

# How to Derive Principles of Interpretability Logic A Toolkit

Joost J. Joosten  
and  
Albert Visser

September 4, 2004

## Abstract

In this paper, we develop a toolkit to derive principles for the interpretability logic of all reasonable theories by fine-tuning **ILM** and **ILP** proofs.

## 1 Introduction

### 1.1 Preludium

When, around 1987, the project of Interpretability Logic started, Dick de Jongh, in collaboration with Frank Veltman, took responsibility for the modal part of the project. The first fruit of this research was the fundamental paper [dJV90]. From this time on, Dick actively worked on the subject. His work can be found in such papers as [dJV91], [AdJH98], [dJV99], [dJJ98]. We always enjoyed his remarkable style of doing modal logic, where visualization of models plays such an important role.

Dick was not so active on the arithmetical side of the project. For this reason, it seemed a good idea to show that arithmetical soundness proofs can be given using modal methods. Perhaps, this will draw him into arithmetical matters again.

In this paper, we use enriched versions of interpretability logic to verify arithmetically sound principles of ordinary interpretability logic. The general methodology of using modal logic to study modal logic has always been part of the strongly self-reflexive provability logic tradition. We just remind the reader of the seminal paper [JJM91] in which the Solovay result itself is proved in a modal logic.

The logics we study arise from well known logics by ‘diversifying’ the modal operators. In another contribution to this Liber, Johan van Benthem discusses the idea of scattering. It seems to us that the auxiliary logics presented in this paper are typical examples of such scattered logics. This either proves the coherence of the Dutch tradition or the naturalness of the idea.

## 1.2 Contents of the paper

In this paper we present and explore a remarkable methodology. We want to derive principles for the interpretability logic of all reasonable theories. This logic is a sublogic both of **ILM** and **ILP**, but it is not the intersection of **ILM** and **ILP**. (See [Vis97].) It turns out that one may derive principles for the logic of all theories by fine-tuning proofs in **ILM** or **ILP**. In fact, *every principle we know* can be derived *both* by fine-tuning an **ILM** and an **ILP** proof.

The fine-tuning procedure is best mediated by auxiliary modal logics that we present in Section 2 (for the **ILM**-case) and in Section 3 (for the **ILP**-case). In Section 4, we put the methods developed in action, and show how to derive a number of principles.

The phenomenon of having two different proofs for the same theorem is always strange. In some cases different proofs reflect different underlying concepts (like the two proofs of the commutativity of addition). In our case the strangeness is increased: we do not just have two proofs, but two natural classes of proofs for the same theorems. No different underlying concepts are in sight. The only clear difference is a difference in *scope*: the **M**-style proofs use sequentiality and the **P**-style proofs do not.

The explicit development of modal logics to analyze the methods known so far to derive interpretability principles, has the following aim. We hope that it will enable us to pin down precisely what principles can be derived using **M**-style methods and which ones in the **P**-style. The question which interpretability principles can be derived in **P**-style and which in **M**-style becomes a modal question open to study with Kripke models. Apart from this aim, we submit that we found exciting and intriguing modal systems.

At the present stage, it is not yet fully clear that our extended systems are definitive. Do we have the right notions? Did we articulate all possible principles for the chosen notions? Do we have to extend or restrict the expressive power? More experimentation is necessary.

## 1.3 Convention

Interpretability will in this paper be *theorems interpretability*, i.o.w.

- $k : U \triangleright V : \iff \forall \phi (\Box_V \phi \rightarrow \Box_U \phi^k)$ .

## 2 A logic for relativization to cuts

In this section, we present a logic that incorporates a number of principles concerning provability predicates relativized to definable cuts. This logic will enable us to fine-tune **ILM**-proofs.

In the present section, theories will be  $\Delta_1^b$ -axiomatized *sequential* theories. We assume that every theory comes equipped with a designated interpretation of  $I\Delta_0 + \Omega_1$ . Quantifiers will range over the numbers given by this interpretation. Arithmetized concepts will be implicitly relativized to this interpretation. For example, suppose our theory is **ZF** and **neumann**

is the von Neumann interpretation of number theory in ZF, then

$$\text{ZF} \vdash \forall x, y, z (x + y = z \rightarrow \Box_{\text{ZF}} x + y = z)$$

means:

$$\text{ZF} \vdash (\forall x, y, z (x + y = z \rightarrow \Box_{\text{ZF}}(x + y = z)^{\text{neumann}}))^{\text{neumann}}.$$

We specify the promised logic, which we will call CuL, and its arithmetical semantics. The logic will have two levels, the inner and the outer. We specify the language. We have a set of propositional variables  $p_0, p_1, \dots$ . The meta-variables  $p, q, r, \dots$  will range over the propositional variables. We have a set of cut-variables,  $I_0, I_1, \dots$ . We have one cut-constant  $\text{id}$ . The meta-variables  $I, J, J', \dots$  will range over the cut-variables and  $\text{id}$ . *Outer formulas* are the smallest class containing the propositional variables, closed under the formation rules corresponding to the propositional connectives (including  $\perp$  and  $\top$ ) and closed under the rule:

- if  $A$  and  $B$  are outer formulas, then so are  $\Box^I A$  and  $A \triangleright B$ .

We will write  $\Box A$  for  $\Box^{\text{id}} A$ . *Inner formulas* are the smallest set containing the propositional variables, closed under the formation rules corresponding to the propositional connectives, such that

- if  $A$  and  $B$  are outer formulas, then  $\Box A$  and  $A \triangleright B$  are inner formulas.

Our logic is specified as an extension of a suitable sequent system for propositional logic for the modal language. The following rules are shared by the inner and the outer system, but for the fact that, in the inner system, formulas are constrained to be inner formulas. Concretely, this means that, for the inner system, we must delete the letter  $I$  in the statement of the principles, as formulated below.

$$\begin{array}{l} (\rightarrow)^J \quad \vdash \Box^I A \rightarrow \Box A \\ \text{L}_1^J \quad \vdash \Box^I (A \rightarrow B) \rightarrow (\Box^I A \rightarrow \Box^I B) \\ \text{L}_2^J \quad \vdash \Box^I A \rightarrow \Box^I \Box^J A \\ \text{L}_3^J \quad \vdash \Box^I (\Box^J A \rightarrow A) \rightarrow \Box^I A \\ \text{J}_1^J \quad \vdash \Box (A \rightarrow B) \rightarrow A \triangleright B \\ \text{J}_2^J \quad \vdash (A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C \\ \text{J}_3^J \quad \vdash (A \triangleright C) \wedge (B \triangleright C) \rightarrow A \vee B \triangleright C \\ \text{J}_4^J \quad \vdash A \triangleright B \rightarrow (\diamond A \rightarrow \diamond B) \\ \text{J}_5^J \quad \vdash \diamond^J A \triangleright A \\ \text{Nec}^J \quad \vdash A \Rightarrow \vdash \Box^I A \end{array}$$

By a well-known trick, we can derive  $\text{L}_2^J$  from  $\text{L}_3^J$ . Moreover, we can derive  $\text{L}_{2,\text{in}}^J$  also from  $\text{J}_5^J$  and  $\text{J}_4^J$ . The further rules are as follows.

$$\begin{array}{l} \text{M}^J \quad \Gamma, (A \wedge \Box^J C \triangleright B \wedge \Box^{J'} C) \vdash_{\text{in}} D \Rightarrow \Gamma, A \triangleright B \vdash_{\text{in}} D \\ \quad \text{Here } J \text{ must be a variable and } J \neq J' \text{ and} \\ \quad J \text{ does not occur in } \Gamma, A, B, D \\ \text{IO}^J \quad \Gamma \vdash_{\text{in}} A \Rightarrow \Gamma \vdash_{\text{out}} A \\ \text{OI}^J \quad \vdash_{\text{out}} A \Rightarrow \vdash_{\text{in}} \Box A \end{array}$$

Note that  $\text{Nec}_{\text{in}}^J$  follows from  $\text{IO}^J$  and  $\text{OI}^{J,1}$ .

