

## A NOTE ON THE STRENGTH OF DISENTANGLED TRUTH-THEORIES

RICHARD KIMBERLY HECK

ABSTRACT. So-called ‘disentangled’ truth-theories are supposed to prevent assumptions about the truth of statements in the object-language from inadvertently strengthening the background syntax. In earlier work, I proved some limitative results in an attempt to show that the strategy works, but those results leave several questions unanswered. We address some of them here. We also discuss a subtlety that has so far been overlooked in discussions of these theories.

Every axiom of Peano Arithmetic is true. The rules of inference preserve truth. So every theorem of PA is true. Since  $0 = 1$  is not true, it is not a theorem of PA. So PA is consistent.

Such ‘soundness arguments’ play an important role in discussion of theories of truth. Because they establish the consistency of the ‘base theory’—PA in this case—any theory in which such an argument can be formalized must be stronger than the base theory, by Gödel’s second incompleteness theorem. But there are some puzzling aspects about the way such arguments are usually formalized.

We start with some arithmetical base theory  $\Sigma$ ; we assume that  $\Sigma$  is strong enough to be able to talk about its own syntax, via Gödel numbering. We then expand the language by adding semantic vocabulary—a truth-predicate and the like—and extend  $\Sigma$  with Tarski-style compositional axioms. Those axioms will allow us to prove the ‘T-sentences’ for the language of arithmetic

$$A \equiv \text{Tr}(\ulcorner A \urcorner)$$

and so to prove that each axiom of  $\Sigma$  is itself true, by a trivial argument:

- (1)  $A$ , since  $A$  is an axiom
- (2)  $A \equiv \text{Tr}(\ulcorner A \urcorner)$ , since the T-sentences are provable
- (3) So,  $\text{Tr}(\ulcorner A \urcorner)$

To carry out the induction that will take us from the truth of the axioms to the truth of the theorems, we also need to extend whatever induction axioms are present in  $\Sigma$  to allow semantic vocabulary to occur in those axioms: The most natural formalization requires  $\Sigma$  to contain induction for  $\Sigma_1$  formulas (Heck, 2015, Theorem 3.20), though it is in fact sufficient

to have induction for  $\Delta_0$  formulas (Łełyk, 2022).<sup>1</sup> Call the result of extending  $\Sigma$  in this way  $\text{CT}[\Sigma]$ .<sup>2</sup> Then, so long as  $\Sigma$  extends  $\text{I}\Delta_0$ ,<sup>3</sup>  $\text{CT}[\Sigma]$  will prove  $\text{Con}(\Sigma)$ .

This works, but it is in some ways unsatisfying. Suppose we want in this way to prove  $\text{Con}(\text{PA})$ . So we work in  $\text{CT}[\text{PA}]$ . Only finitely many of PA's axioms can be used in the proof (and PA is not finitely axiomatizable). So  $\text{Con}(\text{PA})$  must be provable in  $\text{CT}[\text{I}\Sigma_n]$ , for some  $n$ . Moreover, there's a lacuna in the proof sketched earlier. It's clear enough how  $\text{CT}[\text{PA}]$  proves that *each* axiom of PA is true. But how can we prove the general statement that *all* axioms of PA are true?<sup>4</sup> The answer is: by induction. There is, in fact, a single  $\Sigma_1$  formula, in the extended language, from which all the infinitely many induction axioms of PA follow (Heck, 2015, pp. 447–50). What this means is that  $\text{CT}[\text{I}\Sigma_1]$  contains PA. And, as noted,  $\Sigma_1$  induction is enough to carry out the induction at the heart of the soundness argument. So, in fact,  $\text{CT}[\text{I}\Sigma_1]$  proves  $\text{Con}(\text{PA})$  already (Heck, 2015, Theorem 3.20), and in fact  $\text{CT}[\text{I}\Delta_0]$  does so (Łełyk, 2022).

You might have thought that  $\text{CT}[\text{I}\Sigma_1]$  would prove  $\text{Con}(\text{I}\Sigma_1)$ ,  $\text{CT}[\text{I}\Sigma_2]$  would prove  $\text{Con}(\text{I}\Sigma_2)$ , and so forth, and that we'd need  $\text{CT}[\text{PA}]$  to prove  $\text{Con}(\text{PA})$ . But no. If you try to prove  $\text{Con}(\text{I}\Sigma_1)$  via this kind of argument, you end up doing so in a theory already capable of proving  $\text{Con}(\text{PA})$ . That's what's unsatisfying.

