

12 The First Incompleteness Theorem, syntactic version

We now use the same construction of a Gödel sentence as in the previous chapter to show again that PA is incomplete, but this time making only syntactic assumptions. And then we show how to generalize this syntactic version of the incompleteness theorem.

12.1 ω -completeness, ω -consistency

We need to define two key notions. We'll assume in this section that we are dealing with theories whose language includes the language of basic arithmetic. And take all the quantifiers mentioned to run over the natural numbers.¹

First, then,

Defn. 49. *A theory T is ω -incomplete iff, for some open wff $\varphi(x)$, T can prove $\varphi(\bar{n})$ for each natural number n , but T can't go on to prove $\forall x\varphi(x)$.*

We saw in §5.6 that Q is ω -incomplete: that's because it can prove each instance of $0 + \bar{n} = \bar{n}$, but can't prove $\forall x(0 + x = x)$. We added induction to Q hoping to repair as much ω -incompleteness as we could: but, as we'll see, PA remains ω -incomplete, assuming it is consistent.²

Second, we want the following idea:

Defn. 50. *A theory T is ω -inconsistent iff, for some open wff $\varphi(x)$, T can prove each $\varphi(\bar{n})$ and T can also prove $\neg\forall x\varphi(x)$.*

Or, entirely equivalently of course, we could say that T is ω -inconsistent if, for some open wff $\psi(x)$, $T \vdash \exists x\psi(x)$, yet for each number n we have $T \vdash \neg\psi(\bar{n})$.

Note that ω -inconsistency, like ordinary inconsistency, is a syntactically defined property: it is characterized in terms of what wffs can be proved by the theory, not in terms of what the wffs mean. Note too that, in a classical context,

¹If necessary, therefore, read $\forall x\varphi(x)$ as a restricted quantifier $\forall x(Nx \rightarrow \varphi(x))$, where 'N' picks out the numbers from the domain of the theory's native quantifiers (see Defn. 11).

²Why the ' ω ' in ' ω -incomplete'? Because ' ω ' is a standard label for the set of natural numbers (when equipped with their usual ordering).

ω -consistency – defined of course as not being ω -inconsistent! – trivially implies plain consistency. That’s because T ’s being ω -consistent is a matter of its *not* being able to prove a certain combination of wffs, which entails that T can’t prove *all* wffs, hence T can’t be inconsistent.

Now compare and contrast. Suppose T can prove $\varphi(\bar{n})$ for each n . T is ω -incomplete if it can’t prove something we’d then also like it to prove, namely $\forall x\varphi(x)$. While T is ω -inconsistent if it can actually prove the *negation* of what we’d like it to prove, i.e. it can prove $\neg\forall x\varphi(x)$.

So ω -incompleteness in a theory of arithmetic is a regrettable weakness. But ω -inconsistency is a Very Bad Thing (not as bad as outright inconsistency, maybe, but still bad enough). For evidently, a theory T that can prove each of $\varphi(\bar{n})$ and yet also prove $\neg\forall x\varphi(x)$ is just not going to be an acceptable candidate for regimenting arithmetic.

Bring semantic ideas back into play for a moment. Suppose T ’s standard numerals denote the numbers and the quantifier here runs over the natural numbers. Then it can’t be the case that each of $\varphi(\bar{n})$ is true and yet $\neg\forall x\varphi(x)$ is true too. So our ω -inconsistent T can’t be sound.

Given that we want formal arithmetics to have theorems which *are* all true on a standard interpretation, we must therefore want ω -consistent arithmetics. And given that we think e.g. PA *is* sound on its standard interpretation, we are committed to thinking that it *is* ω -consistent.

12.2 The First Theorem for PA – the syntactic version

G is by definition the diagonalization of the open wff $U =_{\text{def}} \forall x\neg\text{Gdl}(x, y)$, i.e. G is the wff $\forall x\neg\text{Gdl}(x, \ulcorner U \urcorner)$ (see §11.1). And now recall Defn. 46: the wff Gdl by hypothesis doesn’t just express Gdl but *captures* it. Using this fact about Gdl , we can again show that PA does not prove G , but this time *without* making the semantic assumption that PA is sound:

Theorem 39. *If PA is consistent, $\text{PA} \not\vdash G$.*

Proof. Suppose that $\text{PA} \vdash G$.

