

# Neural Net Semantics with Modal Operators

This paper is inspired by Hannes Leitgeb's [Nonmonotonic Reasoning by Inhibition Nets](#), which proves completeness for the neuro-symbolic interface suggested by Balkenius and Gärdenfors' [Nonmonotonic Inferences in Neural Networks](#).

Hannes Leitgeb showed that feed-forward neural networks are complete with respect to certain conditional laws of  $\Rightarrow$ . But  $\varphi \Rightarrow \psi$  just reads “ $\psi \subseteq \text{Prop}(\varphi)$ ” (i.e.  $\psi$  is in the propagation of the signal  $\varphi$ ), which we can re-write in modal language as  $\mathbf{T}\varphi \rightarrow \psi$ . In the same way that Hannes shows that feed-forward nets and preferential-conditional models are equivalent, it shouldn't be too hard at all to show that feed-forwards nets and neighborhood models are equivalent. (Note that it is well-known that neighborhood models are a generalization of preferential models.)

I also think I should be able to throw in a **K** modality (graph-reachability) in there, almost for free.

**Why bother with completeness?.** In formal specifications (of AI agents, or otherwise), we're often content with just listing some sound rules or behaviors that the agent will always follow. And it's definitely cool to see that neural networks satisfy some sound logical axioms. But if we want to fundamentally bridge the gap between logic and neural networks, we should set our aim higher: Towards *complete* logical characterizations of neural networks.

A more practical reason: Completeness gives us model-building, i.e. given a specification  $\Gamma$ , we can *build* a neural network  $\mathcal{N}$  satisfying  $\Gamma$ .

**Why bother with this modal language?.** Almost all of the previous work bridging logic and neural networks has focused on neural net models of *conditionals*. In some sense, doing this in modal language is just a re-write of this old work. But this previous work hasn't addressed how *learning* or *update* in neural networks can be cast in logical terms. This is not merely due to circumstance — integrating conditionals with update is a long-standing controversial issue. So instead, we believe that it is more natural to work with modalities (instead of conditionals), because

*Modal language natively supports update.*

In other words, our modal setting sets us up to easily cast update operators (e.g. neural network learning) as modal operators in our logic.

Also this gives me an excuse to title a paper *Neural Network Models à la Mode :-)* (This is a play on both modal logic and also bringing some old work back in style!)

And LOL I can name a section “Learning: The Cherry on Top”

## Related Papers:

### Neural Network Semantics / Semantic Encodings.

**Classic Papers.** [17] [11]

**Conditional Logic (Feedforward Net).** [2], [14], [15], [7] (soundness), [8] (model-building)

[Any other relevant work by the Garcez lab?]

**Description Logic w. Typicality.** [9], [10] [Any other relevant work by the Giordano lab?]

**Modal Logic w. Typicality.** [13]

[Any other big trends I'm missing? See the new survey by Odense + Garcez!]

**Miscellaneous.** [4], [5]

**Surveys.** [18] [1], [20], [12], [16], [3], [21] (the first few sections are a great introduction to Neural Network Semantics)

### Help with Technical Details.

**Neighborhood Models.** [19]

**Temporal Logic Rules.** [6]

# 1 Interpreted Neural Nets

## 1.1 Basic Definitions

DEFINITION 1.1. An **interpreted ANN** (Artificial Neural Network) is a pointed directed graph  $\mathcal{N} = \langle N, E, W, A, I \rangle$ , where

- $N$  is a finite nonempty set (the set of **neurons**)
- $E \subseteq N \times N$  (the set of **excitatory neurons**)
- $W: E \rightarrow \mathbb{R}$  (the **weight** of a given connection)
- $A$  is a function which maps each  $n \in N$  to  $A^{(n)}: \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}$  (the **activation function** for  $n$ , where  $k$  is the indegree of  $n$ )
- $I: \text{propositions} \rightarrow \mathcal{P}(N)$  is an assignment of propositions to sets of neurons (the **interpretation function**).

DEFINITION 1.2. A **BFNN** (Binary Feedforward Neural Network) is an interpreted ANN  $\mathcal{N} = \langle N, E, W, A, I \rangle$  that is

- **Feed-forward**:  $E$  does not contain any cycles
- **Binary**: the output of each neuron is in  $\{0, 1\}$
- $A^{(n)}$  is **zero at zero** in the first parameter:  $A^{(n)}(\vec{0}, \vec{w}) = 0$

DEFINITION 1.3. Given a BFNN  $\mathcal{N}$ ,  $\text{Set} = \mathcal{P}(N) = \{S \mid S \subseteq N\}$

DEFINITION 1.4. For  $S \in \text{Set}$ , let  $\chi_S: N \rightarrow \{0, 1\}$  be given by  $\chi_S = 1$  iff  $n \in S$

We write  $W_{ij}$  to mean  $W(i, j)$  for  $(i, j) \in E$ . To keep the notation from getting really messy, I'll also define:

DEFINITION 1.5. Let  $S \in \text{Set}$ ,  $\vec{m} = m_1, \dots, m_k$  be a sequence where each  $m_i \in N$ , and let  $n \in N$ . Then:

$$\text{Act}_S(\vec{m}, n) = A^{(n)}((\chi_S(m_1), \dots, \chi_S(m_k)); (W(m_1, n), \dots, W(m_k, n)))$$

i.e. the  $m_i \in S$  subsequently “activate”  $n$ .

PROPOSITION 1.6. Let  $S_1, S_2 \in \text{Set}$ ,  $\vec{m} = m_1, \dots, m_k$  be a sequence where each  $m_i \in N$ , and let  $n \in N$ . Suppose that  $S_1$  and  $S_2$  agree on all  $m_i$ , i.e. for all  $1 \leq i \leq k$ ,  $m_i \in S_1$  iff  $m_i \in S_2$ . Then

$$\text{Act}_{S_1}(\vec{m}, n) = \text{Act}_{S_2}(\vec{m}, n)$$

**Proof.** We have:

$$\begin{aligned} \text{Act}_{S_1}(\vec{m}, n) &= A^{(n)}((\chi_{S_1}(m_1), \dots, \chi_{S_1}(m_k)); (W(m_1, n), \dots, W(m_k, n))) \\ &= A^{(n)}((\chi_{S_2}(m_1), \dots, \chi_{S_2}(m_k)); (W(m_1, n), \dots, W(m_k, n))) \\ &= \text{Act}_{S_2}(\vec{m}, n) \end{aligned}$$

□

## 1.2 Prop and Reach

DEFINITION 1.7. (Adapted from [14, Definition 3.4]) Let  $\text{Prop}: \text{Set} \rightarrow \text{Set}$  be defined recursively as follows:  $n \in \text{Prop}(S)$  iff either

**Base Case.**  $n \in S$ , or

**Constructor.** For those  $\vec{m} = m_1, \dots, m_k$  such that  $(m_i, n) \in E$ ,  $\text{Act}_{\text{Prop}(S)}(\vec{m}, n) = 1$ .

DEFINITION 1.8. Let  $\text{Reach}: \text{Set} \rightarrow \text{Set}$  be given by  $\text{Reach}(S) = \{n \mid \exists m \in S \text{ with } E\text{-path from } m \text{ to } n\}$

DEFINITION 1.9. Let  $\text{Reach}^\downarrow: \text{Set} \rightarrow \text{Set}$  be given by  $\text{Reach}^\downarrow(S) = \{m \mid \exists n \in S \text{ with } E\text{-path from } m \text{ to } n\}$

PROPOSITION 1.10. For all  $S_1, \dots, S_k \in \text{Set}$ ,  $\bigcup_i \text{Reach}(S_i) = \text{Reach}(\bigcup_i S_i)$

**Proof.**

$$\begin{aligned}
 n \in \bigcup_i \text{Reach}(S_i) & \quad \text{iff} \quad \exists S_i \text{ with } n \in \text{Reach}(S_i) & \quad (\text{by definition of union}) \\
 & \quad \text{iff} \quad \exists S_i, \exists m \in S_i \text{ with } E\text{-path from } m \text{ to } n & \quad (\text{by definition of Reach}) \\
 & \quad \text{iff} \quad \exists m \in \bigcup_i S_i \text{ with } E\text{-path from } m \text{ to } n & \quad (\text{by definition of union}) \\
 & \quad \text{iff} \quad n \in \text{Reach}\left(\bigcup_i S_i\right) & \quad (\text{by definition of Reach})
 \end{aligned}$$

□

PROPOSITION 1.11. For all  $S \in \text{Set}$ ,  $\text{Reach}^\downarrow(S) = \bigcup_{n \in S} \bigcap_{n \notin \text{Reach}(X)} X^c$

**Proof.** ( $\rightarrow$ ) Suppose  $u \in \text{Reach}^\downarrow(S)$ . So  $\exists n \in S$  with  $E$ -path from  $u$  to  $n$ . Let  $X$  be such that  $n \notin \text{Reach}(X)$ . By definition of  $\text{Reach}$ , there is no  $m \in X$  with an  $E$ -path from  $m$  to  $n$ . But since there *is* such a path from  $u$ , we must have  $u \notin X$ , i.e.  $u \in X^c$ . Since  $X$  was arbitrary,  $u \in \bigcap_{n \notin \text{Reach}(X)} X^c$ . So  $u \in \bigcup_{n \in S} \bigcap_{n \notin \text{Reach}(X)} X^c$ .