Inspection of the verifications of interpretability principles in Section 4, reveals that we only use one cut-variable. So the restriction of our system to one variable, say  $I_0$ , is of some interest. Let us write  $\Delta A$  for  $\Box^{I_0} A$  and  $\nabla A$  for  $\neg\Delta\neg A$ . We get the following bimodal system.<sup>2</sup> (We omitt some superfluous principles.)

$(\rightarrow)^{J,1}$	$\vdash_{\text{out}} \Delta A \rightarrow \Box A$
$\text{L}_1^{J,1}$	$\vdash \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$
	$\vdash_{\text{out}} \Delta(A \rightarrow B) \rightarrow (\Delta A \rightarrow \Delta B)$
$\text{L}_2^{J,1}$	$\vdash \Box A \rightarrow \Box \Delta A$
	$\vdash_{\text{out}} \Delta A \rightarrow \Delta \Delta A$
$\text{L}_3^{J,1}$	$\vdash \Box(\Delta A \rightarrow A) \rightarrow \Box A$
	$\vdash_{\text{out}} \Delta(\Delta A \rightarrow A) \rightarrow \Delta A$
$\text{J}_1^{J,1}$	$\vdash \Box(A \rightarrow B) \rightarrow A \triangleright B$
$\text{J}_2^{J,1}$	$\vdash (A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C$
$\text{J}_3^{J,1}$	$\vdash (A \triangleright C) \wedge (B \triangleright C) \rightarrow A \vee B \triangleright C$
$\text{J}_4^{J,1}$	$\vdash A \triangleright B \rightarrow (\diamond A \rightarrow \diamond B)$
$\text{J}_5^{J,1}$	$\vdash \nabla A \triangleright A$
$\text{M}^{J,1}$	$\Gamma, (A \wedge \Delta C \triangleright B \wedge \Box C) \vdash_{\text{in}} D \Rightarrow \Gamma, A \triangleright B \vdash_{\text{in}} D$
	$\Gamma, A, B, D$ must be $\Delta$ -free
$\text{Nec}^{J,1}$	$\vdash_{\text{out}} A \Rightarrow \vdash_{\text{out}} \Delta A$
$\text{IO}^{J,1}$	$\Gamma \vdash_{\text{in}} A \Rightarrow \Gamma \vdash_{\text{out}} A$
$\text{OI}^{J,1}$	$\vdash_{\text{out}} A \Rightarrow \vdash_{\text{in}} \Box A$

We turn to the ‘arithmetical semantics’. Consider a sequential theory  $T$  and let  $N$  be the designated interpretation of  $I\Delta_0 + \Omega_1$ . For any formula  $\alpha$  of the language of  $T$  having at most  $x$  free, we define:

- $\text{cut}(\alpha)$  is the formula that expresses that  $\alpha$  is an  $N$ -cut, i.e. that  $\{x \mid \alpha\}$  is a subset of  $\{x \mid \delta_N\}$  and that  $\{x \mid \alpha\}$  is closed under the operations  $0, \mathsf{S}, +, \times, \omega_1$  and that  $\{x \mid \alpha\}$  is downwards closed under  $<$ .
- $\mathbf{c}(\alpha) := (\delta_N \wedge (\text{cut}(\alpha) \rightarrow \alpha))$ .

It is easy to see that  $T \vdash \text{cut}(\mathbf{c}(\alpha))$  and  $T \vdash \text{cut}(\alpha) \rightarrow (\alpha \leftrightarrow \mathbf{c}(\alpha))$ . so if  $\alpha$  ranges over all formulas having at most  $x$  free, then  $\mathbf{c}(\alpha)$  ranges modulo  $T$ -provable equivalence over all  $T$ -cuts.

Arithmetical realizations in a theory  $T$  are given by pairs  $\sigma, \tau$ , where  $\sigma$  sends the propositional variables to sentences of the language of  $T$  and  $\tau$  sends the cut-variables to formulas  $\alpha$  of the language of  $T$  having only  $x$  free. We demand that  $\sigma$  and  $\tau$  take for all but finitely many arguments

<sup>1</sup>A possible strengthening of the principle  $\text{M}^J$  that we will not further explore in this paper is:

$$\text{M}^{J+} \quad \Gamma, \{(A_i \wedge \Box^J C_i \triangleright B_i \wedge \Box^{J'} C_i) \mid i < n\} \vdash_{\text{in}} D \Rightarrow \Gamma, \{(A_i \triangleright B_i) \mid i < n\} \vdash_{\text{in}} D$$

Here  $J$  must be a variable and  $J \notin \{J'_0, \dots, J'_{n-1}\}$  and  
 $J$  may not occur in  $\Gamma, A_0, \dots, A_{n-1}, B_0, \dots, B_{n-1}, D$

Alternatively, we could add an operation  $\sqcap$  on cuts and formulate the appropriate axiom.

<sup>2</sup>Remember that  $\Box$  is definable from  $\triangleright$ .

the value  $\top$  resp.  $x = x$ , thus making  $\sigma$  and  $\tau$  into finite objects. The assignments  $\sigma$  and  $\tau$  are extended to the full language in the usual way. The main novelty is that we take, for a cut-variable  $I$ ,

- $(\Box^I A)^{\sigma, \tau} := \Box_T^{c(\tau(I))} A^{\sigma, \tau}$ .

The mapping  $A \mapsto A^{\sigma, \tau}$  can be arithmetized. A consequence of this is that we can treat  $\tau$  in the interpretation of *inner formulas* as a *numerical variable*. This observation makes the definition of *inner validity* below sensible. We define:

- $\Gamma \models_{T, \text{out}} A := \forall \sigma, \tau \ T, \Gamma^{\sigma, \tau} \vdash A^{\sigma, \tau}$ ,
- $\Gamma \models_{T, \text{in}} A := \forall \sigma \ T \vdash \forall \tau \ (\bigwedge \Gamma^{\sigma, \tau} \rightarrow A^{\sigma, \tau})$ .

It is not hard to see that our principles do indeed hold w.r.t. the intended notions of validity and that our rules are admissible. We will briefly consider each principle. The principle  $(\rightarrow)^J$  is a triviality. The principle  $L_1^J$  reflects that concatenation of proofs in a cut, remains within this cut as concatenation is approximately multiplication. The principle  $L_2^J$  is Lemma 2.1. The principle  $L_3^J$  is Löb's theorem with cuts, as proved in Lemma 2.2. The principle  $J_5^J$  follows from the formalized Henkin theorem. The rule  $\text{Nec}_{\text{out}}^J$  follows because, if we have a proof of  $A^{\sigma, \tau}$ , then this proof is standard and, hence,  $T$ -provably in every  $T$ -cut. The rule  $M^J$  follows from Lemma 2.3. The rule  $\text{IO}^J$  is trivial. Since the verification of the validity of  $\vdash_{\text{in}}$  can be formalized in  $I\Delta_0 + \Omega_1$ , we have  $\text{Ol}^J$ .

We write  $\text{Form}_U^J$  for the set of  $U$ -formulas having at most  $x$  free.

