must meet if the translation is to satisfy the Uniformity Condition.
2. We then isolate those values of  $Q$  for which the translation procedure delivers the desired results.

### □ Necessary and Sufficient Conditions for Uniformity<sup>17</sup>

We claim that in the Strong Kleene evaluation scheme, the Uniformity Condition is satisfied by a very narrow class of Generalized Quantifiers, which satisfy a condition of Finite Insensitivity:

<sup>16</sup> Due to the simplicity of the series  $\{<s_i, Tr(s_{i+1})>: i \geq 0\}$ , it is hard to see how any translation could avoid this difficulty. But it would be interesting to settle this question in full generality, something that we do not attempt here.

<sup>17</sup> Thanks to Denis Bonnay for pointing out one error and one omission in an earlier version of this paragraph.

(39) The following conditions are equivalent:

*Uniformity Condition (UC):* For any sentence  $F$  of  $L'$ , in every admissible fixed point  $I^*$  of  $L^*$ , for all  $k$ ,  $k' \geq 0$ , for all  $s \in L'$ ,  $I^*(h_k(f)) = I^*(h_{k'}(F))$ .

*Finite Insensitivity (FI):* For all finite  $i$ ,  $i' \geq 0$ ,  $\underline{Q}^2(\infty, i) = \underline{Q}^2(\infty, i')$  and  $\underline{Q}^2(i, \infty) = \underline{Q}^2(i', \infty)$

### 1) Uniformity Condition $\Rightarrow$ Finite Insensitivity

Suppose that (FI) fails, for instance because there are  $i, i' \geq 0$  such that  $\underline{Q}^2(\infty, i) \neq \underline{Q}^2(\infty, i')$  (the case in which there are  $i, i' \geq 0$  such that  $\underline{Q}^2(i, \infty) \neq \underline{Q}^2(i', \infty)$  is treated in the same way by duality, i.e. by permuting  $1$  and  $0$  in the reasoning below).

**Case 1.**  $\underline{Q}^2(\infty, 0) = 0$ .

Let  $i^*$  be the least  $i$  such that  $\underline{Q}^2(\infty, i) = 1$ . Thus we have:

for each  $i \leq i^* - 1$   $\underline{Q}^2(\infty, i) = 0$ ;  $\underline{Q}^2(\infty, i^*) = 1$ .

Now consider the series  $\{ \langle s_i, \text{Tr}(s_{i+1}) \rangle : i \geq 0 \}$  (this is a series in  $L'$ , and thus the  $s_i$  are constants, not functional terms). For each  $i \geq 0$ ,  $h(\langle s_i, \text{Tr}(s_{i+1}) \rangle) = \{ \langle s_i(k), [\text{Q}k' : k' > 0] \text{Tr}(s_{i+1}(k')) \rangle : k \geq 0 \}$ . It can be seen that the following distribution of truth values yields a coherent valuation for the set of translations.

$\underline{s}_0(\cdot)$	$\underline{s}_1(\cdot)$	$\underline{s}_2(\cdot)$	...	$\underline{s}_n(\cdot)$	...
#	#	#	...	#	...
0	#	#	...	#	...
...	...	...	...	...	...
0	# ( $i^*+1$ times)	#	...	#	...
0	0	#	...	#	...
...	...	...	...	...	...
0	0	# ( $2i^* + 1$ times)	...	#	...
0	0	0	...	#	...
0	0	0	...	...	...
0	0	0	...	# ( $ni^*+1$ times)	...
0	0	0	...	0	...
...	...	...	...	...	...

Note that each cells represents the value of a sentence  $s$  that 'talks about' the truth values of the sentences that come 'under it' in the column immediately to its right. Let us call these sentences the 'followers' of  $s$ . We do not have a full specification of the semantics of  $\underline{Q}$ , but we know enough to determine that each sentence  $s$  is:

-true if  $s$  has infinitely many followers with the value 0 and exactly  $i^*$  followers with the value 1,

-false if  $s$  has infinitely many followers with the value 0 and at most  $i^*-1$  followers with the value 1.

Furthermore, our Strong Kleene semantics entails that  $s$  has the value # if there are two ways to resolve its indeterminate followers as classical, one of which makes  $s$  true and the other one of which makes  $s$  false.

In each column, it can be checked that each indeterminate sentence  $s$  has infinitely many false followers and at least  $i^*$  indeterminate followers. If  $i^*$  of the indeterminate followers are resolved as true while the others are resolved as false,  $s$  will be true; but if all

indeterminate followers are resolved as false,  $s$  will be false. This disagreement shows that  $s$  should indeed be indeterminate. By contrast, each false sentence  $s$  in the table has infinitely many false followers and at most  $i^*-1$  indeterminate followers. No matter how the latter are resolved,  $s$  should indeed be false, as is desired.

**Case 2.**  $Q^2(\infty, 0)=1$ .

Let  $i^*$  be the least  $i$  such that  $Q^2(\infty, i)=0$ . Thus:

for each  $i \leq i^*-1$ ,  $Q^2(\infty, i)=1$

$Q^2(\infty, i^*)=0$

Consider the series  $\{ \langle s_i, \neg \text{Tr}(s_{i+1}) \rangle : i \geq 0 \}$ . Its translation is the set  $\{ \langle s_i(k), [Qk': k' > k] \neg \text{Tr}(s_{i+1}(k')) \rangle : i \geq 0, k \geq 0 \}$ . It can be checked that the following valuation is coherent:

$\underline{s_0}(\cdot)$	$\underline{s_1}(\cdot)$	$\underline{s_2}(\cdot)$	...	$\underline{s_n}(\cdot)$	...
#	#	#	...	#	...
1	#	#	...	#	...
...	...	...	...	...	...
1	# ( $i^*+1$ times)	#	...	#	...
1	1	#	...	#	...
...	...	...	...	...	...
1	1	# ( $2i^* + 1$ times)	...	#	...
1	1	1	...	#	...
1	1	1	...	...	...
1	1	1	...	# ( $ni^*+1$ times)	...
1	1	1	...	1	...
...	...	...	...	...	...

Consider the first column.  $\underline{s_0}(\underline{0})$  is the formula  $[Qk': k' > 0] \neg \text{Tr}(s_1(k'))$ .  $\underline{s_0}(\underline{0})$  has  $i^*$  indeterminate followers and infinitely many true followers. Depending on how the indeterminates are resolved,  $[Qk': k' > 0] \neg \text{Tr}(s_1(k'))$  may be resolved as false (in case all the indeterminate followers are resolved as false, which means that their negations are resolved as true) or as true (otherwise). Thus  $\underline{s_0}(\underline{0})$  should indeed have the value #. By contrast,  $\underline{s_0}(\underline{1})$  has  $i^*-1$  indeterminate followers and infinitely many true followers. No matter how the indeterminates are resolved,  $[Qk': k' > 1] \neg \text{Tr}(s_1(k'))$  will be true. Thus  $\underline{s_0}(\underline{1})$  should indeed have the value 1. More generally, the sentences with the value # have at least  $i^*$  indeterminate followers and infinitely many true followers, as is required by their semantic content. By contrast, the sentences with the value 1 have fewer than  $i^*$  indeterminate followers, as their content also requires.