( $\leftarrow$ ) Suppose  $u \in \bigcup_{n \in S} \bigcap_{n \notin \text{Reach}(X)} X^c$ . Let  $n \in S$  be such that for all  $X$ , if  $n \notin \text{Reach}(X)$  then  $u \in X^c$ . Consider in particular

$$X = \{m \mid \text{there is an } E\text{-path from } m \text{ to } n\}^c$$

Notice that  $n \notin \text{Reach}(X)$  (since  $X$  is the set of nodes where there is *not* a path to  $n$ ). And so  $u \in X^c$ , i.e. there *is* an  $E$ -path from  $u$  to  $n$ . □

PROPOSITION 1.12. Let  $\mathcal{N} \in \text{Net}$ . For all  $S, S_1, S_2 \in \text{Set}$ ,  $n, m \in N$ ,  $\text{Reach}$  is

**(Inclusive).**  $S \subseteq \text{Reach}(S)$

**(Idempotent).**  $\text{Reach}(S) = \text{Reach}(\text{Reach}(S))$

**(Antisymmetric).** If  $m \in \text{Reach}(\{n\})$  and  $n \in \text{Reach}(\{m\})$  then  $n = m$ .

**(Monotonic).** If  $S_1 \subseteq S_2$  then  $\text{Reach}(S_1) \subseteq \text{Reach}(S_2)$

**Proof.** We check each in turn:

**(Inclusive).** If  $n \in S$ , then there is a trivial path from  $n \in S$  to itself. So  $n \in \text{Reach}(S)$ .

**(Idempotent).** The ( $\subseteq$ ) direction is just Inclusion. As for ( $\supseteq$ ), let  $n \in \text{Reach}(\text{Reach}(S))$ . So there is a path from some  $m \in \text{Reach}(S)$  to  $n$ . But since  $m \in \text{Reach}(S)$ , there is a path from some  $u \in S$  to  $m$ . But then we have a path from  $u \in S$  to  $n$ , and so  $n \in \text{Reach}(S)$ .

**(Acyclic).** Suppose  $m \in \text{Reach}(\{n\})$  and  $n \in \text{Reach}(\{m\})$ . By definition of  $\text{Reach}$ , there is an  $E$ -path from  $m$  to  $n$ , and another path from  $n$  to  $m$ . But  $\mathcal{N}$  is feed-forward, i.e.  $E$  contains no cycles! So we must have  $n = m$ .

**(Monotonic).** Let  $n \in \text{Reach}(S_1)$ . So there is a path from some  $m \in S_1$  to  $n$ . Since  $S_1 \subseteq S_2$ ,  $m \in S_2$ . But then we have a path from  $m \in S_2$  to  $n$ . And so  $n \in \text{Reach}(S_2)$ . □

PROPOSITION 1.13. (Adapted from [14, Remark 4]) Let  $\mathcal{N} \in \text{Net}$ . For all  $S, S_1, S_2 \in \text{Set}$ ,  $\text{Prop}$  is

**(Inclusive).**  $S \subseteq \text{Prop}(S)$

**(Idempotent).**  $\text{Prop}(S) = \text{Prop}(\text{Prop}(S))$

**(Contained in Reach).**  $\text{Prop}(S) \subseteq \text{Reach}(S)$

**Proof.** We check each in turn:

**(Inclusive).** Similar to the proof of Inclusion for Reach.

**(Idempotent).** The  $(\subseteq)$  direction is just Inclusion. As for  $(\supseteq)$ , let  $n \in \text{Prop}(\text{Prop}(S))$ , and proceed by induction on  $\text{Prop}(\text{Prop}(S))$ .

**Base Step.**  $n \in \text{Prop}(S)$ , and so we are done.

**Inductive Step.** For those  $\vec{m} = m_1, \dots, m_k$  such that  $(m_i, n) \in E$ ,

$$\text{Act}_{\text{Prop}(\text{Prop}(S))}(\vec{m}, n) = 1$$

By inductive hypothesis,  $m_i \in \text{Prop}(\text{Prop}(S))$  iff  $m_i \in \text{Prop}(S)$ . By Proposition 1.6,  $\text{Act}_{\text{Prop}(S)}(\vec{m}, n) = 1$ , and so  $n \in \text{Prop}(S)$ .

**(Contained in Reach).** Let  $n \in \text{Prop}(S)$ , and proceed by induction on Prop.

**Base Step.**  $n \in S$ . So  $n \in \text{Reach}(S)$ .

**Inductive Step.** For those  $\vec{m} = m_1, \dots, m_k$  such that  $(m_i, n) \in E$ ,

$$\text{Act}_{\text{Prop}(S)}(\vec{m}, n) = 1$$

Since  $A^{(n)}$  is zero at zero, we have  $m_i \in \text{Prop}(S)$  for some  $m = m_i$ . By inductive hypothesis,  $m \in \text{Reach}(S)$ . And since  $(m, n) \in E$ , by definition of Reach,  $n \in \text{Reach}(S)$ .  $\square$

PROPOSITION 1.14. **(Minimal Cause)** For all  $n \in N$ , if  $n \in \text{Reach}^\downarrow(T)$  then the following are equivalent:

1.  $n \in \text{Prop}(S)$
2.  $n \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$

**Proof.** Suppose  $n \in \text{Reach}^\downarrow(T)$ . To show  $(1 \rightarrow 2)$ , let  $n \in \text{Prop}(S)$  and proceed by induction on Prop.

**Base Step.**  $n \in S$ . By the base step of Prop,  $n \in \text{Prop}(S)$ . But we also have  $n \in \text{Reach}^\downarrow(T)$ , and so  $n \in \text{Prop}(S) \cap \text{Reach}^\downarrow(T)$ . By the base step of Prop,  $n \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$ .

**Inductive Step.** For those  $\vec{m} = m_1, \dots, m_k$  such that  $(m_i, n) \in E$ ,

$$\text{Act}_{\text{Prop}(S)}(\vec{m}, n) = 1$$

Since each  $(m_i, n) \in E$ , and  $n \in \text{Reach}^\downarrow(T)$ , by the constructor case of  $\text{Reach}^\downarrow$  each  $m_i \in \text{Reach}^\downarrow(T)$ . So we can apply our inductive hypothesis to each  $m_i$ :  $m_i \in \text{Prop}(S)$  iff  $m_i \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$ .

By Proposition 1.6,  $\text{Act}_{\text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))}(\vec{m}, n) = 1$ , and so  $n \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$ .

As for  $(2 \rightarrow 1)$ , let  $n \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$ , and proceed by induction on the outer Prop.

**Base Step.**  $n \in \text{Prop}(S) \cap \text{Reach}^\downarrow(T)$ . So in particular,  $n \in \text{Prop}(S)$

**Inductive Step.** For that  $\vec{m} = m_1, \dots, m_k$  such that  $(m_i, n) \in E$ ,

$$\text{Act}_{\text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))}(\vec{m}, n) = 1$$

Since each  $(m_i, n) \in E$ , and  $n \in \text{Reach}^\downarrow(T)$ , by the constructor case of  $\text{Reach}^\downarrow$  each  $m_i \in \text{Reach}^\downarrow(T)$ . So we can apply our inductive hypothesis to each  $m_i$ :  $m_i \in \text{Prop}(\text{Prop}(S) \cap \text{Reach}^\downarrow(T))$  iff  $m_i \in \text{Prop}(S)$ .

