

Draft of November 2018—please do not cite without permission.

**The Arithmetization of Syntax and
the New Paradoxes of Self-Reference¹**

T. Parent

In Parent (ms.), I argued that if self-reference is unconstrained in a language L , then L is not classical. Specifically, it will include well-formed sentences that are both true and false. That is so, even if L is “semantically open” (using Tarski’s 1944 idiom), that is, even if L is free of semantic terms like ‘true’ and ‘denotes’ defined on expressions of L .

In conversation, Tim Button worried that this may have dire consequences for Peano Arithmetic (PA), or any other axiomatizations extending Robinson Arithmetic (Q). After all, the method of Gödel numbering allows for something functionally like self-reference. So if unrestricted self-reference enables wff that are both true and false, as the new paradoxes suggest, then Gödel numbering could conceivably be used to demonstrate that the language of arithmetic is non-classical—and in particular, that there are truths in the language which are also false.

I do not believe that the new paradoxes show any such thing. They indeed show that something is awry, but they need not show that the problem lies in the object language. Rather, the problem may well lie in the metalanguage; in particular, it may be that the arithmetization of the object language within the metalanguage enables the paradoxes. If this is correct, then it is not arithmetic itself which is unsound, but rather, any metatheory where Gödel numbering has no restrictions on its use, including Gödel’s (1931) metatheory.

¹ I owe thanks to friends at Notre Dame for excellent feedback on this draft: Rachael Alvir, Tim Bays, Dan Turetsky, and especially Matteo Bianchetti.

1. Preparatory work

The task shall be to present one of the paradoxes from the earlier paper (the “Laputan” paradox), but here, it is to be formulated exclusively in the language of arithmetic; call it “ \mathcal{L} .” For this purpose, we will forego the use of self-referential terms, and instead make use of Gödel numbering. (Since this paper is somewhat more technical, readers may be better served by reading the earlier paper first.)

Expressions of language \mathcal{L} shall be understood as follows. First, its terms are defined recursively as follows:

1. $\underline{0}$ is a term.
2. If τ is a term, then so is $\lceil \tau' \rceil$. (Terms introduced by clauses 1 and 2 are the numerals.)
3. If \underline{n} is a numeral, then $\lceil v_n \rceil$ is a term. (Terms introduced by clause 3 are the variables.)
4. If τ and σ are terms, then so are $\lceil \tau + \sigma \rceil$, and $\lceil \tau \cdot \sigma \rceil$.
5. For any $k > 0$, if $\langle \tau_1, \dots, \tau_k \rangle$ is any sequence of terms, and \underline{n} and \underline{m} are numerals, then $\lceil f_m^n(\tau_1, \dots, \tau_k) \rceil$ is a term.
6. Nothing else is a term.

Assume that terms formed by clauses 1, 2, and 4 have their standard interpretation, where $\underline{0}$ is the numeral denoting 0, ‘ $'$ ’ expresses the successor-function, ‘ $+$ ’ expresses the addition-function, and ‘ \cdot ’ expresses the multiplication-function. Variables from clause 3 have their denotation relative to a variable-assignment, i.e. a selection of a sequence such that the n th member is to be the denotation of the n th variable. Finally, a saturated function-symbol $\lceil f_m^n(\tau_1, \dots, \tau_k) \rceil$ denotes m = the output of the expressed function, given the sequenced denotations of $\langle \tau_1, \dots, \tau_k \rangle$ as input.

The well-formed formulae (wffs) can then be defined by recursion thusly:

1. If $\langle \tau_1, \tau_2 \rangle$ is any pair whose members are terms, then $\lceil =(\tau_1, \tau_2) \rceil$ is a wff.

2. For any $k > 0$, if $\langle \tau_1, \dots, \tau_k \rangle$ is any sequence of terms, and \underline{n} and \underline{m} are numerals, then $\lceil F_m^n(\tau_1, \dots, \tau_k) \rceil$ is a wff.
3. If Φ is a wff, then so is $\sim\Phi$.
4. If Φ and Ψ are wffs, then so is $\lceil \Phi \wedge \Psi \rceil$.
5. If \underline{n} is a numeral, and Φ is a wff with exactly $\lceil v_n \rceil$ free, then $\lceil \forall v_n \Phi \rceil$ is a wff. (A variable $\lceil v_n \rceil$ is free in Φ iff Φ has $\lceil v_n \rceil$ as a part but not $\lceil \forall v_n \rceil$.)
6. Nothing else is a wff.