If G has a proof, then there is some super g.n. m that codes its proof. But G is the diagonalization of the wff with g.n. $\ulcorner U \urcorner$. Hence, by definition, $\text{Gdl}(m, \ulcorner U \urcorner)$.

Since Gdl captures the relation Gdl and we have $\text{Gdl}(m, \ulcorner U \urcorner)$, then by the definition of capturing (i) $\text{PA} \vdash \text{Gdl}(\bar{m}, \ulcorner U \urcorner)$.

But our supposition $\text{PA} \vdash G$, i.e. $\text{PA} \vdash \forall x\neg\text{Gdl}(x, \ulcorner U \urcorner)$, implies (ii) $\text{PA} \vdash \neg\text{Gdl}(\bar{m}, \ulcorner U \urcorner)$, just by instantiating the quantifier.

So, combining (i) and (ii), the supposition that $\text{PA} \vdash G$ entails that PA is inconsistent.

Therefore, if PA is consistent, $\text{PA} \not\vdash G$. □

We’ll now show that PA also can’t prove the *negation* of G , again without assuming PA’s soundness:

Theorem 40. *If PA is ω -consistent, $PA \not\vdash \neg G$.*

Proof. Suppose PA is ω -consistent and $PA \vdash \neg G$. We derive a contradiction, and the theorem follows.

Given PA is ω -consistent, it is consistent. So given PA proves $\neg G$, it can't prove G. So no number m is the super g.n. of a proof for G. But, again, G is the diagonalization of the wff with g.n. $\ulcorner U \urcorner$. Hence, for every number m , $Gdl(m, \ulcorner U \urcorner)$ is false.

Since Gdl captures the relation Gdl and $Gdl(m, \ulcorner U \urcorner)$ is false for each m , by the definition of capturing we have (i) for each m , $PA \vdash \neg Gdl(\overline{m}, \ulcorner U \urcorner)$.

But our supposition $PA \vdash \neg G$ is equivalent to (ii) $PA \vdash \exists x Gdl(x, \ulcorner U \urcorner)$.

Combining (i) and (ii), PA is ω -inconsistent, contradicting our initial supposition. \square

Now recall that G is a Π_1 sentence. That observation put together with what we've shown in this section gives us the following portmanteau result:

Theorem 41. *If PA is consistent, then there is a Π_1 sentence G such that $PA \not\vdash G$, and if PA is ω -consistent $PA \not\vdash \neg G$. Hence, assuming ω -consistency and so consistency, PA is negation incomplete.*

12.3 Two quick corollaries

Theorem 42. *If PA is consistent, it is ω -incomplete.*

Proof. Assume PA is consistent. By Theorem 39, it doesn't prove G.

That means that no number m is the super g.n. of a proof of G. So, exactly as in the proof of the previous theorem, (i) for each m , $PA \vdash \neg Gdl(\overline{m}, \ulcorner U \urcorner)$.

But since G is unprovable, (ii) $PA \not\vdash \forall x \neg Gdl(x, \ulcorner U \urcorner)$.

The combination of (i) and (ii) shows that PA is ω -incomplete. \square

Theorem 43. *If PA is consistent, then $PA + \neg G$ (the theory you get by adding $\neg G$ as an additional axiom) is also consistent but is ω -inconsistent.*

Proof. Assume PA is consistent. If $PA + \neg G$ were inconsistent, then PA would prove G, contrary to Theorem 39. So $PA + \neg G$ is also consistent.

Now, since the expanded theory can prove everything PA can prove, we know as before that (i) for each m , $PA + \neg G \vdash \neg Gdl(\overline{m}, \ulcorner U \urcorner)$.

But just by the definition of $\neg G$, (ii) $PA + \neg G \vdash \exists x Gdl(x, \ulcorner U \urcorner)$.