## 2) Finite Insensitivity $\Rightarrow$ Uniformity Condition

We show something stronger, namely that for *any admissible interpretation* for  $L^*$  (not just fixed points) and for *any series*  $([s]_k)_{k \geq 0}$  (not just those obtained through our definition of  $[s]_k$ ) the entailment holds.

Let  $I^*$  be any interpretation of  $L^*$ . Let us define  $h_k(s) := [Qk': k' \geq k] [s]_{k'}$ .

Let  $f$  be an assignment function. Let us define:

$0_s := |\{ K' \geq 0 : I^*_{f[K' \rightarrow K']}([s]_{K'}) = 0 \}|$



$$1_s := |\{K' \geq 0: I^*_{\{K' \rightarrow K'\}}([s]_{K'}) [K' \rightarrow K'] = 1\}|$$

$$\#_s := |\{K' \geq 0: I^*_{\{K' \rightarrow K'\}}([s]_{K'}) [K' \rightarrow K'] = \#\}|$$

**Case 1.**  $\#_s$  is finite

**1a)** If  $0_s = 1_s = \infty$ , it immediately follows that for all  $k \geq 0$ ,  $I^*(h_k(s)) = \underline{Q}^2(\infty, \infty)$ , and the Uniformity Condition is satisfied.

**1b)** Otherwise, exactly one of  $0_s, 1_s$  is  $\infty$ . Suppose it is  $0_k$ . It follows from Condition 2 that for all  $k \geq 0$ ,  $I^*(h_k(s)) = \underline{Q}^2(\infty, f)$  for arbitrary finite  $f$ , and hence the Uniformity Condition is satisfied. The reasoning is parallel if  $1_s = \infty$ .

**Case 2.**  $\#_s = \infty$

**2a)** If  $0_s = 1_s = \infty$ , the Uniformity Condition follows as in Case 1a).

**2b)** Suppose  $0_s$  and  $1_s$  are both finite.

-If for some  $i, i' \geq 0$  (including  $\infty$ )  $\underline{Q}^2(\infty, i) \neq \underline{Q}^2(i', \infty)$ , then all the sentences in the series have the value  $\#$ , and the Uniformity Condition follows.

-If for all  $i, i' \geq 0$  (including  $\infty$ )  $\underline{Q}^2(\infty, i) = \underline{Q}^2(i', \infty) = a$ , the Uniformity Condition follows (all sentences have the value  $a$ ).

**2c)** Suppose that exactly one of  $0_s, 1_s$  is  $\infty$  - say, that  $0_k = \infty$  and that  $1_s \neq \infty$  (the opposite case is parallel).

-If  $\underline{Q}^2(\infty, f) = \underline{Q}^2(\infty, \infty) = a$ , the Uniformity Condition follows (all sentences have the value  $a$ )

-If  $\underline{Q}^2(\infty, f) \neq \underline{Q}^2(\infty, \infty)$ , all sentences in the series have the value  $\#$ , and the Uniformity Condition follows again.

#### □ *Characterization of the Generalized Quantifiers that can be used in the Translation*

Clearly, a necessary condition for the translation to be successful is that  $Q$  should satisfy the Uniformity Condition, and thus Finite Insensitivity. In addition, we may observe that even if the Uniformity Condition is satisfied, the translation will fail in case  $\underline{Q}^2(\infty, i) = 1$  or  $\underline{Q}^2(i, \infty) = 0$ . To see this, let us consider the translation of the Truth-Teller  $\langle s_i, \text{Tr}(s_i) \rangle$ . Given our translation scheme, we have that  $h(\langle s_i, \text{Tr}(s_i) \rangle) = \{ \langle s_i(\mathbf{k}), [Qk': k' > \mathbf{k}] \text{Tr}(s_i(k')) \rangle : k \geq 0 \}$ . However, if  $\underline{Q}^2(\infty, i) = 1$ , the translations cannot be assigned the value 0, while if  $\underline{Q}^2(i, \infty) = 0$ , the translations cannot be assigned the value 1. This means that the translation will fail, since the Truth-Teller can equally coherently be assigned the values 0 and 1.

Conversely, if  $\underline{Q}^2(\infty, i) = 0$  and  $\underline{Q}^2(i, \infty) = 1$  and  $Q$  satisfies the Uniformity Condition / Finite Insensitivity, the translation will work as desired. The reasoning is as follows:

1. Property 1 just the Uniformity Condition, which is equivalent to Finite Insensitivity.
2. Property 2 is proved as in Section 3.3, where it was noted that the only properties of  $Q$  that matter is that (Q1)  $Q$  satisfies the Uniformity Condition, and (Q2) when  $F$  contains no bound variables,  $[Qk': k' > i] F$  has the same value as  $F$ . (Q2) is guaranteed to hold by the requirement that  $\underline{Q}^2(\infty, 0) = 0$  and  $\underline{Q}^2(0, \infty) = 1$ , and thus the proof of Section 3.3 can be replicated here. The results we have obtained so far can finally be summarized as an Adequacy Condition on Generalized Quantifiers that are used in the translation:

#### (40) Adequacy Condition

A Generalized Quantifier  $Q$  can be used in the translation scheme defined earlier if and only if:

for all finite  $i \geq 0$ ,  $\underline{Q}^2(\infty, i) = 0$  and  $\underline{Q}^2(i, \infty) = 1$

If we are only interested in the behavior of the quantifier on infinite domains, there are only two cases to consider:

Case 1.  $Q^2(\infty, i) = Q^2(\infty, \infty) = 0$  and  $Q^2(i, \infty) = 1$

This defines a quantifier that behaves like *all but finitely many* on infinite domains.

Case 2.  $Q^2(\infty, i) = 0$  and  $Q^2(i, \infty) = Q^2(\infty, \infty) = 1$

This defines a quantifier that behaves like *infinitely many* on infinite domains.

#### □ *Comparison between the Uniformity Property and the Uniformity Condition*

In Section 2 we asked in which cases a Yablo-series of the form  $S_Q = \{ \langle s(i), [Qk: k > i] \text{Tr}(s_i(k)) \rangle \}$  satisfies the Uniformity Property, which requires that for any fixed point  $I'$ , all the members of  $S_Q$  have the same value according to  $I'$ . As was mentioned earlier,  $S_Q$  is simply the translation of the Truth-Teller obtained when the generalized quantifier is  $Q$ . Comparing the general case to this special case, we see that

-if we only require that the translations of the Truth-Teller have a uniform truth value, *any* quantifier can be used in the translation (in some cases this won't give an adequate translation, for instance if the quantifier is *No*, which incorrectly translates the Truth-Teller as an  $\exists$ -Liar; still, Uniformity is satisfied by this incorrect translation).

-by contrast, when we wish to obtain translations that have a uniform truth value for *any* sentence of the original language, we can only use for  $Q$  quantifiers that satisfy Finite Insensitivity, i.e.  $Q^2(\infty, i) = Q^2(\infty, i')$  and  $Q^2(i, \infty) = Q^2(i', \infty)$  for all finite  $i, i'$ .

