

Review of: "Nonlinearity and Illfoundedness in the Hierarchy of Large Cardinal Consistency Strength"

Sakaé Fuchino¹

¹ Kobe University

Potential competing interests: No potential competing interests to declare.

Joel Hamkins' NONLINEARITY AND ILLFOUNDEDNESS IN THE HIERARCHY OF LARGE CARDINAL CONSISTENCY STRENGTH is really an inspiring read. As Peter Holly also writes in his review, the second part of the paper from section 7 on is highly recommended for a broad audience in mathematics and philosophy of science who are interested in the role of modern research of set theory, in connection with the foundational questions in mathematics in particular.

This being said, I have to confess my enormous difficulty in the technical details of the first part of the paper. At the moment the difficulty is not yet totally resolved. The following is thus merely a primary report of my struggle with this paper. I shall update this report later with more comments when I obtain better understanding of the material. Meanwhile I hope this version of the report is helpful for those who have difficulties similar to mine in this paper.

My understanding of statements about relative consistency in set theory is that they should be (at least in the end) formulated in metamathematics.

In meta-mathematics there is no semantics. We can only talk about provability which is expressed in statements like “there is/is not a (concretely given) proof p such that $\phi \vdash \ulcorner p \urcorner$ ” where ϕ here is some base theory (e.g. some fragment of a second order extension of PA or some weak set theory) in which we would like to code the logic. ϕ may be chosen such that we can also formulate the notion of models, model relation etc., and prove Completeness and Incompleteness Theorems (as theorems in ϕ). In such theory ϕ , we have at least four different notions of “truth”: ① : “ ϕ holds” ($\phi \vdash \phi$), ② : “(ϕ thinks) ϕ is provable in ϕ ” ($\phi \vdash \ulcorner \ulcorner \phi \urcorner \urcorner \vdash \ulcorner \phi \urcorner$), ③ : “(ϕ thinks) ϕ holds in a model \mathcal{M} ” ($\phi \vdash \ulcorner \ulcorner \phi \urcorner \urcorner \vdash \ulcorner \phi \urcorner$), ④ : “(ϕ thinks) thinks ϕ is provable in ϕ ” ($\phi \vdash \ulcorner \ulcorner \ulcorner \phi \urcorner \urcorner \urcorner \vdash \ulcorner \phi \urcorner$). Actually we can easily extend this list ad infinitum by adding “(ϕ thinks) a model \mathcal{M} thinks that ϕ holds in a model N_{ϕ} ” ($\phi \vdash \ulcorner \ulcorner \ulcorner \ulcorner \phi \urcorner \urcorner \urcorner \vdash \ulcorner \phi \urcorner$),... etc.

Cohen's result saying “ZFC + \neg CH is consistent”, for example, is to be understood as a meta-mathematical statement: “If ZF is consistent then ZFC + \neg CH is also consistent”. More precisely, this should mean that there is a mechanical procedure with which a given proof p of inconsistency in ZFC + \neg CH (if such a p ever exists) can be recast into a proof p' of inconsistency in ZF. Not all experts of logic are aware of this. I remember that now almost 20 years ago when I gave an introductory tutorial on forcing, a prominent Japanese proof theorist, who was among the audience, could not believe this statement and I had to explain how to see it in length.

Now, a statement like

“ $PA + Con(S) + \neg Con(T)$ is consistent”

should mean in metamathematics that there is no (concretely given) proof for the (concretely given) theory $PA + Con(\ulcorner \ulcorner \urcorner) + \neg Con(\ulcorner \ulcorner \urcorner)$ where \ulcorner and \urcorner are concretely given theories, both of which extend a concretely given theory \ulcorner for which the consistency of $\ulcorner + Con(\ulcorner \urcorner)$, or sometimes even the consistency of $\ulcorner + Con(\ulcorner + Con(\ulcorner \urcorner))$ is assumed and from which at least the consistence of $PA + Con(\ulcorner \urcorner)$ and that of $PA + \neg Con(\ulcorner \urcorner)$ follow.