By Proposition 1.6,  $\text{Act}_{\text{Prop}(S)}(\vec{m}, n) = 1$ , and so  $n \in \text{Prop}(S)$ .  $\square$

PROPOSITION 1.15. The Cumulative and Loop properties from [14] [The KLM Cumulative & Loop properties, actually], i.e.

**(Cumulative).** If  $S_1 \subseteq S_2 \subseteq \text{Prop}(S_1)$  then  $\text{Prop}(S_1) \subseteq \text{Prop}(S_2)$

**(Loop).** If  $S_1 \subseteq \text{Prop}(S_0), \dots, S_n \subseteq \text{Prop}(S_{n-1})$  and  $S_0 \subseteq \text{Prop}(S_n)$ ,  
then  $\text{Prop}(S_i) = \text{Prop}(S_j)$  for all  $i, j \in \{0, \dots, n\}$

follow from the properties of Prop and Reach above.

**Proof.** [Todo – note that (Cumulative) actually follows from (Loop). Use acyclic property of Reach to get (Loop)]  $\square$

### 1.3 Neural Network Semantics

DEFINITION 1.16. Formulas of our language  $\mathcal{L}$  are given by

$$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{K}\varphi \mid \mathbf{T}\varphi$$

where  $p$  is any propositional variable, and  $i$  is any nominal (denoting a neuron). Material implication  $\varphi \rightarrow \psi$  is defined as  $\neg\varphi \vee \psi$ . We define  $\perp, \vee, \leftrightarrow, \Leftrightarrow$ , and the dual operators  $\langle \mathbf{K} \rangle, \langle \mathbf{T} \rangle$  in the usual way.

DEFINITION 1.17. Let  $\mathcal{N} \in \text{Net}$ . The semantics  $\llbracket \cdot \rrbracket : \mathcal{L} \rightarrow \text{Set}$  for  $\mathcal{L}$  are defined recursively as follows:

$\llbracket p \rrbracket$	$= I(p) \in \text{Set}$
$\llbracket \neg\varphi \rrbracket$	$= \llbracket \varphi \rrbracket^c$
$\llbracket \varphi \wedge \psi \rrbracket$	$= \llbracket \varphi \rrbracket \cap \llbracket \psi \rrbracket$
$\llbracket \langle \mathbf{K} \rangle \varphi \rrbracket$	$= \text{Reach}(\llbracket \varphi \rrbracket)$
$\llbracket \langle \mathbf{K}^\downarrow \rangle \varphi \rrbracket$	$= \text{Reach}^\downarrow(\llbracket \varphi \rrbracket)$
$\llbracket \langle \mathbf{T} \rangle \varphi \rrbracket$	$= \text{Prop}(\llbracket \varphi \rrbracket)$

DEFINITION 1.18. **(Truth at a neuron)**  $\mathcal{N}, n \models \varphi$  iff  $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}$ .

DEFINITION 1.19. **(Truth in a net)**  $\mathcal{N} \models \varphi$  iff  $\mathcal{N}, n \models \varphi$  for all  $n \in N$ .

DEFINITION 1.20. **(Entailment)**  $\Gamma \models_{\text{BFNN}} \varphi$  if for all BFNNs  $\mathcal{N}$  for all neurons  $n \in N$ , if  $\mathcal{N}, n \models \Gamma$  then  $\mathcal{N}, n \models \varphi$ .

## 2 Neighborhood Models

### 2.1 Basic Definitions

DEFINITION 2.1. [19, Definition 1.9] A **neighborhood frame** is a pair  $\mathcal{F} = \langle W, f \rangle$ , where  $W$  is a non-empty set of **worlds** and  $f: W \rightarrow \mathcal{P}(\mathcal{P}(W))$  is a **neighborhood function**. A **multi-frame** may have more than one neighborhood function, but to keep things simple I won't distinguish between frames and multi-frames.

DEFINITION 2.2. [19, Section 1.1] Let  $\mathcal{F} = \langle W, f \rangle$  be a neighborhood frame, and let  $w \in W$ . The set  $\bigcap_{X \in f(w)} X$  is called the **core of  $f(w)$** , abbreviated  $\cap f(w)$ . If  $X \subseteq W$ , the set  $\bigcup_{w \in X} \cap f(w)$  is called the **core of  $f(X)$** , abbreviated  $\cap f(X)$ .

DEFINITION 2.3. [19, Definition 1.4] Let  $\mathcal{F} = \langle W, f \rangle$  be a frame.  $\mathcal{F}$  is a **proper filter** iff:

- $f$  is **closed under finite intersections**: for all  $w \in W$ , if  $X_1, \dots, X_n \in f(w)$  then their intersection  $\bigcap_{i=1}^n X_i \in f(w)$
- $f$  is **closed under supersets**: for all  $w \in W$ , if  $X \in f(w)$  and  $X \subseteq Y \subseteq W$ , then  $Y \in f(w)$

- $f$  **contains the unit**: iff  $W \in f(w)$

PROPOSITION 2.4. [19, Corollary 1.1] If  $\mathcal{F} = \langle W, f \rangle$  is a filter, and  $W$  is finite, then  $\mathcal{F}$  contains its core.

DEFINITION 2.5.  $\mathcal{F} = \langle W, f, g \rangle$  is a **preferential filter** iff:

- $W$  is finite
- $\langle W, f \rangle$  forms a proper filter, and  $g$  contains the unit
- $f$  is **antisymmetric**: for all  $u, v \in W$ , if  $u \in \cap f(v)$  and  $v \in \cap f(u)$  then  $u = v$ .
- $f, g$  are **reflexive**: for all  $w \in W$ ,  $w \in \cap f(w)$  (similarly for  $g$ )
- $f, g$  are **transitive**: for all  $w \in W$ , if  $X \in f(w)$  then  $\{u \mid X \in f(u)\} \in f(w)$  (similarly for  $g$ )
- $g$  **contains**  $f$ : for all  $w \in W$ , if  $X \in f(w)$  then  $X \in g(w)$
- $f$  is a **skeleton** of  $g$ : for all  $w \in W$  and  $Y \subseteq W$  such that  $w \in \cap f(Y)$ ,

$$X \in g(w) \quad \text{iff} \quad \{u \mid X \in g(u)\} \cup (\cap f(Y))^c \in g(w)$$

PROPOSITION 2.6. Let  $\mathcal{F} = \langle W, f, g \rangle$  be a preferential filter. For all  $w \in W$ , we have in particular:

$$X \in g(w) \quad \text{iff} \quad \{u \mid X \in g(u)\} \cup (\cap f(w))^c \in g(w)$$

**Proof.** This is just the  $Y = \{w\}$  instance of the ‘skeleton’ property. Note that for  $Y = \{w\}$ ,

$$\cap f(Y) = \bigcup_{u \in \{w\}} f(u) = \cap f(w)$$

Since  $f$  is reflexive,  $w \in \cap f(w)$ . Since  $f$  is the skeleton of  $g$ , we have our conclusion.  $\square$

## 2.2 Neighborhood Semantics

DEFINITION 2.7. [19, Definition 1.11] Let  $\mathcal{F} = \langle W, f, g \rangle$  be a neighborhood frame. A **neighborhood model** based on  $\mathcal{F}$  is  $\mathcal{M} = \langle W, f, g, V \rangle$ , where  $V: \mathcal{L} \rightarrow \mathcal{P}(W)$  is a valuation function.