The reason we get this result is that the 'base theory' is playing two different roles. Suppose we naïvely pursue the mentioned strategy in an attempt to prove  $\text{Con}(\text{I}\Sigma_1)$ . So we work in  $\text{CT}[\text{I}\Sigma_1]$ . Then, on the one hand,  $\text{I}\Sigma_1$  is the theory whose consistency we are proving. In that role, it is what allows us to prove the truth of  $\text{I}\Sigma_1$ 's axioms, via the trivial argument mentioned above.<sup>5</sup> On the other hand, an extension of  $\text{I}\Sigma_1$  is also the theory *in which* we are proving consistency: It's what allows us to reason about syntax and semantics. But when we extend the induction axioms of  $\text{I}\Sigma_1$  *qua* syntax and semantics—which we need to do in order to carry out the induction that takes us from the truth of the axioms to

<sup>1</sup>A formula is  $\Delta_0$  if all quantifiers occurring in it are bounded, i.e., are of the form  $\forall x < t$  or  $\exists x < t$ ;  $\Sigma_1$  if it is of the form  $\exists v_1 \dots \exists v_n(\phi)$ , where  $\phi$  is  $\Delta_0$ ;  $\Pi_1$ , of the form  $\forall v_1 \dots \forall v_n(\phi)$ , where  $\phi$  is  $\Delta_0$ ;  $\Sigma_n$ , of the form  $\exists v_1 \dots \exists v_n(\phi)$ , where  $\phi$  is  $\Pi_{n-1}$ ;  $\Pi_n$ , of the form  $\forall v_1 \dots \forall v_n(\phi)$ , where  $\phi$  is  $\Sigma_{n-1}$ .

<sup>2</sup>We'll not be interested here in the weaker theory  $\text{CT}^-[\Sigma]$ , in which that has not been done. As is well known,  $\text{CT}^-[\Sigma]$  is a conservative extension of  $\Sigma$  and so does not prove  $\text{Con}(\Sigma)$ . That does not mean, however, that  $\text{CT}^-[\Sigma]$  is no stronger than  $\Sigma$  itself. In fact, if  $\Sigma$  is finitely axiomatizable, it is stronger (Heck, 2015, §3.2).

<sup>3</sup> $\text{I}\Delta_0$  is Robinson arithmetic, Q, plus the induction axioms for  $\Delta_0$  formulas. Similarly, for  $\text{I}\Sigma_n$ .

<sup>4</sup>Compare: For each  $n$ , PA proves that  $n$  does not code a proof of  $0 = 1$ . But PA does not prove the general statement: For all  $n$ ,  $n$  does not code a proof of  $0 = 1$ , since that just is  $\text{Con}(\text{PA})$ .

<sup>5</sup> $\text{I}\Sigma_1$  is finitely axiomatizable. So, if we use some finite axiomatization, we do not have to worry about how to get from 'each' to 'all'.

the truth of the theorems—we simultaneously give ourselves the ability to prove the truth of *much more* than just the axioms of  $I\Sigma_1$ . In fact, we are now able to prove the truth of all the axioms of PA. The semantic induction axioms, as it were, inadvertently strengthen the base theory itself.

In earlier work (Heck, 2009, 2015), suggested that we might address this problem by following Tarski (1956) and ‘disentangling’ the two roles played by the base theory.<sup>6</sup> We now work in a multi-sorted framework. One sort contains numbers: the objects that our ‘target’ theory is about.<sup>7</sup> Another sort contains syntactic objects—terms, formulas, proofs, and the like, from the language of the target theory—and we have a syntactic theory that allows us to reason about these objects. It might, for example, contain axioms that allow us to carry out induction on the complexity of expressions, or on the length of proofs. A third sort contains assignments: functions from variables to their values. We then have a semantic theory that relates the syntactic objects to the numbers. Officially, we might think of the syntactic theory as formulated in a language whose sole primitive is a symbol for concatenation:  $x \frown y$ , so that the theory really is a theory of syntax.<sup>8</sup> For convenience and familiarity, however, we may take the syntactic theory also to be an arithmetical theory. Still, we want to think of the syntactic language  $\mathcal{S}$  as distinct from the language of arithmetic,  $\mathcal{A}$ . Perhaps it is “written in boldface, or something of the sort” (Heck, 2015, p. 451).