**Lemma 2.1.** *Consider any theories  $U, V$  with designated natural numbers. ( $U$  and  $V$  need not be sequential.) We have, for any  $V$ -sentence  $\alpha$ ,*

$$I\Delta_0 + \Omega_1 \vdash \forall \beta \in \text{Form}_V^1 (\Box_U \alpha \rightarrow \Box_V \Box_U^{c(\beta)} \alpha).$$

*Proof.* Reason in  $I\Delta_0 + \Omega_1$ . We assume that  $\text{proof}_U(p, \alpha)$  for some  $p$ . Since the sentence  $\text{proof}_U(p, \alpha)$  is in  $\exists \Sigma_1^b$  we get, by verifiable  $\exists \Sigma_1^b$ -completeness, that, for some  $p'$ ,  $\text{proof}_V(p', \text{proof}_U(p, \alpha))$ . By the OBIS-principle<sup>3</sup>, we have  $\Box_V(p \in c(\beta))$ . Ergo  $\Box_V \Box_U^{c(\beta)} \alpha$ .  $\dashv$

We note that  $L_{2, \text{out}}^J$ , follows because we have  $I\Delta_0 + \Omega_1$  on  $T$ -cut  $I$ . We take  $U := T$  and  $V := T$  and we specialize  $\beta$  to the standard formula for  $J$ , noting that its standard code will be  $T$ -provably in  $I$ . We get  $L_{2, \text{in}}^J$  by noting that we have  $I\Delta_0 + \Omega_1$  in  $N$ , the designated numbers of  $T$ .

**Lemma 2.2.** *Consider any theory  $U$  with designated natural numbers. ( $U$  need not be sequential.) We have, for any  $U$ -sentence  $\alpha$ ,*

$$I\Delta_0 + \Omega_1 \vdash \forall \beta \in \text{Form}_U^1 (\Box_U (\Box_U^{c(\beta)} \alpha \rightarrow \alpha) \rightarrow \Box_U \alpha).$$

*Proof.* We can do this in two ways. We can simply adapt the proof of Löb's Theorem, or, alternatively, we can use Löb's Theorem. We follow

---

<sup>3</sup>OBIS stands for *Outside Big Inside Small*. The proof of the OBIS-principle consists in a  $p$ -time transformation of a proof  $q$  of  $\text{cut}(c(\beta))$  and of  $p$  into a proof of  $(p \in c(\beta))$ . We note that the possible 'non-standardness' of  $\beta$  does not matter here.

the second route here. Reason in  $I\Delta_0 + \Omega_1$ . Fix  $\beta \in \text{Form}_U^1$ . We have:

$$\begin{aligned} \Box_U(\Box_U^{c(\beta)}\alpha \rightarrow \alpha) &\rightarrow \Box_U\Box_U^{c(\beta)}(\Box_U^{c(\beta)}\alpha \rightarrow \alpha) \\ &\rightarrow \Box_U\Box_U^{c(\beta)}\alpha \\ &\rightarrow \Box_U\alpha. \end{aligned}$$

⊢

The derivation of  $\mathbf{L}_{3,\text{out}}^J$  and  $\mathbf{L}_{3,\text{in}}^J$  is immediate. Inspecting the proof of Lemma 2.2, we see that  $\mathbf{L}_{3,\text{in}}^J$  follows from  $\mathbf{L}_{3,\text{out}}^J$  plus the other principles and rules.

**Lemma 2.3.** *Suppose  $U$  and  $V$  are theories with designated numbers. Suppose further that  $U$  is sequential. For any  $\Sigma_0^1$ -sentence  $\sigma$ , we have:*

$$I\Delta_0 + \Omega_1 \vdash \forall \beta \in \text{Form}_V^1 (U \triangleright V \rightarrow \exists \gamma \in \text{Form}_U^1 (U + \sigma^{c(\gamma)}) \triangleright (V + \sigma^{c(\beta)})).$$

(Note our slightly sloppy notation: the function  $c$  depends on the designated natural numbers. So, it is not the same in both occurrences.)

*Proof.* This is a direct consequence of Pudlák's Lemma. We proceed to reason in  $I\Delta_0 + \Omega_1$  and suppose  $j : U \triangleright V$  and  $\beta \in \text{Form}_V^1$ . We change the designated numbers of  $V$  to those given by  $c(\beta)$ . Pudlák's Lemma provides us with a  $U$ -cut  $c(\gamma)$  and an isomorphism  $h$  between  $c(\gamma)$  and an initial segment of  $c(\beta)$ . Now, let  $\sigma$  be  $\exists y \sigma_0(y)$  with  $\sigma_0 \in \Delta_0$ . We get that  $\Box_U \forall y \in c(\gamma) (\sigma_0(y) \leftrightarrow \sigma_0^{j \circ c(\beta)}(h(y)))$  and thus certainly

$$\Box_U(\sigma^{c(\gamma)} \rightarrow \sigma^{j \circ c(\beta)}). \quad (1)$$

The desired fact is now immediate. ⊢

Using the modal principles given here, many interesting facts can be derived. With  $A \equiv B$  we shall denote that  $A$  and  $B$  are mutually interpretable. That is,  $(A \triangleright B) \ \& \ (B \triangleright A)$ .

**Lemma 2.4.** *We have:  $\text{CuL} \vdash_{\text{in}} A \equiv A \wedge \Box^I \neg A \equiv A \vee \Diamond^J A$ .*

*Proof.* Just copy the proofs from **IL**, replacing some regular principles with the new principles relativized to a cut. Note that for the first mutual interpretability, we do need to detour via the outer system. ⊢

**Lemma 2.5.**  $\text{CuL} \vdash_{\text{in}} \neg(A \triangleright \neg C) \rightarrow \Diamond(A \wedge \Box^J C)$ .

*Proof.* Reason internally in **CuL**. By contraposition we get that (sloppy notation):  $\Box(A \rightarrow \Diamond^J \neg C) \rightarrow A \triangleright \Diamond^J \neg C \triangleright \neg C$ . ⊢

### 3 The approximating theory

In this section we develop a logic for fine-tuning **ILP**-proofs. Theories in this section will be  $\Delta_1^b$ -axiomatized theories with a designated interpretation of  $\mathbf{S}_2^1$ . (Since  $\mathbf{S}_2^1$  is finitely axiomatizable, it is in this context more convenient to use than  $I\Delta_0 + \Omega_1$ . Of course, nothing substantial is entailed by this choice.)

### 3.1 The approximating theory defined

For finitely axiomatized theories  $V$ , we have:  $S_2^1 \vdash U \triangleright V \rightarrow \Box_{S_2^1}(U \triangleright V)$ , because  $U \triangleright V$  is a  $\exists\Sigma_1^b$ -sentence. To mimic the P-style behavior for an arbitrary theory  $V$ , we will modify  $V$  to a new theory  $V'$  approximating  $V$  to obtain  $\vdash U \triangleright V \rightarrow \Box_{S_2^1}(U \triangleright V')$ . Of course,  $V'$  should be sufficiently like  $V$  to be useful. Thus, we define a theory  $V'$  that is extensionally the same as  $V$ , but for which  $U \triangleright V'$  is so simple that we can easily infer  $\Box_{S_2^1}(U \triangleright V')$ .

The idea is as follows. Let us *define* the set of axioms  $V'$  as consisting of just those axioms  $\phi$  of  $V$  such that  $U \vdash \phi^k$ . Note that, if  $k : U \triangleright V$ , then  $V$  and  $V'$  have the same axioms. However, we cannot take this insight with us when we proceed to reason inside a box. This idea works modulo some trifling details. Firstly, the definition of the new axiom set does not have the right complexity. Secondly, if the argument is not set up in a careful way, we may seem to need both  $\Sigma_1$ -collection and  $\text{exp}$ . We shall use a variation of Craig's trick so that our axiom set will be  $\Delta_1^b$ -definable. The same trick makes the use of strong principles, like  $\Sigma_1$ -collection and  $\text{exp}$ , superfluous.

