

Interpretability in PRA

Marta Bilkova[†], Dick de Jongh^{*}, and Joost J. Joosten^{*},

*Institute for Logic Language and Computation
University of Amsterdam
and

†Department of Logic
Charles University; Prague

14th July 2007

- ▶ We all use the notion $T \triangleright S$: T interprets S

- ▶ We all use the notion $T \triangleright S$: T interprets S
- ▶ $T \triangleright S$ means (modulo some technical details)

- ▶ We all use the notion $T \triangleright S$: T interprets S
- ▶ $T \triangleright S$ means (modulo some technical details)
- ▶ $\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{ Proof}_T(p, \ulcorner \varphi^j \urcorner))$

- ▶ We are interested in the structural behavior of the notion of interpretability.

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for
- ▶ $T + \varphi \triangleright T + \psi$

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for
- ▶ $T + \varphi \triangleright T + \psi$
- ▶ We are interested in the interpretability logic of a theory T :

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for
- ▶ $T + \varphi \triangleright T + \psi$
- ▶ We are interested in the interpretability logic of a theory T :
- ▶ The set of all model propositional logical formulas in the language \square, \triangleright which are true regardless how you interpret the variables as arithmetical sentences

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for
- ▶ $T + \varphi \triangleright T + \psi$
- ▶ We are interested in the interpretability logic of a theory T :
- ▶ The set of all model propositional logical formulas in the language \square, \triangleright which are true regardless how you interpret the variables as arithmetical sentences
- ▶ Of course, reading \triangleright as \triangleright_T , etc.

- ▶ We are interested in the structural behavior of the notion of interpretability.
- ▶ Interpretability can easily be formalized/arithmetized.
- ▶ We shall consider sentential extensions of a base theory
- ▶ $\varphi \triangleright_T \psi$ stands for
- ▶ $T + \varphi \triangleright T + \psi$
- ▶ We are interested in the interpretability logic of a theory T :
- ▶ The set of all model propositional logical formulas in the language \square, \triangleright which are true regardless how you interpret the variables as arithmetical sentences
- ▶ Of course, reading \triangleright as \triangleright_T , etc.
- ▶ Example: $(\varphi \triangleright \psi) \wedge (\psi \triangleright \chi) \rightarrow (\varphi \triangleright \chi)$

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that
- ▶ For example, for theories with full induction, we have that *Montagna's Axiom* holds

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that
- ▶ For example, for theories with full induction, we have that *Montagna's Axiom* holds

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ It turns out that precisely ILM is, e.g. $IL(PA)$ (Shavrukov 1988; Berarducci 1990)

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that
- ▶ For example, for theories with full induction, we have that *Montagna's Axiom* holds

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ It turns out that precisely ILM is, e.g. $IL(PA)$ (Shavrukov 1988; Berarducci 1990)
- ▶ Likewise, the interpretability logic for finitely axiomatized theories is known

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that
- ▶ For example, for theories with full induction, we have that *Montagna's Axiom* holds

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ It turns out that precisely ILM is, e.g. $IL(PA)$ (Shavrukov 1988; Berarducci 1990)
- ▶ Likewise, the interpretability logic for finitely axiomatized theories is known
- ▶ And no other!

- ▶ For all theories T , $IL(T)$ contains some sort of core logic IL
- ▶ $IL(T)$ is characterized by some additional axiom schemes on top of that
- ▶ For example, for theories with full induction, we have that *Montagna's Axiom* holds

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ It turns out that precisely ILM is, e.g. $IL(PA)$ (Shavrukov 1988; Berarducci 1990)
- ▶ Likewise, the interpretability logic for finitely axiomatized theories is known
- ▶ And no other!
- ▶ That's where PRA comes in

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3
- ▶ When S has finitely many axioms, then Σ_1

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3
- ▶ When S has finitely many axioms, then Σ_1
- ▶ When T is reflexive, then Π_2 . (Orey-Hájek).

