

## Lecture 26: Consistency statements

Tin Lok Wong

15 November, 2018

We conclude the course by pointing out a connection between proofs of finite consistency statements and complexity theory.

From the Incompleteness Theorems, one learns that consistency statements are in general hard to prove. In the case of the Second Incompleteness Theorem, ‘hard to prove’ simply means ‘unprovable’, while in the case of the Finite Incompleteness Theorem, ‘hard to prove’ means ‘not polynomial-length provable’. One naturally asks where the limit of this incompleteness phenomenon is. For instance, can one improve ‘polynomial’ to ‘exponential’ in the Finite Incompleteness Theorem? As mentioned briefly in the previous lecture, the answer is in general no. Let us state this more properly here. Recall the definition of  $\text{PA}^-$  from Lemma 11.1.

**Theorem 26.1** (Pudlák). Let  $T$  be a finite  $\mathcal{L}_A(\text{exp})$  theory extending  $\text{PA}^-$ . Then one can construct a  $\Delta_0(\text{exp})$  formula  $\Box^{\leq}(w, y)$  with some  $k \in \mathbb{N}$  such that

- (1)  $\{(n, m) \in \mathbb{N}^2 : \mathbb{N} \models \Box^{\leq}(n, m)\} = \{(n, \ulcorner \theta \urcorner) \in \mathbb{N}^2 : \theta \text{ is an } \mathcal{L}_A(\text{exp}) \text{ formula and } T \vdash^n \theta\}$ ;
- (2)  $\Box^{\leq}(w, y)$  satisfies the derivability conditions (M), (N), (IN), and ( $\Box$ D) listed in the statement of Theorem 25.7 for all  $n \in \mathbb{N}$  and all  $\mathcal{L}_A(\text{exp})$  sentences  $\sigma, \tau$ ; and
- (3)  $T \vdash^{\frac{2^{k|n|}+k}{|n|}} \neg \Box^{\leq}(n, \perp)$  for all  $n \in \mathbb{N}$ . □

Although Theorem 26.1 says every finite theory  $T \supseteq \text{PA}^-$  has exponential-length proofs of its own finite consistencies, it does not exclude the existence of such a theory  $T$  which does not have exponential-length proofs of finite consistencies of a *stronger theory*  $T^*$ . In view of the general impression that consistency statements are hard to prove, one may actually expect every consistent such theory  $T$  to *not* have exponential-length proofs of finite consistencies of some sufficiently strong theory  $T^* \supseteq T$ . Given a finite consistent theory  $T$ , the most prominent example of a strictly stronger (but not necessarily consistent) theory is  $T^* = T + \text{Con}(T)$  by the Second Incompleteness Theorem, where  $\text{Con}(T)$  is some sentence expressing the consistency of  $T$  reasonably. This leads to the following conjecture. To avoid introducing new notation, we formulate this conjecture entirely in terms of the finite provability predicate  $\Box^{\leq}(w, y)$  for  $T$ . In particular, we write  $\text{Con}(T) = \forall w \neg \Box^{\leq}(w, \perp)$ , and  $T + \text{Con}(T) \not\vdash \perp$  is rewritten as  $T \not\vdash \neg \text{Con}(T)$  via (RAA') and ( $\perp$ ).

**Conjecture 26.2** (Pudlák). Fix a finite consistent  $\mathcal{L}_A(\text{exp})$  theory  $T \supseteq \text{PA}^-$ . Let  $\Box^{\leq}(w, y)$  be an  $\mathcal{L}_A(\text{exp})$  formula satisfying both (1) and (2) in the statement of Theorem 26.1. Then for every  $k \in \mathbb{N}$ , there is  $n \in \mathbb{N}$  such that

$$T \vdash^{\frac{2^{k|n|}+k}{|n|}} \neg \Box^{\leq}(n, \ulcorner \forall w \neg \Box^{\leq}(w, \perp) \urcorner).$$

As Pudlák himself observed, his conjecture implies  $\text{NP} \neq \text{co-NP}$ , and thus also  $\text{P} \neq \text{NP}$ . Since it is generally believed that radically new techniques have to be developed to prove these inequations in complexity theory, one does not expect to see a proof of Conjecture 26.2 in the near future. Therefore, realistically, to make progress, one considers either weaker conjectures which are easier to prove, or stronger conjectures which are easier to refute. The rest of this lecture is devoted to an example of the latter. More precisely, we claim that if  $\text{Con}(T) = \forall w \neg \Box^{\leq}(w, \perp)$  is replaced by an arbitrary  $T$ -unprovable sentence, then the conjecture becomes false in general.