**Definition 3.1.** Let  $k$  be a translation of the right type. We define  $V^{[k]}$  as follows.

$$\begin{aligned} \text{axioms}_{V^{[k]}}(x) \quad &:= \quad \exists p (x = \ulcorner \varphi \wedge (\underline{p} = \underline{p}) \urcorner \wedge \\ &\quad \text{axioms}_V(\varphi) \wedge \text{proof}_U(p, \varphi^k)). \end{aligned}$$

It is clear that  $\text{axioms}_{V^{[k]}}(x)$  is in poly-time decidable if  $\text{axioms}_V(x)$  and  $\text{axioms}_U(x)$  are. Note that it is essential here that we work with efficient numerals  $\underline{p}$ .

**Lemma 3.2.**

1.  $S_2^1 \vdash \forall k (\text{id} : V \triangleright V^{[k]})$ .
2.  $S_2^1 \vdash \forall k (k : U \triangleright V \rightarrow \text{id} : V^{[k]} \triangleright V)$ .

We see that  $S_2^1$  verifies that  $k : U \triangleright V$  implies that  $V$  and  $V^{[k]}$  are extensionally equal.

*Proof.* Ad (1). Reason in  $S_2^1$ . We have to show:  $\Box_{V^{[k]}}\varphi \rightarrow \Box_V\varphi$ . This is easily seen to be true, since we can replace every axiom  $\varphi \wedge (\underline{p} = \underline{p})$  of  $V^{[k]}$  by a proof of  $\varphi \wedge (\underline{p} = \underline{p})$  from the  $V$ -axiom  $\varphi$ . The resulting transformation is clearly p-time.

Ad (2). Reason in  $S_2^1$ . Suppose  $k : U \triangleright V$  and  $\Box_V\varphi$ . We have a proof  $p$  of  $\varphi$  from axioms,  $\tau_0, \dots, \tau_n$ . Let  $\tau$  be the conjunction of these axioms. Note that  $\tau$  is bounded by  $p$ . Since, clearly,  $\Box_V\tau$ , we may find, using  $k : U \triangleright V$ , a  $U$ -proof  $q$  of  $\tau^k$ . We may use  $q$  to obtain  $U$ -proofs of  $q_i$  of  $\tau_i^k$ . Clearly,  $|q_i|$  is bounded by a term of order  $|q|^2$ . We can now replace every axiom occurrence of  $\tau_i$  in  $p$  by

$$\frac{\tau_i \wedge (q_i = q_i)}{\tau_i} \wedge E, l$$

and obtain a  $V^{[k]}$ -proof  $r$  of  $\varphi$ . We find that  $|r|$  is bounded by a term of order  $|p| \cdot |q|^2$ . So  $r$  can indeed be found in p-time from the given  $p$  and  $q$ .  $\dashv$

Note that, although we do have  $\Box_{S_2^1}(\Box_{V^{[k]}}\varphi \rightarrow \Box_V\varphi)$  we shall, in general, not have  $\Box_{S_2^1}(\Box_V\varphi \rightarrow \Box_{V^{[k]}}\varphi)$ .

**Lemma 3.3.**  $S_2^1 \vdash \forall k (k : U \triangleright V^{[k]})$ .

*Proof.* Reason in  $S_2^1$ . Suppose  $p$  is a  $V^{[k]}$ -proof of  $\phi$ . We want to construct a  $U$ -proof of  $\phi^k$ . As a first step we transform  $p$  into a  $V$ -proof  $p'$  as we did in the proof of Lemma 3.2,(1). Next we transform  $p'$ , using  $k$ , into a predicate logical proof  $q$  of  $\phi^k$  from assumptions  $\tau^k$ , where  $\tau$  is a  $V$ -axiom. It is well known that this transformation is p-time. Finally each axiom  $\tau$  extracted from  $p$ , comes from a  $V^{[k]}$ -axiom  $\tau \wedge (\underline{x} = \underline{x})$ , where  $r$  is a  $U$ -proof of  $\tau^k$ . So our final step is to extend  $q$  to a  $U$ -proof  $q'$  by prepending the  $U$ -proofs  $r$  above the corresponding  $\tau^k$ . This extension will at most double the number of symbols of  $q$ , so  $q' \approx q^2$ .  $\dashv$

### 3.2 A modal logic for approximation

We proceed to articulate modal principles reflecting facts about approximations. We will call our modal system AtL. We first specify the language. We have propositional variables  $p_0, p_1, p_2, \dots$ , we will use  $p, q, r, \dots$ , and we have interpretation variables  $k_0, k_1, k_2, \dots$ . We have one interpretation constant  $\text{id}$ . the meta-variables  $k, \ell, m, \dots$  will range over the interpretation variables and  $\text{id}$ . The modal language is the smallest language containing the propositional variables, closed under the propositional connectives, including  $\top$  and  $\perp$ , and closed under the following rule.

- If  $A, B$  are in the language and  $k$  is an interpretation term, then  $\Box^{[k]}A$  and  $A \triangleright^{[k]}B$  are in the language.

We will write  $\Box$  for  $\Box^{[\text{id}]}$  and  $\triangleright$  for  $\triangleright^{[\text{id}]}$ .

$$\begin{array}{ll}
(\rightarrow \Box)^k & \vdash \Box^{[k]}A \rightarrow \Box A \\
(\rightarrow \triangleright)^k & \vdash A \triangleright B \rightarrow A \triangleright^{[k]}B \\
L_1^k & \vdash \Box^{[k]}(A \rightarrow B) \rightarrow (\Box^{[k]}A \rightarrow \Box^{[k]}B) \\
L_2^k & \vdash \Box^{[\ell]}A \rightarrow \Box^{[k]}\Box^{[\ell]}A \\
L_3^k & \vdash \Box^{[k]}(\Box^{[k]}A \rightarrow A) \rightarrow \Box^{[k]}A \\
J_1^k & \vdash \Box^{[k]}(A \rightarrow B) \rightarrow A \triangleright^{[k]}B \\
J_2^k \text{a} & \vdash (A \triangleright B) \wedge (B \triangleright^{[k]}C) \rightarrow A \triangleright^{[k]}C \\
J_2^k \text{b} & \vdash (A \triangleright^{[k]}B) \wedge \Box^{[k]}(B \rightarrow C) \rightarrow A \triangleright^{[k]}C \\
J_3^k & \vdash (A \triangleright^{[k]}C) \wedge (B \triangleright^{[k]}C) \rightarrow A \vee B \triangleright^{[k]}C \\
J_4^k & \vdash A \triangleright^{[k]}B \rightarrow (\Diamond A \rightarrow \Diamond^{[k]}B) \\
J_5^k & \vdash A \triangleright^{[k]} \Diamond^{[n]}B \rightarrow A \triangleright^{[n]}B \\
P^k & \Gamma, \Delta, \Box(A \triangleright^{[k]}B) \vdash C \Rightarrow \Gamma, A \triangleright B \vdash C \\
\text{Nec}^k & \vdash A \Rightarrow \vdash \Box^{[k]}A
\end{array}$$

Here  $P^k$  is subject to the following conditions<sup>4</sup> :

<sup>4</sup>We realize that this formulation of  $P^k$  is not going to win a beauty contest. However, the primary focus of this paper is to formulate systems that (i) are arithmetically correct



1.  $k$  is an interpretation variable;
2.  $k$  does not occur in  $\Gamma, A, B, C$ ;
3.  $\Delta$  consists of formulas of the form  $(E \triangleright^{[k]} F \rightarrow E \triangleright F)$  and  $(\Box E \rightarrow \Box^{[k]} E)$ .