We thus see that the only quantifiers that can be used in the translation are those that guarantee that any two sentences in the series will have exactly the same semantic content, in the sense that *even* in interpretations that are not fixed points the Uniformity Condition will be satisfied. This is because our proof of the entailment from Finite Insensitivity to the Uniformity Condition did not rely on the assumption that the interpretation was a fixed point. Since all the translations of a given sentence have exactly the same semantic content, we might be tempted to claim that our translation procedure yields sentences that are self-referential in a broader sense: they refer to other sentences that have exactly the same semantic content as them. It would certainly be interesting to *define* in precise terms the broader notion of self-reference that underlies this criticism. But I believe that this objection would eventually fail for a minor variant of our translation. To see this, suppose that we had opted for a slightly more sophisticated translation, defined as follows:

- (41) a. *Translation:* For each positive integer  $i$ ,  $h_i(F) = [Qk': k' > i] [F]_{k'} \wedge (\text{Tr}(\underline{i} = \underline{i}) \leftrightarrow (i = i))$   
 b. *Denotation:*  $s$  denotes  $F$  according to  $N'$  iff  $s(i)$  denotes  $h_i(F)$  according to  $N^*$ .

It is clear that in every fixed point for  $L^*$  the additional conjunct  $(\text{Tr}(\underline{i} = \underline{i}) \leftrightarrow (i = i))$  is true (since  $\underline{i} = \underline{i}$  is classical), and hence innocuous. But in interpretations that are not fixed points this need not be the case. In this sense this modified translation procedure does *not* give the same semantic content to all the translations of a given sentence of  $L'$ . We could easily add different 'empirical' conjuncts to  $h_i(F)$  (different for various values of  $i$ ) to further ensure that the various translations of a given sentence  $F$  are not equivalent to each other (but simply turn out to have the same value in a given interpretation  $I$  and in a given fixed point  $I^*$  extending  $I$ ).

### 4.3 Another Translation Procedure

The translation scheme we adopted was particularly simple, but it might be interesting to explore alternatives to it. It appears that the same generalized quantifiers that could be used in the procedure outlined above can also be used in a different translation scheme, which is defined as follows:

- (42) a. *Translation*: For each positive integer  $i$ ,  $g_i(F) = \{F\}_i$ , as defined below.  
 b. *Denotation*:  $s$  denotes  $F$  according to  $N'$  iff  $s(i)$  denotes  $h_i(F)$  according to  $N^*$ .

The translation procedure  $\{.\}_{k'}$  is defined recursively as:

- (43) a.  $\{P^n_i(o_1, \dots, o_n)\}_{k'} = P^n_i(o_1, \dots, o_n)$   
 b.  $\{Tr(s)\}_{k'} = [Qk': k' > k] Tr(s(k'))$   
 c.  $\{\neg F\}_{k'} = \neg \{F\}_{k'}$   
 d.  $\{(F_1 \wedge F_2)\}_{k'} = (\{F_1\}_{k'} \wedge \{F_2\}_{k'})$   
 e.  $\{(F_1 \vee F_2)\}_{k'} = (\{F_1\}_{k'} \vee \{F_2\}_{k'})$

(If  $F$  is a formula, we will write  $g_k(F)$  instead of  $\{F\}_k$ ).

#### □ Examples

To get a feel for the translation, let us look at a few simple examples.

- (44) Translation of *Tr*-free formulas  
 $g(\langle c_1, P^0_1 \rangle) = \{\langle c_1(\mathbf{k}), P^0_1 \rangle : k \geq 0\}$   
 This translation is obviously adequate.
- (45) Translation of sentences that 'talk about' *Tr*-free formulas  
 With  $c_1$  as in (44), we consider:  
 $g(\langle c_2, Tr(c_1) \rangle) = \{\langle c_2(\mathbf{k}), [Qk': k' > \mathbf{k}] Tr(c_1(k')) \rangle : k \geq 0\}$   
 This is precisely the result we obtained according to the 'old' translation scheme  $h$ :  
 $g(\langle c_2, Tr(c_1) \rangle) = h(\langle c_2, Tr(c_1) \rangle)$
- (46) Translation of the Liar  
 $g(\langle c_3, \neg Tr(c_3) \rangle) = \{\langle c_3(\mathbf{k}), \neg [Qk': k' > \mathbf{k}] Tr(c_3(k')) \rangle : k \geq 0\}$   
 $\{\langle c_3, \neg Tr(c_3) \rangle\}$  is the simple Liar.  $\{\langle c_3(\mathbf{k}), \neg [Qk': k' > \mathbf{k}] Tr(c_3(k')) \rangle : k \geq 0\}$  is an infinite Liar. From the semantics of  $Q$ , the members of  $\{[Qk': k' > \mathbf{k}] Tr(c_3(k'))\}_{k \geq 0}$  must have a constant value. Therefore all the sentences in the series have the same value, and it follows that this value can only be #.
- (47) Translation of the Truth-Teller  
 $g(\langle c_4, Tr(c_4) \rangle) = \{\langle c_4(\mathbf{k}), [Qk': k' > \mathbf{k}] Tr(c_4(k')) \rangle : k \geq 0\}$   
 This is exactly the same result we obtained in the 'old' translation scheme:  
 $g(\langle c_4, Tr(c_4) \rangle) = h(\langle c_4, Tr(c_4) \rangle)$

#### □ Necessary and Sufficient Conditions for Uniformity

This 'new' translation has essentially the same properties as the 'old' one: the same quantifiers satisfy the Uniformity Condition, and the same quantifiers can be used in the translation.

### 1) Uniformity Condition $\Rightarrow$ Finite Insensitivity

Suppose Condition 2 fails, for instance because there are  $i, i' \geq 0$  such that  $\underline{Q}^2(\infty, i) \neq \underline{Q}^2(\infty, i')$  (the case in which there are  $i, i' \geq 0$  such that  $\underline{Q}^2(i, \infty) \neq \underline{Q}^2(i', \infty)$  is treated in the same way by duality, i.e. by permuting 1 and 0 in the reasoning below).

**Case 1.**  $\underline{Q}^2(\infty, i) = 0$ . Let  $i^*$  be the least  $i$  such that  $\underline{Q}^2(\infty, i) = 1$ .

The non-uniform valuation that was defined earlier for the translation of the series  $\{ \langle s_i, \text{Tr}(s_{i+1}) \rangle : i \geq 0 \}$  will work just as well in the present context, since the 'old' and the 'new' translations agree on sentences of the form  $\text{Tr}(c)$ .

**Case 2.**  $\underline{Q}^2(\infty, 0) = 1$ . Let  $i^*$  be the least  $i$  such that  $\underline{Q}^2(\infty, i) = 0$ . Thus:

for each  $i \leq i^* - 1$ ,  $\underline{Q}^2(\infty, i) = 1$

$\underline{Q}^2(\infty, i^*) = 0$

Consider the series  $\{ \langle s_i, \neg \text{Tr}(s_{i+1}) \rangle : i \geq 0 \}$ . Its translation is the set  $\{ \langle s_i(k), \neg [Qk': k' > k] \text{Tr}(s_{i+1}(k')) \rangle : i \geq 0, k \geq 0 \}$ . It can be checked that the following valuation is coherent:

$\underline{s}_0(\cdot)$	$\underline{s}_1(\cdot)$	$\underline{s}_2(\cdot)$	...	$\underline{s}_n(\cdot)$	...
#	#	#	...	#	...
0	#	#	...	#	...
...	...	...	...	...	...
0	# ( $i^*+1$ times)	#	...	#	...
0	0	#	...	#	...
...	...	...	...	...	...
0	0	# ( $2i^* + 1$ times)		#	...
0	0	0	...	#	...
0	0	0	...	...	...
0	0	0	...	# ( $ni^*+1$ times)	...
0	0	0	...	0	...
...	...	...	...	...	...