In addition to the derivability conditions listed in Theorem 25.7, we will employ three properties of the usual finite provability predicate. The first one is Theorem 26.1. The second and the third ones are the following duals of (N) and (IN).

**Lemma 26.3.** Provided  $T \vdash \text{I}\Delta_0(\text{exp})$ , one can add to the list of conditions in Theorem 26.1 the following clause: there are  $p(X), q(X) \in \mathbb{N}[X]$  such that for all  $n \in \mathbb{N}$  and all  $\mathcal{L}_A(\text{exp})$  formulas  $\theta$ ,

- if  $T \vdash^n \theta$ , then  $T \vdash^{\underline{p(2^n)}} \neg \square^{\leq}(\underline{n}, \theta)$ ; and
- $T \vdash^{\underline{q(n+\text{len}(\theta))}} \neg \square^{\leq}(\underline{n}, \theta) \rightarrow \square^{\leq}(\underline{p(2^n)}, \neg \square^{\leq}(\underline{n}, \theta))$ .

*Proof sketch.* Very carefully study the actual shape of the  $\Delta_0(\text{exp})$  formula  $\square^{\leq}(w, y)$  and the proof of formalized  $\Sigma_1$  completeness, i.e., Theorem 12.6. The  $2^n$  comes from the exponential number of proofs of length less than  $n$ . The details of the verification are cumbersome and so are omitted here.  $\square$

To establish the claim, we will need to produce a sentence  $\sigma$  such that  $T + \sigma$  is strictly stronger than the consistent theory  $T$  we are given. Recall from Theorem 26.1 that  $T$  has exponential-length proofs of its own finite consistencies. If  $T + \sigma$  is not much stronger than  $T$ , then intuitively the finite consistencies of  $T + \sigma$  are not much harder to prove than those of  $T$ , and thus they should have exponential-length proofs in  $T$  too. This, in a sense, explains why the witness comparison method is employed in the construction of our  $\sigma$  below. The following theorem can be seen as a formalization of First Incompleteness Theorem with explicit bounds on proof lengths. We modify the self-referential sentence slightly so as to keep the lengths of the required proofs exponential. Note that a polynomial in  $n$  is an exponential in  $|n|$ .

**Notation.** Write  $T \vdash_*^n S(n)$  and  $T \vdash_*^{2^n} S(n)$  respectively for

$$\exists t(X) \in \mathbb{N}[X] \quad \forall n \in \mathbb{N} \quad T \vdash^{t(n)} S(n) \quad \text{and} \quad \exists t(X) \in \mathbb{N}[X] \quad \forall n \in \mathbb{N} \quad T \vdash^{t(2^n)} S(n).$$

**Theorem 26.4** (Hrubeš). Fix an  $\mathcal{L}_A(\text{exp})$  theory  $T \vdash \text{I}\Delta_0(\text{exp})$ . Let  $\square^{\leq}(w, y)$  be an  $\mathcal{L}_A(\text{exp})$  formula and  $p(X), q(X) \in \mathbb{N}[X]$  such that for all  $n \in \mathbb{N}$  and all  $\mathcal{L}_A(\text{exp})$  sentences  $\sigma, \tau$ ,

- (M)  $T \vdash \forall w, y (\square^{\leq}(w, y) \rightarrow \forall w' \geq w \square^{\leq}(w', y))$ ;
- (N) If  $T \vdash^n \sigma$ , then  $T \vdash^{\underline{p(n)}} \square^{\leq}(\underline{n}, \sigma)$ ;
- (IN)  $T \vdash^{\underline{q(|n|+\text{len}(\sigma))}} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(\underline{p(n)}, \square^{\leq}(\underline{n}, \sigma))$ ;
- (□D)  $T \vdash^{\underline{q(|n|+\text{len}(\sigma)+\text{len}(\tau))}} \square^{\leq}(\underline{n}, \sigma \rightarrow \tau) \rightarrow (\square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(\underline{p(n)}, \tau))$ ;
- (C)  $T \vdash^{\underline{p(n)}} \neg \square^{\leq}(\underline{n}, \perp)$ ;
- (¬N) If  $T \vdash^n \sigma$ , then  $T \vdash^{\underline{p(2^n)}} \neg \square^{\leq}(\underline{n}, \sigma)$ ; and
- (¬IN)  $T \vdash^{\underline{q(n+\text{len}(\sigma))}} \neg \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(\underline{p(2^n)}, \neg \square^{\leq}(\underline{n}, \sigma))$ .