If I interpret it correctly, the same statement in Hamkins' narrations means “ $\mathfrak{o}' \vdash Con(\ulcorner PA \urcorner + Con(\ulcorner S \urcorner) + \neg Con(\ulcorner T \urcorner))$ ” where \mathfrak{o}' is some strong enough extension of the base theory which may be assumed to be consistent (I say “assumed” since I can not say any more due to the Second Incompleteness Theorem). Working in \mathfrak{o}' he then uses the Completeness Theorem (as a theorem in \mathfrak{o}') and take a model \mathfrak{M} of $\ulcorner PA \urcorner + Con(\ulcorner S \urcorner) + \neg Con(\ulcorner T \urcorner)$ in \mathfrak{o}' (note that $\ulcorner PA \urcorner$ here is not the PA in \mathfrak{M} but what \mathfrak{o}' think is $\ulcorner PA \urcorner$).

It seems, when a conclusion like “ \ast is consistent” is obtained in Hamkins' setting, what is actually attained is sometimes not the statement “ \ast is consistent” in metamathematics but rather “ $\mathfrak{o}' \vdash Con(\ulcorner \ast \urcorner)$ ”. In many cases, this creates no problem for metamathematical consideration because of the following:

Lemma 1. If $\mathfrak{o} \vdash Con(\ulcorner \mathfrak{1} \urcorner)$, then $\mathfrak{1}$ is consistent (in metamathematics) assuming that \mathfrak{o} is consistent.

Proof. Suppose $\mathfrak{1}$ is not consistent. This means that there is a proof $\ulcorner \ulcorner \urcorner$ such that $\mathfrak{o} \vdash \ulcorner \ulcorner \urcorner \equiv 1$. This can be translated to $\mathfrak{o} \vdash \ulcorner \ulcorner \urcorner \vdash \ulcorner \urcorner \equiv 1$.

By the assumption of $\mathfrak{o} \vdash Con(\ulcorner \mathfrak{1} \urcorner)$, it follows that $\mathfrak{o} \vdash \ulcorner \ulcorner \urcorner \equiv 1$. This is a contradiction to the assumption of consistency of \mathfrak{o} . \square

Even though we can reinterpret Hamkins' arguments as a corresponding metamathematical statement via Lemma 1 above, some details of his proofs remain extremely difficult to understand for me. To explain my difficulty, let me try to analyze the following paragraph from the proof of Theorem 2:

[The last but one paragraph of the proof of Theorem 2.]: “Since σ is not refutable, it follows that $PA + Con(PA) + \sigma$ is consistent, and so it is also consistent with the assertion of its own inconsistency $\neg Con(PA + Con(PA) + \sigma)$. In any model of this combined theory, σ is refutable in $PA + Con(PA)$, but since also σ is true there, there must not be any smaller refutation of τ . Since this syntactic situation will be provable in PA, it follows in light of what the sentences assert that the model thinks that PA proves that σ is true and τ is false. So from $Con(PA)$ it follows both that $Con(PA + \sigma)$ and $\neg Con(PA + \tau)$ in this model.”

[My reading of the paragraph]: We work in \mathfrak{o} . Where I assume that \mathfrak{o} implies $Con(PA + Con(PA))$. Since σ is not refutable in $PA + Con(PA)$ (i.e. $PA + Con(PA) \not\vdash \neg \sigma$), $PA + Con(PA) + \sigma$ is consistent. By the Second Incompleteness Theorem (formulated in \mathfrak{o}), it follows that $\ast := PA + Con(PA) + \sigma + \neg Con(PA + Con(PA) + \sigma)$ is consistent. By the Completeness

Theorem, there is a model \models^* . By $\models \neg \text{Con}(\text{PA} + \text{Con}(\text{PA}) + \sigma)$, we have $\models \text{PA} + \text{Con}(\text{PA}) \not\models \sigma$. By the choice of σ (and since $\models \text{PA}$), this means

$$\models \text{PA} + \text{Con}(\text{PA}) \not\models \exists p(\text{PA} + \text{Con}(\text{PA}) \not\models p \rightarrow \tau \wedge \forall q(p(\text{PA} + \text{Con}(\text{PA}) \not\models q \rightarrow \sigma)) \dots$$

My reading has been stuck at this point and I could go any further from here. First when I learned from Taishi Kurahashi that corresponding formal proof (in meta-mathematics) uses the formalized Σ_1 -completeness, I came to the following alternative proof. Here the Σ_1 -completeness is the following statement.