Consider the first column.  $\underline{s}_0(\mathbf{0})$  is the formula  $\neg [Qk': k' > \mathbf{0}] \text{Tr}(s_1(k'))$ .  $\underline{s}_0(\mathbf{0})$  has  $i^*$  indeterminate followers and infinitely many false followers. Depending on how the indeterminates are resolved,  $[Qk': k' > \mathbf{0}] \text{Tr}(s_1(k'))$  may be resolved as false (in case all indeterminates are resolved as true, since  $\underline{Q}^2(\infty, i^*) = 0$ ) or as true (otherwise). Thus  $\underline{s}_0(\mathbf{0})$  should indeed have the value #. By contrast,  $\underline{s}_0(\mathbf{1})$  has  $i^*-1$  indeterminate followers and infinitely many false followers. No matter how the indeterminates are resolved,  $[Qk': k' > \mathbf{1}] \text{Tr}(s_1(k'))$  will be true, and hence  $\underline{s}_0(\mathbf{1})$  will be false. Thus  $\underline{s}_0(\mathbf{1})$  should indeed have the value 0. More generally, any sentence with value # have exactly at least  $i^*$  indeterminate followers and infinitely many false followers; this guarantees that it should have the value #. By contrast, any sentences with value 0 has at most  $i^*-1$  indeterminate followers, which guarantees that it should be false.

## 2) Finite Insensitivity $\Rightarrow$ Uniformity Condition

In brief, we give a proof by induction on the construction of formulas of  $L'$  which shows that for each subformula  $G$  of a formula  $g_k(F)$  of  $g(L')$ , for any assignment function  $f$ , for any interpretation (and *a fortiori* for any fixed point)  $I^*$  of  $L^*$ , for all  $k$ ,  $k' \geq 0$ ,  $I_f^*(g_k(G)) = I_f^*(g_{k'}(G))$ .

The key clause in the induction proof is the translation of subformulas of the form  $Tr(c)$ , which get translated as  $[Qk': k' > k]Tr(c(k'))$ . Crucially, because of the choice of the quantifier  $Q$ , the value of  $[Qk': k' > k]Tr(c(k'))$  does not depend on  $k$ , hence the desired result.

The proof is by induction on the construction of formulas of  $L'$ .

-If  $F$  is atomic and does not contain the predicate  $Tr$ , for each  $k \geq 0$ ,  $g_k(F) = F$ , hence the desired result.

-If  $F = Tr(c)$ , for each  $k \geq 0$ ,  $g_k(F) = [Qk': k' > k] Tr(c(k'))$ . Uniformity follows from our proof of *Finite Insensitivity  $\Rightarrow$  Uniformity Condition* for the 'old' translation scheme.

-If  $F = (F_1 \wedge F_2)$  and Uniformity holds of  $F_1$  and  $F_2$ , it holds of  $F$  as well.

-If  $F = \neg F'$  and Uniformity holds of  $F'$ , it holds of  $F$  as well.

It can be checked the Adequacy Condition holds of the new translation without any modifications:  $Q$  can be used in the translation just in case for all finite  $i \geq 0$ ,  $Q^2(\infty, i) = 0$  and  $Q^2(i, \infty) = 1$ .

## 5 Elimination of Self-Reference in a Language With Quantifiers (Sketch<sup>18</sup>)

So far we have assumed that our base language does not contain quantifiers. But we might well want to eliminate Self-Reference from a language that contains quantifiers - say, First-Order Quantifiers. We will briefly sketch a strategy to extend our construction to this case.

We assume that the base language contains sorted variables  $x_i, \dots$ , and  $y_i, \dots$ , ranging over objects and sentences respectively. Similarly there are quantifiers  $\exists$  and  $\forall$ , which can bind object- or sentence-denoting variables, as the case may be. For simplicity we assume that predicate symbols are sorted, in the sense that in any given argument position they take either object-denoting or sentence-denoting terms, but not both (in other words, a given predicate may take both object-denoting and sentence-denoting terms, but not in the same argument positions). The target language is identical to the initial language, except that (i) it contains arithmetic vocabulary (as was the case in our earlier examples), (ii) any sentence-denoting constant  $c$  of the initial language is a function symbol in the target language (it takes a number-denoting term as argument to form a sentence-denoting term), and (iii) it includes a rank predicate  $rk$ , which takes as arguments a number-denoting term and a sentence-denoting term. The rank predicate will be used in the translations to ensure that quantification over sentences never involves self-reference; in essence, a translation with rank  $i$  will only involve quantification over sentences with rank higher than  $i$ .

Let us now turn to a tentative translation procedure. In our earlier examples, the translation involved both a translation in the narrow sense and a specification of the denotation of the sentence-denoting terms of the target language. This will also be the case here, but in addition we will have to modify the interpretation of predicates that take

---

<sup>18</sup> This section will be refined in future drafts.

sentence-denoting terms as arguments. Why? Suppose we consider in the original language a sentence that says that *some beautiful sentence is true* [=  $\exists y_1 (B(y_1) \wedge \text{Tr}(y_1))$ ]. Of course the interpretation will have to specify which sentences of the original language are indeed beautiful. In fact, that very sentence might well be the only beautiful one, which should clearly make the sentence self-referential. But what about the translation? If we do not adapt the interpretation of *beautiful*, we will be faced with two problems: (i) The translations may fail to be 'semantically autonomous'; in other words, we might have to 'look' at properties of sentences that are *not* found in the set of translations to determine what the truth value of some of the translations is. (ii) We may also fail to find an adequate translation. The problem is that if the sentence mentioned above is the *only* beautiful one, and if all we can do in our translation procedure is provide a guarantee about the *translations*, it will *fail* to be the case that one of the translations is also beautiful (since by assumption the only beautiful sentence is the one we mentioned above, which does not have the right form to be a translation of anything).

The solution is to give the translation a bit more leeway. As was the case in our earlier endeavors, we will keep the interpretation of the 'non-linguistic' vocabulary constant, but we will allow ourselves to modify the interpretation of the 'linguistic' vocabulary. In our earlier constructions, the only linguistic vocabulary we had apart from the Truth predicate was a set of sentence-denoting constants, which was duly transformed into a set of a sentence-denoting functional terms in the target language. But in the case at hand we also have at our disposal in the initial language a set of predicates some of which may take sentence-denoting terms (whether constants or variables) as arguments. What shall we do with these? We will modify their interpretations by stipulating that a formula  $F$  lies in the extension of a predicate  $P$  according to  $I'$  just in case its translations  $h_0(F)$ ,  $h_1(F)$ , ... lie in the extension of  $P$  according to  $I^*$ . We will then restrict attention to *admissible interpretations* of  $L^*$ , which will have to be built on this modification of the initial classical interpretation  $I$ , and satisfy the conditions stated below on the interpretation of the rank predicate  $rk$  and of the sentence-denoting terms.