And (i) and (ii) together imply that $PA + \neg G$ is ω -inconsistent. \square

12.4 Generalizing the proof

The proof for Theorem 41 evidently generalizes. Suppose T is a p.r. axiomatized theory which (perhaps after introducing some new vocabulary by definitions)

contains \mathcal{Q} – in the obvious sense that the language of T includes the language of basic arithmetic, and T can prove every \mathcal{Q} -theorem. Then, assuming we are working with normal scheme for Gödel-numbering wffs of T , the relation $Gdl_T(m, n)$ which holds when m numbers a T -proof of the diagonalization of the wff with number n will be primitive recursive again.

Since T can prove everything \mathcal{Q} proves, T will be able to capture the p.r. relation Gdl_T by a Σ_1 wff Gld_T . Just as we did for PA, we'll be able to construct the corresponding Π_1 wff G_T . And, exactly the same arguments as before will then show, more generally,

Theorem 44. *If T is a consistent p.r. axiomatized theory which contains \mathcal{Q} , then there will be a Π_1 sentence G_T such that $T \not\vdash G_T$, and if T is ω -consistent, $T \not\vdash \neg G_T$. Hence, assuming ω -consistency and so consistency, T is negation incomplete.*

And note, this gives us another incompleteness theorem: if we keep chucking more and more additional axioms at our theory T , it will still remain negation incomplete, unless it stops ω -consistent or stops being p.r. axiomatized.

When people refer to the *First Incompleteness Theorem* (without qualification), they typically mean something like our last Theorem, deriving incompleteness from *syntactic* assumptions. Let's re-emphasize that last point. Being p.r. axiomatized is a syntactic property; containing \mathcal{Q} is a matter of \mathcal{Q} 's axioms again being adopted as axioms or being derivable, a syntactic property; being consistent here is a matter of no contradictory pair $\varphi, \neg\varphi$ being derivable; being ω -consistent is another syntactic property as we stressed before. The chains of argument that lead to our Theorem depend just on the given syntactic assumptions, via e.g. the proof that \mathcal{Q} can capture all p.r. functions – another claim about a syntactically definable property. That is why we are calling this the *syntactic* incompleteness theorem.

Of course, we are *interested* in these various syntactically definable properties because of their semantic relevance: for example, we care about the idea of capturing p.r. functions because we are interested in what an interpreted theory might be able to prove in the sense of establish-as-true. But it is one thing for us to have a semantic motivation for being interested in a certain concept, it is another thing for that concept to have semantic content.

12.5 Comparing old and new syntactic incompleteness theorems

Compare Theorem 44 with our initially announced

Theorem 2. *Suppose T is a formal axiomatized theory whose language contains the language of basic arithmetic. Then, if T is consistent and can prove a certain modest amount of arithmetic (and has a certain additional property that any sensible formalized arithmetic will share), there will be a sentence G_T of basic arithmetic such that $T \not\vdash G_T$ and $T \not\vdash \neg G_T$, so T is negation incomplete.*

Our new theorem fills out the old one in various respects. It fixes the ‘modest amount of arithmetic’ that T is assumed to contain and it also spells out the ‘additional desirable property’ of ω -consistency which we previously left mysterious. Further it tells us more about the undecidable Gödel sentence – namely it has minimal quantifier complexity, i.e. it is a Π_1 sentence of arithmetic. Our new theorem is weaker, however, as it only applies to p.r. axiomatized theories, not to effectively axiomatized theories more generally. But we’ve already noted, that that’s not much loss. (And again, if we insist, we can in fact go on to make up the shortfall: see the next Interlude.)

12.6 Gödel's own Theorem

As we said, Theorem 44, or something like it, is what people usually mean when they speak without qualification of ‘The First Incompleteness Theorem’. But since the stated theorem refers to Robinson Arithmetic Q (developed by Robinson in 1950), and Gödel didn’t originally know about that (in 1931), our version can’t be quite what Gödel originally proved. But it is a near miss. Let’s explain.