We will call the licence to use  $(\Box E \rightarrow \Box^{[k]} E)$  provided by  $\mathbf{P}^k$ :  $(\mathbf{E}\Box)^k$ , and we will call the licence to use  $(E \triangleright^{[k]} F \rightarrow E \triangleright F)$ :  $(\mathbf{E}\triangleright)^k$ .<sup>5</sup> In our verification of interpretability principles in Section 4, we use only one interpretation variable, say  $k_0$ . Thus, it could be a good idea to study the bimodal system obtained by restricting AtL to one variable. We leave the obvious formulation of the bimodal case to the reader.

We proceed to formulate the desired notion of arithmetical validity. Let  $T$  be any theory with a designated interpretation, say  $N$ , of  $\mathbf{S}_2^1$ . Let  $\alpha^*$  be a conjunction of  $T$ -axioms that implies  $(\mathbf{S}_2^1)^N$ .

Define, for any translation  $k$  of the language of  $T$  to the language of  $T$ :

- $\mathbf{s}(k) := k\langle(\alpha^*)^k\rangle\text{id}$ .

Thus  $\mathbf{s}(k)$  is the interpretation that is equal to  $k$  in case  $\alpha^*$  and that is the identity interpretation otherwise. We have, verifiably in  $\mathbf{S}_2^1$ , for any  $k$ ,

- $T \vdash (\alpha^*)^{\mathbf{s}(k)}$ ,
- $T \vdash (\alpha^*)^k \rightarrow (\phi^{\mathbf{s}(k)} \leftrightarrow \phi^k)$ ,

Thus, modulo  $T$ -provable equivalence,  $\mathbf{s}(k)$  ranges precisely over all interpretations of  $\alpha^*$ . We have:

**Lemma 3.4.**  $\mathbf{S}_2^1 \vdash \forall k \Box_{T^{\mathbf{s}(k)}} (\mathbf{S}_2^1)^N$ .

*Proof.* Reason in  $\mathbf{S}_2^1$ . Consider any  $k$ . We have a proof in  $T$  of  $(\alpha^*)^{\mathbf{s}(k)}$ . Hence, we have proofs of  $\alpha^{\mathbf{s}(k)}$ , for standardly finitely many  $T$ -axioms  $\alpha$  together implying  $(\mathbf{S}_2^1)^N$ . Clearly,  $T^{[\mathbf{s}(k)]}$  will imply each of these  $\alpha$ , and hence,  $(\mathbf{S}_2^1)^N$ .  $\dashv$

We will take  $N$  as the designated interpretation of  $\mathbf{S}_2^1$  in the  $T^{[k]}$ .

Our assignments for the arithmetical interpretation are pairs  $\sigma, \tau$ , where  $\sigma$  maps the propositional variables to  $T$ -sentences and  $\tau$  maps the interpretation variables to translations from the language of  $T$  to the language of  $T$ . We stipulate that the  $\sigma$  are  $\top$  for all but finitely many arguments and that the  $\tau$  are  $\text{id}$  for all but finitely many arguments. The  $\sigma, \tau$  are lifted to the arithmetical language in the obvious way, taking:

- $(\Box^{[k]} A)^{\sigma, \tau} := \Box_{T^{[\mathbf{s}(k)]}} A^{\sigma, \tau}$ ,

---

for the given interpretation, (ii) enable us to formalize the desired reasoning and (iii) are as parsimonious as possible in expressive power. There is definitely work to do to obtain formalizations of the systems fitting a good proof-theoretical format.

<sup>5</sup>We might wish to consider the following possible strengthening of  $\mathbf{P}^k$ .

$$\mathbf{P}^{k^+} \quad \Gamma, \Delta, \{\Box(A_i \triangleright^{[k]} B_i \mid i < n + 1)\} \vdash C \Rightarrow \Gamma, \{A_i \triangleright B_i \mid i < n + 1\} \vdash C$$

We put the obvious conditions on occurrences of  $k$  and on  $\Delta$ .

- $(A \triangleright^{[k]} B)^{\sigma, \tau} := (T + A^{\sigma, \tau}) \triangleright (T^{[s(k)]} + B^{\sigma, \tau})$ .

Note that the interpretation  $\mathfrak{s}(k)$  is applied only at locations where it is inside a provability. Thus, we can arithmetize its use in  $T$ . For this reason, the following definition makes sense.

- $\Gamma \models_T A :\Leftrightarrow \forall \sigma \ S_2^1 \vdash \forall \tau (\bigwedge \Gamma^{\sigma, \tau} \rightarrow A^{\sigma, \tau})$ .

The iterated modalities make sense because of Lemma 3.4. Note that, if we drop the superscripts in  $J_5^k$ , we get a formula that is equivalent over **J1**, **J2** to the ordinary version of **J5**. Moreover, in our system, we can derive:  $\vdash \diamond^{[k]} A \triangleright^{[k]} A$ .

We turn to checking the validity of our principles and rules. The principles  $(\rightarrow \square)^k$  and  $(\rightarrow \triangleright)^k$  are immediate from Lemma 3.2. The principles  $L_1^k$  to  $J_4^k$  are simple. The validity of  $J_5^k$  follows from the observation that:

$$S_2^1 \vdash \forall k, n (T^{[s(k)]} + \text{con}(T^{[s(n)]} + B)) \triangleright (T^{[s(n)]} + B).$$

This follows by the usual formalization of Henkin's Theorem.

We now consider  $P^k$ . Reason in  $S_2^1$ . Suppose  $k^* : (T + \alpha) \triangleright (T + \beta)$ . Let  $k' := k^*(\alpha)\text{id}$  be the translation that acts like  $k^*$  if  $\alpha$  and like  $\text{id}$  if  $\neg\alpha$ . Let  $k := \mathfrak{s}(k')$ . (This last move is only of an administrative nature, since, in the present context,  $k'$  and  $k$  will be the same in their behaviour as interpretations.) Then, we have both  $k : T \triangleright T$  and  $k : (T + \alpha) \triangleright (T + \beta)$ . By Lemma 3.3, we have  $\square_T(k : T \triangleright T^{[k]})$ . Also we have  $k : (T + \alpha) \triangleright \beta$ , and, hence,  $\square_T(k : (T + \alpha) \triangleright \beta)$ . Combining, we find:  $\square_T(k : (T + \alpha) \triangleright (T^{[k]} + \beta))$ .

By Lemma 3.2, we find that  $\text{id} : T \equiv T^{[k]}$ . So, from  $\square_T \delta$  we will get  $\square_{T^{[k]}} \delta$ . Moreover, Suppose  $m : (T + \delta) \triangleright (T^{[k]} + \varepsilon)$ . It follows that:

$$T + \varepsilon \xrightarrow{\text{id}} T^{[k]} + \varepsilon \xrightarrow{m} T + \delta.$$

So,  $m : (T + \delta) \triangleright (T + \varepsilon)$ .

Finally,  $\text{Nec}^k$  is evident.

## 4 Arithmetical soundness results

In this section, we shall give arithmetical soundness proofs for interpretability principles that hold in all reasonable arithmetical theories. These principles should thus certainly hold in any finitely axiomatizable and in any essentially reflexive theory. This means that the principles should be provable both in **ILP** and **ILM**. We shall see that these two modal proofs give rise to two different arithmetical soundness proofs. The M-style proofs use definable cuts and find place in the modal system **CuL**. The P-style proofs are based on the use of approximating theories. This behavior is captured in the system **AtL**.