Let  $\text{CTD}[\Sigma]$  be the semantic theory for the language of arithmetic built on the syntactic theory  $\Sigma$ .<sup>9</sup> As before, a natural formalization of the soundness argument requires  $\Sigma_1$  induction, so we assume that  $\Sigma$  contains  $I\Sigma_1$ . If we want to prove the consistency of some arithmetical theory  $\Theta$  in this way, then we also need to assume that  $\Theta$ ’s axioms are true. But, if we do, then we can indeed prove  $\text{Con}(\Theta)$ . That is:<sup>10</sup>

**Theorem 1** (Heck, 2015, Theorem 4.11).  $\text{CTD}[I\Sigma_1] + \text{Tr}(\Theta)$  *proves*  $\text{Con}(\Theta)$ . *Moreover, if  $\Theta$  is finitely axiomatized, then  $\text{CTD}[I\Sigma_1] + \Theta$  proves  $\text{Con}(\Theta)$ .*

So the disentangled framework still allows us to carry out the soundness argument.

<sup>6</sup>Leigh and Nicolai (2013) also develop this approach.

<sup>7</sup>An additional advantage of this way of proceeding is that the target theory could also be ZF, say, without that affecting the syntactic theory in any dramatic way. The syntax will talk about ‘ $\epsilon$ ’ instead of ‘0’ and ‘ $S$ ’, but it will otherwise be unchanged.

<sup>8</sup>Halbach and Leigh (2024) discuss such theories at length.

<sup>9</sup>I borrow this notation from Leigh and Nicolai (2013).

<sup>10</sup> $\text{Tr}(\Theta)$  is the formalization of “All of  $\Theta$ ’s axioms are true”. There will be many possible formalizations, in fact, depending upon how the set of  $\Theta$ ’s axioms is represented. If  $\Theta$  is finitely axiomatizable, then there is a canonical choice:  $x = A_1 \vee \dots \vee x = A_n$ . But if it is not, then there famously is no canonical choice (Feferman, 1960). So problems of intensionality will arise here, as they do in the case of the second incompleteness theorem.

This result allows us to illustrate another advantage of ‘disentangling’. The theory  $\text{CT}[\text{I}\Sigma_1 + \neg\text{Con}(\text{PA})]$  is inconsistent, because  $\text{CT}[\text{I}\Sigma_1]$  proves  $\text{Con}(\text{PA})$  already. So we cannot, in that way, formalize a (putative) soundness proof for  $\text{I}\Sigma_1 + \neg\text{Con}(\text{PA})$ . But we can do so in the disentangled framework:  $\text{Con}(\text{I}\Sigma_1 + \neg\text{Con}(\text{PA}))$  is provable in  $\text{CTD}[\text{I}\Sigma_1] + (\text{I}\Sigma_1 + \neg\text{Con}(\text{PA}))$ , which is consistent—since, by Theorem 2 below, it is mutually interpretable with  $\text{I}\Sigma_1 + \text{Con}(\text{I}\Sigma_1 + \neg\text{Con}(\text{PA}))$ . And how do we know that  $\text{I}\Sigma_1 + \text{Con}(\text{I}\Sigma_1 + \neg\text{Con}(\text{PA}))$  is consistent? Because it is true!<sup>11</sup>

Does the disentangled framework solve the problems that motivated this framework in the first place? Heck (2015) does prove some limitative results that are meant to address this question. For our purposes, the most important of these is:

**Theorem 2** (Heck, 2015, Corollary 4.14). *CTD $[\text{I}\Sigma_n] + \text{Tr}(\Theta)$  is mutually interpretable with  $\text{I}\Sigma_n + \text{Con}(\Theta)$ . Moreover, if  $\Theta$  is finitely axiomatized, then  $\text{CTD}[\text{I}\Sigma_n] + \Theta$  is mutually interpretable with  $\text{I}\Sigma_n + \text{Con}(\Theta)$ .*