Assume that '=', ' \sim ', ' \wedge ', and quantifier-expressions $\lceil \forall v_n \rceil$ have their standard interpretation, and that an n -place predicate has a set of n -tuples as its extension. Per usual, a sentence is a wff with no free variables; let us also stipulate that a *designator* is any term which has no variable as a (proper or improper) part. Also, n.b., Arabic numerals will be used below for subscripts and superscripts, so to reduce clutter; although '1' as a subscript or superscript will often just be omitted.

Let us now define Gödel numbers for expressions of the language. The coding scheme for the basic symbols of \mathcal{L} is as follows (where $n > 0$ and $m > 0$):

Symbol:	0	'	()	,	\sim	\wedge	\forall	=	+	·	$\lceil v_n \rceil$	$\lceil F_m^n \rceil$	$\lceil f_m^n \rceil$
Code:	3	5	7	9	11	13	15	17	19	21	23	$2 \cdot 5^n$	$2^2 \cdot 3^n \cdot 5^m$	$2^3 \cdot 3^n \cdot 5^m$

The compound expressions have a unique Gödel number determined in the usual way, exploiting Gaussian prime decomposition: The codes of the n basic parts are first assigned (in order) as exponents to the first n members of the sequence of primes $\langle 2, 3, 5, 7, \dots \rangle$. The compound expression's Gödel number is then the product of the exponentiated primes. Thus, the sentence ' $\lceil (0, 0) \rceil$ ' will have a Gödel number equal to $2^{19} \cdot 3^7 \cdot 5^3 \cdot 7^{11} \cdot 11^3 \cdot 13^9$. (Further measures are needed to code proofs as sequences of sentences, but this is unnecessary for the remarks below.)

Another preliminary is necessary. Given a sentence of the form $\ulcorner \forall v_3 \forall v_2 \forall v \Phi(v, v_2, v_3) \urcorner$, let C be a “coordinated” substitution instance (or for short, a “CSI”) of the sentence iff:

$$C = \ulcorner \forall v \Phi(v, n, \delta) \urcorner$$

In this, δ is a designator with Gödel number n whose numeral \underline{n} has replaced ‘ v_2 ’. (The numeral in question we shall call the *Gödel-numeral* of δ .)

Thus, a CSI of a sentence $\ulcorner \forall v_3 \forall v_2 \forall v \Phi(v, v_2, v_3) \urcorner$ will replace the second variable with the Gödel-numeral for δ , and replace the third variable with δ itself. Since it is decidable whether \underline{n} is the Gödel-numeral for δ , and since this type of sentence is otherwise defined by its form, it is decidable whether a sentence is of this type. (This can be justified via Church’s thesis).

2. Paradox with Gödel numbers

Consider now the following sentence of \mathcal{L} :²

$$(\dagger) \quad \forall v_3 \forall v_2 \forall v (=f(v), v_2) \supset (F(v_3) \equiv (=f(v), 0''') \wedge =(0, v_3)))$$

We have yet to define the predicate ‘ $F(v_3)$ ’ and the function-symbol ‘ $f(v)$ ’—regardless, we can still decide whether a sentence is a CSI of (\dagger) . For instance, the following sentence is such a CSI:

$$(1) \quad \forall v (=f(v), 0''') \supset (F(0) \equiv (=f(v), 0''') \wedge =(0, 0)))$$

The reason is that ‘ v_2 ’ in (\dagger) is replaced with the Gödel-numeral for the designator that replaces ‘ v_3 ’. In particular, ‘ v_2 ’ is replaced by $\underline{0''''}$, and ‘ v_3 ’ is replaced by $\underline{0}$.

² For concision’s sake, ‘ \supset ’ and ‘ \equiv ’ are here used as if they are part of the object language, even though they were not mentioned in specifying \mathcal{L} . But if preferred, one could revise formulae of the form $\ulcorner \Phi \supset \Psi \urcorner$ to $\ulcorner \sim(\Phi \wedge \sim\Psi) \urcorner$, and revise formulae of the form $\ulcorner \Phi \equiv \Psi \urcorner$ to $\ulcorner \sim(\Phi \wedge \sim\Psi) \wedge \sim(\Psi \wedge \sim\Phi) \urcorner$.