With these conditions in place, we give an example of a translation procedure that would seem to deliver the desired results.

- (48) a. *Translation*: For each positive integer  $i$ ,  $h_i(F) = [Qk': k' > i] [F]_{k'}$ , where  $k$  and  $k'$  are 'fresh' number-denoting variables and where  $[F]_{k'}$  is the result of replacing:
- i) each occurrence of the form  $\text{Tr}(s)$  [where  $c$  is a *constant*] with  $\text{Tr}(s(k'))$
  - ii) each subformula of the form  $\exists y_i \_$  and  $\forall y_i \_$  with  $[\exists y_i: rk(y_i, k')] \_$  and  $[\forall y_i: rk(y_i, k')] \_$  (or an unrestricted quantificational version of these, i.e.  $\exists y_i (rk(y_i, k') \wedge \_)$  and  $\forall y_i (rk(y_i, k') \rightarrow \_)$  respectively).
- b. *Denotation*: If  $s$  is a sentence-denoting constant of the initial language  $L'$ ,  $s$  denotes  $F$  according to  $N'$  iff  $s(i)$  denotes  $h_i(F)$  according to  $N^*$ .
- c. *Rank*:  $\langle F^*, i \rangle$  is in the extension of  $rk$  according to  $N^*$  iff for some natural number  $i$ , for some formula  $F$  of  $L'$ ,  $F^* = h_i(F)$ .
- d. *Interpretation of predicates that take sentence-denoting terms as arguments*:  
If  $P$  is a predicate of  $L$ , if  $F_0 \dots F_k$  are formulas of  $L'$ , and if  $d_0, \dots, d_{i_0}, \dots, d_{i_{k+1}}$  are objects,  $\langle d_0, \dots, d_{i_0}, F_0, d_{i_0+1}, \dots, d_{i_1}, F_1, d_{i_1+1}, \dots, d_{i_2}, \dots, F_k, d_{i_k+1}, \dots, d_{i_{k+1}} \rangle \in I'(P)$  if and only if for all  $i \geq 0$ ,  $\langle d_0, \dots, d_{i_0}, h_i(F_0), d_{i_0+1}, \dots, d_{i_1}, h_i(F_1), d_{i_1+1}, \dots, d_{i_2}, \dots, h_i(F_k), d_{i_k+1}, \dots, d_{i_{k+1}} \rangle \in I^*(P)$
- e. *Admissible interpretations*  
An admissible interpretation of  $L'$  is one that extends  $I$  and is compatible with  $N'$ .  
An admissible interpretation of  $L^*$  is one that extends the modification of  $I$  defined by d. and is compatible with b. and c.

Let us immediately turn to some examples.

(49) 'Some sentence is true'

$$h(<s_1, \exists y_1 \text{Tr}(y_1)>) = \{<s_1(\mathbf{i}), [Qk': k'>\mathbf{i}] [\exists y_1: \text{rk}(y_1, k')] \text{Tr}(y_1)>: i \geq 0\}$$

It is clear that both the original and its translation are true in any admissible interpretations.

(50) 'Some beautiful sentence is true'

$$h(<s_1, \exists y_1 (B(y_1) \wedge \text{Tr}(y_1))>) = \{<s_1(\mathbf{i}), [Qk': k'>\mathbf{i}] [\exists y_1: \text{rk}(y_1, k')] (B(y_1) \wedge \text{Tr}(y_1))>: i \geq 0\}$$

If  $I'(B) = \{s_1\}$ ,  $I^*(B) = \{s_1(\mathbf{i}): i \geq 0\}$ , and both the original and its translations are true.

More generally,  $I^*(B) = \{h_i(s): s \in I'(B) \wedge i \geq 0\}$ , which guarantees that the originally and its translations have the same value (which turns out to be the same in all admissible interpretations).

Can we ascertain that this translation scheme satisfies both the Uniformity Condition and the Isomorphism Condition? A full proof is left for future research, but here is a brief sketch of a positive argument.

1. The Uniformity Condition is satisfied because all the translations of a sentence  $F$  are of the form  $h_i(F) = [Qk': k'>\mathbf{i}] [F]_k$ , where  $Q$  is in effect the quantifier 'all but finitely many'.

2. To study the Isomorphism Condition, we use the same notations as in Section 3.3, with the difference that the admissible fixed points of  $L^*$  are defined as extensions of a modification of the initial interpretation  $I$ , as defined above.

Proof (Sketch): We write that  $J(I', [I^*]) = \text{just in case } I' \text{ and } I^* \text{ are admissible fixed points for } L' \text{ and } L^* \text{ respectively, and } I^{*+}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I^{*+}(\text{Tr})\}$ ,  $I^{*-}(\text{Tr}) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I^{*-}(\text{Tr})\}$ .

1) Let  $I'$  be an admissible fixed point for  $L'$ . We show that there is exactly one equivalence class of admissible fixed points  $[I^*]$  for  $L^*$  satisfying  $J(I', [I^*])$ .

'At most one': given  $N^*$  and  $I$ , the truth value of any member of  $h(L')$  is fixed by the restriction of the interpretation of  $Tr$  to  $h(L')$ . As a result, once  $I^{*+}(\text{Tr}) \cap h(L')$  and  $I^{*-}(\text{Tr}) \cap h(L')$  are fixed, so is the value of each of the members of  $h(L')$ .

'At least one': we show how to construct an admissible fixed point  $I^*$  for  $L^*$  which satisfies  $J(I', [I^*])$ .

(i)  $I^*_0$  is defined by:

$$I^{*+}_0(\text{Tr}) = \{h_k(F): k \geq 0 \text{ and } F \in I^{*+}(\text{Tr})\}$$

$$I^{*-}_0(\text{Tr}) = \{h_k(F): k \geq 0 \text{ and } F \in I^{*-}(\text{Tr})\}$$

$I^*_0$  can be shown to be a fixed point of  $h(L')$  because for each sentence  $F$  of  $L'$ ,

$$h_i(F) \in I^{*+}_0(\text{Tr}) \text{ (resp. } I^{*-}_0(\text{Tr})) \text{ iff } F \in I^{*+}(\text{Tr}) \text{ (resp. } I^{*-}(\text{Tr}))$$

$$\text{iff } I'(F) = 1 \text{ (resp. } = 0) \text{ [because } I' \text{ is a fixed point]}$$

$$\text{iff for each } k \geq 0, I^*_0([F]_k) = 1 \text{ (resp. } = 0), \text{ where } [F]_k \text{ is obtained}$$

from  $F$  by replacing (a) each occurrence of the form  $Tr(c)$  with  $Tr(c(\mathbf{k}))$ , and (b) each formula of the form  $\exists y_i \text{ ---}$  and  $\forall y_i \text{ ---}$  with  $[\exists y_i: \text{rk}(y_i, \mathbf{k})] \text{ ---}$  and  $[\forall y_i: \text{rk}(y_i, \mathbf{k})] \text{ ---}$

The latter equivalence is proven by induction on the construction of formulas of  $L'$ . Specifically, we show that for all  $k \geq 0$ ,  $I'_s(F) = I^*_{0 \ s^*}([F]_k)$ , where for all  $i$   $s^*(y_i) = h_k(s(y_i))$ .