We now come to the soundness proofs of the following principles.

$$\begin{array}{ll} \mathbf{W} & \vdash A \triangleright B \rightarrow A \triangleright (B \wedge \square \neg A) \\ \mathbf{M}_0 & \vdash A \triangleright B \rightarrow (\diamond A \wedge \square C) \triangleright (B \wedge \square C) \\ \mathbf{W}^* & \vdash A \triangleright B \rightarrow (B \wedge \square C) \triangleright (B \wedge \square C \wedge \square \neg A) \\ \mathbf{P}_0 & \vdash A \triangleright \diamond B \rightarrow \square(A \triangleright B) \\ \mathbf{R} & \vdash A \triangleright B \rightarrow \neg(A \triangleright \neg C) \triangleright B \wedge \square C \end{array}$$

The principles  $M_0$  and  $P_0$  both follow from  $R$  and  $W^*$  follows from  $M_0$  and  $W$ . So it would be sufficient<sup>6</sup> to just prove the soundness of  $R$  and  $W$ . However, we have decided to give short proofs for all principles. Like this, the close match between the modal systems comes better to the fore. For every principle we shall give a proof in **ILP** and in **ILM**. These proofs can then be copied almost literally to yield arithmetical soundness proofs.

In **CuL** and **AtL** we will reason in an informal way, as if  $M^J$  and  $P^k$  were formulated with an existential quantifier. It is easy to see how to convert this reasoning to the official format.

## 4.1 The principle $W$

We start with the **ILP**-proof of  $W$ .

**Fact 4.1.** **ILP**  $\vdash W$ .

*Proof.* We reason in **ILP**. Suppose  $A \triangleright B$ . Then,  $\Box(A \triangleright B)$ . Hence,  $(*) \Box(\Diamond A \rightarrow \Diamond B)$ , and, thus,  $(**) \Box(\Box \neg B \rightarrow \Box \neg A)$ .

Moreover, from  $A \triangleright B$ , we have  $A \triangleright (B \wedge \Box \neg A) \vee (B \wedge \Diamond A)$ . So it is sufficient to show:  $B \wedge \Diamond A \triangleright B \wedge \Box \neg A$ . We have, by  $(*)$ ,

$$\begin{array}{lll} B \wedge \Diamond A & \triangleright & \Diamond B & \text{by } L_3 \\ & \triangleright & \Diamond(B \wedge \Box \neg B) & \\ & \triangleright & B \wedge \Box \neg B & \text{by } (**) \\ & \triangleright & B \wedge \Box \neg A. & \end{array}$$

□

Note that the proof of Fact 4.1, already works in **ILP<sub>R</sub>**, where:

$$P_R \vdash A \triangleright B \rightarrow \Box(\Diamond A \rightarrow \Diamond B).$$

We turn to the **ILM**-proof of  $W$ .

**Fact 4.2.** **ILM**  $\vdash W$ .

*Proof.* We reason in **ILM**. Suppose  $A \triangleright B$ . We find  $A \wedge \Box \neg A \triangleright B \wedge \Box \neg A$ . But  $A \triangleright A \wedge \Box \neg A$ , whence  $A \triangleright B \wedge \Box \neg A$ . □

**P-style soundness proof of  $W$**  We just follow the modal proof of  $W$  in **ILP**. At some places, axioms are replaced by their counterparts that deal with approximations.

Reason in **AtL**. Suppose that  $A \triangleright B$ . By  $P^k$  we have that, for some  $k$ ,  $\Box(A \triangleright^{[k]} B)$ . Hence, by  $J_4^k$ , we have  $(*) \Box(\Diamond A \rightarrow \Diamond^{[k]} B)$  and, so,  $(**) \Box(\Box^{[k]} \neg B \rightarrow \Box \neg A)$ .

Moreover, from  $A \triangleright B$ , we have  $A \triangleright (B \wedge \Box \neg A) \vee (B \wedge \Diamond A)$ . So it is sufficient to show  $B \wedge \Diamond A \triangleright B \wedge \Box \neg A$ . We have, by  $(*)$ ,

$$\begin{array}{lll} B \wedge \Diamond A & \triangleright & \Diamond^{[k]} B & \text{by } L_3^k \\ & \triangleright & \Diamond^{[k]}(B \wedge \Box^{[k]} \neg B) & \text{by } J_5^k \text{ and } (E\triangleright)^k \\ & \triangleright & B \wedge \Box^{[k]} \neg B & \text{by } (** ) \\ & \triangleright & B \wedge \Box \neg A. & \end{array}$$

<sup>6</sup>Originally in [JG04], a slightly different, but equivalent version of  $R$  was given. In the contribution of Joosten and Goris to this volume, a new principle is given that is precisely  $W$  and  $R$  together.

**M-style soundness proof of  $W$**  We reason in inner CuL. We assume  $A \triangleright B$ . By  $M^J$ , we may find a  $J$  such that  $A \wedge \Box^J \neg A \triangleright B \wedge \Box \neg A$ . By Lemma 2.4,  $A \triangleright A \wedge \Box^J \neg A$ , whence  $A \triangleright B \wedge \Box \neg A$ . We may conclude  $A \triangleright B \wedge \Box \neg A$ .

## 4.2 The principle $M_0$

We start with the **ILP**-proof of  $M_0$ .

**Fact 4.3.** **ILP**  $\vdash M_0$ .

*Proof.* Reason in **ILP**. Suppose  $A \triangleright B$ . By **P**, we have:

$$\begin{array}{ll}
\Box(A \triangleright B) & \rightarrow \\
\Box(\Diamond A \rightarrow \Diamond B) & \rightarrow \\
\Box(\Diamond A \wedge \Box C \rightarrow \Diamond B \wedge \Box C) & \rightarrow \\
\Diamond A \wedge \Box C \triangleright \Diamond B \wedge \Box C & \rightarrow \\
\Diamond A \wedge \Box C \triangleright \Diamond(B \wedge \Box C) & \rightarrow \\
\Diamond A \wedge \Box C \triangleright B \wedge \Box C & \rightarrow
\end{array}$$

□

Note that the proof of Fact 4.3, already works in **ILP<sub>R</sub>**. We proceed with the **ILM**-proof of  $M_0$ .

**Fact 4.4.** **ILM**  $\vdash M_0$ .

*Proof.* Reason in **ILM**. Suppose  $A \triangleright B$ . We find  $A \wedge \Box C \triangleright B \wedge \Box C$ . But,  $\Diamond A \wedge \Box C \triangleright \Diamond(A \wedge \Box C) \triangleright A \wedge \Box C$ , whence  $\Diamond A \wedge \Box C \triangleright B \wedge \Box C$ . □

**P-style soundness proof of  $M_0$**  Reason in AtL. Suppose  $A \triangleright B$ . By **P<sup>k</sup>**, we have, for some  $k$ ,

$$\begin{array}{ll}
\Box(A \triangleright^{[k]} B) & \rightarrow J_4^k \\
\Box(\Diamond A \rightarrow \Diamond^{[k]} B) & \rightarrow \\
\Box(\Diamond A \wedge \Box C \rightarrow \Diamond^{[k]} B \wedge \Box C) & \rightarrow \\
\Diamond A \wedge \Box C \triangleright \Diamond^{[k]} B \wedge \Box C & \rightarrow \text{a.o. by } L_2^k \\
\Diamond A \wedge \Box C \triangleright \Diamond^{[k]}(B \wedge \Box C) & \rightarrow \text{by } J_5^k \text{ and } (E\triangleright)^k \\
\Diamond A \wedge \Box C \triangleright B \wedge \Box C &
\end{array}$$

**M-style soundness proof of  $M_0$**  We reason in inner CuL. Suppose  $A \triangleright B$ . We have, for some  $J$ ,  $A \wedge \Box^J C \triangleright B \wedge \Box C$ . By  $L_2^J$ , and necessitation, we also have  $\Box(\Diamond A \wedge \Box C \rightarrow \Diamond A \wedge \Box \Box^J C)$ , whence

$$\begin{array}{ll}
\Diamond A \wedge \Box C & \triangleright \Diamond A \wedge \Box \Box^J C \\
& \triangleright \Diamond(A \wedge \Box^J C) \\
& \triangleright A \wedge \Box^J C \\
& \triangleright B \wedge \Box C
\end{array}$$

### 4.3 The principle $W^*$

We start with the **ILP**-proof of  $W^*$ .