Looking again at our analysis of the syntactic argument for incompleteness, we see that we are interested in theories which extend Q *because we are interested in theories which can capture p.r. relations like Gdl* . It’s being able to capture Gdl that is the crucial condition for our proof to go through. So let’s say

Defn. 51. *A theory T is p.r. adequate if it can capture all primitive recursive functions and relations.*

Then, instead of mentioning Q , let’s instead explicitly write in the requirement of p.r. adequacy. So, by just the same arguments,

Theorem 45. *If T is a p.r. adequate, p.r. axiomatized theory whose language includes L_A , then there is Π_1 sentence φ such that, if T is consistent then $T \not\vdash \varphi$, and if T is ω -consistent then $T \not\vdash \neg\varphi$.*

And *this* is pretty much Gödel’s own general version of the incompleteness result. I suppose that it has as much historical right as any to be called *Gödel’s First Theorem*.

(‘Hold on! If *that’s* the original First Theorem, we don’t need to do the hard work showing that Q and PA are p.r. adequate, do we?’ Well, yes and no. No, proving *this* original version of the Theorem of course doesn’t depend on proving that any particular theory is p.r. adequate. But yes, showing that this Theorem has real bite, showing that it does actually apply to familiar arithmetics, does depend on proving that these arithmetics are indeed p.r. adequate.)

Thus, in his 1931 paper, Gödel first proves his Theorem VI which, with a bit of help from his Theorem VIII shows that the formal system P – which is his

12 The First Incompleteness Theorem, syntactic version

simplified version of the hierarchical type-theory of *Principia Mathematica* – has a formally undecidable Π_1 sentence.³ Then he immediately generalizes:

In the proof of Theorem VI no properties of the system P were used besides the following:

1. The class of axioms and the rules of inference (that is, the relation ‘immediate consequence’) are [primitive] recursively definable (as soon as we replace the primitive signs in some way by the natural numbers).
2. Every [primitive] recursive relation is definable [i.e., in our terms, is ‘capturable’] in the system P .

Therefore, in every formal system that satisfies the assumptions 1 and 2 and is ω -consistent, there are undecidable propositions of the form $[\forall xF(x)]$, where F is a [primitive] recursively defined property of natural numbers, and likewise in every extension of such a system by a recursively definable ω -consistent class of axioms.

Which gives us our Theorem 45.

³Or as Gödel put it, the undecidable sentence is ‘of Goldbach type’. The allusion here is to Goldbach’s conjecture that every even number other than two is the sum of two primes. The claim that the particular number n is an even number other than two and is the sum of two primes is expressible by a Δ_0 wff (why?). So Goldbach’s conjecture, the universal quantification of this claim about n , is expressible by a Π_1 wff.

13 The Diagonalization Lemma, and Rosser's Theorem

In the previous chapter, we proved Theorem 44, a version of the syntactic First Theorem. In this chapter, we start by proving this incompleteness theorem again by a slightly different route – this time highlighting the so-called *Diagonalization Lemma*.

The same important Lemma can be used in proving further important theorems. We will look at two. The first of these is due to J. Barkley Rosser, and shows us how to tweak the Gödelian syntactic theorem so that we no longer need the assumption of ω -consistency. Then, in the short following chapter, we show how we can very easily derive an important two-part theorem about truth usually ascribed to Alfred Tarski (though known to Gödel).

13.1 The Diagonalization Lemma

First, it's useful to recall the definition of capturing a function:

Defn. 37. *The theory T captures the one-place function f by the open wff $\psi(x, y)$ iff, for any m, n ,*

- i. if $f(m) = n$, then $T \vdash \psi(\bar{m}, \bar{n})$,*
- ii. if $f(m) \neq n$, then $T \vdash \neg\psi(\bar{m}, \bar{n})$.*
- iii. $T \vdash \exists!y\psi(\bar{m}, y)$.*

And this time, for use in just a moment, let's note explicitly that (i) and (iii) together imply