This implies, in particular, that  $\text{CTD}[\text{I}\Sigma_1] + \text{I}\Sigma_1$  is mutually interpretable with  $\text{I}\Sigma_1 + \text{Con}(\text{I}\Sigma_1)$ , which is a sub-theory not only of PA but of  $\text{I}\Sigma_2$ . So if you want to prove  $\text{Con}(\text{I}\Sigma_2)$ , you do have to work in a stronger theory, namely  $\text{CTD}[\text{I}\Sigma_1] + \text{I}\Sigma_2$ . The ‘semantic’ induction axioms in  $\text{CTD}[\text{I}\Sigma_1]$  thus do not inadvertently strengthen the object-language theory.

However, this does not, by itself, answer all the questions one might want to ask here. One might wonder not just whether semantic induction strengthens the target theory but also whether our assumptions about the semantics of the target theory somehow infect the syntactic theory. Consider, for example,  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$ . We can, of course, reason about the language of arithmetic in PA itself, via Gödel numbering, and prove facts about that language in PA that we cannot prove in  $\text{I}\Sigma_1$ . But the semantic theory establishes a mapping from numerals to numbers: Each numeral  $\bar{n}$  denotes the corresponding number  $n$ ; so the open term “ $\text{den}(\bar{x})$ ”—the denotation of the  $x^{\text{th}}$  numeral—describes a function that embeds the numbers in the syntactic sort, the ‘*snumbers*’, into the numbers in the target sort. Might this somehow allow us to transfer facts about the numbers in the target sort back into the syntax? If we can prove that all *numbers* have some (coded) syntactic property, and if the *snumbers* are provably mapped onto an initial segment of the numbers, might that imply that the *snumbers* have the (coded) syntactic property, too?

Theorem 2 does tell us that  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  is mutually interpretable with  $\text{I}\Sigma_1 + \text{Con}(\text{PA})$ . But *that* theory does interpret PA, by the arithmetized

<sup>11</sup>Feferman (1960, Theorem 6.6) showed that PA interprets  $\text{PA} + \neg\text{Con}(\text{PA})$ . His proof relies upon the reflexivity of PA, but Visser (2014) has proved that  $\text{I}\Sigma_n$  interprets  $\text{I}\Sigma_n + \neg\text{Con}(\text{I}\Sigma_n)$ —and a more general result still. In light of remarks below,  $\text{I}\Sigma_n$  will therefore prove  $\text{Con}(\text{I}\Sigma_n) \equiv \text{Con}(\text{I}\Sigma_n + \neg\text{Con}(\text{I}\Sigma_n))$ .

completeness theorem,<sup>12</sup> so it isn't clear, at this point, whether we might be able to prove all of the *purely syntactic* induction axioms (as opposed to the ones with semantic vocabulary). The worry isn't, then, that  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  is no different from  $\text{CTD}[\text{PA}] + \text{Tr}(\text{PA})$ . It seems unlikely, to say the least, that we'll be able to prove all the *extended* induction axioms, since our target theory has no semantic vocabulary. But the non-extended induction axioms—that is, the induction axioms in the purely syntactic language  $\mathcal{S}$ —are a different matter. The question is whether we can prove those. We *can*, after all, prove their analogue in the object-theory.

As we'll see, this worry can be addressed.

We'll start with a relatively simple case. Suppose we could prove all the induction axioms for  $\mathcal{S}$  in  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$ . Then, one might reason, since we can also prove  $\text{Con}(\text{PA})$ , and since  $\text{Con}(\text{PA})$  is itself sentence of  $\mathcal{S}$ ,  $\text{PA} + \text{Con}(\text{PA})$  would be a sub-theory of  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$ . So it will be enough to show that  $\text{I}\Sigma_1 + \text{Con}(\text{PA})$  does not interpret  $\text{PA} + \text{Con}(\text{PA})$ . And this is easy to show.<sup>13</sup> Since  $\text{PA}$  is essentially reflexive, so is  $\text{PA} + \text{Con}(\text{PA})$ , which means that  $\text{PA} + \text{Con}(\text{PA})$  proves the consistency of every one of its finite sub-theories. But  $\text{I}\Sigma_1 + \text{Con}(\text{PA})$  is such a sub-theory, so  $\text{PA} + \text{Con}(\text{PA})$  proves  $\text{Con}(\text{I}\Sigma_1 + \text{Con}(\text{PA}))$ . Since no (consistent) theory can interpret any theory that proves its consistency,<sup>14</sup> it thus follows that  $\text{I}\Sigma_1 + \text{Con}(\text{PA})$  does not interpret  $\text{PA} + \text{Con}(\text{PA})$ .