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3
- ▶ When S has finitely many axioms, then Σ_1
- ▶ When T is reflexive, then Π_2 . (Orey-Hájek).
- ▶ When T is reflexive, we have access to Montagna's Principle:

$$(T \triangleright S) \rightarrow ((T \wedge \Box \gamma) \triangleright (S \wedge \Box \gamma))$$

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3
- ▶ When S has finitely many axioms, then Σ_1
- ▶ When T is reflexive, then Π_2 . (Orey-Hájek).
- ▶ When T is reflexive, we have access to Montagna's Principle:

$$(T \triangleright S) \rightarrow ((T \wedge \Box \gamma) \triangleright (S \wedge \Box \gamma))$$

- ▶ Every extension of PRA by Σ_2 sentences is reflexive (Parsons, Beklemishev, etc)

- ▶ Consider again

$$\exists j \forall \varphi (\text{Axiom}_S(\varphi) \rightarrow \exists p \text{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

- ▶ Certainly Σ_3
- ▶ When S has finitely many axioms, then Σ_1
- ▶ When T is reflexive, then Π_2 . (Orey-Hájek).
- ▶ When T is reflexive, we have access to Montagna's Principle:

$$(T \triangleright S) \rightarrow ((T \wedge \Box \gamma) \triangleright (S \wedge \Box \gamma))$$

- ▶ Every extension of PRA by Σ_2 sentences is reflexive (Parsons, Beklemishev, etc)
- ▶ $(\alpha \triangleright_{\text{PRA}} \beta) \rightarrow ((\alpha \wedge \Box \gamma) \triangleright_{\text{PRA}} (\beta \wedge \Box \gamma))$
 whenever $\alpha \in \Sigma_2$

► $B := (A \triangleright B) \rightarrow (A \wedge \Box C) \triangleright (B \wedge \Box C)$ for $A \in \text{ES}_2$

- ▶ $B := (A \triangleright B) \rightarrow (A \wedge \Box C) \triangleright (B \wedge \Box C)$ for $A \in \text{ES}_2$
- ▶ where

- ▶ $B := (A \triangleright B) \rightarrow (A \wedge \Box C) \triangleright (B \wedge \Box C)$ for $A \in \text{ES}_2$
- ▶ where
- ▶

$$\text{ES}_2 := \Box A \mid \neg \Box A \mid \text{ES}_2 \wedge \text{ES}_2 \mid \text{ES}_2 \vee \text{ES}_2 \mid \neg(\text{ES}_2 \triangleright A)$$

- ▶ If T and S are Π_2 axiomatized theories with

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$
- ▶ then, $T \equiv_1 (T \cup S)$

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$
- ▶ then, $T \equiv_1 (T \cup S)$
- ▶ So,

$$(\alpha \triangleright \beta) \wedge (\beta \triangleright \alpha) \rightarrow (\alpha \triangleright (\alpha \wedge \beta))$$

whenever,

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$
- ▶ then, $T \equiv_1 (T \cup S)$
- ▶ So,

$$(\alpha \triangleright \beta) \wedge (\beta \triangleright \alpha) \rightarrow (\alpha \triangleright (\alpha \wedge \beta))$$

whenever,

- ▶ $\alpha, \beta \in \Sigma_2$, and

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$
- ▶ then, $T \equiv_1 (T \cup S)$
- ▶ So,

$$(\alpha \triangleright \beta) \wedge (\beta \triangleright \alpha) \rightarrow (\alpha \triangleright (\alpha \wedge \beta))$$

whenever,

- ▶ $\alpha, \beta \in \Sigma_2$, and
- ▶ $\alpha, \beta \in \Pi_2$.

- ▶ If T and S are Π_2 axiomatized theories with
- ▶ $T \equiv_1 S$
- ▶ then, $T \equiv_1 (T \cup S)$
- ▶ So,

$$(\alpha \triangleright \beta) \wedge (\beta \triangleright \alpha) \rightarrow (\alpha \triangleright (\alpha \wedge \beta))$$

whenever,

- ▶ $\alpha, \beta \in \Sigma_2$, and
- ▶ $\alpha, \beta \in \Pi_2$.
- ▶ In other words: $\alpha, \beta \in \Delta_2$

- ▶ $Z \quad (A \triangleright B) \wedge (B \triangleright A) \rightarrow (A \triangleright (A \wedge B))$ for A and B in ED_2

▶ Z $(A \triangleright B) \wedge (B \triangleright A) \rightarrow (A \triangleright (A \wedge B))$ for A and B in ED_2

▶

$$ED_2 := \Box A \mid \neg ED_2 \mid ED_2 \wedge ED_2 \mid ED_2 \vee ED_2$$

▶ $Z \quad (A \triangleright B) \wedge (B \triangleright A) \rightarrow (A \triangleright (A \wedge B))$ for A and B in ED_2



$$ED_2 := \Box A \mid \neg ED_2 \mid ED_2 \wedge ED_2 \mid ED_2 \vee ED_2$$

▶ Is this all?

