

# Principles of interpretability logic in the intersection of ILP and ILM

---

Vicent Navarro Arroyo

joint work with

Joost J. Joosten

July 7th-11th

Logic Colloquium 2025, Vienna, Austria

Department of Philosophy,  
University of Barcelona

# Motivation

---

# Motivation

In Provability Logic, for a fixed theory  $T$  (PL)  $\Box A$  reads as

“ $A$ ” is provable in  $T$ .

Interpretability Logic (IL) extends PL adding  $A \triangleright B$  which means

$T + A$  interprets  $T + B$

We say that  $S$  interprets  $T$  –  $S \triangleright T$  – if there exists a mapping

$$j: \text{Form}_T \rightarrow \text{Form}_S$$

that *preserves structure*, for example, if  $\circ$  is a binary logical connective, then  $(\varphi \circ \psi)^j = \varphi^j \circ \psi^j$  such that moreover

$$\forall \varphi \left( \Box_T \varphi \rightarrow \Box_S \varphi^j \right).$$

## Example

Natural numbers can be interpreted as sets.

Gödel's Second Incompleteness Theorem is modally expressed as

$$\Diamond T \rightarrow \neg \Box \Diamond T.$$

In interpretability logic it can be generalized to

$$\Diamond T \rightarrow \neg(T \triangleright \Diamond T). \quad (\text{Feferman})$$

We can define the interpretability logic of a theory  $T$ .

$$\text{IL}(T) := \{A \mid \forall * \ T \vdash A^*\},$$

where  $A$  is a formula in the language  $L_{\Box, \triangleright}$

$$F := \perp \mid \text{Prop} \mid F \rightarrow F \mid \Box F \mid F \triangleright F,$$

and  $*$  is a translation sending propositional variables to arithmetical sentences.

# Motivation

The axioms of the basic interpretability **IL** are

L1 $\Box(A \rightarrow B) \rightarrow (\Box A \rightarrow \Box B)$	J2 $(A \triangleright B) \wedge (B \triangleright C) \rightarrow A \triangleright C$
L2 $\Box A \rightarrow \Box \Box A$	J3 $A \triangleright C \wedge B \triangleright C \rightarrow A \vee B \triangleright C$
L3 $\Box(\Box A \rightarrow A) \rightarrow \Box A$	J4 $A \triangleright B \rightarrow (\Diamond A \rightarrow \Diamond B)$
J1 $\Box(A \rightarrow B) \rightarrow A \triangleright B$	J5 $\Diamond A \triangleright A$

## Remark

- J1 tells us that the identity translation yields an interpretation.
- J5 represents Henkin's completeness theorem formalised.

# Motivation

There are some interesting principles of interpretability.

Namely,

$$M := A \triangleright B \rightarrow A \wedge \Box C \triangleright B \wedge \Box C \quad (\text{Montagna})$$

$$P := A \triangleright B \rightarrow \Box(A \triangleright B) \quad (\text{Persistence})$$

It is known that

$$IL(PA) := ILM \quad (\text{Full induction})$$

and

$$IL(I\Sigma_1) := ILP \quad (\text{Finitely Axiomatized}).$$

ILM and ILP motivate the characterisation of  $\text{IL}(\text{All})$ .

$$\text{IL}(\text{All}) := \{A \mid \forall T \supseteq \text{I}\Delta_0 + \text{Exp} \ \forall * T \vdash A^*\},$$

the interpretability logic of all “reasonable” arithmetical theories.

## Remark

$$\text{IL}(\text{All}) \subsetneq \text{ILM} \cap \text{ILP}$$

We present some advances on its modal characterization.



# Semantics and intersections

---

# Semantics and intersections

In interpretability logic, models are 4-tuples

$$\mathcal{M} := \langle W, R, \{S_x\}_{x \in W}, V \rangle$$

where

- $W \neq \emptyset$
- $R \subseteq W \times W$
- $S_x \subseteq x \upharpoonright \times x \upharpoonright$
- $V: \text{Prop} \rightarrow \mathcal{P}(W)$

$$x \upharpoonright := \{y \mid xRy\}.$$

$R$  transitive and conversely well-founded;

$S_x$  is reflexive transitive and contains  $R$  on  $x \upharpoonright$ .