In fact, however, there is a subtlety here that we have overlooked. Volker Halbach once observed that we need to distinguish the consistency statement  $\text{Con}_{\mathcal{S}}(\Theta)$  in the *syntactic* language from the (coded) consistency statement  $\text{Con}_{\mathcal{A}}(\Theta)$  in the target language. What  $\text{CTD}[\Sigma] + \Theta$  proves is  $\text{Con}_{\mathcal{S}}(\Theta)$ . Halbach additionally observed that  $\text{CTD}[\Sigma] + \Theta$  is always a conservative extension of  $\Theta$ , so  $\text{CTD}[\Sigma] + \Theta$  does *not* prove  $\text{Con}_{\mathcal{A}}(\Theta)$ , unless  $\text{CTD}[\Sigma] + \Theta$  is inconsistent (Leigh and Nicolai, 2013, §3.2; Heck, 2018, §5). So we cannot transfer results from the syntax to the target.

Similarly, then, the statement  $\text{Con}_{\mathcal{S}}(\text{PA})$  that  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  proves is a statement about the *arithmetical* version of  $\text{PA}$ . So it is really  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{A}})$  that  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  proves, *not*  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ , which asserts the consistency of the *syntactic* version of  $\text{PA}$ . Now, on reflection, this does not affect the argument two paragraphs back: That depended only upon the fact that  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{A}})$  belongs to the language  $\mathcal{S}$ . But it does raise the question whether  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$  can also be proven. This might seem trivial and unimportant. But we'll need to know below that  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$

<sup>12</sup>Indeed,  $\text{Q} + \text{Con}(\mathcal{T})$  always interprets  $\mathcal{T}$ .

<sup>13</sup>Thanks to Carlo Nicolai for this observation.

<sup>14</sup>In fact, something stronger is true: Even if  $\Theta$  only proves  $\text{Con}(\Sigma)$  on a cut,  $\Sigma$  cannot interpret  $\Theta$ . See Heck (2015, §2.5) for a proof of this folklorish result, which is implicit in Pudlák (1985).

proves  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ , and the issue is more obviously important, and not at all trivial, if we take our syntactic theory to be, as it ideally should be, a theory of concatenation. Then  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$  does not even make sense, since PA is not a theory of concatenation. Similar issues arise when the object-language is, say, the language of set-theory.  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{ZF})$  is mutually interpretable with  $\text{I}\Sigma_1 + \text{Con}(\text{ZF})$ . It may not be immediately obvious whether this theory proves  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ .

However,  $\text{PA}^{\mathcal{S}}$  and  $\text{PA}^{\mathcal{A}}$  are (trivially) mutually interpretable. And  $\text{I}\Sigma_1$  is strong enough both (i) to show that  $\text{PA}^{\mathcal{S}}$  is interpretable in  $\text{PA}^{\mathcal{A}}$  and (ii) to prove that interpretability implies relative consistency, i.e., to prove that  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{A}})$  implies  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ . So  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  does prove  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ , as wanted. This does not depend upon the triviality of the interpretation. Generally speaking,  $\text{I}\Sigma_1$  is strong enough to establish the existence of interpretations, where they exist, since this is just a matter of the existence of certain proofs.<sup>15</sup> So  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{ZF})$  does prove  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ , by showing that  $\text{PA}^{\mathcal{S}}$  is interpretable in ZF and so that  $\text{Con}(\text{ZF})$  implies  $\text{Con}_{\mathcal{S}}(\text{PA}^{\mathcal{S}})$ . Similar results will be available when the syntax is the syntactic analogue of  $\text{I}\Sigma_1$ :<sup>16</sup> Such a theory will be strong enough both to establish results about interpretability and to show that interpretability implies relative consistency. The argument that  $\text{I}\Sigma_1 + \text{Con}(\text{PA})$  does not interpret  $\text{PA} + \text{Con}(\text{PA})$  will thus generalize in the right way.