Suppose  $\sigma$  is an  $\mathcal{L}_A(\text{exp})$  sentence which satisfies

$$T \vdash \sigma \leftrightarrow \neg \exists w (\square^{\leq}(w, \sigma) \wedge \neg \square^{\leq}(2^w, \neg \sigma)).$$

Then there are non-constant  $r(X), s(X) \in \mathbb{N}[X]$  such that for all  $n \in \mathbb{N}$ , the following hold.

- (1) If  $T \vdash^{\underline{r(2^{2^n})}} \perp$ , then  $T \vdash^n \sigma$ .
- (2) If  $T \vdash^{\underline{r(n)}} \perp$ , then  $T \vdash^n \neg \sigma$ .
- (3)  $T \vdash^{\underline{s(2^n)}} \neg \square^{\leq}(\underline{r(2^{2^n})}, \perp) \rightarrow \neg \square^{\leq}(\underline{n}, \sigma)$ .

$$(4) T \frac{|s(n)|}{*} \neg \square^{\leq}(r(n), \perp) \rightarrow \neg \square^{\leq}(\underline{n}, \neg \sigma).$$

*Proof.* Let us introduce the function  $x \mapsto |x|$  to  $\mathcal{L}_A(\text{exp})$  in a way similar to how we introduced the function  $\text{len}$  in Lecture 10. The usual properties of the function  $x \mapsto |x|$  can routinely be proved within  $\text{ID}_0(\text{exp})$ ; cf. Lecture 11. It is a straightforward exercise involving numerous applications of (M), (N) and ( $\square$ D) to show that, by choosing a larger  $p(X)$  if necessary, one can make the following hold.

- (i)  $T \frac{|n|}{*} \square^{\leq}(p(n), \theta(\underline{n})) \wedge \square^{\leq}(p^2(n), \neg \theta(\underline{n})) \rightarrow \square^{\leq}(p^3(n), \perp)$  for all  $\mathcal{L}_A(\text{exp})$  formulas  $\theta(x)$ .
- (ii)  $T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \rightarrow \square^{\leq}(p^2(n), \exists w (\square^{\leq}(w, \sigma) \wedge \neg \square^{\leq}(2^w, \neg \sigma)))$ .
- (iii)  $T \frac{|n|}{*} \square^{\leq}(p(n), \exists w (\square^{\leq}(w, \sigma) \wedge \neg \square^{\leq}(2^w, \neg \sigma))) \wedge \square^{\leq}(p^2(n), \square^{\leq}(\underline{n}, \neg \sigma))$   
 $\rightarrow \square^{\leq}(p^3(n), \square^{\leq}(|n|, \sigma))$ .

The proofs of (3) and (4) are essentially formalizations of those of (1) and (2). So we concentrate on (3) and (4).

(3) On the one hand,

$$\begin{aligned} & T \frac{|n|}{*} \square^{\leq}(\underline{n}, \sigma) \wedge \square^{\leq}(2^n, \neg \sigma) \rightarrow \square^{\leq}(p^3(2^n), \perp) \quad \text{by (i);} \\ \therefore & T \frac{2^n}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \neg \square^{\leq}(2^n, \neg \sigma) \quad \text{as } T \frac{p^4(2^n)}{*} \neg \square^{\leq}(p^3(2^n), \perp) \text{ by (C);} \\ \therefore & T \frac{2^n}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(p(2^{2^n}), \neg \square^{\leq}(2^n, \neg \sigma)) \quad \text{by } (\neg\text{IN}). \end{aligned}$$

On the other hand, condition (IN) implies

$$T \frac{|n|}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(p(n), \square^{\leq}(\underline{n}, \sigma)).$$