$\mathcal{F} = \langle W, R, \{S_x\}_{x \in W} \rangle$  denotes a frame.

Sometimes we denote models as  $\mathcal{M} = \langle \mathcal{F}, V \rangle$ .

Propositions, implications and *falsum* ( $\perp$ ) are forced as usual.

The forcing of formulas  $\Box A$  is

$$\mathcal{M}, x \Vdash \Box A: \iff \forall y (xRy \rightarrow \mathcal{M}, y \Vdash A).$$

The forcing of formulas  $A \triangleright B$  is

$$\mathcal{M}, x \Vdash A \triangleright B: \iff \forall y (xRy \wedge \mathcal{M}, y \Vdash A \rightarrow \exists z: yS_x z \wedge \mathcal{M}, z \Vdash B).$$

# Semantics and intersections

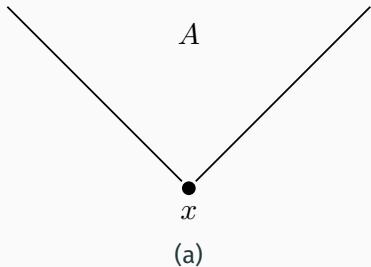
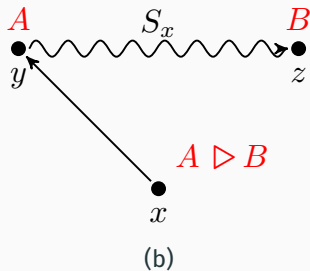


Figure 1: (a)  $\Box A$  is forced at  $x$



(b)  $A \triangleright B$  is forced at  $x$

Validity on models and frames is defined as follows.

## Validity

- **Validity of a formula on a model:**

$$\mathcal{M} \models \varphi \text{ iff } \mathcal{M}, w \Vdash \varphi, \text{ for all } w \in W.$$

- **Validity of a formula on a frame:**

$$\mathcal{F} \models \varphi \text{ iff } \forall V \langle \mathcal{F}, V \rangle \models \varphi.$$

- **Validity of a scheme:** A model or a frame validates a scheme  $X$  ( $\mathcal{M} \models X$  and  $\mathcal{F} \models X$ , respectively) iff it validates all  $X$ 's instances.

# Semantics and intersections

The **frame condition** of a scheme  $X$  is a first (or higher) order predicate formula  $\mathcal{C}$  such that

$$\forall \mathcal{F} (\mathcal{F} \models \mathcal{C} \iff \mathcal{F} \models X).$$

## Example

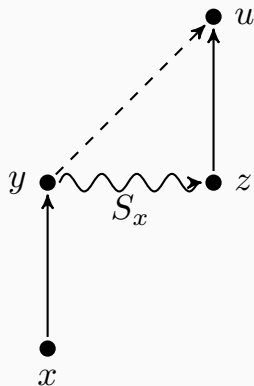
$$\mathcal{F} \models \Box A \rightarrow \Box \Box A \iff \mathcal{F} \models \forall x, y, z \left( xRy \wedge yRz \rightarrow xRz \right)$$

Frame conditions of ILM and ILP.

$$\mathcal{F} \models M \iff \mathcal{F} \models xRyS_xzRu \rightarrow yRu.$$

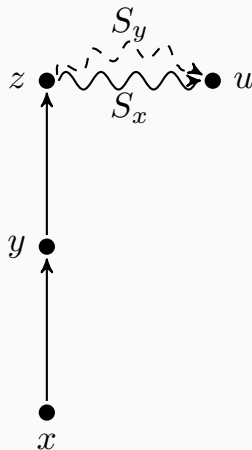
$$\mathcal{F} \models P \iff \mathcal{F} \models xRyRzS_xu \rightarrow zS_yu.$$

# Semantics and intersections



(a)