- iv. if $f(m) = n$, then $T \vdash \forall x(\psi(\bar{m}, x) \leftrightarrow x = \bar{n})$.*

Next, let's restate Theorem 32 about the *diag* function, and apply Theorems 27 and 28 to add an important clause about expressing and capturing it (we assume as usual that we have normal Gödel-numbering scheme in place for the relevant theory T):

Theorem 46. *If T is a p.r. axiomatized theory which contains Q , there is a p.r. function $\text{diag}_T(n)$ which, when applied to a number n which is the g.n. of some T -wff with one free variable, yields the g.n. of that wff's diagonalization, and*

13 The Diagonalization Lemma, and Rosser's Theorem

yields n otherwise. And, as with any p.r. function, T can express and capture this function by a Σ_1 wff $\text{Diag}_T(x, y)$.

We can now officially state and prove the two-part *Diagonalization Lemma* (Rudolf Carnap noted a version of part (i): he is often, but I think wrongly, attributed (ii) as well):

Theorem 47. *If T is a p.r. axiomatized theory which contains \mathbf{Q} , and φ is a one-place open sentence of T 's language, then there is sentence δ such that (i) $\delta \leftrightarrow \varphi(\overline{\overline{\delta}})$ is true, and moreover (ii) $T \vdash \delta \leftrightarrow \varphi(\overline{\overline{\delta}})$.*

To avoid unsightly rashes of subscripts, let's drop subscript ' T 's. Then we can argue as follows. (The proofs *look* complicated at first glance; but after the initial construction of δ , we are just applying various definitions and easy consequences.)

Proof for (i). Put $\alpha =_{\text{def}} \forall z(\text{Diag}(y, z) \rightarrow \varphi(z))$, and let δ be the diagonalization of α . So by definition δ is $\forall z(\text{Diag}(\overline{\overline{\alpha}}, z) \rightarrow \varphi(z))$.

Because diagonalizing α yields δ , by definition $\text{diag}(\ulcorner \alpha \urcorner) = \ulcorner \delta \urcorner$. So since Diag expresses *diag*, we know that $\text{Diag}(\overline{\overline{\alpha}}, \overline{\overline{\delta}})$; indeed, $\text{Diag}(\overline{\overline{\alpha}}, z)$ is *only* satisfied by $\ulcorner \delta \urcorner$.

So $\forall z(\text{Diag}(\overline{\overline{\alpha}}, z) \rightarrow \varphi(z))$ is true if and only if $\ulcorner \delta \urcorner$ satisfies $\varphi(z)$. In other words, δ is true if and only if $\varphi(\overline{\overline{\delta}})$ is true. \square

Proof for (ii). Since by hypothesis Diag captures *diag* in T , (iv) from Defn. 37 just above tells us that *if $\text{diag}(\ulcorner \alpha \urcorner) = \ulcorner \delta \urcorner$, then $T \vdash \forall z(\text{Diag}(\overline{\overline{\alpha}}, z) \leftrightarrow z = \overline{\overline{\delta}})$.*

But by definition $\text{diag}(\ulcorner \alpha \urcorner) = \ulcorner \delta \urcorner$. Hence $T \vdash \forall z(\text{Diag}(\overline{\overline{\alpha}}, z) \leftrightarrow z = \overline{\overline{\delta}})$.

Since T can prove the equivalence of $\text{Diag}(\overline{\overline{\alpha}}, z)$ and $z = \overline{\overline{\delta}}$, T can also prove the equivalence of $\forall z(\text{Diag}(\overline{\overline{\alpha}}, z) \rightarrow \varphi(z))$ and $\forall z(z = \overline{\overline{\delta}} \rightarrow \varphi(z))$.

In other words, $T \vdash \delta \leftrightarrow \forall z(z = \overline{\overline{\delta}} \rightarrow \varphi(z))$. Which trivially gives us $T \vdash \delta \leftrightarrow \varphi(\overline{\overline{\delta}})$. \square

A bit of jargon: by a mild abuse of mathematical terminology, we say

Defn. 52. *If δ is such that $T \vdash \delta \leftrightarrow \varphi(\overline{\overline{\delta}})$, then it is said to be a fixed point for φ .*

So the Diagonalization Lemma – or rather, part (ii) of it – is often called the Fixed Point Theorem: for appropriate theories T , every one-place open sentence has a fixed point.