The crucial induction step is that for which  $F = \exists y_i G$ . We then have that  $[F]_k = [\exists y_i: rk(y_i, k)] [G]_k$ , and

$I'_s(F)=1$	iff	for some sentence $d$ of $L'$ , $I'_{0 \ s[y_i \rightarrow d]}(G)=1$
	iff	for some sentence $d$ of $L'$ , $I^*_{0 \ (s[y_i \rightarrow d])^*} [G]_k=1$ (induction hypothesis)
	iff	for some sentence $d$ of $L'$ , $I^*_{0 \ s^*[y_i \rightarrow h_k(d)]}([G]_k)=1$
	iff	$I^*_{0 \ s^*}([\exists y_i: rk(y_i, k)] [G]_k)=1$
	iff	$I^*_{0 \ s^*}([F]_k)=1$
$I'_s(F)=0$	iff	for every sentence $d$ of $L'$ , $I'_{s[y_i \rightarrow d]}(G)=0$
	iff	for every sentence $d$ of $L'$ , $I^*_{0 \ (s[y_i \rightarrow d])^*} [G]_k=0$ (induction hypothesis)
	iff	for every sentence $d$ of $L'$ , $I^*_{0 \ s^*[y_i \rightarrow h_k(d)]}([G]_k)=0$
	iff	$I^*_{0 \ s^*}([\exists y_i: rk(y_i, k)] [G]_k)=0$
	iff	$I^*_{0 \ s^*}([F]_k)=0$

In the special case of sentences, we have that  $I'(F) = I^*_0([F]_k)$ , as desired. We can finish the proof:

$h_i(F) \in I^{*+}_0(Tr)$  (resp.  $I^{*-}(Tr)$ ) iff for each  $k \geq 0$ ,  $I^*_0([F]_k)=1$  (resp.  $=0$ ),  
iff  $I^*_0([Qk': k' \geq i][F]_{k'})=1$  (resp.  $=0$ ) [this follows because (i)  $Q$  satisfies property (Q2): when  $F$  contains no bound variables,  $[Qk': k' > i] F$  has the same value as  $F$ , and (ii) for all  $k', k'' \geq 0$ ,  $I^*_0([F]_{k'}) = I^*_0([F]_{k''})$ ]  
iff  $I^*_0(h_i(F))=1$  (resp.  $=0$ )

(ii) By the Extension Lemma, this local fixed point can be extended to a global fixed point.

2) Let  $I^*$  be an admissible fixed point for  $L^*$ . We show that there is exactly one admissible fixed point  $I'$  for  $L'$  satisfying  $J(I', [I^*])$ .

Given the Uniformity Property, for all  $k, k' \geq 0$ ,  $I^*(h_k(s)) = I^*(h_{k'}(s))$ . Given  $I$  and  $N'$ , we can thus define an interpretation  $I'$  by  $I'^+(Tr) = \{s: \text{for some } k \geq 0, h_k(s) \in I^{*+}\}$  and  $I'^-(Tr) = \{s: \text{for some } k \geq 0, h_k(s) \in I^{*-}\}$ . It is then immediate that  $I'^+(Tr) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^+(Tr)\}$ ,  $I'^-(Tr) \cap h(L') = \{h_k(s): k \geq 0 \text{ and } s \in I'^-(Tr)\}$ . All that remains to be shown is that  $I'$  is a fixed point.

2a. From the Uniformity Condition and the definition of an admissible interpretation for  $L^*$ , it follows that for any formula  $F$  of  $L'$ ,  $I^*(h_i(F)) = I^*(h_i(F)/_{0/k'})$ , where  $h_i(F)/_{0/k'}$  is obtained from  $h_i(F)$  by replacing:

every formula of the form  $Tr(c(k'))$  with  $Tr(c(\mathbf{0}))$

every formula of the form  $[\exists y_i: rk(y_i, k')] \_$  with  $[\exists y_i: rk(y_i, \mathbf{0})] \_$

Therefore for all  $i \geq 0$ ,  $I^*(h_i(F)) = I^*([Qk': k' > k] / [F]_{k'} /_{0/k'})$

$= I^*([F]_0)$  [because quantification is vacuous, and  $Q$  satisfies Property (Q2)]

$= I'(F)$  [this could be proven in a proof by induction].

2b. We can now reason as follows:

$F \in I'^+(Tr)$  (resp.  $I'^-(Tr)$ ) iff for each  $i \geq 0$ ,  $h_i(F) \in I^{*+}(Tr)$  (resp.  $I^{*-}(Tr)$ )



fixed point]                      iff      for each  $i \geq 0$ ,  $I^*(h_i(F))=1$  (resp.  $=0$ )    [because  $I^*$  is a  
    iff       $I'(F)=1$  (resp.  $=0$ )                      [from 2a].

Taken together 1) and 2) show that  $J$  is a 1-1, onto function from the admissible fixed points of  $L'$  to the equivalence classes of admissible fixed points of  $L^*$ . We henceforth write  $[I^*]=j(I')$  for  $J(I', [I^*])$ . It is immediate from the meaning of  $J$  that  $I'_1 \leq I'_2$  iff  $j(I'_1) \leq_{h(L^*)} j(I'_2)$ .

## 6 Perspectives

The main results of this investigation can be summarized as follows. (1) Yablo-Series satisfy the Uniformity Property no matter which Generalized Quantifier is used. (2) Self-reference can be eliminated from *every* sentence of a given language by generalizing the procedure at work in Yablo-Series. The latter come out as the translations of the Truth-Teller for various values of the Generalized Quantifier  $Q$  used in the translation. However the procedure can work in the general case only when  $Q$  can guarantee that in any interpretation (not just fixed points) all the translations of a given sentence  $F$  share the same truth value. Over infinite domains,  $Q$  should be one of two quantifiers: *all but finitely many* or *infinitely many*.

We hope that the present attempt will have convinced the reader that it might be fruitful to study *general* results about the elimination of self-reference rather than to simply attempt -piecemeal, so to speak- to replicate Yablo's results on a variety of semantic phenomena.

## Appendix. A Sufficient Condition of Non-Self-Reference

### □ Basic ideas

It has sometimes been argued that Yablo's Paradox is in fact 'covertly' referential. A large part of the difficulty is that there is no accepted criterion of what it means for a quantificational sentence to 'refer' to anything. We will not solve the problem in full generality, but we will develop a plausible criterion of *non*-self-reference which is strong enough to show that Yablo's Paradox does not involve self-reference.

### • Reference of Quantified Sentences

To start with the simplest idea, we could determine that an unrestricted quantifier 'refers' to all objects in its range. But this would not be a particularly helpful analysis, since it would entail that any statement that contains any quantifier 'refers' to every object in the domain.

But we can develop a more discriminating analysis by taking a hint from natural language, which includes the device of *restricted quantification*. *Every student is sick* is understood to 'refer' in a broad sense to the students, but certainly not to anyone else. And similarly for *Some student is sick*. There is a logical reason for this intuition: natural language quantifiers have two properties, called 'extension' and 'conservativity', which ensure that in evaluating any sentence of the form *Q student is sick* we can safely disregard those individuals that are neither sick nor students (this is a consequence of 'extension'), and that furthermore among those we may restrict attention to the students, disregarding the sick people who are not students (this is guaranteed by 'conservativity').