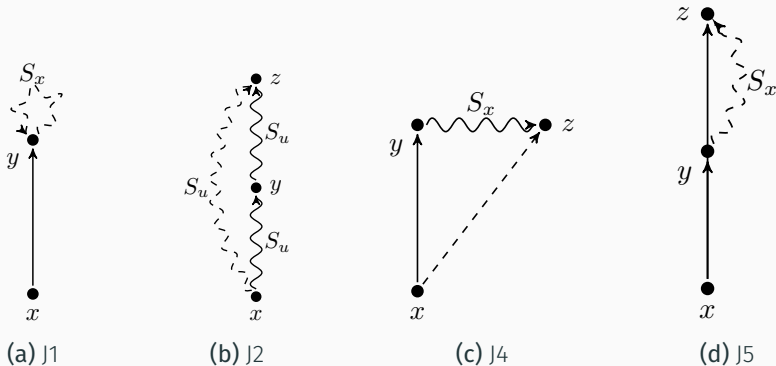
Figure 2: Frame condition of M (a)



(b)

Frame condition of P (b)

# Semantics and intersections



**Figure 3:** Frame definition reflecting axioms  $\Box(A \rightarrow B) \rightarrow A \triangleright B$  (J1),  $A \triangleright B \wedge B \triangleright C \rightarrow A \triangleright C$  (J2),  $A \triangleright B \rightarrow (\Diamond A \rightarrow \Diamond B)$  (J4) and  $\Diamond A \triangleright A$  (J5)



# Semantics and intersections

Sometimes we need to close on the frame properties.

## Closure

The closure of a (proto-) frame  $\mathcal{F} := \langle W, R, \{S_x\}_{x \in W} \rangle$  under some principle  $X$  is the smallest structure

$\overline{\mathcal{F}}^X := \langle W, \overline{R}^X, \{\overline{S}_x^X\}_{x \in W} \rangle$  satisfying  $X$  such that  $R \subseteq \overline{R}^X$  and  $S_x \subseteq \overline{S}_x^X$ , for every  $x \in W$ .

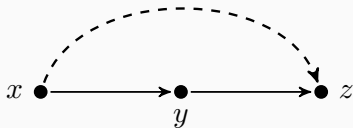


Figure 4: Transitive closure

# Semantics and intersection

## Frame operator

If  $L = \{\phi_i\}_i$  is a set of atomic predicates (like  $xRy$  or  $yS_xz$ , etc.), we define the **IL-frame induced by  $L$** ,  $\overline{\mathcal{F}(\bigwedge_i \phi_i)}^{\text{IL}}$ , as the universal closure of the smallest proto-frame that satisfies all atomic predicates.

For brevity, we will write  $\mathcal{F}(\bigwedge_i \phi_i)$ .

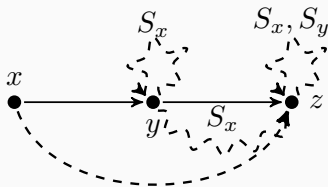


Figure 5: Closure of  $\{xRy, yRz\}$  under IL frame requirements.

# Semantics and intersection

Let  $\mathfrak{F}$  be a class of **IL**-frames. We define the interpretability logic corresponding to  $\mathfrak{F}$ .

$$\mathbf{IL}[\mathfrak{F}] := \{A : \text{for all } \mathcal{F} \in \mathfrak{F}, \mathcal{F} \models A\}.$$

Let  $F(x, y, z)$  denote any first or higher order formula where the only free variables are  $x, y, z$ . We now define the following class of conditions.

$$\mathcal{C}_{\mathbf{ILP} \cap_S \mathbf{ILM}} :=$$

$$\{F(x, y, z) \rightarrow xS_yz : \mathbf{ILP} \models F(x, y, z) \rightarrow xS_yz \wedge \mathbf{ILM} \models F(x, y, z) \rightarrow xS_yz\}.$$

Also, we define the class

$$\mathfrak{A} \mathfrak{I} \mathfrak{I} := \{\mathcal{F} \models \mathbf{ILW} : \forall C \in \mathcal{C}_{\mathbf{ILP} \cap_S \mathbf{ILM}}, \mathcal{F} \models C\}.$$

The principle  $W$  is

$$W := A \triangleright B \rightarrow A \triangleright (B \wedge \Box \neg A)$$

and its frame condition is that there are no  $S_x; R$  infinite chains.