13.2 Incompleteness from the Diagonalization Lemma

Suppose as usual that $\text{Prf}_T(m, n)$ is the relation which holds just if m numbers a T proof of a sentence with g.n. n (we continue to assume, of course, that we have a normal Gödel-numbering scheme in place). If T is p.r. axiomatized and contains \mathbf{Q} , this relation is p.r. decidable and can be captured in T by a Σ_1 wff $\text{Prf}_T(x, y)$. And now we pick up again an idea we first met in §3.6:

Defn. 16. Put $\text{Prov}_T(y) =_{\text{def}} \exists x \text{Prf}_T(x, y)$ (where the quantifier runs over all the numbers in the domain). Then $\text{Prov}_T(\bar{n})$ says that some number Gödel-numbers a T -proof of the wff with Gödel-number n , i.e. the wff with code number n is a T -theorem. So $\text{Prov}_T(x)$ is naturally called a provability predicate.

Hence $\text{Prov}_T(\overline{\neg\varphi})$ is true just when φ is a theorem. Note, by the way, that since Prov_T is built by existentially quantifying a Σ_1 wff, it is also Σ_1 .

And here now is a general result about fixed points for the *negation* of such a provability predicate:

Theorem 48. *Suppose T is p.r. axiomatized, contains Q , and some sentence γ is a fixed point for $\neg\text{Prov}_T$; in other words, suppose $T \vdash \gamma \leftrightarrow \neg\text{Prov}_T(\overline{\neg\gamma})$. Then (i) if T is consistent, $T \not\vdash \gamma$. And (ii) if T is ω -consistent, $T \not\vdash \neg\gamma$.*

Again, to avoid unsightly rashes of subscripts, let's drop subscript ' T 's. Then we can argue like this:

Proof. (i) Suppose $T \vdash \gamma$. Then, since $T \vdash \gamma \leftrightarrow \neg\text{Prov}_T(\overline{\neg\gamma})$, we have $T \vdash \neg\text{Prov}(\overline{\neg\gamma})$. But if there *is* a proof of γ , then for some m , $\text{Prf}(m, \overline{\neg\gamma})$, so $T \vdash \text{Prf}(\bar{m}, \overline{\neg\gamma})$, since T captures Prf by Prf . Hence $T \vdash \exists x \text{Prf}(x, \overline{\neg\gamma})$, i.e. we also have $T \vdash \text{Prov}(\overline{\neg\gamma})$, making T inconsistent. So if T is consistent, $T \not\vdash \gamma$.

(ii) Suppose $T \vdash \neg\gamma$. Then, since $T \vdash \gamma \leftrightarrow \neg\text{Prov}_T(\overline{\neg\gamma})$, we have $T \vdash \text{Prov}(\overline{\neg\gamma})$, i.e. $T \vdash \exists x \text{Prf}(x, \overline{\neg\gamma})$. Given T is consistent and proves $\neg\gamma$, there is no proof of γ , i.e. for every m , not- $\text{Prf}(m, \overline{\neg\gamma})$, whence for every m , $T \vdash \neg\text{Prf}(\bar{m}, \overline{\neg\gamma})$. So we have a $\psi(x)$ such that T proves $\exists x \psi(x)$ while it refutes each instance $\psi(\bar{m})$. Hence, if it is consistent, T is ω -inconsistent. So if T is ω -consistent (and hence consistent), $T \not\vdash \neg\gamma$. \square

But of course the general Diagonalization Lemma tells us that, as a special case,

Theorem 49. *If T is p.r. axiomatized, contains Q , then there exists a sentence γ such that $T \vdash \gamma \leftrightarrow \neg\text{Prov}_T(\overline{\neg\gamma})$.*

Moreover, since Prov_T is Σ_1 , $\neg\text{Prov}_T$ is Π_1 , and the diagonalization construction produces a Π_1 fixed point γ (see if you can work out why).