DEFINITION 2.8. [19, Definition 1.12] Let  $\mathcal{M} = \langle W, f, g, V \rangle$  be a model based on  $\mathcal{F} = \langle W, f, g \rangle$ . The (neighborhood) semantics for  $\mathcal{L}$  are defined recursively as follows:

$\mathcal{M}, w \Vdash p$	iff	$w \in V(p)$
$\mathcal{M}, w \Vdash \neg \varphi$	iff	$\mathcal{M}, w \not\Vdash \varphi$
$\mathcal{M}, w \Vdash \varphi \wedge \psi$	iff	$\mathcal{M}, w \Vdash \varphi$ and $\mathcal{M}, w \Vdash \psi$
$\mathcal{M}, w \Vdash \mathbf{K}\varphi$	iff	$\{u \mid \mathcal{M}, u \Vdash \varphi\} \in f(w)$
$\mathcal{M}, w \Vdash \mathbf{K}^\downarrow \varphi$	iff	$\forall u \in W$ , if $w \in \cap f(u)$ then $\mathcal{M}, u \Vdash \varphi$
$\mathcal{M}, w \Vdash \mathbf{T}\varphi$	iff	$\{u \mid \mathcal{M}, u \Vdash \varphi\} \in g(w)$

In neighborhood semantics, the operators  $\mathbf{K}$ , and  $\mathbf{T}$  are more natural to interpret. But when we gave our neural semantics, we instead interpreted the *duals*  $\langle \mathbf{K} \rangle$ , and  $\langle \mathbf{T} \rangle$ . Since we need to relate the two, I'll write the explicit neighborhood semantics for the duals here:

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{K} \rangle \varphi & \quad \text{iff} \quad \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin f(w) \\ \mathcal{M}, w \Vdash \langle \mathbf{K}^\downarrow \rangle \varphi & \quad \text{iff} \quad \exists u \in W \text{ such that } w \in \cap f(u) \text{ and } \mathcal{M}, u \not\Vdash \varphi \\ \mathcal{M}, w \Vdash \langle \mathbf{T} \rangle \varphi & \quad \text{iff} \quad \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin g(w) \end{aligned}$$

DEFINITION 2.9. [19, Definition 1.13] (**Truth in a model**)  $\mathcal{M} \models \varphi$  iff  $\mathcal{M}, w \Vdash \varphi$  for all  $w \in W$ .

DEFINITION 2.10. [19, Definition 2.32] (**Entailment**) Let  $F$  be a collection of neighborhood frames.  $\Gamma \models_F \varphi$  if for all models  $\mathcal{M}$  based on a frame  $\mathcal{F} \in F$  and for all worlds  $w \in W$ , if  $\mathcal{M}, w \Vdash \Gamma$  then  $\mathcal{M}, w \Vdash \varphi$ .

**Note.** This is the *local* consequence relation in modal logic.

### 3 From Nets to Frames

**This is the easy (“soundness”) direction!**

DEFINITION 3.1. Given a BFNN  $\mathcal{N}$ , its **simulation frame**  $\mathcal{F}^\bullet = \langle W, f, g \rangle$  is given by:

- $W = N$
- $f(w) = \{S \subseteq W \mid w \notin \text{Reach}(S^c)\}$
- $g(w) = \{S \subseteq W \mid w \notin \text{Prop}(S^c)\}$

Moreover, the **simulation model**  $\mathcal{M}^\bullet = \langle W, f, g, V \rangle$  based on  $\mathcal{F}^\bullet$  has:

- $V(p) = I(p)$

THEOREM 3.2. Let  $\mathcal{N}$  be a BFNN, and let  $\mathcal{M}^\bullet$  be the simulation model based on  $\mathcal{F}^\bullet$ . Then for all  $w \in W$ ,

$$\mathcal{M}^\bullet, w \Vdash \varphi \quad \text{iff} \quad \mathcal{N}, w \Vdash \varphi$$

**Proof.** By induction on  $\varphi$ . The propositional,  $\neg\varphi$ , and  $\varphi \wedge \psi$  cases are trivial.

**$\langle \mathbf{K} \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}^\bullet, w \Vdash \langle \mathbf{K} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}^\bullet, u \not\Vdash \varphi\} \notin f(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket\} \notin f(w) \text{ (IH)} \\ & \text{ iff } \llbracket \varphi \rrbracket^c \notin f(w) \\ & \text{ iff } w \in \text{Reach}(\llbracket (\varphi^c)^c \rrbracket) \text{ (by choice of } f) \\ & \text{ iff } w \in \text{Reach}(\llbracket \varphi \rrbracket) \\ & \text{ iff } w \in \llbracket \langle \mathbf{K} \rangle \varphi \rrbracket \text{ (by definition)} \\ & \text{ iff } \mathcal{N}, w \Vdash \langle \mathbf{K} \rangle \varphi \text{ (by definition)} \end{aligned}$$

**$\langle \mathbf{K}^\downarrow \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}^\bullet, w \Vdash \langle \mathbf{K}^\downarrow \rangle \varphi & \text{ iff } \exists u \text{ such that } w \in \cap f(u) \text{ and } \mathcal{M}^\bullet, u \Vdash \varphi \text{ (by definition)} \\ & \text{ iff } \exists u \text{ such that } w \in \cap f(u) \text{ and } u \in \llbracket \varphi \rrbracket \text{ (IH)} \\ & \text{ iff } \exists u \in \llbracket \varphi \rrbracket \text{ such that } w \in \bigcap_{X \in f(u)} X \\ & \text{ iff } \exists u \in \llbracket \varphi \rrbracket \text{ such that } w \in \bigcap_{u \notin \text{Reach}(X^c)} X \text{ (by choice of } f) \\ & \text{ iff } w \in \text{Reach}^\downarrow(\llbracket \varphi \rrbracket) \text{ (by definition)} \\ & \text{ iff } \mathcal{N}, w \Vdash \langle \mathbf{K}^\downarrow \rangle \varphi \text{ (by definition)} \end{aligned}$$

**$\langle \mathbf{T} \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}^\bullet, w \Vdash \langle \mathbf{T} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}^\bullet, u \not\Vdash \varphi\} \notin g(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket\} \notin g(w) \text{ (IH)} \\ & \text{ iff } \llbracket \varphi \rrbracket^c \notin g(w) \\ & \text{ iff } w \in \text{Prop}(\llbracket (\varphi^c)^c \rrbracket) \text{ (by choice of } g) \\ & \text{ iff } w \in \text{Prop}(\llbracket \varphi \rrbracket) \\ & \text{ iff } w \in \llbracket \langle \mathbf{T} \rangle \varphi \rrbracket \text{ (by definition)} \\ & \text{ iff } \mathcal{N}, w \Vdash \langle \mathbf{T} \rangle \varphi \text{ (by definition)} \end{aligned}$$



□

COROLLARY 3.3.  $\mathcal{M}^\bullet \models \varphi$  iff  $\mathcal{N} \models \varphi$ .

THEOREM 3.4.  $\mathcal{F}^\bullet$  is a preferential filter.

**Proof.** We show each in turn:

**W is finite.** This holds because our BFNN is finite.

**$f$  is closed under finite intersection.** Suppose  $X_1, \dots, X_n \in f(w)$ . By definition of  $f$ ,  $w \notin \bigcup_i \text{Reach}(X_i^c)$  for all  $i$ . By Proposition 1.10 we have  $\bigcup_i \text{Reach}(X_i^c) = \text{Reach}(\bigcup_i X_i^c) = \text{Reach}((\bigcap_i X_i)^c)$  (note that this is where we use the fact that  $\text{Reach}$  is monotonic). So  $w \notin \text{Reach}((\bigcap_i X_i)^c)$ . But this means that  $\bigcap_i X_i \in f(w)$ .

**$f$  is closed under superset.** Suppose  $X \in f(w)$ ,  $X \subseteq Y$ . By definition of  $f$ ,  $w \notin \text{Reach}(X^c)$ . Note that  $Y^c \subseteq X^c$ , and so by monotonicity of  $\text{Reach}$  we have  $w \notin \text{Reach}(Y^c)$ . But this means  $Y \in f(w)$ , so we are done.

**$f$  contains the unit.** Note that for all  $w \in W$ ,  $w \notin \text{Reach}(\emptyset) = \text{Reach}(W^c)$ . So  $W \in f(w)$ .

**$g$  contains the unit.** Same as the proof for  $f$ , except that we use the fact that for all  $w$ ,  $w \notin \text{Prop}(\emptyset)$ .

**$f$  is antisymmetric.** Suppose  $u \in \cap f(v)$  and  $v \in \cap f(u)$ . Expanding the definition of core,  $u \in \bigcap_{X \in f(v)} X$ , and  $v \in \bigcap_{X \in f(u)} X$ . By definition of  $f$ ,  $u \in \bigcap_{v \notin \text{Reach}(X^c)} X$  and  $v \in \bigcap_{u \notin \text{Reach}(X^c)} X$ . Substituting  $X^c$  for  $X$  we get  $u \in \bigcap_{v \notin \text{Reach}(X)} X^c$  and  $v \in \bigcap_{u \notin \text{Reach}(X)} X^c$ . By Proposition 1.11,  $u \in \text{Reach}^\downarrow(\{v\})$  and  $v \in \text{Reach}^\downarrow(\{u\})$ . But by the definition of  $\text{Reach}^\downarrow$ , this just means there is a path from  $u$  to  $v$  and a path from  $v$  to  $u$  in  $E$ , i.e.  $v \in \text{Reach}(\{u\})$  and  $u \in \text{Reach}(\{v\})$  by Proposition ?. So  $u = v$  by Antisymmetry of  $\text{Reach}$ .