**Conjecture 1 (Goris, Joosten 2020)**

$$IL(All) = IL[\mathfrak{M}].$$

**Recall**

$$IL(All) := \{A \mid \forall T \supseteq I\Delta_0 + \text{Exp} \ \forall * T \vdash A^*\}.$$

## $M \cap P$ -closure

Given a proto-frame  $\mathcal{F} = \langle W, R, S \rangle$ , its  $M \cap P$ -closure is  $\overline{\mathcal{F}}^{M \cap P} := \overline{\mathcal{F}}^M \cap \overline{\mathcal{F}}^P = \langle W, \overline{R}^M \cap \overline{R}^P, \overline{S}^M \cap \overline{S}^P \rangle$ .

As an example, consider the principle  $M_0$

$$M_0 := A \triangleright B \rightarrow \Diamond A \wedge \Box C \triangleright B \wedge \Box C,$$

whose frame condition is

$$\forall x, y, z, u, v \left( xRyRzS_xuRv \rightarrow yRv \right).$$

# Semantics and intersection

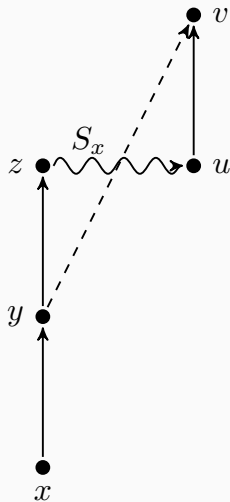
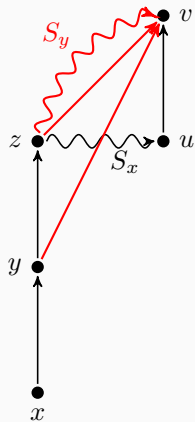


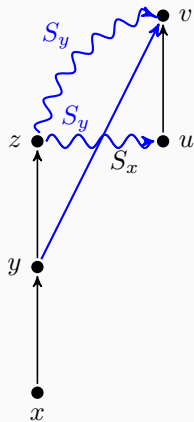
Figure 6:  $M_0$

# Semantics and intersection



(a)

Figure 7: (a) M closure and



(b)

(b) P closure

# Semantics and intersection

## $M \cap_{\mathcal{F}} P$ -clause set

We define the  $M \cap_{\mathcal{F}} P$ -clause set as

$$\bigwedge_i \phi_i \rightarrow \varphi : \in M \cap_{\mathcal{F}} P \text{ iff } \overline{\mathcal{F}(\bigwedge_i \phi_i)}^{M \cap P} \models \varphi$$

whenever  $\{\phi_i\}_i \cup \{\varphi\}$  is a set of atomic predicates so that  $\mathcal{F}(\bigwedge_i \phi_i)$  defines a proto-frame.

## Remark

$\bigwedge_i \phi_i \rightarrow \varphi$  is a Horn clause.

Non-empty since the  $M_0$  frame condition belongs to it.

It is known that the *Broad* series and the *Slim* hierarchy belong to it.