So putting those last two theorems together, we immediately recover Theorem 44.¹

13.3 Proving our old G_T is in fact a fixed point for $\neg\text{Prov}_T$

How does our new proof of the syntactic incompleteness theorem relate to the old one?

¹Warning. Some authors call any fixed point γ for $\neg\text{Prov}_T$ a Gödel sentence for T . That's fine, as long as you are alert to the fact that not everything that is true of Gödel sentences in the narrow sense we introduced in the preceding chapters is true of Gödel sentences in this wider sense. There's more about this in *IGT2*.

13 The Diagonalization Lemma, and Rosser's Theorem

Our canonical Gödel sentence G_T , recall, was so constructed that it is true if and only if unprovable-in T . This fact can now be *expressed* inside T itself, by the wff $G_T \leftrightarrow \neg \text{Prov}_T(\ulcorner G_T \urcorner)$. But T doesn't just express this fact but can quite easily *prove* it too, i.e. we have

Theorem 50. *If T is p.r. axiomatized and contains Q , $T \vdash G_T \leftrightarrow \neg \text{Prov}_T(\ulcorner G_T \urcorner)$.*

In other words, G_T is one of those fixed points for $\neg \text{Prov}_T$ which are all undecidable by T . And note that (notation apart) this gives us the key result we simply announced way back in §3.6. The proof is a bit fiddly but elementary. Dropping subscripts once more:

Proof. By definition, $Gdl(m, n)$ holds iff $\text{Prf}(m, \text{diag}(n))$. We can therefore fix on the following canonical definition:

$$Gdl(x, y) =_{\text{def}} \exists z (\text{Prf}(x, z) \wedge \text{Diag}(y, z)).$$

And now let's do some elementary manipulations:

$$\begin{aligned} G &=_{\text{def}} \forall x \neg Gdl(x, \ulcorner U \urcorner) \\ &\leftrightarrow \forall x \neg \exists z (\text{Prf}(x, z) \wedge \text{Diag}(\ulcorner U \urcorner, z)) && \text{(definition of Gdl)} \\ &\leftrightarrow \forall x \forall z \neg (\text{Prf}(x, z) \wedge \text{Diag}(\ulcorner U \urcorner, z)) && \text{(pushing in the negation)} \\ &\leftrightarrow \forall z \forall x \neg (\text{Prf}(x, z) \wedge \text{Diag}(\ulcorner U \urcorner, z)) && \text{(swapping quantifiers)} \\ &\leftrightarrow \forall z (\text{Diag}(\ulcorner U \urcorner, z) \rightarrow \neg \exists x \text{Prf}(x, z)) && \text{(rearranging after '}\forall z\text{'}) \\ &=_{\text{def}} \forall z (\text{Diag}(\ulcorner U \urcorner, z) \rightarrow \neg \text{Prov}(z)) && \text{(definition of Prov)} \end{aligned}$$

Since this is proved by simple logical manipulations, that means we can prove the equivalence inside the formal first-order logic built into Q and hence in T . Therefore we have

$$T \vdash G \leftrightarrow \forall z (\text{Diag}(\ulcorner U \urcorner, z) \rightarrow \neg \text{Prov}(z)).$$

Now, diagonalizing U yields G . Hence, just by the definition of *diag*, we have $\text{diag}(\ulcorner U \urcorner) = \ulcorner G \urcorner$. Since by hypothesis Diag captures *diag*, it follows that

$$T \vdash \forall z (\text{Diag}(\ulcorner U \urcorner, z) \leftrightarrow z = \ulcorner G \urcorner).$$

Putting those two results together, we immediately get

$$T \vdash G \leftrightarrow \forall z (z = \ulcorner G \urcorner \rightarrow \neg \text{Prov}(z)).$$

But the right-hand side of that biconditional is trivially equivalent to $\neg \text{Prov}(\ulcorner G \urcorner)$. So we've proved the desired result. \square

We should note a simple corollary of that last theorem. Suppose Prov_T not only expresses but *captures* the property of being a T -theorem. Then, by the definition of capturing, if φ is a non-theorem, then we'd have $T \vdash \neg \text{Prov}_T(\ulcorner \varphi \urcorner)$. And in particular, since G_T is a non-theorem, $T \vdash \neg \text{Prov}_T(\ulcorner G_T \urcorner)$. But, given Theorem 50, that would imply $T \vdash G_T$ which we know is false.

