Why and how study interpretability
Proof theoretic characteristics of PRA
Modal matters

# 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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

- ▶ 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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
**Modal matters**

**Interpretations**
Interpretability logics

- ► 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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

- ▶ 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
- ▶ $\varphi \rhd_T \psi$ stands for
- ▶ $T + \varphi \rhd T + \psi$

**Why and how study interpretability**
Proof theoretic characteristics of PRA
**Modal matters**

**Interpretations**
Interpretability logics

- ▶ 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
- ▶ $\varphi \rhd_T \psi$ stands for
- ▶ $T + \varphi \rhd T + \psi$
- ▶ We are interested in the interpretability logic of a theory $T$:

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

**Interpretations**
Interpretability logics

- ▶ 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
- ▶ $\varphi \rhd_T \psi$ stands for
- ▶ $T + \varphi \rhd T + \psi$
- ▶ We are interested in the interpretability logic of a theory $T$:
- ▶ The set of all model propositional logical formulas in the language $\Box, \rhd$ which are true regardless how you interpret the variables as arithmetical sentences

**Why and how study interpretability**
Proof theoretic characteristics of PRA
**Modal matters**

**Interpretations**
Interpretability logics

- ▶ 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
- ▶ $\varphi \rhd_T \psi$ stands for
- ▶ $T + \varphi \rhd T + \psi$
- ▶ We are interested in the interpretability logic of a theory $T$:
- ▶ The set of all model propositional logical formulas in the language $\Box, \rhd$ which are true regardless how you interpret the variables as arithmetical sentences
- ▶ Of course, reading $\rhd$ as $\rhd_T$, etc.

**Why and how study interpretability**
Proof theoretic characteristics of PRA
**Modal matters**

**Interpretations**
Interpretability logics

- ▶ 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
- ▶ $\varphi \rhd_T \psi$ stands for
- ▶ $T + \varphi \rhd T + \psi$
- ▶ We are interested in the interpretability logic of a theory $T$:
- ▶ The set of all model propositional logical formulas in the language $\Box, \rhd$ which are true regardless how you interpret the variables as arithmetical sentences
- ▶ Of course, reading $\rhd$ as $\rhd_T$, etc.
- ▶ Example: $(\varphi \rhd \psi) \wedge (\psi \rhd \chi) \rightarrow (\varphi \rhd \chi)$

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

Interpretations
**Interpretability logics**

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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

Interpretations
**Interpretability logics**

- 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

**Why and how study interpretability**
Proof theoretic characteristics of PRA
Modal matters

Interpretations
**Interpretability logics**

- 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

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

**Why and how study interpretability**
**Proof theoretic characteristics of PRA**
**Modal matters**

Interpretations
**Interpretability logics**

- 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)

**Why and how study interpretability**
**Proof theoretic characteristics of PRA**
**Modal matters**

Interpretations
**Interpretability logics**

- ▶ 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

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

- ▶ 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

**Why and how study interpretability**
**Proof theoretic characteristics of PRA**
**Modal matters**

Interpretations
**Interpretability logics**

▶ 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

$$(A \rhd B) \to ((A \wedge \Box C) \rhd (B \wedge \Box C))$$

▶ 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!

**Why and how study interpretability**
Proof theoretic characteristics of PRA
**Modal matters**

Interpretations
**Interpretability logics**

- ▶ 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

$$(A \rhd B) \to ((A \land \Box C) \rhd (B \land \Box C))$$

- ▶ 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 were PRA comes in

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

► Consider again

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

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

▶ Consider again

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

▶ Certainly $\Sigma_3$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

► Consider again

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

► Certainly $\Sigma_3$

► When $S$ has finitely many axioms, then $\Sigma_1$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

▶ Consider again

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

▶ Certainly $\Sigma_3$
▶ When $S$ has finitely many axioms, then $\Sigma_1$
▶ When $T$ is reflexive, then $\Pi_2$. (Orey-Hájek).

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

► Consider again

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

► Certainly $\Sigma_3$

► When $S$ has finitely many axioms, then $\Sigma_1$

► When $T$ is reflexive, then $\Pi_2$. (Orey-Hájek).