# Semantics and intersection

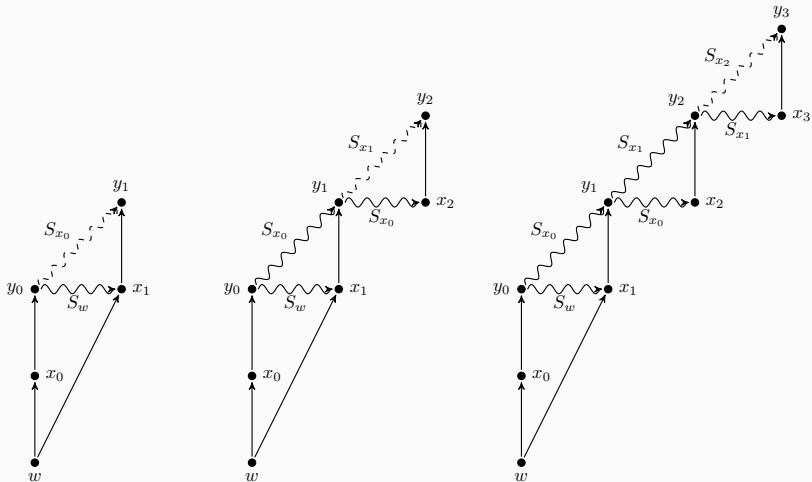


Figure 8: Slim (or Staircase) hierarchy

# Semantics and intersection

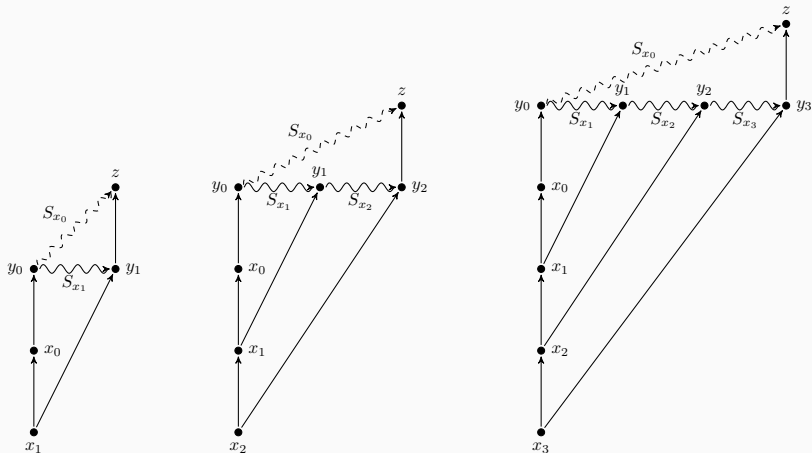


Figure 9: Broad series

# Semantics and intersection

$M \cap_{\mathcal{F}} P$  defines a fragment of  $IL[\mathfrak{M}]$ .

Let us define the lower-case class of **IL**-frames

$$\mathfrak{all} := \{\mathcal{F} \models ILW : \forall C \in M \cap_{\mathcal{F}} P, \mathcal{F} \models C\}.$$

## Theorem

$$IL[\mathfrak{all}] \subseteq IL[\mathfrak{M}].$$

## Remark

- It is unknown if  $IL[\mathfrak{all}] \subset IL[\mathfrak{M}]$ .
- $IL[\mathfrak{all}]$  entails the frame conditions of *Broad* and *Slim*.

It is natural to conjecture that

## Conjecture 2

$$IL[\mathfrak{all}] = IL(All).$$

This new conjecture strengthens the old conjecture.

## Conjecture 1 (Goris, Joosten 2020)

$$IL(All) = IL[\mathfrak{all}].$$

How can we get a grip on  $M \bigcap_{\mathcal{F}} P$ ?

One may try to focus on the clauses that imply an  $R$ -pair and conjecture that

## Conjecture 3

Consider an **IL**-frame  $\mathcal{F} = \langle W, R, S \rangle$ . Then, for any  $x, y \in W$ , we have that  $x\bar{R}^M y \wedge x\bar{R}^P y \wedge \neg(xRy) \rightarrow x\bar{R}^{M_0} y$ .

Nonetheless, this is disproven by the...

# Pencil frame

---

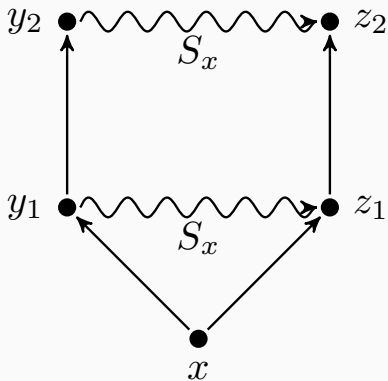


Figure 10: Pencil frame.

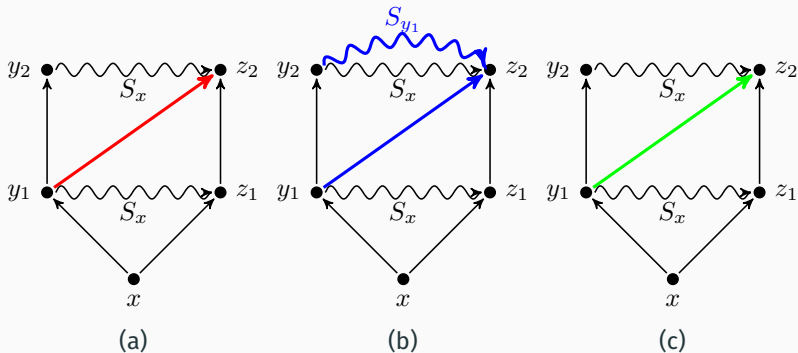


Figure 11: (a) M-closure (b) P-closure (c) Intersection.