**Fact 4.5.**  $\mathbf{ILP} \vdash W^*$ .

*Proof.* We reason in **ILP**: Suppose  $A \triangleright B$ . Then:

$$\Box(\Box\neg B \rightarrow \Box\neg A) \quad (2)$$

and

$$\Box(\Diamond A \wedge \Box C \rightarrow \Diamond B \wedge \Box C) \quad (3)$$

Moreover,  $B \wedge \Box C \triangleright (B \wedge \Box C \wedge \Box\neg A) \vee (B \wedge \Box C \wedge \Diamond A)$ . Thus, it is sufficient to show  $B \wedge \Box C \wedge \Diamond A \triangleright B \wedge \Box C \wedge \Box\neg A$ . We have:

$$\begin{aligned} B \wedge \Box C \wedge \Diamond A &\triangleright \Diamond A \wedge \Box C && \text{by (3)} \\ &\triangleright \Diamond B \wedge \Box C && \text{by } L_3 \\ &\triangleright \Diamond(B \wedge \Box\neg B) \wedge \Box C && \text{by } L_2 \\ &\triangleright \Diamond(B \wedge \Box C \wedge \Box\neg B) && \text{by } J_5 \\ &\triangleright B \wedge \Box C \wedge \Box\neg B && \text{by (2)} \\ &\triangleright B \wedge \Box C \wedge \Box\neg A \end{aligned}$$

⊢

Note that the proof of Fact 4.5, already works in  $\mathbf{ILP}_R$ , We proceed with the **ILM**-proof of  $W^*$ .

**Fact 4.6.**  $\mathbf{ILM} \vdash W^*$ .

*Proof.* We reason in **ILM**. Suppose  $A \triangleright B$ . We have:

$$B \wedge \Box C \triangleright (B \wedge \Box C \wedge \Box\neg A) \vee (B \wedge \Box C \wedge \Diamond A).$$

So it sufficient to show,  $B \wedge \Box C \wedge \Diamond A \triangleright B \wedge \Box C \wedge \Box\neg A$ . We have:

$$\begin{aligned} B \wedge \Box C \wedge \Diamond A &\triangleright \Diamond A \wedge \Box C && \text{by } L_3 \\ &\triangleright \Diamond(A \wedge \Box\neg A) \wedge \Box C && \text{by } L_2 \\ &\triangleright \Diamond(A \wedge \Box C \wedge \Box\neg A) && \text{by } J_5 \\ &\triangleright A \wedge \Box C \wedge \Box\neg A && \text{by M and } A \triangleright B \\ &\triangleright B \wedge \Box C \wedge \Box\neg A \end{aligned}$$

⊢

**P-style soundness proof of  $W^*$**  Reason in **AtL**. Suppose  $A \triangleright B$ . By  $P^k$  we obtain a  $k$ , such that we have:

$$\begin{aligned} \Box(A \triangleright^{[k]} B) &\rightarrow \text{by } J_4^k \\ \Box(\Diamond A \rightarrow \Diamond^{[k]} B) &\rightarrow \\ \Box(\Box^{[k]} \neg B \rightarrow \Box\neg A) &\rightarrow (*) \\ \Box(\Diamond A \wedge \Box C \rightarrow \Diamond^{[k]} B \wedge \Box C) &\rightarrow (**) \end{aligned}$$

By the usual reasoning, we only need to show:

$$B \wedge \Box C \wedge \Diamond A \triangleright B \wedge \Box C \wedge \Box\neg A.$$

We have:

$$\begin{array}{lll}
B \wedge \Box C \wedge \Diamond A & \triangleright & \Diamond A \wedge \Box C & \text{by } (**) \\
& \triangleright & \Diamond^{[k]} B \wedge \Box C & \text{by } \mathbf{L}_3^k \\
& \triangleright & \Diamond^{[k]} (B \wedge \Box^{[k]} \neg B) \wedge \Box C & \text{by } \mathbf{L}_2^k \\
& \triangleright & \Diamond^{[k]} (B \wedge \Box C \wedge \Box^{[k]} \neg B) & \text{by } \mathbf{J}_5^k \text{ and } (\mathbf{E}\triangleright)^k \\
& \triangleright & B \wedge \Box C \wedge \Box^{[k]} \neg B & \text{by } (*) \\
& \triangleright & B \wedge \Box C \wedge \Box \neg A & 
\end{array}$$

**M-style soundness proof of  $W^*$**  We reason in inner CuL. Suppose  $A \triangleright B$ . By  $\mathbf{M}^J$ , we have, for some  $J$ ,  $A \wedge \Box^J (C \wedge \neg A) \triangleright B \wedge \Box (C \wedge \neg A)$ , and, hence, by elementary reasoning (via the outer system),

$$A \wedge \Box^J C \wedge \Box^J \neg A \triangleright B \wedge \Box C \wedge \Box \neg A.$$

It is sufficient to show:  $B \wedge \Box C \wedge \Diamond A \triangleright B \wedge \Box C \wedge \Box \neg A$ . We have:

$$\begin{array}{lll}
B \wedge \Box C \wedge \Diamond A & \triangleright & \Diamond A \wedge \Box C & \text{by } \mathbf{L}_3^J \\
& \triangleright & \Diamond (A \wedge \Box^J \neg A) \wedge \Box C & \text{by } \mathbf{L}_2^J \\
& \triangleright & \Diamond (A \wedge \Box^J C \wedge \Box^J \neg A) & \text{by } \mathbf{J}_5^J \\
& \triangleright & A \wedge \Box^J C \wedge \Box^J \neg A & \\
& \triangleright & B \wedge \Box C \wedge \Box \neg A & 
\end{array}$$

#### 4.4 The principle $P_0$

We start with the **ILP**-proof of  $P_0$ .

**Fact 4.7.** **ILP**  $\vdash P_0$ .

*Proof.* Reason in **ILP**. Suppose  $A \triangleright \Diamond B$ . Then,  $\Box(A \triangleright \Diamond B)$  and, so,  $\Box(A \triangleright B)$ .  $\dashv$

**Fact 4.8.** **ILP<sub>R</sub>**  $\not\vdash P_0$ .

*Proof.* It is easy to see that frames satisfying  $uRxRyS_u z \rightarrow xRz$  are sound for **ILP<sub>R</sub>**. And it is equally easy to provide such a model on which  $P_0$  does not hold.  $\dashv$

Fact 4.8 nicely reflects that the frame condition for  $P_0$  essentially involves new  $S$ -transitions. We proceed with the **ILM**-proof of  $P_0$ .