**$f$  is reflexive.** We want to show that  $w \in \cap f(w)$ . Well, suppose  $X \in f(w)$ , i.e.  $w \notin \text{Reach}(X^c)$  (by definition of  $f$ ). Since for all  $S$ ,  $S \subseteq \text{Reach}(S)$ , we have  $w \notin X^c$ . But this means  $w \in X$ , and we are done.

**$g$  is reflexive.** Same as the proof for  $f$ , except we use the fact that for all  $S$ ,  $S \subseteq \text{Prop}(S)$ .

**$f$  is transitive.** Suppose  $X \in f(w)$ , i.e.  $w \notin \text{Reach}(X^c)$ . Well,

$$\begin{aligned} \text{Reach}(X^c) &= \text{Reach}(\text{Reach}(X^c)) && \text{(by Idempotence of Reach)} \\ &= \text{Reach}(\{u \mid u \in \text{Reach}(X^c)\}) \\ &= \text{Reach}(\{u \mid u \notin \text{Reach}(X^c)\}^c) \\ &= \text{Reach}(\{u \mid X \in f(u)\}^c) && \text{(by definition of } f) \end{aligned}$$

So by definition of  $f$ ,  $\{u \mid X \in f(u)\} \in f(w)$ .

**$g$  is transitive.** Same as the proof for  $f$ , except we use the fact that  $\text{Prop}$  is idempotent.

**$g$  contains  $f$ .** Suppose  $X \in f(w)$ , i.e.  $w \notin \text{Reach}(X^c)$ . Since for all  $S$ ,  $\text{Prop}(S) \subseteq \text{Reach}(S)$ , we have  $w \notin \text{Prop}(X^c)$ . And so  $X \in g(w)$ , and we are done.

**$f$  is the skeleton of  $g$ .** Suppose  $w \in \cap f(Y)$ . We will show the  $(\leftarrow)$  direction; the other direction is similar. Suppose  $\{u \mid X \in g(u)\} \cup (\cap f(Y))^c \in g(w)$ . By choice of  $g$ ,  $w \notin \text{Prop}([\{u \mid X \in g(u)\} \cup (\cap f(Y))^c]^c)$ . Distributing the outer complement, we have  $w \notin \text{Prop}(\{u \mid X \notin g(u)\} \cap (\cap f(Y)))$ . Again by choice of  $g$ ,  $w \notin \text{Prop}(\{u \mid u \in \text{Prop}(X^c)\} \cap (\cap f(Y))) = \text{Prop}(\text{Prop}(X^c) \cap (\cap f(Y)))$ . By choice of  $f$ ,  $w \notin \text{Prop}(\text{Prop}(X^c) \cap \bigcup_{w \in Y} (\bigcap_{w \notin \text{Reach}(Y^c)} Y))$ . Substituting  $Y^c$  for  $Y$ , we get  $w \notin \text{Prop}(\text{Prop}(X^c) \cap \bigcup_{w \in Y^c} (\bigcap_{w \notin \text{Reach}(Y)} Y^c))$ . By our alternative characterization of  $\text{Reach}^\downarrow$  (Proposition 1.11),  $w \notin \text{Prop}(\text{Prop}(X^c) \cap \text{Reach}^\downarrow(Y^c))$ .

Similarly,  $w \in \cap f(Y)$  gives us  $w \in \text{Reach}^\downarrow(Y^c)$ , the precondition of Minimal Cause (Proposition 1.14). By Minimal Cause, we conclude that  $w \notin \text{Prop}(X^c)$ , i.e.  $X \in g(w)$ . □



## 4 From Frames to Nets

**This is the harder (“completeness”) direction!**

DEFINITION 4.1. Let  $\mathcal{M}$  be a model based on preferential filter  $\mathcal{F} = \langle W, f, g \rangle$ . Its **simulation net**  $\mathcal{N}^\bullet = \langle N, E, W, A, I \rangle$  is the BFNN given by:

- $N = W$
- $(u, v) \in E$  iff  $u \in \cap f(v)$

Now let  $m_1, \dots, m_k$  list those nodes such that  $(m_i, n) \in E$ .

- $W(m_i, n) = \text{[Does not matter; arbitrary]}$
- $A^{(n)}(\vec{x}, \vec{w}) = 1$  iff  $\{m_i | x_i = 1\}^c \notin g(n)$
- $I(p) = V(p)$

CLAIM 4.2.  $\mathcal{N}^\bullet$  is a BFNN.

**Proof.** Clearly  $\mathcal{N}^\bullet$  is a binary ANN. We check the rest of the conditions:

**$\mathcal{N}^\bullet$  is feed-forward.** Suppose that  $E$  contains a cycle, i.e.  $n_1, \dots, n_k \in N$  such that  $n_1 E n_2, \dots, n_{k-1} E n_k, n_k E n_1$ . We show that each  $n_i = n_j$  by induction on  $k$ .

**Base Case.**  $k = 1$ , and so trivially  $n_1 = n_k$ .

**Inductive Case.** Let  $k \geq 1$ . By our inductive hypothesis,  $n_1 = \dots = n_{k-1}$ . We will show in particular that  $n_1 = n_k$  (the other cases are similar). Since  $n_{k-1} E n_k$  and  $n_1 = n_{k-1}$ , we have  $n_1 E n_k$ . From our earlier assumptions, we also have  $n_k E n_1$ . By definition of  $E$ ,  $n_1 \in \cap f(n_k)$ , and  $n_k \in \cap f(n_1)$ . Since  $f$  is antisymmetric,  $n_1 = n_k$ .

**$O^{(n)} \circ A^{(n)}$  is zero at zero.** Suppose for contradiction that  $A^{(v)}(\vec{0}, \vec{w}) = 1$ . Then  $\emptyset^c = W \notin g(v)$ , which contradicts the fact that  $f$  contains the unit.  $\square$

LEMMA 4.3.  $\text{Reach}_{\mathcal{N}^\bullet}(S) = \{n | S^c \notin f(n)\}$

**Proof.** For the  $(\supseteq)$  direction, suppose  $S^c \notin f(n)$ . We claim that  $\cap f(n) \not\subseteq S^c$ . Why not? If  $\cap f(n) \subseteq S^c$ , we have  $\cap f(n) \in f(n)$  (since  $f$  is closed under finite intersection) and so  $S^c \in f(n)$  (since  $f$  is closed under superset). This would contradict  $S^c \notin f(n)$ .

So  $\cap f(n) \not\subseteq S^c$ . This means that there is some  $m \in \cap f(n)$  such that  $m \notin S^c$ . That is,  $(m, n) \in E$  and  $m \in S$ . But then  $m \in \text{Reach}_{\mathcal{N}^\bullet}(S)$ . So we have  $m \in \text{Reach}_{\mathcal{N}^\bullet}(S)$  and a path  $(m, n) \in E$  from  $m$  to  $n$ , i.e.  $n \in \text{Reach}_{\mathcal{N}^\bullet}(S)$ .

Now for the  $(\subseteq)$  direction. Suppose  $n \in \text{Reach}_{\mathcal{N}^\bullet}(S)$ . So there is a path from some  $m \in S$  to  $n$ . We proceed by induction on the length  $l$  of this path.

**Base step.**  $l = 0$ , i.e.  $m = n$ , which gives us  $n \in S$ . Suppose for contradiction that  $S^c \in f(n)$ . By definition of core,  $\cap f(n) \subseteq S^c$ . But since  $\mathcal{F}$  is reflexive,  $n \in \cap f(n)$ . So  $n \in S^c$ , which contradicts  $n \in S$ .