## Remark

Observe the **green arrow** is not in the  $M_0$ -closure.



# Pencil frame

We observe the Pencil frame is **not** modally definable.

## Frame definability

Given a first or higher order predicate formula  $\mathcal{C}$ . The class of frames that make true  $\mathcal{C}$  is **modally definable** if

$$\exists A \in L_{\Box, \triangleright} \forall \mathcal{F} (\mathcal{F} \models \mathcal{C} \iff \mathcal{F} \models A).$$

## Example

The class of transitive frames is defined by  $\Box A \rightarrow \Box \Box A$ .

## Remark

Consider the formula  $\mathcal{C}_P := xRy_1S_xz_1Rz_2 \wedge y_1Ry_2S_xz_2 \rightarrow y_1Rz_2$ . Notice that  $y_1Rz_2$  is precisely the **green arrow**.

# Pencil frame

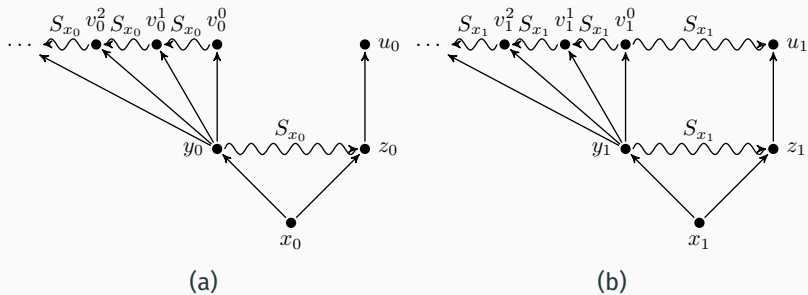


Figure 12: (a)  $\mathcal{F}_0$  satisfies  $\mathcal{C}_P$

(b)  $\mathcal{F}_1$  does not satisfy  $\mathcal{C}_P$

## Theorem: Pencil frame is not modally definable

- By *Reductio ad Absurdum*, assume it is, that is,

$$\exists A \in L_{\Box, \triangleright} \forall \mathcal{F} (\mathcal{F} \models \mathcal{C}_P \iff \mathcal{F} \models A).$$

- Consider  $\mathcal{F}_0$  and  $\mathcal{F}_1$ . Notice  $\mathcal{F}_0 \models \mathcal{C}_P$  whereas  $\mathcal{F}_1 \not\models \mathcal{C}_P$ .
- Then, by hypothesis,  $\mathcal{F}_0 \models A$  and  $\mathcal{F}_1 \not\models A$ .
- **Claim:**  $\forall V_1 \exists V_0: \langle \mathcal{F}_1, V_1 \rangle \sim_{\text{bisimilar}} \langle \mathcal{F}_0, V_0 \rangle$ .
- Bisimilar image-finite models prove the same modal formulas (**Hennessy–Milner**).
- Thus,  $\exists V_0 \langle \mathcal{F}_0, V_0 \rangle \not\models A$ . **Contradiction!** ( $\mathcal{F}_0 \models A$ )

New series

---

Given that the Pencil frame is not modally definable and its frame condition is in  $M \bigcap_{\mathcal{F}} P$  and not induced by neither *Broad* nor *Slim*, a natural question arises:

**Is there a class of modally definable frames whose frame condition is in  $M \bigcap_{\mathcal{F}} P$  but it is not induced by *Slim* nor *Broad*?**

We found out that the answer is positive 😊

## New series

We will inductively define a series of schemes.

Firstly, we inductively define the following series of formulas.

$$\begin{aligned}\varphi^0 &:= \Diamond ((D \triangleright D_0) \wedge \Diamond \neg(A \triangleright \neg C)), \\ \varphi^n &:= \Diamond ((D_{n-2} \triangleright D_{n-1}) \wedge \varphi^{n-1}).\end{aligned}\quad (n \geq 1)$$