Proposition $\aleph 2$. (Formalized Σ_1 -Completeness) For an arithmetical Σ_1 -sentence θ , we have

$$\text{PA}' \vdash \ulcorner \ulcorner \neg \urcorner \urcorner \vdash \ulcorner \theta \urcorner \rightarrow (\ulcorner \ulcorner \neg \urcorner \urcorner \vdash \ulcorner \theta \urcorner). \quad \square$$

Note that the above can be reformulated as

$$\text{PA}' \vdash \ulcorner \ulcorner \neg \urcorner \urcorner \vdash \ulcorner \theta \urcorner \rightarrow \neg \text{Con}(\ulcorner \ulcorner \neg \urcorner \urcorner + \ulcorner \neg \theta \urcorner).$$

PA' denotes here a theory which might be slightly stronger than PA with a second-order part which is strong enough accommodate a reasonable model theory. It is often assumed that PA satisfies this Σ_1 -Completeness. But for me this assumption seems to entail a bit more than the very strict variant finitary standpoint in metamathematics. Proposition $\aleph 2$ is a theorem for a weak second-order arithmetical system which has the same Σ_1 part as PA. However, this Σ_1 equivalence is established using model theoretic method which seems to exceed the very strict version of finitary standpoint. This is one of the things I am not yet quite sure at the moment and on which I am in an on-going discussion with Hiroshi Sakai. For simplicity, however we assume in the following $\text{PA}' = \text{PA}$ also satisfies Proposition $\aleph 2$.

[An alternative to my reading of the paragraph]: We consider

$$* := \text{PA} + \text{Con}(\text{PA}) + \tau + \neg \text{Con}(\text{PA} + \text{Con}(\text{PA}) + \tau)$$

(Note that we use τ in place of σ). Working in $*$, $\neg \text{Con}(\text{PA} + \text{Con}(\text{PA}) + \tau)$ means $\text{PA} + \text{Con}(\text{PA}) \not\models \neg \tau$. Note that $\neg \tau$ is a Σ_1 -formula by definition of τ given in the paper. By the Formal Σ_1 -Completeness, it follows that $\text{PA} + \text{Con}(\text{PA}) \not\models \neg \text{Con}(\text{PA} + \tau)$ (If we switch in the proof of the Formal Σ_1 -Completeness here, a model-theoretic argument is deployed at this point).

The previous paragraphs in Hamkins' proof of Theorem 2 translates to

$\text{PA} + \text{Con}(\text{PA} + \text{Con}(\text{PA})) \not\models \text{Con}(\text{PA} + \sigma)$. Thus, by Lemma $\aleph 1$ above, we can conclude that $\text{PA} + \text{Con}(\text{PA} + \sigma) + \neg \text{Con}(\text{PA} + \tau)$ is consistent. \square

Theorems 3,4 also should be able to be treated in this manner.

At the moment I cannot yet correctly work out Theorem 18 which deals with what is called "cautious enumeration" of a

theory. In the proof of the theorem, finite subsets of a rearranged enumeration of a theory are treated. If we talk about a finite fragment of a theory T , the finiteness here can have several different meanings. It can mean concretely given finite collection of formulas S in the sense of metamathematics but it can also mean finite set of the set $\ulcorner S \urcorner$ where $\ulcorner S \urcorner$ is the set defined in \mathcal{M} by the same recursive definition as the one given in metamathematics to decide which formula belongs to S .

The finiteness here is what the base theory \mathcal{M} thinks is finite. If we are working in the base theory \mathcal{M} and consider a model \mathcal{M}' of some set theory there, then $\ulcorner S \urcorner^{\mathcal{M}'}$ and $\omega^{\mathcal{M}'}$ can be totally different from what \mathcal{M} thinks are $\ulcorner S \urcorner$ and ω , since \mathcal{M}' can contain nonstandard numbers and formulas.

This subtle distinction seems to be quite relevant here since it seems that Lévy-Montague Reflection Theorem is applied in the proof. The Lévy-Montague Reflection Theorem is famous for easily producing apparently correct proof of the inconsistency of the set theory if we are sloppy with the fine difference between these notions of finiteness.

There are still many other issues I want talk about in connection with this paper but I shall write them in the next update of this review.