**Inductive step.** Let  $l \geq 0$ . Let  $u$  immediately precede  $n$  on this path i.e. there is a path from  $m$  to  $u$  and  $(u, n) \in E$  (and so  $u \in \cap f(n)$ ). Note that  $u \in \text{Reach}_{\mathcal{N}^\bullet}(S)$ , and so by our inductive hypothesis  $S^c \notin f(u)$ . Now suppose for contradiction that  $S^c \in f(n)$ . Since  $f$  is transitive,  $\{t | S^c \in f(t)\} \in f(n)$ . By definition of core,  $\cap f(n) \subseteq \{t | S^c \in f(t)\}$ . Since  $u \in \cap f(n)$ ,  $S^c \in f(u)$ . But this contradicts  $S^c \notin f(u)$ !  $\square$

LEMMA 4.4.  $\text{Prop}_{\mathcal{N}^\bullet}(S) = \{n | S^c \notin g(n)\}$

**Proof.** First, let's consider the  $(\supseteq)$  direction. Since  $\mathcal{N}^\bullet$  is feed-forward (i.e. acyclic), we can perform a topological sort on its nodes to get a well-ordering  $<$  over  $N$  such that for all  $n \neq m \in N$ ,

$$\text{If } (m, n) \in E \text{ then } m < n$$

Let  $n \in N$  be such that  $S^c \notin g(n)$ . We proceed by induction on the ordering  $<$ .

**Base Step.** There are no nodes  $m$  such that  $m < n$ , and hence no nodes  $m \neq n$  with  $(m, n) \in E$ . We also have  $S^c \notin g(n)$  from before.

CLAIM.  $n \in S$  (and so we can conclude that  $n \in \text{Prop}_{N^\bullet}(S)$  by the base case of Prop)

**Proof.**  $S^c \notin g(n)$ , and so  $S^c \notin f(n)$  (since  $g$  contains  $f$ ). Since  $f$  is closed under superset and finite intersection,  $\cap f(n) \not\subseteq S^c$ , i.e. there is some  $u \in \cap f(n)$  such that  $u \in S$ .  $u \in \cap f(n)$  gives us  $(u, n) \in E$  (by choice of  $E$ ), but there are no  $u \neq n$  with  $(u, n) \in E$ . But that means that  $u = n$  — and so  $n \in S$ !  $\square$

**Inductive Step.** Let  $\vec{m} = m_1, \dots, m_k$  be all those nodes  $m_i \neq n$  such that  $(m_i, n) \in E$ . In particular, each  $m_i < n$ , and so we can apply our Inductive Hypothesis to each  $m_i$ :

**Inductive Hypothesis.**  $m_i \in \text{Prop}_{N^\bullet}(S)$  iff  $S^c \notin g(m_i)$

CLAIM 4.5.  $\{m_i \mid m_i \in \text{Prop}_{N^\bullet}(S)\} = \{u \mid S^c \notin g(u)\} \cap (\cap f(n))$

**Proof.**

$$\begin{aligned} \{m_i \mid m_i \in \text{Prop}_{N^\bullet}(S)\} &= \{m_i \mid S^c \notin g(m_i)\} && \text{(by Inductive Hypothesis)} \\ &= \{u \mid S^c \notin g(u) \text{ and } (u, n) \in E\} && \text{(by choice of } m_1, \dots, m_k) \\ &= \{u \mid S^c \notin g(u) \text{ and } u \in \cap f(n)\} && \text{(by choice of } E) \\ &= \{u \mid S^c \notin g(u)\} \cap (\cap f(n)) \end{aligned} \quad \square$$

CLAIM 4.6.  $\text{Act}_S(\vec{m}, n) = 1$  iff  $\{m_i \mid m_i \in S\}^c \notin g(n)$

**Proof.**  $\text{Act}_S(\vec{m}, n) = 1$  iff:

$$\begin{aligned} &A^{(n)}((\chi_S(m_1), \dots, \chi_S(m_k)); (W(m_1, n), \dots, W(m_k, n))) = 1 && \text{(by definition of Act)} \\ \text{iff} &\{m_i \mid \chi_S(m_i) = 1\}^c \notin g(n) && \text{(by our choice of } A^{(n)}) \quad \square \\ \text{iff} &\{m_i \mid m_i \in S\}^c \notin g(n) && \text{(by definition of } \chi_S) \end{aligned}$$

Putting everything together, we have:

$$\begin{aligned} S^c \notin g(n) &\rightarrow \{u \mid S^c \in g(u)\} \cup (\cap f(n))^c \notin g(n) && \text{(by Proposition 2.6, since } f \text{ is a skeleton of } g) \\ &\rightarrow [\{u \mid S^c \notin g(u)\} \cap (\cap f(n))]^c && \text{(distributing the complement)} \\ &\rightarrow \{m_i \mid m_i \in \text{Prop}_{N^\bullet}(S)\}^c \notin g(n) && \text{(by Claim 4.5 above)} \\ &\rightarrow \text{Act}_{\text{Prop}_{N^\bullet}(S)}(\vec{m}, n) = 1 && \text{(by Claim 4.6 above)} \\ &\rightarrow n \in \text{Prop}_{N^\bullet}(S) && \text{(by the constructor of Prop)} \end{aligned}$$

As for the  $(\subseteq)$  direction, suppose  $n \in \text{Prop}_{N^\bullet}(S)$ , and proceed by induction on Prop.

**Base step.**  $n \in S$ . Suppose for contradiction that  $S^c \in g(n)$ . Since  $g$  is reflexive,  $n \in \cap g(n)$ . By definition of core, we have  $\cap g(n) \subseteq S^c$ . But then  $n \in \cap g(n) \subseteq S^c$ , i.e.  $n \in S^c$ , which contradicts  $n \in S$ .

**Inductive step.** Let  $\vec{m} = m_1, \dots, m_k$  list those nodes such that  $(u_i, v) \in E$ . We have

$$\text{Act}_{\text{Prop}_{N^\bullet}(S)}(\vec{m}, n) = 1$$

Re-using Claim 4.6 above, this means that  $\{m_i \mid m_i \in \text{Prop}_{N^\bullet}(S)\}^c \notin g(n)$ . But by our inductive hypothesis,  $\{m_i \mid m_i \in \text{Prop}_{N^\bullet}(S)\} = \{m_i \mid S^c \notin g(m_i)\}$ . For convenience, let  $T$  be this latter set, i.e.  $T = \{m_i \mid S^c \notin g(m_i)\}$ . So we have  $T^c \notin g(n)$ .

We would like to show that  $S^c \notin g(n)$ . Suppose for contradiction that  $S^c \in g(n)$ . Notice that, by definition of  $T$ ,  $T^c = \{u_i \mid S^c \in g(u_i)\}$ . Since  $S^c \in g(v)$  and  $\mathcal{G}$  is transitive,  $T^c \in g(v)$ , which contradicts  $T^c \notin g(v)$ .

□

**THEOREM 4.7.** Let  $\mathcal{M}$  be a model based on a preferential filter  $\mathcal{F}$ , and let  $\mathcal{N}^\bullet$  be the corresponding simulation net. We have, for all  $w \in W$ ,

$$\mathcal{M}, w \Vdash \varphi \quad \text{iff} \quad \mathcal{N}^\bullet, w \Vdash \varphi$$

**Proof.** By induction on  $\varphi$ . Again, the propositional,  $\neg\varphi$ , and  $\varphi \wedge \psi$  cases are trivial.

**$\langle \mathbf{K} \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{K} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin f(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}\} \notin f(w) \text{ (Inductive Hypothesis)} \\ & \text{ iff } \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}^c \notin g(w) \\ & \text{ iff } w \in \text{Reach}_{\mathcal{N}^\bullet}(\llbracket \varphi \rrbracket) \text{ (by Lemma 4.3)} \\ & \text{ iff } w \in \llbracket \langle \mathbf{K} \rangle \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ (by definition)} \\ & \text{ iff } \mathcal{N}^\bullet, w \Vdash \langle \mathbf{K} \rangle \varphi \text{ (by definition)} \end{aligned}$$

**$\langle \mathbf{K}^\downarrow \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{K}^\downarrow \rangle \varphi & \text{ iff } \exists u \text{ such that } w \in \cap f(u) \text{ and } \mathcal{M}, u \Vdash \varphi \text{ (by definition)} \\ & \text{ iff } \exists u \text{ such that } w \in \cap f(u) \text{ and } u \in \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ (IH)} \\ & \text{ iff } \exists u \in \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ such that } w \in \bigcap_{X \in f(u)} X \\ & \quad \exists u \in \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ such that } w \in \bigcap_{u \notin \text{Reach}_{\mathcal{N}^\bullet}(X^c)} X \text{ (by Lemma 4.3)} \\ & \quad w \in \text{Reach}^\downarrow(\llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}) \text{ (by Proposition 1.11)} \\ & \text{ iff } \mathcal{N}^\bullet, w \Vdash \langle \mathbf{K}^\downarrow \rangle \varphi \text{ (by definition)} \end{aligned}$$