Then, we inductively define  $V$  as the series of all the principles  $V^n$ , for any  $n \in \mathbb{N}$ , where

$$\begin{aligned}V^0 &:= A \triangleright B \rightarrow ((D_0 \triangleright \Diamond D_1) \wedge \varphi^0) \triangleright B \wedge \Box C \wedge (D \triangleright D_1), \\ V^{n+1} &:= V^n[\varphi^n / \varphi^{n+1}; \\ &\quad D_n \triangleright \Diamond D_{n+1} / D_{n+1} \triangleright \Diamond D_{n+2}; \\ &\quad D \triangleright D_{n+1} / D \triangleright D_{n+2}]\end{aligned}$$

For example,

$v^0$  :=

$$A \triangleright B \rightarrow ((D_0 \triangleright \Diamond D_1) \wedge \Diamond ((D \triangleright D_0) \wedge \Diamond \neg(A \triangleright \neg C))) \triangleright B \wedge C \wedge (D \triangleright D_1),$$

$v^1$  :=

$$A \triangleright B \rightarrow ((D_1 \triangleright \Diamond D_2) \wedge \Diamond ((D_0 \triangleright D_1) \wedge \Diamond ((D \triangleright D_0) \wedge \Diamond \neg(A \triangleright \neg C)))) \triangleright B \wedge C \wedge (D \triangleright D_2).$$

Their frame conditions are, respectively, ...

# New series

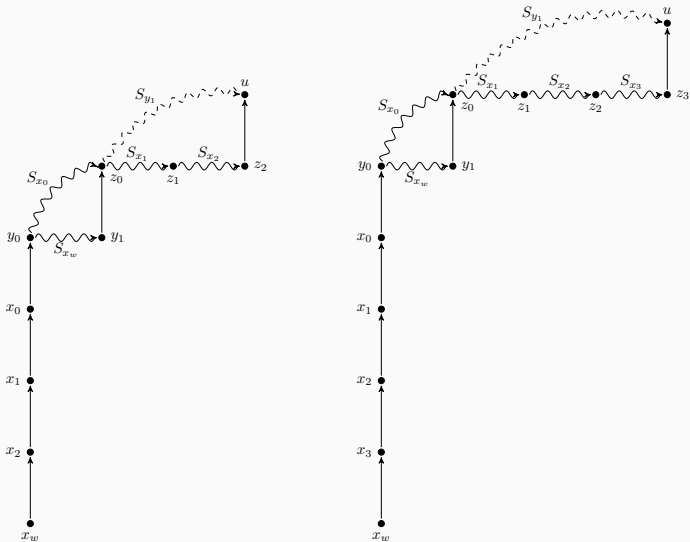
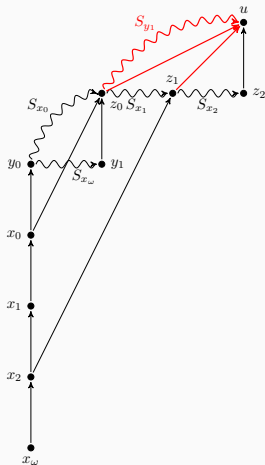


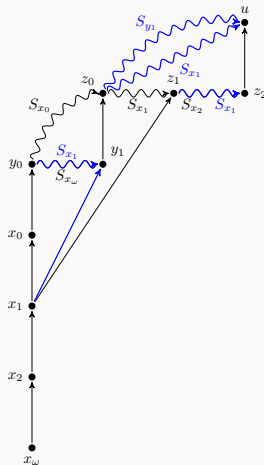
Figure 13: (Left) Frame condition of  $V_0$ . (Right) Frame condition of  $V_1$  35



# New series



(a)



(b)

Figure 14: (a) M-closure

(b) P-closure

## Remark

The V series is **not** a hierarchy.

Wrapping up:

- The classes of frames that satisfy the frame conditions of the V series are modally definable.
- The frame conditions of the V series belong to  $M \bigcap_{\mathcal{F}} P$ .
- It can be shown that neither the *Broad* series nor the *Slim* hierarchy induce the V series.

Also, these principles of the V series are arithmetically valid through **arithmetical definable cuts**.

## Summary (Summer-e)

1. We strengthen the old conjecture by focusing on Horn clauses;
2. We show that all known principles fall in this class;
3. We show that some frame properties are modally undefinable;
4. We found a new series of principles;
5. We have proven the new principles to be arithmetically sound;
6. Thus the conjecture still stands;
7. Preprint and paper coming out 'soon'.

Thank you! Danke!