We now turn to the question whether even the  $\Sigma_2$  induction axioms for  $\mathcal{S}$  are provable in  $\text{CTD}[\text{I}\Sigma_1] + \text{Tr}(\text{PA})$  (that is, whether  $\text{I}\Sigma_2^{\mathcal{S}}$  is a sub-theory thereof). We shall prove something stronger and more general:

**Theorem 3.**  $\text{CTD}[\text{I}\Sigma_n^{\mathcal{S}}] + \text{I}\Sigma_{n+1}^{\mathcal{A}}$  does not contain  $\text{I}\Sigma_{n+1}^{\mathcal{S}}$  as a sub-theory (for  $n > 0$ ).

The beginning of the argument is the same as before. Suppose otherwise. Then, since  $\text{CTD}[\text{I}\Sigma_n^{\mathcal{S}}] + \text{I}\Sigma_{n+1}^{\mathcal{A}}$  proves  $\text{Con}_{\mathcal{S}}(\text{I}\Sigma_{n+1}^{\mathcal{A}})$ , it also proves  $\text{Con}_{\mathcal{S}}(\text{I}\Sigma_{n+1}^{\mathcal{S}})$ , so it contains  $\text{I}\Sigma_{n+1}^{\mathcal{S}} + \text{Con}_{\mathcal{S}}(\text{I}\Sigma_{n+1}^{\mathcal{S}})$  as a sub-theory. Since  $\text{CTD}[\text{I}\Sigma_n^{\mathcal{S}}] + \text{I}\Sigma_{n+1}^{\mathcal{A}}$  is mutually interpretable with  $\text{I}\Sigma_n + \text{Con}(\text{I}\Sigma_{n+1})$ , it will suffice to establish:

<sup>15</sup>When the interpreted theory is not finitely axiomatizable, this becomes a bit more subtle, and may depend upon how the set of axioms is described, as Feferman (1960) famously shows. I'll ignore such subtleties here.

<sup>16</sup>It's most natural to take 'bounded' in a theory of concatenation to be defined in terms of substrings, as Halbach and Leigh (2024, §8.4) do. We can then define  $\Sigma_n$ , etc, in the usual way, and show that the syntactic theory with  $\Sigma_n$  induction is mutually interpretable with  $\text{I}\Sigma_n$ . This is because the most natural interpretation of arithmetic in syntax interprets numbers by strings like  $|$ ,  $\|$ , etc, and the usual coding of syntax in arithmetic has it that, if  $a$  is a substring of  $b$ , then the code of  $a$  is less than the code of  $b$ . (This is the much discussed 'monotonicity' of the usual codings (Heck, 2007; Halbach and Visser, 2014a,b; Grabmayr and Visser, 2021).)

**Lemma 4.**  $|\Sigma_{n+1} + \text{Con}(|\Sigma_{n+1}|)$  is not interpretable in  $|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)$ .

This is true despite the fact that  $|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)$  does interpret  $|\Sigma_{n+1}|$  itself, by the arithmetized completeness theorem. So, in a sense, this is best possible.

For the proof, we need the following important fact:

**Theorem 5** (Hájek and Pudlák, 1993, p. 106, Theorem 4.33).  $|\Sigma_n$  proves that the theory whose (arithmetical) axioms are the true  $\Pi_{n+1}$  sentences is consistent.