► When $T$ is reflexive, we have access to Montagna's Principle:

$$(T \rhd S) \to ((T \wedge \Box \gamma) \rhd (S \wedge \Box \gamma))$$

Why and how study interpretability of PRA
Proof theoretic characteristics of PRA
Modal matters

Beklemishev's principle
Zambella's Principle

▶ Consider again

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

▶ Certainly $\Sigma_3$

▶ When $S$ has finitely many axioms, then $\Sigma_1$

▶ When $T$ is reflexive, then $\Pi_2$. (Orey-Hájek).

▶ When $T$ is reflexive, we have access to Montagna's Principle:

$$(T \rhd S) \to ((T \wedge \Box \gamma) \rhd (S \wedge \Box \gamma))$$

▶ Every extension of PRA by $\Sigma_2$ sentences is reflexive (Parsons, Beklemishev, etc)

Why and how study interpretability of PRA
Proof theoretic characteristics of PRA
Modal matters

Beklemishev's principle
Zambella's Principle

▶ Consider again

$$\exists j \; \forall \varphi (\mathsf{Axiom}_S(\varphi) \to \exists p \; \mathsf{Proof}_T(p, \ulcorner \varphi^j \urcorner))$$

▶ Certainly $\Sigma_3$

▶ When $S$ has finitely many axioms, then $\Sigma_1$

▶ When $T$ is reflexive, then $\Pi_2$. (Orey-Hájek).

▶ When $T$ is reflexive, we have access to Montagna's Principle:

$$(T \rhd S) \to ((T \wedge \Box \gamma) \rhd (S \wedge \Box \gamma))$$

▶ Every extension of PRA by $\Sigma_2$ sentences is reflexive (Parsons, Beklemishev, etc)

▶ $(\alpha \rhd_{\mathsf{PRA}} \beta) \to ((\alpha \wedge \Box \gamma) \rhd_{\mathsf{PRA}} (\beta \wedge \Box \gamma))$
whenever $\alpha \in \Sigma_2$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

▶ $B := (A \rhd B) \to (A \wedge \Box C) \rhd (B \wedge \Box C)$    for $A \in \mathsf{ES}_2$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

- $B := (A \rhd B) \to (A \wedge \Box C) \rhd (B \wedge \Box C)$     for $A \in \mathsf{ES}_2$
- where

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

**Beklemishev's principle**
Zambella's Principle

- $B := (A \rhd B) \rightarrow (A \wedge \Box C) \rhd (B \wedge \Box C)$     for $A \in ES_2$
- where
-

$$ES_2 \quad := \quad \Box \mathcal{A} \mid \neg \Box \mathcal{A} \mid ES_2 \wedge ES_2 \mid ES_2 \vee ES_2 \mid \neg(ES_2 \rhd \mathcal{A})$$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

► If $T$ and $S$ are $\Pi_2$ axiomatized theories with

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$
- then, $T \equiv_1 (T \cup S)$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
Zambella's Principle

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$
- then, $T \equiv_1 (T \cup S)$
- So,

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

whenever,

Marta Bilkova$^{\dagger}$, Dick de Jongh$^{*}$, and Joost J. Joosten$^{*}$,     Interpretability in PRA

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$
- then, $T \equiv_1 (T \cup S)$
- So,

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

  whenever,

- $\alpha, \ \beta \in \Sigma_2$, and

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$
- then, $T \equiv_1 (T \cup S)$
- So,

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

whenever,

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

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- If $T$ and $S$ are $\Pi_2$ axiomatized theories with
- $T \equiv_1 S$
- then, $T \equiv_1 (T \cup S)$
- So,

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

  whenever,

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

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

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

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

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

$$ED_2 \quad := \quad \Box \mathcal{A} \mid \neg ED_2 \mid ED_2 \wedge ED_2 \mid ED_2 \vee ED_2$$

Why and how study interpretability
**Proof theoretic characteristics of PRA**
Modal matters

Beklemishev's principle
**Zambella's Principle**

- Z $\quad (A \rhd B) \wedge (B \rhd A) \rightarrow (A \rhd (A \wedge B))$ $\quad$ for $A$ and $B$ in $ED_2$
-
$$ED_2 \quad := \quad \Box \mathcal{A} \mid \neg ED_2 \mid ED_2 \wedge ED_2 \mid ED_2 \vee ED_2$$

- Is this all?