**$\langle \mathbf{T} \rangle \varphi$  case:**

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{T} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin g(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}\} \notin g(w) \text{ (Inductive Hypothesis)} \\ & \text{ iff } \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}^c \notin g(w) \\ & \text{ iff } w \in \text{Prop}_{\mathcal{N}^\bullet}(\llbracket \varphi \rrbracket) \text{ (by Lemma 4.4)} \\ & \text{ iff } w \in \llbracket \langle \mathbf{T} \rangle \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ (by definition)} \\ & \text{ iff } \mathcal{N}^\bullet, w \Vdash \langle \mathbf{T} \rangle \varphi \text{ (by definition)} \end{aligned}$$

□

**COROLLARY 4.8.**  $\mathcal{M} \models \varphi$  iff  $\mathcal{N}^\bullet \models \varphi$ .

## 5 Completeness

### 5.1 The Base Modal Logic

**DEFINITION 5.1.** Our logic **L** is the smallest set of formulas in  $\mathcal{L}$  containing the axioms

- (DISTR<sub>K</sub>)  $\mathbf{K}(\varphi \rightarrow \psi) \rightarrow (\mathbf{K}\varphi \rightarrow \mathbf{K}\psi)$
- (REFL<sub>K</sub>)  $\mathbf{K}\varphi \rightarrow \varphi$
- (TRANS<sub>K</sub>)  $\mathbf{K}\varphi \rightarrow \mathbf{K}\mathbf{K}\varphi$
- (GRZ<sub>K</sub>)  $\mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \rightarrow \varphi$

- (DISTR<sub>K</sub>)  $\mathbf{K}^\downarrow(\varphi \rightarrow \psi) \rightarrow (\mathbf{K}^\downarrow\varphi \rightarrow \mathbf{K}^\downarrow\psi)$
- (BACK)  $\varphi \rightarrow \mathbf{K}\langle\mathbf{K}^\downarrow\rangle\varphi$
- (FORTH)  $\varphi \rightarrow \mathbf{K}^\downarrow\langle\mathbf{K}\rangle\varphi$
- (REFL<sub>T</sub>)  $\mathbf{T}\varphi \rightarrow \varphi$
- (TRANS<sub>T</sub>)  $\mathbf{T}\varphi \rightarrow \mathbf{T}\mathbf{T}\varphi$
- (INCL)  $\mathbf{K}\varphi \rightarrow \mathbf{T}\varphi$
- (SKEL)  $\langle\mathbf{K}^\downarrow\rangle\psi \rightarrow (\mathbf{T}\varphi \leftrightarrow \mathbf{T}(\langle\mathbf{K}^\downarrow\rangle\psi \rightarrow \mathbf{T}\varphi))$

that is closed under:

- (NEC) If  $\varphi \in \mathbf{L}$  then  $\Box\varphi \in \mathbf{L}$  for  $\Box \in \{\mathbf{K}, \mathbf{K}^\downarrow, \mathbf{T}\}$

DEFINITION 5.2. [19, Definition 2.30] (**Deduction for L**)  $\vdash\varphi$  iff either  $\varphi$  is an axiom, or  $\varphi$  follows from some previously obtained formula by one of the inference rules. If  $\Gamma \subseteq \mathcal{L}$  is a set of formulas,  $\Gamma \vdash \varphi$  whenever there are finitely many  $\psi_1, \dots, \psi_k \in \Gamma$  such that  $\vdash \psi_1 \wedge \dots \wedge \psi_k \rightarrow \varphi$ .

DEFINITION 5.3. [19, Definition 2.36]  $\Gamma$  is **consistent** iff  $\Gamma \not\vdash \perp$ .  $\Gamma$  is **maximally consistent** if  $\Gamma$  is consistent and for all  $\varphi \in \mathcal{L}$  either  $\varphi \in \Gamma$  or  $\varphi \notin \Gamma$ .

LEMMA 5.4. [19, Lemma 2.19] (“Lindenbaum's Lemma”) We can extend any consistent set  $\Gamma$  to a maximally consistent set  $\Delta \supseteq \Gamma$ .

DEFINITION 5.5. [19, Definition 2.36] (**Proof Set**)  $|\varphi|_{\mathbf{L}} = \{\Delta \mid \Delta \text{ is maximally consistent and } \varphi \in \Delta\}$

PROPOSITION 5.6. Let  $\Delta$  be maximally consistent, and let  $\Box \in \{\mathbf{K}, \mathbf{K}^\downarrow, \mathbf{T}\}$ . We have  $\Box\varphi \in \Delta$  iff

$$\forall \Sigma \text{ maximally consistent, if } \forall \psi, \Box\psi \in \Delta \text{ implies } \psi \in \Sigma, \text{ then } \varphi \in \Sigma$$

**Proof.** The  $(\rightarrow)$  direction is straightforward. As for the  $(\leftarrow)$  direction, suppose contrapositively that  $\Box\varphi \notin \Delta$ , and let  $\Sigma = \{\psi \mid \Box\psi \in \Delta\} \cup \{\neg\varphi\}$ .

First, we need to check that  $\Sigma$  is consistent. Suppose for contradiction that  $\perp \in \Sigma$ . By definition of  $\Sigma$ , either  $\Box\perp \in \Delta$  or  $\perp \equiv \neg\varphi$  (i.e.  $\varphi$  is  $\top$ ). In the former case, the appropriate instances of (REFL) give us  $\perp \in \Delta$ , which contradicts the consistency of  $\Delta$ . In the latter case,  $\varphi$  is  $\top$ , and so our assumption that  $\Box\varphi \notin \Delta$  gives us  $\Box\top \notin \Delta$ . But  $\Box\top$  follows from our axioms (specifically, (NEC)).

So by Lindenbaum's Lemma, we can extend  $\Sigma$  to maximally consistent  $\Sigma^*$ . From here, we can apply our hypothesis. By construction, for all  $\psi$ ,  $\Box\psi \in \Delta$  implies  $\psi \in \Sigma \subseteq \Sigma^*$ , but  $\varphi \notin \Sigma^*$  (since  $\neg\varphi \in \Sigma \subseteq \Sigma^*$ ).  $\square$

LEMMA 5.7. Let  $\Sigma, \Delta$  be maximally consistent. The following are equivalent:

1.  $\mathbf{K}\varphi \in \Sigma$  implies  $\varphi \in \Delta$
2.  $\mathbf{K}^\downarrow\varphi \in \Delta$  implies  $\varphi \in \Sigma$

**Proof.** Suppose (1) holds, and suppose  $\mathbf{K}^\downarrow\varphi \in \Delta$ . For contradiction, suppose  $\varphi \notin \Sigma$ . Since  $\Sigma$  is maximally consistent,  $\neg\varphi \in \Sigma$ . Applying the (BACK) axiom, we get  $\mathbf{K}\langle\mathbf{K}^\downarrow\rangle\neg\varphi \in \Sigma$ , i.e.  $\mathbf{K}\neg\mathbf{K}^\downarrow\varphi \in \Sigma$ . By (1),  $\neg\mathbf{K}^\downarrow\varphi \in \Delta$ , i.e.  $\mathbf{K}^\downarrow\varphi \notin \Delta$ . But this contradicts  $\mathbf{K}^\downarrow\varphi \in \Delta$ !