Here, the true  $\Pi_{n+1}$  sentences are characterized using a ‘partial’ truth-predicate for  $\Pi_{n+1}$  sentences: a formula  $\text{Tr}_{\Pi_{n+1}}(x)$  such that, if  $A$  is (no more complex than)  $\Pi_{n+1}$ ,  $|\Sigma_n$  (in fact,  $|\Sigma_1$ ) proves  $A \equiv \text{Tr}_{\Pi_{n+1}}(\ulcorner A \urcorner)$ .<sup>17</sup>

Theorem 5 is what allows us to show that  $|\Sigma_{n+1}$  proves  $\text{Con}(|\Sigma_n|)$ , the reason being that  $|\Sigma_n$  has a finite  $\Pi_{n+2}$  axiomatization.<sup>18</sup> And  $|\Sigma_{n+1}$  will prove that all of those axioms are true by the sort of trivial argument given above: Let  $A$  be one of the axioms; then  $A$  (since  $|\Sigma_n$  is a sub-theory of  $|\Sigma_{n+1}|$ ); but also  $A \equiv \text{Tr}_{\Pi_{n+1}}(\ulcorner A \urcorner)$ ; so  $\text{Tr}_{\Pi_{n+1}}(\ulcorner A \urcorner)$ . So  $|\Sigma_{n+1}$  proves that  $|\Sigma_n$  is a sub-theory of  $\text{Tr}(\Pi_{n+2})$ . But even  $\text{Q}$  knows that a sub-theory of a consistent theory is consistent. So  $|\Sigma_{n+1}$  proves that  $|\Sigma_n$  is consistent if  $\text{Tr}(\Pi_{n+2})$  is and so proves that  $|\Sigma_n$  is consistent.

The proof of Lemma 4 just extends this argument.

*Proof of Lemma 4.*  $\text{Con}(|\Sigma_{n+1}|)$  is a  $\Pi_1$  sentence, and  $|\Sigma_{n+1} + \text{Con}(|\Sigma_{n+1}|)$  proves that it is a true one, by the trivial argument. So  $|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)$  also has a  $\Pi_{n+2}$  axiomatization, and  $|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)$  proves that all those axioms are true. So, by the same reasoning just rehearsed,  $|\Sigma_{n+1} + \text{Con}(|\Sigma_{n+1}|)$  proves  $\text{Con}(|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)|)$ . So  $|\Sigma_n + \text{Con}(|\Sigma_{n+1}|)$  is not interpretable in  $|\Sigma_{n+1} + \text{Con}(|\Sigma_{n+1}|)$ .<sup>19</sup>  $\square$

So that completes the proof of Theorem 3.

It’s important to note that the theory mentioned in Theorem 5 is the one containing as axioms those  $\Pi_{n+1}$  sentences that  $|\Sigma_n$  *thinks* are true, not necessarily the ones that actually are true—though, since  $|\Sigma_n$  is sound, the sentences it thinks are true really are true.<sup>20</sup> To put it more

<sup>17</sup>For the details, see Hájek and Pudlák (1993, pp. 50–61).

<sup>18</sup>This is because the usual, non-finite axiomatization is  $\Pi_{n+2}$ —just check the complexity of the induction axioms—and the finitely many axioms can be chosen from among the usual axioms (Hájek and Pudlák, 1993, p. 78, Theorem 2.52). (In my courses on these matters, I often assign the following exercise: If some infinitely axiomatized theory  $\mathcal{T}$  has a finite axiomatization, then it has a finite axiomatization all of whose axioms are among those of  $\mathcal{T}$ .)

<sup>19</sup>This is really a special case of Hájek and Pudlák’s Corollary 4.34(ii). But it is worth filling in some of the details they omit.

<sup>20</sup>Moreover, in the standard model, the extension of  $\text{Tr}_{\Pi_n}(x)$  is exactly the set of true  $\Pi_n$  sentences.

formally, what Theorem 5 says is that  $\mathcal{I}\Sigma_n$  proves: There is no sequence of formulas, ending with  $0 = 1$ , each of which is either a true  $\Pi_{n+2}$  sentence or a logical axiom or else follows from earlier members of the sequence by one of the rules of inference. The proof of Lemma 4 depends upon the fact that  $\mathcal{I}\Sigma_{n+1} + \text{Con}(\mathcal{I}\Sigma_{n+1})$  proves that  $\mathcal{I}\Sigma_n + \text{Con}(\mathcal{I}\Sigma_{n+1})$  is a sub-theory of the theory whose axioms are the true  $\Pi_{n+2}$  sentences. It is irrelevant whether that claim is itself true.