Combining the two, we deduce that

$$\begin{aligned} & T \frac{2^n}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(r(2^{2^n}), \square^{\leq}(\underline{n}, \sigma) \wedge \neg \square^{\leq}(2^n, \neg \sigma)) \quad \text{for some } r(X) \in \mathbb{N}[X]; \\ \therefore & T \frac{2^n}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(r(2^{2^n}), \neg \sigma) \quad \text{for some } r(X) \in \mathbb{N}[X], \\ & \quad \text{by the choice of } \sigma, \text{ plus (N) and } (\square\text{D}); \\ \therefore & T \frac{2^n}{*} \square^{\leq}(\underline{n}, \sigma) \rightarrow \square^{\leq}(r(2^{2^n}), \perp) \quad \text{for some } r(X) \in \mathbb{N}[X], \text{ by (i).} \end{aligned}$$

(4) By (IN), we know  $T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \rightarrow \square^{\leq}(p(n), \square^{\leq}(\underline{n}, \neg \sigma))$ . Combining this with (ii) using (iii), we deduce that

$$T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \rightarrow \square^{\leq}(p^4(n), \square^{\leq}(|n|, \sigma)).$$

Consequently, since  $T \frac{|n|}{*} \neg \square^{\leq}(|n|, \sigma) \rightarrow \square^{\leq}(p(n), \neg \square^{\leq}(|n|, \sigma))$  by ( $\neg$ IN),

$$\begin{aligned} & T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \wedge \neg \square^{\leq}(|n|, \sigma) \rightarrow \square^{\leq}(p^6(n), \perp) \quad \text{by (i);} \\ \therefore & T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \rightarrow \square^{\leq}(|n|, \sigma) \quad \text{as } T \frac{p^7(n)}{*} \neg \square^{\leq}(p^6(n), \perp) \text{ by (C);} \\ \therefore & T \frac{|n|}{*} \square^{\leq}(\underline{n}, \neg \sigma) \rightarrow \square^{\leq}(p^3(n), \perp) \quad \text{by (i).} \quad \square \end{aligned}$$

From Theorem 26.4, one can easily derive what we claimed on page 126.

**Corollary 26.5** (Hrubeš). Let  $T$  be any consistent  $\mathcal{L}_A(\text{exp})$  theory extending  $\text{ID}_0(\text{exp})$ . Given an  $\mathcal{L}_A(\text{exp})$  formula  $\square^{\leq}(w, y)$  and  $p(X), q(X) \in \mathbb{N}[X]$  satisfying (M)–( $\neg$ IN) in the statement of Theorem 26.4 for all  $n \in \mathbb{N}$  and all  $\mathcal{L}_A(\text{exp})$  sentences  $\sigma, \tau$ , one can construct an  $\mathcal{L}_A(\text{exp})$  sentence  $\sigma$  such that

- $T \not\vdash \sigma$  and  $T \not\vdash \neg\sigma$ ; and
- $\exists k \in \mathbb{N} \forall n \in \mathbb{N} T \vdash_{2^{k|n|+k}} \neg \Box^{\leq}(n, \neg\sigma)$ .

*Proof.* Apply the Diagonal Lemma to find an  $\mathcal{L}_A(\text{exp})$  sentence  $\sigma$  such that

$$T \vdash \sigma \leftrightarrow \neg \exists w (\Box^{\leq}(w, \sigma) \wedge \neg \Box^{\leq}(2^w, \neg\sigma)).$$

Let  $r(X), s(X) \in \mathbb{N}[X]$  given by an application of Theorem 26.4. As  $T$  is consistent, we see from Theorem 26.4(1,2) that  $T \not\vdash \sigma$  and  $T \not\vdash \neg\sigma$ . Moreover, since (C) implies  $T \vdash_{p(r(n))} \neg \Box^{\leq}(r(n), \perp)$  for all  $n \in \mathbb{N}$ , we know  $T \vdash_{\frac{n}{*}} \neg \Box^{\leq}(n, \neg\sigma)$  by Theorem 26.4(4).  $\square$

Although Corollary 26.5 produces a sentence  $\sigma$  such that  $T$  has exponential-length proofs of the finite consistencies of  $T + \sigma$ , it does not guarantee the same for  $T + \neg\sigma$ . In fact, Theorem 26.4 only gives a triply-exponential length-bound for  $T$ -proofs of the finite consistencies of  $T + \neg\sigma$ . Using a different variation of the witness comparison method, Jeřábek managed to produce another such  $\sigma$  such that  $T$  also has exponential-length proofs of the finite consistencies of  $T + \neg\sigma$ , thus restoring the symmetry between  $T + \sigma$  and  $T + \neg\sigma$  in the usual proof of the First Incompleteness Theorem.