The logic IL

$$\text{L1: } \Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$$

$$\text{L2: } \Box A \rightarrow \Box \Box A$$

$$\text{L3: } \Box(\Box A \rightarrow A) \rightarrow \Box A$$

$$\text{J1: } \Box(A \rightarrow B) \rightarrow A \triangleright B$$

$$\text{J2: } (A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C$$

$$\text{J3: } (A \triangleright C) \wedge (B \triangleright C) \rightarrow A \vee B \triangleright C$$

$$\text{J4: } A \triangleright B \rightarrow (\Diamond A \rightarrow \Diamond B)$$

$$\text{J5: } \Diamond A \triangleright A$$

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.
- ▶ R is conversely well-founded and transitive

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.
- ▶ R is conversely well-founded and transitive
- ▶ $yS_xz \rightarrow xRy \wedge xRz$

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.
- ▶ R is conversely well-founded and transitive
- ▶ $yS_xz \rightarrow xRy \wedge xRz$
- ▶ $xRyRz \rightarrow yS_xz$

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.
- ▶ R is conversely well-founded and transitive
- ▶ $yS_xz \rightarrow xRy \wedge xRz$
- ▶ $xRyRz \rightarrow yS_xz$
- ▶ S_x is transitive and reflexive for each x

- ▶ A Veltman frame $F = \langle W, R, S \rangle$,
 $R \subseteq W \times W$,
 $S_w \subseteq W \times W$ for each $w \in W$.
- ▶ R is conversely well-founded and transitive
- ▶ $yS_xz \rightarrow xRy \wedge xRz$
- ▶ $xRyRz \rightarrow yS_xz$
- ▶ S_x is transitive and reflexive for each x

A model $M = \langle W, R, S, \Vdash \rangle$,
 $\Vdash \subseteq W \times \text{Prop}$

▶ $w \not\Vdash \perp$

A model $M = \langle W, R, S, \Vdash \rangle$,
 $\Vdash \subseteq W \times \text{Prop}$

- ▶ $w \not\Vdash \perp$
- ▶ $w \Vdash A \rightarrow B$ iff $w \not\Vdash A$ or $w \Vdash B$

A model $M = \langle W, R, S, \Vdash \rangle$,
 $\Vdash \subseteq W \times \text{Prop}$

- ▶ $w \not\Vdash \perp$
- ▶ $w \Vdash A \rightarrow B$ iff $w \not\Vdash A$ or $w \Vdash B$
- ▶ $w \Vdash \Box A$ iff $\forall v (wRv \Rightarrow v \Vdash A)$

A model $M = \langle W, R, S, \Vdash \rangle$,
 $\Vdash \subseteq W \times \text{Prop}$

- ▶ $w \not\Vdash \perp$
- ▶ $w \Vdash A \rightarrow B$ iff $w \not\Vdash A$ or $w \Vdash B$
- ▶ $w \Vdash \Box A$ iff $\forall v (wRv \Rightarrow v \Vdash A)$
- ▶ $w \Vdash A \triangleright B$ iff $\forall u (wRu \wedge u \Vdash A \Rightarrow \exists v (uS_w v \Vdash B))$

- ▶ Montagna has a nice frame condition

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ Montagna has a nice frame condition

$$(A \triangleright B) \rightarrow ((A \wedge \Box C) \triangleright (B \wedge \Box C))$$

- ▶ Beklemishev is somewhat similar

A B-simulation on a frame is a binary relation \mathcal{S} for which the following holds.