**Fact 4.9.** **ILM**  $\vdash P_0$ .

*Proof.* Reason in **ILM**.

$$\begin{array}{ll}
A \triangleright \Diamond B & \rightarrow A \wedge \Box \neg B \triangleright \perp \\
& \rightarrow \Box(A \rightarrow \Diamond B) \\
& \rightarrow \Box\Box(A \rightarrow \Diamond B) \\
& \rightarrow \Box(A \triangleright \Diamond B) \\
& \rightarrow \Box(A \triangleright B)
\end{array}$$

$\dashv$

**P-style soundness proof of  $P_0$**  Reason in AtL. Suppose  $A \triangleright \Diamond B$ . We have, for some  $k$ ,  $\Box(A \triangleright^{[k]} \Diamond B)$ . Hence, by  $\mathbf{J}_5^k$ ,  $A \triangleright \Diamond B \rightarrow \Box(A \triangleright B)$ .

**M-style soundness proof of  $P_0$**  We reason in the inner system of CuL. Suppose  $A \triangleright \diamond B$ . By  $M^J$ , we have, for some  $J$ ,  $A \wedge \square^J \neg B \triangleright \perp$ . It follows that  $\square(A \rightarrow \diamond^J B)$  and, hence,  $\square \square(A \rightarrow \diamond^J B)$ . We may conclude  $\square(A \triangleright \diamond^J B)$ . So,  $\square(A \triangleright B)$ .

Note: the principle  $A \triangleright \diamond B \rightarrow \square(A \triangleright \diamond B)$  is also provable in both **ILM** and **ILP**. In [Vis97] it is shown that this principle is not valid in **PRA**. It is nice to see where proof-attempts of this principle in our systems fail.

## 4.5 The principle R

Before we see that **ILP**  $\vdash$  R, we first prove an auxiliary lemma.

**Lemma 4.10.** **IL**  $\vdash \neg(A \triangleright \neg C) \wedge (A \triangleright B) \rightarrow \diamond(B \wedge \square C)$ .

*Proof.* We prove the logical equivalent  $(A \triangleright B) \wedge \square(B \rightarrow \diamond \neg C) \rightarrow A \triangleright \neg C$  in **IL**. But this is clear, as  $(A \triangleright B) \wedge \square(B \rightarrow \diamond \neg C) \rightarrow A \triangleright \diamond \neg C$  and  $\diamond \neg C \triangleright \neg C$ .  $\dashv$

**Fact 4.11.** **ILP**  $\vdash$  R.

*Proof.* We reason in **ILP**. Suppose  $A \triangleright B$ . It follows that  $\square(A \triangleright B)$ . Using this together with Lemma 4.10 we get:

$$\begin{aligned} \neg(A \triangleright \neg C) &\triangleright \neg(A \triangleright \neg C) \wedge (A \triangleright B) \\ &\triangleright \diamond(B \wedge \square C) \\ &\triangleright B \wedge \square C \end{aligned}$$

$\dashv$

**Fact 4.12.** **ILP<sub>R</sub>**  $\not\vdash$  R.

*Proof.* By providing a countermodel as in the proof of Fact 4.8.  $\dashv$

**Fact 4.13.** **ILM**  $\vdash$  R.

*Proof.* In **IL**, it is easy to see that  $\neg(A \triangleright \neg C) \rightarrow \diamond(A \wedge \square C)$ . Reason in **ILM**. Suppose  $A \triangleright B$ . Then,

$$\begin{aligned} \neg(A \triangleright \neg C) &\triangleright \diamond(A \wedge \square C) \\ &\triangleright A \wedge \square C \\ &\triangleright B \wedge \square C \end{aligned}$$

$\dashv$

**P-style soundness proof of R** Reason in AtL. We first show that  $(A \triangleright^{[k]} B) \wedge \neg(A \triangleright \neg C) \rightarrow \diamond^{[k]}(B \wedge \Box C)$ . Suppose that  $A \triangleright^{[k]} B$  and  $\Box^{[k]}(B \rightarrow \diamond \neg C)$ , then, by  $J_2^k$ ,  $A \triangleright^{[k]} \diamond \neg C$ . Thus, by  $J_5^k$ , we find  $A \triangleright \neg C$ . By necessitation,

$$\Box((A \triangleright^{[k]} B) \wedge \neg(A \triangleright \neg C) \rightarrow \diamond^{[k]}(B \wedge \Box C)). \quad (4)$$

We now turn to the main proof. Suppose  $A \triangleright B$ . Then, for some  $k$ , we have  $\Box(A \triangleright^{[k]} B)$  and, thus,

$$\begin{aligned} \neg(A \triangleright \neg C) &\triangleright \neg(A \triangleright \neg C) \wedge (A \triangleright^{[k]} B) && \text{by (4)} \\ &\triangleright \diamond^{[k]}(B \wedge \Box C) && \text{by } J_5^k \text{ and } (E\triangleright)^k \\ &\triangleright B \wedge \Box C. \end{aligned}$$

**M-style soundness proof of R** Reason in the inner system of CuL. Suppose that  $A \triangleright B$ . Then, for some  $J$ , we have  $A \wedge \Box^J C \triangleright B \wedge \Box C$ . By Lemma 2.5, we find that, without assumptions:  $\neg(A \triangleright \neg C) \rightarrow \diamond(A \wedge \Box^J C)$ . Hence,  $\Box(\neg(A \triangleright \neg C) \rightarrow \diamond(A \wedge \Box^J C))$  and, so,  $\neg(A \triangleright \neg C) \triangleright \diamond(A \wedge \Box^J C)$ . We have:

$$\begin{aligned} \neg(A \triangleright \neg C) &\triangleright \diamond(A \wedge \Box^J C) \\ &\triangleright A \wedge \Box^J C \\ &\triangleright B \wedge \Box C. \end{aligned}$$

## References

- [AdJH98] C. Areces, D. de Jongh, and E. Hoogland. The interpolation theorem for IL and ILP. In *Proceedings of AiML98. Advances in Modal Logic*, Uppsala, Sweden, October 1998. Uppsala University.
- [dJJ98] D. de Jongh and G. Japaridze. The Logic of Provability. In S.R. Buss, editor, *Handbook of Proof Theory*. Studies in Logic and the Foundations of Mathematics, Vol.137., pages 475–546. Elsevier, Amsterdam, 1998.
- [dJV90] D.H.J. de Jongh and F. Veltman. Provability logics for relative interpretability. In *[Pet90]*, pages 31–42, 1990.
- [dJV91] D.H.J. de Jongh and A. Visser. Explicit fixed points in interpretability logic. *Studia Logica*, 50:39–50, 1991.
- [dJV99] Dick de Jongh and Frank Veltman. The modal completeness of ILW. In Gerbrandy et al. [GMdRV99], page /contribs/jongh/. ISBN 90 5629 104 1.
- [GMdRV99] Jelle Gerbrandy, Maarten Marx, Maarten de Rijke, and Yde Venema, editors. *JFAK. Essays Dedicated to Johan van Benthem on the Occasion of his 50th Birthday*, Vossiuspers, Amsterdam, 1999. Amsterdam University Press. ISBN 90 5629 104 1.
- [JG04] J.J. Joosten and E. Goris. Modal matters in interpretability logic. Logic Group Preprint Series 226, Department of Philosophy, Utrecht University, Heidelberglaan 8, 3584 CS Utrecht, March 2004.



- [JJM91] de D. Jongh, M. Jumelet, and F. Montagna. On the proof of Solovay's theorem. *Studia Logica*, 50:51–70, 1991.
- [Pet90] P.P. Petkov, editor. *Mathematical logic, Proceedings of the Heyting 1988 summer school in Varna, Bulgaria*. Plenum Press, Boston, 1990.
- [Vis97] A. Visser. An overview of interpretability logic. In M. Kracht, M. de Rijke, and H. Wansing, editors, *Advances in modal logic '96*, pages 307–359. CSLI Publications, Stanford, CA, 1997.