Hence Prov_T does *not* capture the property of being a T -theorem.

13.4 Rosser's Theorem

One half of the syntactic First Theorem requires the assumption that we are dealing with a theory T which is not only consistent but is ω -consistent. But we can improve on this. Following Barkley Rosser, we can construct a Rosser sentence R_T such that we only need to assume T is plain consistent in order to show that R_T is formally undecidable.

So how does Rosser do the trick? Essentially, where Gödel constructs a sentence G_T that indirectly says 'I am unprovable in T ', Rosser constructs R_T which indirectly says 'if I am provable in T , then my negation is already provable' (i.e. it says that if there is a proof of R_T with super g.n. n , then there is a proof of $\neg R_T$ with a smaller code number).

How do we formally implement this idea? Start by considering the relation $\overline{Prf}_T(m, n)$ which holds when m numbers a T -proof of the *negation* of the wff with number n . This relation is obviously p.r. given that Prf_T is; so assuming T has the usual properties it will be captured by a wff $\overline{Prf}_T(x, y)$. So we can construct *the Rosser provability predicate* defined as follows:

Defn. 53. $RProv_T(x) =_{\text{def}} \exists v(Prf_T(v, x) \wedge (\forall w \leq v)\neg\overline{Prf}_T(w, x))$.

Then a sentence is Rosser-provable in T – its g.n. satisfies the Rosser provability predicate – if it has a proof (in the ordinary sense) and there's no 'smaller' proof of its negation.

Now we apply the Diagonalization Lemma, not to the negation of a regular provability predicate (which is what we just did to get Gödel's First Theorem again), but to the negation of the Rosser provability predicate. The Lemma then tells us,

Theorem 51. *Given that T is p.r. axiomatized and contains Q , then there is a sentence R_T such that $T \vdash R_T \leftrightarrow \neg RProv_T(\overline{R_T})$.*

We call such a sentence R_T a Rosser sentence for T . We can now show that

Theorem 52. *Let T be consistent p.r. axiomatized theory which contains Q , and let ρ be any fixed point for $\neg RProv_T(x)$. Then $T \not\vdash \rho$ and $T \not\vdash \neg\rho$.*

And since the Diagonalization Lemma tells us that there is such a fixed point, it follows that T has an undecidable sentence R_T , without now requiring ω -consistency. Moreover, with extra effort, we can get ourselves a Π_1 sentence which is undecidable if T is consistent. Which gives us Rosser's Theorem in the form:

Theorem 53. *Let T be a consistent p.r. axiomatized theory which contains Q , then there is Π_1 sentence φ such that $T \not\vdash \varphi$ and $T \not\vdash \neg\varphi$.*

But what we gain on the swings (a slightly stronger result) we lose on the roundabouts (in the complications in getting there). There's no avoiding it; the

13 The Diagonalization Lemma, and Rosser's Theorem

proof of Theorem 52 is decidedly unpretty. We then have to do more work to beef up that proof idea to show that in fact (as with Gödel's original proof) we can find a Π_1 sentence which is undecidable so long as T is consistent.² But you do not need to know these proofs! – it is just good to know that they exist.

Note, though, that once more Rosser's Theorem gives us an incompleteness result. Suppose T is a consistent p.r. axiomatized theory which contains Q . It is incomplete. Add more axioms. It will remain incomplete unless it becomes inconsistent or stops being p.r. axiomatized.

²Masochists can check out the proof of Theorem 25.4 and further details in *IGT2*.