This allows us to establish a generalization of Lemma 4 that is worth stating separately:

**Corollary 6.** *Let  $A$  be a  $\Pi_{n+2}$  sentence that is consistent with  $\mathcal{I}\Sigma_n$ . Then  $\mathcal{I}\Sigma_{n+1} + A$  is not interpretable in  $\mathcal{I}\Sigma_n + A$ .*

*Proof.*  $\mathcal{I}\Sigma_n + A$  is  $\Pi_{n+2}$  axiomatized, and  $\mathcal{I}\Sigma_{n+1} + A$  proves that all its axioms are true. So  $\mathcal{I}\Sigma_{n+1} + A$  proves  $\text{Con}(\mathcal{I}\Sigma_{n+1} + A)$ .  $\square$

We can be confident, then, that ‘disentangling’ does what it is supposed to do: Not only do the semantic induction axioms not inadvertently strengthen the target theory, but assuming the truth of the induction axioms in the target theory does not ‘transfer’ back into the syntactic theory, either. That, however, seems about as much as we are likely to be able to say:  $\text{CTD}[\mathcal{I}\Sigma_n] + \mathcal{I}\Sigma_n$  is not a conservative extension of  $\mathcal{I}\Sigma_n$ , or even of  $\text{CTD}[\mathcal{I}\Sigma_n]$ , not even for  $\Pi_1$  sentences: The whole point is that  $\text{CTD}[\mathcal{I}\Sigma_n] + \mathcal{I}\Sigma_n$  proves the purely syntactic  $\Pi_1$  sentence  $\text{Con}_{\mathcal{I}}(\mathcal{I}\Sigma_{n+1})$ , which is not provable in  $\mathcal{I}\Sigma_n$  itself (unless, of course,  $\mathcal{I}\Sigma_n$  is inconsistent).<sup>21</sup>

#### REFERENCES

- Feferman, S. (1960). ‘Arithmetization of metamathematics in a general setting’, *Fundamenta Mathematicae* 49: 35–92.
- Grabmayr, B. and Visser, A. (2021). ‘Self-reference upfront: A study of self-referential Gödel numberings’, *The Review of Symbolic Logic*. Forthcoming.
- Hájek, P. and Pudlák, P. (1993). *Metamathematics of First-order Arithmetic*. New York, Springer-Verlag.
- Halbach, V. and Leigh, G. E. (2024). *The Road to Paradox: A Guide to Syntax, Truth, and Modality*. Cambridge, Cambridge University Press. Forthcoming.
- Halbach, V. and Visser, A. (2014a). ‘Self-reference in arithmetic I’, *The Review of Symbolic Logic* 7: 671–91.
- (2014b). ‘Self-reference in arithmetic II’, *The Review of Symbolic Logic* 7: 692–712.
- Heck, R. K. (2007). ‘Self-reference and the languages of arithmetic’, *Philosophia Mathematica* 15: 1–29. Originally published under the name “Richard G. Heck, Jr”.

<sup>21</sup>Thanks to Volker Halbach, Graham Leigh, and Carlo Nicolai for helping me sort all this out, and to Matteo Zicchetti for inspiring me to do so.



- (2009). ‘The strength of truth-theories’. Originally published under the name “Richard G. Heck, Jr”.
- (2015). ‘Consistency and the theory of truth’, *Review of Symbolic Logic* 8: 424–66. Originally published under the name “Richard G. Heck, Jr”.
- (2018). ‘The logical strength of compositional principles’, *Notre Dame Journal of Formal Logic* 59: 1–33. Originally published under the name “Richard G. Heck, Jr”.
- Leigh, G. E. and Nicolai, C. (2013). ‘Axiomatic truth, syntax and metatheoretic reasoning’, *The Review of Symbolic Logic* 6: 613–36.
- Łełyk, M. Z. (2022). ‘Model theory and proof theory of the global reflection principle’, *The Journal of Symbolic Logic* 88: 738–79.
- Pudlák, P. (1985). ‘Cuts, consistency statements and interpretations’, *Journal of Symbolic Logic* 50: 423–41.
- Tarski, A. (1956). ‘The concept of truth in formalized languages’, in J. Corcoran (ed.), *Logic, Semantics, and Metamathematics*. Indianapolis, Hackett, 152–278.
- Visser, A. (2014). ‘The interpretability of inconsistency: Feferman’s theorem and related results’. Preprint.