Let us next define ' $f(v)$ ' as denoting the function which maps the CSI of (\dagger) coded by v onto the Gödel code of its final designator. "The final designator" is the designator replacing ' v_3 ', i.e., the final variable in (\dagger).³ Thus, where $g(\delta)$ is the Gödel number of a term τ (with $g^{-1}(v)$ as the inverse function), the suggestion is to define ' $f(v)$ ' as follows:

$$f(v) = \begin{array}{ll} n & \text{if } g^{-1}(v) = \text{the CSI of } (\dagger) \text{ whose final designator is } \delta, \text{ where } g(\delta) = n. \\ \uparrow & \text{otherwise.} \end{array}$$

Again, one can roughly think of $f(v)$ as mapping the CSI of (\dagger) coded by v onto the Gödel code for the CSI's final designator. (If v does not code such a CSI, then the function is undefined, although the undefined cases will have no bearing on the paradox below.)

Suppose now that the predicate 'F' is defined by (\dagger). Then, (1) is an instance of the definition which makes explicit a condition on which 0 is F. Basically, it says that if you take the CSI of (\dagger) whose final designator is coded by 3, then 'F' is satisfied by 0 iff the final designator of that CSI is indeed coded by 3 and $0 = 0$. Notice, then, that this last clause is true. Therefore, the definition indicates that 0 is F. (N.B., (1) is itself the CSI of (\dagger) which has its final designator coded by 3. So (1) defines 0 being F with reference to features of (1) itself.)

But the paradox is that we can also show that 0 is not F. Suppose here that $f_2(v) = 0$, for any v (i.e., it is the constantly zero function). And observe that the designator ' $f_2(0')$ ' (i.e., the function-symbol with $0'$ as the instantiating constant) is coded by $2^{240} \cdot 3^7 \cdot 5^3 \cdot 7^5 \cdot 11^9$. For short, let us say that this is a number h with numeral \underline{h} . Then, another CSI of (\dagger) would be the following:

$$(2) \quad \forall v (=(f(v), h) \supset (F(f_2(0')) \equiv (=(f(v), 0''') \wedge =(0, f_2(0'))))))$$

³ Since there are two occurrences of ' v_3 ' in (\dagger), "the final designator" could be what replaces the second occurrence of ' v_3 ', or it could be the expression-type replacing both occurrences. Either precisification is fine for our purposes.

This counts as a CSI of (\dagger) given that ' v_2 ' in (\dagger) is replaced by the Gödel-numeral for the designator that replaces ' v_3 '.

Consider, then, (2) also provides a condition on which 0 is F, given that $f_2(0) = 0$. It indicates that, where v codes a CSI of (\dagger) whose final designator is coded by h , 'F' is satisfied by 0 iff that final designator is coded by 3 and $0 = f_2(0)$. Now in the antecedent of (2), the formula ' $(=f(v), h)$ ' is satisfied when v is the code of (2) itself. After all, (2) is itself the CSI of (\dagger) whose final designator has code h , given that its final designator is ' $f_2(0)$ '. Yet, contra the final clause of (2), it is false that its final designator is also coded by 3. Thus, (2) reveals that 0 is not F.

And so, the predicate as defined by (\dagger) determines that 'F(0)' is both true and false.

3. Objections and replies

Objection 1: The first objection is that my symbol ' $f(v)$ ' is ill-defined, for its definition refers to instances of (\dagger) , and such instances contain the very symbol being defined. Such circularity is thought to be dubious.

It is correct that the function-symbol defined is with reference to the sentence (\dagger) , and (\dagger) indeed has the function-symbol as a part. However, when (\dagger) is first identified, the function-symbol is thus far treated as uninterpreted. So it is not as if the function-symbol had to be interpreted before one could interpret the function-symbol. If that were so, that may be a dubious kind of circularity. Rather, the symbol just needs to exist, in order to define the symbol. (This is hardly unusual—one always needs the symbol to exist before one can define the symbol.)

It is a bit odd, however, that the function-symbol is defined with reference to a string that includes the function-symbol itself. Yet that just is a type of self-reference.

Accordingly, if one dislikes the self-reference in how the symbol defined, then this is already to accept the lesson of the paradox. The lesson is *precisely* that we should forbid certain kinds of self-reference in a classical setting.

This I take to be a significant and novel result. After all, self-reference is ubiquitous in allegedly classical languages, and it has had free reign. One example is with the Henkin construction in proving the completeness of predicate logic with identity: Henkin (1949) has each constant denote itself in the constructing the appropriate model. As a more humdrum example, suppose ‘ Wx ’ is defined in predicate logic to mean “ x is wff”: The definition will start by listing the singular terms and predicates—and the list of predicates will include ‘ Wx ’ itself. In this case, just like with ‘ $f(v)$ ’, the symbol is defined in part by referring to that very symbol. This sort of thing has been done without any qualms whatsoever.