1. $\mathcal{S}(x, x') \rightarrow x\uparrow = x'\uparrow$
2. $\mathcal{S}(x, x') \ \& \ xRy \rightarrow \exists y'(yS_x y' \wedge \mathcal{S}(y, y') \wedge y'S_{x'}\uparrow \subseteq yS_x\uparrow)$

$F \models \mathcal{C}_B$ if and only if there is a B-simulation \mathcal{S} on F such that for all x and y ,

$$xRy \rightarrow \exists y'(yS_x y' \wedge \mathcal{S}(y, y') \wedge \forall d, e (y'S_{x'} d R e \rightarrow yRd)).$$

$$\begin{aligned}
 & ES_2^0 & := & ED_2 \\
 \blacktriangleright \quad & ES_2^{n+1} & := & ES_2^n \mid ES_2^{n+1} \wedge ES_2^{n+1} \mid ES_2^{n+1} \vee ES_2^{n+1} \mid \\
 & & & \neg(ES_2^n \triangleright \text{Form})
 \end{aligned}$$

- $ES_2^0 := ED_2$
- $ES_2^{n+1} := ES_2^n \mid ES_2^{n+1} \wedge ES_2^{n+1} \mid ES_2^{n+1} \vee ES_2^{n+1} \mid$
 $\neg(ES_2^n \triangleright \text{Form})$
- $S_0(b, u) := b \uparrow = u \uparrow$
- $S_{n+1}(b, u) := S_n(b, u) \wedge$
 $\forall c (bRc \rightarrow \exists c' (uRc' \wedge S_n(c, c') \wedge$
 $cS_b c' \wedge c' S_u \uparrow \subseteq cS_b \uparrow))$

- $$ES_2^0 := ED_2$$
- $$\text{▶ } ES_2^{n+1} := ES_2^n \mid ES_2^{n+1} \wedge ES_2^{n+1} \mid ES_2^{n+1} \vee ES_2^{n+1} \mid \neg(ES_2^n \triangleright \text{Form})$$
- $$S_0(b, u) := b\uparrow = u\uparrow$$
- $$\text{▶ } S_{n+1}(b, u) := S_n(b, u) \wedge \forall c (bRc \rightarrow \exists c' (uRc' \wedge S_n(c, c') \wedge cS_b c' \wedge c'S_u \uparrow \subseteq cS_b \uparrow))$$
- $$\text{▶ For every } i \text{ we define the frame condition } \mathcal{C}_i \text{ to be } \forall a, b (aRb \rightarrow \exists u (bS_a u \wedge S_i(b, u) \wedge \forall d, e (uS_a d R e \rightarrow bR e))).$$

- $$ES_2^0 := ED_2$$
- $$\text{▶ } ES_2^{n+1} := ES_2^n \mid ES_2^{n+1} \wedge ES_2^{n+1} \mid ES_2^{n+1} \vee ES_2^{n+1} \mid \neg(ES_2^n \triangleright \text{Form})$$
- $$S_0(b, u) := b\uparrow = u\uparrow$$
- $$\text{▶ } S_{n+1}(b, u) := S_n(b, u) \wedge \forall c (bRc \rightarrow \exists c' (uRc' \wedge S_n(c, c') \wedge cS_b c' \wedge c'S_u \uparrow \subseteq cS_b \uparrow))$$
- $$\text{▶ For every } i \text{ we define the frame condition } \mathcal{C}_i \text{ to be } \forall a, b (aRb \rightarrow \exists u (bS_a u \wedge S_i(b, u) \wedge \forall d, e (uS_a d R e \rightarrow bR e))).$$

- $ES_2^0 := ED_2$
- $ES_2^{n+1} := ES_2^n \mid ES_2^{n+1} \wedge ES_2^{n+1} \mid ES_2^{n+1} \vee ES_2^{n+1} \mid \neg(ES_2^n \triangleright \text{Form})$
- $S_0(b, u) := b\uparrow = u\uparrow$
- $S_{n+1}(b, u) := S_n(b, u) \wedge \forall c (bRc \rightarrow \exists c' (uRc' \wedge S_n(c, c') \wedge cS_b c' \wedge c'S_u \uparrow \subseteq cS_b \uparrow))$
- For every i we define the frame condition C_i to be

$$\forall a, b (aRb \rightarrow \exists u (bS_a u \wedge S_i(b, u) \wedge \forall d, e (uS_a d R e \rightarrow bR e)))$$

► **Theorem**

A finite frame F validates all instances of Beklemishev's principle if and only if $\forall i F \models C_i$.

► $B \vdash Z$

- ▶ $B \vdash Z$
- ▶ $B \models Z$

- ▶ $B \vdash Z$
- ▶ $B \models Z$
- ▶ Frame condition Zambella?