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

# The logic IL

L1: $\Box(A \to B) \to (\Box A \to \Box B)$

L2: $\Box A \to \Box\Box A$

L3: $\Box(\Box A \to A) \to \Box A$

J1: $\Box(A \to B) \to A \rhd B$

J2: $(A \rhd B) \wedge (B \rhd C) \to A \rhd C$

J3: $(A \rhd C) \wedge (B \rhd C) \to A \vee B \rhd C$

J4: $A \rhd B \to (\Diamond A \to \Diamond B)$

J5: $\Diamond A \rhd A$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

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

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

- 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

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

- 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
- $yS_xz \rightarrow xRy \wedge xRz$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

- 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
- $yS_xz \rightarrow xRy \wedge xRz$
- $xRyRz \rightarrow yS_xz$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

- 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 $x$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

- 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
- $y S_x z \rightarrow x R y \wedge x R z$
- $x R y R z \rightarrow y S_x z$
- $S_x$ is transitive and reflexive for each $x$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

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

- $w \nVdash \bot$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

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

- $w \nVdash \bot$
- $w \Vdash A \rightarrow B$ iff $w \nVdash A$ or $w \Vdash B$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

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

- $w \nVdash \bot$
- $w \Vdash A \to B$ iff $w \nVdash A$ or $w \Vdash B$
- $w \Vdash \Box A$ iff $\forall v \; (wRv \Rightarrow v \Vdash A)$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

**The basics**
Frame conditions

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

- $w \nVdash \perp$
- $w \Vdash A \rightarrow B$ iff $w \nVdash A$ or $w \Vdash B$
- $w \Vdash \Box A$ iff $\forall v \ (wRv \Rightarrow v \Vdash A)$
- $w \Vdash A \triangleright B$ iff $\forall u \ (wRu \wedge u \Vdash A \Rightarrow \exists v(uS_w v \Vdash B))$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

▶ Montagna has a nice frame condition

$$(A \rhd B) \to ((A \land \Box C) \rhd (B \land \Box C))$$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
Frame conditions

▶ Montagna has a nice frame condition

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

▶ Beklemishev is somewhat similar

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

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

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

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

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
Frame conditions

$$
\begin{array}{lcl}
\mathrm{ES}_2^0 & := & \mathrm{ED}_2 \\
\blacktriangleright \quad \mathrm{ES}_2^{n+1} & := & \mathrm{ES}_2^n \mid \mathrm{ES}_2^{n+1} \wedge \mathrm{ES}_2^{n+1} \mid \mathrm{ES}_2^{n+1} \vee \mathrm{ES}_2^{n+1} \mid \\
& & \neg(\mathrm{ES}_2^n \rhd \mathrm{Form})
\end{array}
$$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
Frame conditions

$$\begin{aligned}
\text{ES}_2^0 &:= \text{ED}_2 \\
\text{ES}_2^{n+1} &:= \text{ES}_2^n \mid \text{ES}_2^{n+1} \wedge \text{ES}_2^{n+1} \mid \text{ES}_2^{n+1} \vee \text{ES}_2^{n+1} \mid \\
&\quad \neg(\text{ES}_2^n \rhd \text{Form})
\end{aligned}$$

$$\begin{aligned}
\mathcal{S}_0(b, u) &:= b{\uparrow} = u{\uparrow} \\
\mathcal{S}_{n+1}(b, u) &:= \mathcal{S}_n(b, u) \wedge \\
&\quad \forall c\, (bRc \rightarrow \exists c'\, (uRc' \wedge \mathcal{S}_n(c, c') \wedge \\
&\quad cS_b c' \wedge c' S_u{\uparrow} \subseteq cS_b{\uparrow}))
\end{aligned}$$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

$$\begin{aligned}
\mathsf{ES}_2^0 &:= \mathsf{ED}_2 \\
\blacktriangleright \quad \mathsf{ES}_2^{n+1} &:= \mathsf{ES}_2^n \mid \mathsf{ES}_2^{n+1} \wedge \mathsf{ES}_2^{n+1} \mid \mathsf{ES}_2^{n+1} \vee \mathsf{ES}_2^{n+1} \mid \\
&\qquad \neg(\mathsf{ES}_2^n \rhd \mathsf{Form})
\end{aligned}$$