This is not to say that a classical language must be free of all self-reference. Nor does the paradox show, specifically, that such a language must not talk about which strings are wffs. But it does indicate that a classical language must be cautious with self-reference, even though caution has largely been absent.

Objection 2: The second objection is whether (†) is a legitimate means to define the predicate ‘ $F(v_3)$ ’. The issue concerns the fact that ‘ v ’ and ‘ v_2 ’ are universally quantified, suggesting that the truth-conditions for ‘ $F(0)$ ’ are less straightforward than what is indicated by (1) alone.

The best way to illustrate the concern is to suppose first (contrary to fact) that $\underline{0}$ is the only designator in \mathcal{L} for 0—thus, compound designators for 0 like ‘ $f_2(0')$ ’ are assumed not to exist. (Imagine, if you prefer, that we are dealing with a fragment of \mathcal{L} .) Regardless, the truth-

condition of ‘F(0)’ would not be determined by (1) alone. Consider, after all, the following instance of (†):

$$(3) \quad \forall v (= (f(v), h) \supset (F(0) \equiv (= (f(v), 0''') \wedge = (0, 0))))$$

Note well that (3) is *not* a CSI of (†). After all, \underline{h} replaces ‘ v_2 ’, and \underline{h} is not the Gödel-numeral of the final designator in (3). Yet since (†) is fully general, (3) still defines a truth-condition for ‘F(0)’. Thus, (1) is not the only formula that determines the truth-value for ‘F(0)’.

But a sentence may well have its truth-value determined by multiple formulae. Thus, for a given set A , the truth-value of ‘ A is an ordered pair’ is determined by several definitions, including:

(Kuratowski) A is an ordered pair iff, for some a and b , $A = \{\{a\}, \{a, b\}\}$.

(K-reverse) A is an ordered pair iff, for some a and b , $A = \{\{b\}, \{a, b\}\}$.

There is no problem here, since the different conditions are equivalent.

Yet perhaps the problem with ‘F(0)’ is precisely that the different formulae determining its truth-value are non-equivalent. Indeed, notice that the embedded conjunction in (3) appears to be false. After all, when $f(v) = h$, then it cannot also be that $f(v) = 3$. Specifically, *in the case of* (2), if the code for that CSI is assigned to ‘ v ’, then the antecedent of (3) will be true, but its embedded conjunction will be false. This reveals that (3) gives the opposite verdict of (1). Namely, (3) determines that ‘F(0)’ is false.

But in the end, this just seems to be an independent demonstration for how the predicate ‘F(v_3)’ is paradoxical.⁴ And my claim all along has been that ‘F(0)’ is both true and false. Furthermore, I *agree* that this means we must somehow exclude such a predicate from a classical

⁴ The argument just given is, in fact, entirely parallel to Jay Newhard’s argument for one of the new paradoxes of self-reference. See sections 3 and 4 of Parent (ms.) for more on Newhard’s argument.

language. However, no provisions against such a predicate have ever been given. It may be possible to correct for that, but the lesson here is that it needs correcting.

4. Closing remarks

The suggestion that Gödel's (1931) metatheory contains paradoxical sentences can cause dramatic reactions. But the most revolutionary conclusions do not follow. In particular, it does not follow that the Gödel-sentence, in particular, is both true and false. Nor does it follow that Peano Arithmetic is complete. Indeed, there is at least one alternate proof of incompleteness due to Kripke, which is reported by Putnam (2000).⁵ Regardless, restrictions on Gödel-numbering are necessary in a classical language, even though this has hitherto been unacknowledged.

⁵ Notably, Kripke's proof turns on a number-theoretic statement which, according to Putnam, "is not at all 'self-referring'" (p. 55). Nonetheless, it utilizes Gödel-numbering without restrictions, and so may still be vulnerable to an objection based on the considerations offered here. I hope to examine this further in future work.

References

- Gödel, K. (1931). Über formal unentscheidbare Sätze der *Principia Mathematica* und verwandter Systeme, I. *Monatshefte für Mathematik und Physik* 38: 173–98.
- Henkin, L. (1949). The completeness of the first-order functional calculus. *Journal of Symbolic Logic* 14: 159–66.
- Parent, T. (ms.). Paradox with just self-reference. Available at <http://tparent.net/d.pdf>.
- Putnam, H. (2000). Nonstandard models and Kripke's proof of Gödel's theorem. *Notre Dame Journal of Formal Logic* 41(1): 53–58.
- Tarski, A. (1944). The semantic conception of truth and the foundations of semantics. *Philosophy and Phenomenological Research* 4(3): 341–76.