Let us say that a sentence *s* is in *Restricted Quantifier Notation* if:

- (i) every quantifier *Q* appears with a restrictor, and has thus the form  $[Qx_i: F]F'$ ; in addition no other occurrence of a quantifier introduces the variable  $x_i$  in *s*.
- (ii) no function symbol (including constants) appears except in subformulas of the form  $x=c$  that occur in a restrictor, and
- (iii) all restrictors are classical, i.e. they do not contain the predicate *Tr*.

To give a simple example,  $\exists x(P(x) \wedge Q(x))$  is equivalent to several formulas that are in restricted quantifier notation:  $[\exists x: P(x)] Q(x)$ ,  $[\exists x: Q(x)] P(x)$ ,  $[\exists x: x=x] P(x)$  (the list is not exhaustive). Similarly,  $\forall x(P(x) \Rightarrow Q(x))$  is equivalent to the formula in restricted quantifier notation  $[\forall x: P(x)] Q(x)$ . And by the same token  $P(c)$  is equivalent to  $[\exists x: x=c] P(x)$  or to  $[\forall x: x=c] P(x)$ , which are both in restricted quantifier notation.

Obviously the restricted quantifier notation is particularly well-adapted for 'binary generalized quantifiers' that are found in natural language, such as *most* in *Most students passed the exam*, or for that matter *infinitely many* as in *Infinitely many numbers are prime*. As we observed earlier, the  $\exists\forall$ -Liar' can be rewritten as the set  $\{ \langle s(i), [Qk': k'>i] \neg \text{True}(s(sk')) \rangle : i \geq 0 \}$ , where *Q* is a binary generalized quantifier meaning 'all but finitely many'. We will make use of this notation again later in the paper.

- *Transitive Reference*

The restricted quantifier notation will not solve all our problems, however. Even in the simplest cases, which involve constants and no quantifiers, we need some notion of 'referential closure'. Consider for instance the following set:  $\{ \langle c_1, \text{Tr}(c_2) \rangle, \langle c_2, \text{Tr}(c_1) \rangle \}$ . Intuitively,  $c_1$  doesn't 'directly' refer to itself, but it still refers to itself 'indirectly', because it refers to  $c_2$ , which in turn refers to  $c_1$ . In order to develop an adequate notion of reference, then, we will need to look at the transitive closure of the notion *refers directly to* in order to obtain the generalized (indirect) notion which is of interest for paradoxes. Clearly, it would be of no particular interest to show that Yablo's paradox does not involve any 'direct' self-reference, since much simpler constructions that only involve indirect self-reference are known to generate Liar-like phenomena. This is in particular the case of the Circular Liar:  $\{ \langle c_1, \text{Tr}(c_2) \rangle, \langle c_2, \neg \text{Tr}(c_1) \rangle \}$ .

□ *Definitions and Examples*

- *Primary Reference*

We start with the definition of direct reference for a formula in which all quantifiers appear in restricted quantifier notation. The definitions we give are supposed to be very liberal, in the sense that an element  $d$  will be 'referred to' if there is *some chance* that  $d$  might affect the value of the formula. Since we are after a criterion of *non*-reference (specifically: of non-self-reference), it is of course judicious to make the criterion reference as liberal as we possibly can (so that the criterion of non-self-reference will be as stringent as possible).

The definitions we give hold for formulas in Restricted Quantifier Notation in a trivalent logic satisfying Reasonableness (defined formally in (14)), which can be seen as a trivalent generalization of the Tree of Numbers. We start with the definition of a Reference Set of a formula  $F$  relative to a set  $S$  of assignment function. Intuitively it can be seen as the set of all assignment functions which will be 'accessed' in the evaluation of  $F$  relative to a member of  $F$ . We write  $e$  for the null assignment function. As usual,  $s[x_{i_k} \rightarrow d]$  is the assignment function which is identical to  $s$  with the possible exception that it assigns  $d$  to  $x_{i_k}$  (note in particular that  $e[x_{i_k} \rightarrow d]$  is simply the partial function that assigns  $d$  to  $x_{i_k}$ , which we also write as  $[x_{i_k} \rightarrow d]$ ).

(51) Reference Set of  $F$  relative to  $S$

- If  $F$  is atomic,  $R(F, S) = S$
- If  $F = \neg F'$ ,  $R(F, S) = R(F', S)$
- If  $F = (F' \wedge F'')$  or  $F = (F' \vee F'')$ ,  $R(F, S) = R(F', S) \cup R(F'', S)$
- If  $F = [Qx_i: F'] F''$ ,  $R(F, S) = R(F'', \{s[x_i \rightarrow d]: s \in S \wedge d \in D \wedge I_{[x_i \rightarrow d]}(F') = 1\})$

If  $F$  is a closed formula, we define its Primary Reference  $PR(F)$  to be:

$$(52) \quad PR(F) = \bigcup_{s \in R(F, \{e\})} \text{rng}(s)$$

In other words,  $PR(F)$  is the set of all objects that are in the range of assignment functions that are in the Reference Set of  $F$  relative to  $\{e\}$ .

### Examples

- (53)  $F := [\forall x_1: P(x_1)]Q(x_1)$   
 $R(F, \{e\}) = R(Q(x_1), \{[x_1 \rightarrow d]: d \in D \wedge I_{[x_1 \rightarrow d]}(P(x_1)) = 1\})$   
 $= R(Q(x_1), \{[x_1 \rightarrow d]: d \in D \wedge d \in I^+(P)\})$   
 $= \{[x_1 \rightarrow d]: d \in D \wedge d \in I^+(P)\}$   
 $PR(F) = \bigcup_{s \in R(F, \{e\})} \text{rng}(s) = \{d \in D: d \in I^+(P)\}$
- (54)  $F' := [\forall x_1: P(x_1)][\forall x_2: P'(x_2)]Q(x_1)$   
 $R(F, \{e\}) = R([\forall x_2: P'(x_2)]Q(x_1), \{[x_1 \rightarrow d]: d \in D \wedge I_{[x_1 \rightarrow d]}(P(x_1)) = 1\})$   
 $= R([\forall x_2: P'(x_2)]Q(x_1), \{[x_1 \rightarrow d]: d \in D \wedge d \in I^+(P)\})$   
 $= R(Q(x_1), \{[x_1 \rightarrow d][x_2 \rightarrow d']: d \in D \wedge d' \in D \wedge d \in I^+(P) \wedge d' \in I^+(P')\})$   
 $= \{[x_1 \rightarrow d][x_2 \rightarrow d']: d \in D \wedge d' \in D \wedge d \in I^+(P) \wedge d' \in I^+(P')\}$   
 $PR(F) = \bigcup_{s \in R(F, \{e\})} \text{rng}(s) = \{d \in D: d \in I^+(P)\}$
- $= R(Q(x_1), \{[x_1 \rightarrow d]: d \in D \wedge d \in I^+(P)\})$   
 $= \{[x_1 \rightarrow d]: d \in D \wedge d \in I^+(P)\}$   
 $PR(F') = \bigcup_{s \in R(F, \{e\})} \text{rng}(s) = \{d \in D: d \in I^+(P) \vee d \in I^+(P')\}$

- *Transitive Reference*

We are now in a position to define the Reference Set of Level  $k$  and then simply the Reference Set of a formula. These notions are intended to apply to a fragment whose sentences are all in Restricted Quantifier Notation.