$$\begin{aligned}
\mathcal{S}_0(b, u) &:= b{\uparrow} = u{\uparrow} \\
\blacktriangleright \quad \mathcal{S}_{n+1}(b, u) &:= \mathcal{S}_n(b, u) \wedge \\
&\qquad \forall c \, (bRc \rightarrow \exists c' \, (uRc' \wedge \mathcal{S}_n(c, c') \wedge \\
&\qquad cS_b c' \wedge c' S_u{\uparrow} \subseteq cS_b{\uparrow}))
\end{aligned}$$

► For every $i$ we define the frame condition $\mathcal{C}_i$ to be

$\forall a, b \, (aRb \rightarrow \exists u \, (bS_a u \wedge \mathcal{S}_i(b, u) \wedge \forall d, e \, (uS_a dRe \rightarrow bRe)))$.

Why and how study interpretability in PRA
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

$$\begin{aligned}
\mathsf{ES}_2^0 &:= \mathsf{ED}_2 \\
\blacktriangleright \quad \mathsf{ES}_2^{n+1} &:= \mathsf{ES}_2^n \mid \mathsf{ES}_2^{n+1} \wedge \mathsf{ES}_2^{n+1} \mid \mathsf{ES}_2^{n+1} \vee \mathsf{ES}_2^{n+1} \mid \\
&\qquad \neg(\mathsf{ES}_2^n \rhd \mathsf{Form})
\end{aligned}$$

$$\begin{aligned}
\mathcal{S}_0(b, u) &:= b{\uparrow} = u{\uparrow} \\
\blacktriangleright \quad \mathcal{S}_{n+1}(b, u) &:= \mathcal{S}_n(b, u) \wedge \\
&\qquad \forall c \, (bRc \to \exists c' \, (uRc' \wedge \mathcal{S}_n(c, c') \wedge \\
&\qquad cS_bc' \wedge c'S_u{\uparrow} \subseteq cS_b{\uparrow}))
\end{aligned}$$

► For every $i$ we define the frame condition $\mathcal{C}_i$ to be
  $\forall a, b \, (aRb \to \exists u \, (bS_a u \wedge \mathcal{S}_i(b, u) \wedge \forall d, e \, (uS_a dRe \to bRe)))$.

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

$$\begin{aligned}
\mathrm{ES}_2^0 &:= \mathrm{ED}_2
\end{aligned}$$

▶ $\mathrm{ES}_2^{n+1} := \mathrm{ES}_2^n \mid \mathrm{ES}_2^{n+1} \wedge \mathrm{ES}_2^{n+1} \mid \mathrm{ES}_2^{n+1} \vee \mathrm{ES}_2^{n+1} \mid$
$\qquad\qquad\qquad \neg(\mathrm{ES}_2^n \rhd \mathrm{Form})$

$$\begin{aligned}
\mathcal{S}_0(b, u) &:= b{\uparrow}=u{\uparrow}
\end{aligned}$$

▶ $\mathcal{S}_{n+1}(b, u) := \mathcal{S}_n(b, u) \wedge$
$\qquad\qquad\qquad \forall c\, (bRc \rightarrow \exists c'\, (uRc' \wedge \mathcal{S}_n(c, c') \wedge$
$\qquad\qquad\qquad cS_bc' \wedge c'S_u{\uparrow} \subseteq cS_b{\uparrow}))$

▶ For every $i$ we define the frame condition $\mathcal{C}_i$ to be
$\forall a, b\, (aRb \rightarrow \exists u\, (bS_au \wedge \mathcal{S}_i(b, u) \wedge \forall d, e\, (uS_adRe \rightarrow bRe)))$.

▶ Theorem
*A finite frame F validates all instances of Beklemishev's principle if and only if $\forall i\ F \models \mathcal{C}_i$.*

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
Frame conditions

► B ⊢ Z

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
Frame conditions

- $B \vdash Z$
- $B \models Z$

Why and how study interpretability
Proof theoretic characteristics of PRA
**Modal matters**

The basics
**Frame conditions**

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