# Gödel's Second Incompleteness Theorem Explained in Words of One Syllable

## GEORGE BOOLOS

First of all, when I say "proved", what I will mean is "proved with the aid of the whole of math". Now then: two plus two is four, as you well know. And, of course, *it can be proved* that two plus two is four (proved, that is, with the aid of the whole of math, as I said, though in the case of two plus two, of course we do not need the *whole* of math to prove that it is four). And, as may not be quite so clear, it can be proved that it can be proved that two plus two is four, as well. And it can be proved that it can be proved that it can be proved that two plus two is four. And so on. In fact, if a claim can be proved, then it can be proved that the claim can be proved. And *that* too can be proved.

Now: two plus two is not five. And it can be proved that two plus two is not five. And it can be proved that it can be proved that two plus two is not five, and so on.

Thus: it can be proved that two plus two is not five. Can it be proved as well that two plus two *is* five? It would be a real blow to math, to say the least, if it could. If it could be proved that two plus two is five, then it could be proved that five is not five, and then there would be *no* claim that could *not* be proved, and math would be a lot of bunk.

So, we now want to ask, can it be *proved* that it can't be proved that two plus two is five? Here's the shock: no, it can't. Or to hedge a bit: *if* it can be proved that it can't be proved that two plus two is five, *then* it can be proved as well that two plus two is five, and math is a lot of bunk. In fact, if math is not a lot of bunk, then no claim of the form "claim *X* can't be proved" can be proved.

So, if math is not a lot of bunk, then, though it can't be proved that two plus two is five, it can't be proved *that* it can't be proved that two plus two is five.

By the way, in case you'd like to know: yes, it *can* be proved that if it can be proved that it can't be proved that two plus two is five, then it can be proved that two plus two is five.

## "I wish he would explain his explanation"

If, as we shall assume, the whole of mathematics can be formalized as a formal theory of the usual sort (no small assumption), then there is a formula Proof $(x,y)$

of the language (of that theory) obtainable from a suitable description of the theory ("as" a formal theory) that meets the following three conditions:

(i)   if ⊢ $p$, then ⊢ □ $p$,

(ii)   ⊢ (□ ($p \rightarrow q$) → (□ $p \rightarrow$ □ $q$)), and

(iii)   ⊢ (□ $p \rightarrow$ □ □ $p$)

for all sentences $p$, $q$ of the language. We have written: □ $p$ to abbreviate: ∃xProof(x,˹$p$˺), where ˹$p$˺ is a standard representation in the language for the sentence $p$. (˹$p$˺ might be the numeral for the Gödel number of $p$.) "⊢ " is a preposed verb phrase (of our language) meaning "is provable in the theory". "Proof(x,y)" is a noun phrase (of our language) denoting a formula (of the theory's language) whose construction parallels any standard definition of "… is a proof of ___ in the theory". Thus, for any sentence $p$ of the language, □ $p$ is another sentence of the language that may be regarded as saying that $p$ is provable in the theory.[1]

Conditions (i), (ii), and (iii) are called the Hilbert-Bernays-Löb derivability conditions; they are satisfied by all reasonable formal theories in which a certain small amount of arithmetic can be proved.

Since the theory is standard, all tautologies in its language are provable in the theory, and all logical consequences in its language of provable statements are provable.

It follows that for all sentences $p$, $q$,

(iv)   if ⊢ ($p \rightarrow q$), then ⊢ (□ $p \rightarrow$ □ $q$).

For: if ⊢ ($p \rightarrow q$), then by (i) ⊢ □ ($p \rightarrow q$); but by (ii), ⊢ (□ ($p \rightarrow q$) → (□ $p \rightarrow$ □ $q$)), and then ⊢ (□ $p \rightarrow$ □ $q$) by modus ponens.

⊥ is the zero-place truth-functional connective that is always evaluated as false. Of course ⊥ is a contradiction. We shall need to observe later that (¬$q \rightarrow$ ($q \rightarrow$ ⊥)) is a tautology. If ⊥ is not one of the primitive symbols of the language, it may be defined as any refutable sentence, e.g., one expressing that two plus two is five.

With the aid of ⊥, there is an easy way to say that the theory is consistent: ⊬ ⊥, i.e., ⊥ is not provable in the theory. The sentence of the language stating that the theory is consistent can thus be taken to be ¬□ ⊥, which is identical with ¬∃xProof(x,˹⊥˺).

We may prove Gödel's second incompleteness theorem, as well as the theorem that the second incompleteness theorem is provable in the theory ("the formalized second incompleteness theorem"), as follows.

Via the technique of diagonalization, introduced by Gödel (1931) in "On formally undecidable propositions… ", a sentence $p$ can be found that is equivalent in the theory to the statement that $p$ is unprovable in the theory, i.e. a sentence such that

1.   ⊢ $p \leftrightarrow \neg$□ $p$

2.   ⊢ $p \rightarrow \neg$□ $p$           truth-functionally from 1

---

[1] For an extended account of the application of modal logic to the concept of provability in formal theories, see Boolos (1993).

3.   $\vdash \Box\, p \to \Box\, \neg\, \Box\, p$                 by (iv) from 2

4.   $\vdash \Box\, p \to \Box\, \Box\, p$                 by (iii)

5.   $\vdash \neg\, \Box\, p \to (\Box\, p \to\!\perp)$         a tautology

6.   $\vdash \Box\, \neg\, \Box\, p \to \Box\, (\Box\, p \to\!\perp)$      by (iv) from 5

7.   $\vdash \Box\, (\Box\, p \to\perp) \to (\Box\, \Box\, p \to\Box\, \perp)$   by (ii)

8.   $\vdash \Box\, p \to\Box\, \perp$               truth-functionally from 3, 6, 7, and 4

9.   $\vdash \neg\, \Box\, \perp \to p$                truth-functionally from 8 and 1

10. $\vdash \Box\, \neg\, \Box\, \perp \to \Box\, p$           by (iv) from 9

11. $\vdash \neg\, \Box\, \perp \to \neg\, \Box\, \neg\, \Box\, \perp$       truth functionally from 8 and 10.

(We have omitted outermost parentheses in (1) through (11).)

     Thus if $\vdash \neg\, \Box\, \perp$, then both $\vdash \neg\, \Box\, \neg\, \Box\, \perp$, by (11), and $\vdash \Box\, \neg\, \Box\, \perp$, by (i), whence $\vdash \perp$, by the propositional calculus. So if $\nvdash \perp$, then $\nvdash \neg\, \Box\, \perp$.

*Department of Linguistics and Philosophy*             GEORGE BOOLOS
*M.I.T.*
*Cambridge, MA 02139*
*USA*

## REFERENCES

Boolos, George 1993: *The Logic of Provability*. Cambridge: Cambridge University Press.

Gödel, Kurt 1931: "Über formal unentscheidbare Sätze der *Principia Mathematica* und verwandter Systeme I". *Monatshefte für Mathematik und Physik* 38, pp. 173-198, translated in his *Collected Works*, Volume I, ed. Solomon Feferman et al. Oxford: Oxford University Press, 1986, pp. 145-95.