- (55) We define the Reference Set of Level  $k$  and the Reference Set of  $F$  as follows:

a. Reference Set of Level  $k$

$$\text{Ref}_1(F) = PR(F)$$

$$\text{Ref}_{k+1}(F) = \bigcup \{PR(d): d \in \text{Ref}_k(F) \wedge d \text{ is a closed formula}\}$$

b. Reference Set

$$\text{Ref}(F) = \bigcup_{k \geq 1} \text{Ref}_k(F)$$

### Examples

- (56)  $F := [\forall x_1: P(x_1)]Q(x_1)$   
 Suppose that  $I^+(P) = \{d, F\}$ , where  $d$  is not a sentence.  
 As shown in (53),  $PR(F) = \{d \in D: d \in I^+(P)\}$ . Hence  
 $\text{Ref}_1(F) = PR(F) = \{d \in D: d \in I^+(P)\} = \{d, F\}$   
 $\text{Ref}_2(F) = \bigcup \{PR(d): d \in \text{Ref}_1(F) \wedge d \text{ is a closed formula}\}$   
 $= \bigcup \{PR(F)\}$   
 $= \{d, F\}$   
 $\text{Ref}(F) = \bigcup_{k \geq 1} \text{Ref}_k(F) = \{d, F\}$   
 $F$  is clearly self-referential, since it is contained in its own Reference Set.

$$(57) F' := [\forall x_1: P'(x_1)]Q(x_1)$$

Suppose that  $I(P) = \{d', F\}$ , where  $d'$  is not a formula and  $F$  is as in (56)

Replacing  $P$  with  $P'$  in (56), we obtain:  $PR(F') = \{d', F\}$

$Ref_1(F') = PR(F') = \{d', F\}$

$Ref_2(F') = \cup \{PR(d): d \in Ref_1(F') \wedge d \text{ is a closed formula}\}$

$= \cup \{PR(F)\}$

$= \{d, F\}$

It is clear that:

$Ref(F) = \{d', d, F\}$

$F'$  is not self-referential, since it is not contained in its own Reference Set. However it involves self-reference, in the sense that it refers to a formula, namely  $F$ , which is itself self-referential.

We can finally state (i) a condition that guarantees that a formula is *not self-referential*, and -more interestingly (ii) a condition that guarantees that a formula does *not involve self-reference*, in the sense that it does not refer to any formula that is self-referential. (Note that the second condition entails the first. If  $F$  is self-referential, it refers to  $F$ , and hence it refers to a formula that is self-referential. Failure of the first condition entails failure of the second, and so by contraposition the second condition entails the first).

(58) Sufficient Conditions of Non-Self-Reference

a.  $F$  is guaranteed not to be self-referential if  $F \notin Ref(F)$

b.  $F$  is guaranteed not to involve self-reference if  $(\forall d (d \text{ is a sentence} \wedge d \in Ref(F) \rightarrow d \notin Ref(d)))$

Note: If  $F$  refers (directly or indirectly) to a formula  $F'$  which is *not* in Restricted Quantification Form, the criterion for determining what  $F'$  refers to will not work (since our criterion is not even defined in that case). Still, by inspecting  $Ref(F)$  we will be in a position to tell that our criterion of reference is not reliable, since we will find a formula, namely  $F'$ , which is in  $Ref(F)$  and which does not have the right syntactic form. This limitation will not have any consequences in what follows, since all the formulas we will be considering will only include in their Reference Set formulas that are in Restricted Quantifier Notation.

□ Application to Yablo's Paradox

We now apply the notions we defined to Yablo's paradox. It will be seen that each of Yablo's sentences satisfies our Sufficient Conditions of Non-Self-Reference, and does not involve self-reference. For simplicity we only consider the Universal Liar, which we put in Restricted Quantifier Notation. The language in which it is stated is the target language of our translation schemes; it is sorted, with variables  $k_i$  ranging over non-negative integers and variables  $x_i$  ranging over objects. We assume that  $S_V$  is evaluated with respect to an interpretation  $I$  which is compatible with the naming relation that  $S_V$  defines.

$$(59) S_V := \{ \langle s(\mathbf{k}), [\forall k_1: k_1 > \mathbf{k}] [\forall x_2: x_2 = s(k_1)] \neg Tr(x_2) \rangle : i \geq 0 \}$$

For each  $k \geq 0$ ,

$R(s(\mathbf{k}), \{e\}) = \{ [k_1 \rightarrow n] [x_2 \rightarrow d] : n \in \mathbb{N} \wedge d \in D \wedge n > k \wedge d = I(s)(n) \}$

$PR(s(\mathbf{k})) = \{ n : n > k \} \cup \{ s(\mathbf{n}) : n > k \}$

We note that for each  $k \geq 0$ ,  $\text{PR}(s(\mathbf{k}+1)) \subseteq \text{PR}(s(\mathbf{k}))$   
 $\text{Ref}_1(s(\mathbf{k})) = \text{PR}(s(\mathbf{k}))$   
 $\text{Ref}_{i+1}(s(\mathbf{k})) = \bigcup \{ \text{PR}(d) : d \in \text{Ref}_i(s(\mathbf{k})) \wedge d \text{ is a closed formula} \}$   
 DETAILS TO BE ADDED HERE.  
 In the end  $\text{Ref}(s(\mathbf{k})) = \{n : n > k\} \cup \{s(\mathbf{n}) : n > k\}$   
 We see that for no  $k' \geq k$  is  $s(\mathbf{k}')$  a member of  $\text{Ref}(s(\mathbf{k}))$ .  
 Thus  $s(\mathbf{k})$  does not involve self-reference.

The same methods can be used to show that none of the translations we have used in this paper involves self-reference either.

## Partial References

- Cook, R. 2004. Patterns of Paradox, *Journal of Symbolic Logic*, 69, 3, 767-774  
 Egré, P. 2005. The Knower Paradox in the Light of Provability Interpretations of Modal Logic, *Journal of Logic, Language and Information*, **14**: 13–48  
 Keenan, E. 1996. The Semantics of Determiners, in Lappin, S. (ed.) *The Handbook of Contemporary Semantic Theory*  
 Kripke, S. 1975. Outline of a Theory of Truth, *Journal of Philosophy* 72: 690-716  
 van Benthem, J. 1986. *Essays in Logical Semantics*, Reidel, Dordrecht  
 Schlenker, P. 2005. How to Eliminate Self-Reference: A Précis. Ms., UCLA & IJN.  
 Simmons, K. 1993. *Universality and the Liar*, Cambridge University Press.  
 Tarski, A. 1944. The Semantic Conception of Truth and the Foundations of Semantics, *Philosophy and Phenomenological Research* 4, 3: 341-376  
 Visser, A. 1989. Semantics and the Liar Paradox, in *Handbook of Philosophical Logic* (4)  
 Yablo, S. 1993.  
 Yablo, S. 2004. Circularity and Paradox, in *Self-Reference*, CSLI