Now suppose (2) holds, and suppose  $\mathbf{K}\varphi \in \Sigma$ . For contradiction, suppose  $\varphi \notin \Delta$ . Since  $\Delta$  is maximally consistent,  $\neg\varphi \in \Delta$ . Applying the (FORTH) axiom, we get  $\mathbf{K}^\downarrow\langle\mathbf{K}\rangle\neg\varphi \in \Delta$ , i.e.  $\mathbf{K}^\downarrow\neg\mathbf{K}\varphi \in \Delta$ . By (2),  $\neg\mathbf{K}\varphi \in \Sigma$ , i.e.  $\mathbf{K}\varphi \notin \Sigma$ . But this contradicts  $\mathbf{K}\varphi \in \Sigma$ !  $\square$

## 5.2 Soundness

**THEOREM 5.8. (Soundness)** If  $\Gamma \vdash \varphi$  then  $\Gamma \models_{\text{BFNN}} \varphi$

**Proof.** Suppose  $\Gamma \vdash \varphi$ , and let  $\mathcal{N}, n \models \Gamma$ . We just need to check that each of the axioms and rules of inference are sound, from which we can conclude that  $\mathcal{N}, n \models \varphi$ . We can do this either by the semantics of BFNNs, or instead by checking them in an equivalent preferential frame  $\mathcal{M}^* = \langle W, f, g, V \rangle$ :

To show soundness of:	Use:	Alternative:
(DISTR <sub>K</sub> )	Monotonicity of Reach	$\langle W, f \rangle$ forms a filter
(REFL <sub>K</sub> )	Inclusion of Reach	Reflexivity of $f$
(TRANS <sub>K</sub> )	Idempotence of Reach	Transitivity of $f$
(GRZ <sub>K</sub> )	Antisymmetry of Reach	$f$ is antisymmetric
(DISTR <sub>K</sub> )	Definition of Reach <sup>↓</sup>	Definition of $\mathbf{K}^\downarrow$
(BACK)	Monotonicity of Reach	$\langle W, f \rangle$ forms a filter
(FORTH)	Monotonicity of Reach	$\langle W, f \rangle$ forms a filter
(REFL <sub>T</sub> )	Inclusion of Prop	Reflexivity of $g$
(TRANS <sub>T</sub> )	Idempotence of Prop	Transitivity of $g$
(INCL)	Reach contains Prop	$g$ contains $f$
(SKEL)	Minimal Cause	$f$ is a skeleton of $g$
(NEC)	$\forall w, w \notin \text{Reach}(\emptyset), \text{Prop}(\emptyset)$	$f, g$ contain the unit

□

## 5.3 Model Building

Given a set  $\Gamma \subseteq \mathcal{L}$ , I will show that we can build a net  $\mathcal{N}$  that models  $\Gamma$ . Since preferential filters are equivalent to BFNNs (over  $\mathcal{L}$ ), I will focus instead on building a preferential filter  $\mathcal{F}$ . This is the same strategy taken by [14], who constructs KLM cumulative-ordered models in order to build a neural net.

The following are the standard canonical construction and facts for neighborhood models (see Eric Pacuit's book). Adapting these to our logic of  $\mathbf{K}, \mathbf{K}^\downarrow, \mathbf{T}$  is a straightforward exercise in modal logic.

**LEMMA 5.9.** [19, Lemma 2.12 & Definition 2.37] We can build a **canonical** neighborhood model for  $\mathbf{L}$ , i.e. a model  $\mathcal{M}^C = \langle W^C, f^C, g^C, V^C \rangle$  such that:

- $W^C = \{\Delta \mid \Delta \text{ is maximally consistent}\}$
- For each  $\Delta \in W^C$  and each  $\varphi \in \mathcal{L}$ ,  $|\varphi|_{\mathbf{L}} \in f^C(\Delta)$  iff  $\mathbf{K}\varphi \in \Delta$
- For each  $\Delta \in W^C$  and each  $\varphi \in \mathcal{L}$ ,  $|\varphi|_{\mathbf{L}} \in g^C(\Delta)$  iff  $\mathbf{T}\varphi \in \Delta$
- $V^C(p) = |p|_{\mathbf{L}}$

**Note.** This is where the Necessitation rules come into play — we need them in order to guarantee that we can actually build this model!

**LEMMA 5.10.** [19, Lemma 2.13] (**Truth Lemma**) We have, for canonical model  $\mathcal{M}^C$ ,

$$\{\Delta \mid \mathcal{M}^C, \Delta \Vdash \varphi\} = |\varphi|_{\mathbf{L}}$$

**Proof.** By induction on  $\varphi$ . The propositional, and boolean cases are straightforward.

**K case.**

$$\begin{aligned}
 \mathcal{M}^C, \Delta \Vdash \mathbf{K}\varphi & \quad \text{iff} \quad \{u \mid \mathcal{M}^C, \Sigma \Vdash \varphi\} \in f^C(\Delta) \quad (\text{by definition}) \\
 & \quad \text{iff} \quad |\varphi|_{\mathbf{L}} \in f^C(\Delta) \quad (\text{by IH}) \\
 & \quad \text{iff} \quad \mathbf{K}\varphi \in \Delta \quad (\text{since } \mathcal{M}^C \text{ is canonical}) \\
 & \quad \text{iff} \quad \Delta \in |\mathbf{K}\varphi|_{\mathbf{L}} \quad (\text{by definition})
 \end{aligned}$$

**K<sup>↓</sup> case.**

$$\begin{array}{llll}
\mathcal{M}^C, \Delta \Vdash \mathbf{K}^\downarrow \varphi & \text{iff} & \forall \Sigma \in W^C, \text{ if } \Delta \in \cap f^C(\Sigma) \text{ then } \mathcal{M}, \Sigma \Vdash \varphi & \text{(by definition)} \\
& \text{iff} & \forall \Sigma \in W^C, \text{ if } \Delta \in \cap f^C(\Sigma) \text{ then } \Sigma \in |\varphi|_L & \text{(by IH)} \\
& \text{iff} & \forall \Sigma \in W^C, \text{ if } \Delta \in \cap f^C(\Sigma) \text{ then } \varphi \in \Sigma & \\
& \text{iff} & \forall \Sigma \in W^C, \text{ if } (|\psi|_L \in f^C(\Sigma) \text{ implies } \Delta \in |\psi|_L) \text{ then } \varphi \in \Sigma & \text{(by definition of core)} \\
& \text{iff} & \forall \Sigma \in W^C, \text{ if } (\forall \psi, \mathbf{K}\psi \in \Sigma \text{ implies } \psi \in \Delta) \text{ then } \varphi \in \Sigma & \text{(since } \mathcal{M}^C \text{ is canonical)} \\
& \text{iff} & \forall \Sigma \in W^C, \text{ if } (\forall \psi, \mathbf{K}^\downarrow \psi \in \Delta \text{ implies } \psi \in \Sigma) \text{ then } \varphi \in \Sigma & \text{(by Lemma 5.7)} \\
& \text{iff} & \mathbf{K}^\downarrow \varphi \in \Delta & \text{(by Proposition 5.6)} \\
& \text{iff} & \Delta \in |\mathbf{K}^\downarrow \varphi|_L & \text{(by definition)}
\end{array}$$

**T case.**

$$\begin{array}{llll}
\mathcal{M}^C, \Delta \Vdash \mathbf{T}\varphi & \text{iff} & \{u | \mathcal{M}^C, \Sigma \Vdash \varphi\} \in g^C(\Delta) & \text{(by definition)} \\
& \text{iff} & |\varphi|_L \in g^C(\Delta) & \text{(by IH)} \\
& \text{iff} & \mathbf{T}\varphi \in \Delta & \text{(since } \mathcal{M}^C \text{ is canonical)} \\
& \text{iff} & \Delta \in |\mathbf{T}\varphi|_L & \text{(by definition)}
\end{array}$$

□

**THEOREM 5.11.** [State that our logic has the finite model property]

**Proof.** [Prove it by the usual filtration construction — the fact that the filtration is closed under  $\cap$ ,  $\subseteq$ , reflexive, and transitive are all shown in Pacuit's book. So I just need to show that the same is true of the acyclic & skeleton properties.] □

**PROPOSITION 5.12.** If  $\mathcal{M}$  is finite and satisfies the Truth Lemma, then  $\mathcal{M}$  is a preferential filter.

**Proof.** We check each property in turn:

**$W^C$  is finite.** Holds by assumption.

**$f^C$  is closed under finite intersection.** It's enough to show that  $f^C$  is closed under binary intersections.  $L$  contains all instances of (DISTR<sub>K</sub>), from which we can derive all instances of:

$$(\text{CONJUNCTION-CLOSURE}_K) \quad \mathbf{K}\varphi \wedge \mathbf{K}\psi \rightarrow \mathbf{K}(\varphi \wedge \psi)$$

Suppose  $|\varphi|_L, |\psi|_L \in f^C(\Delta)$ . By definition of  $f^C$ ,  $\mathbf{K}\varphi \in \Delta$  and  $\mathbf{K}\psi \in \Delta$ . So  $\mathbf{K}\varphi \wedge \mathbf{K}\psi \in \Delta$ . Applying (CONJUNCTION-CLOSURE<sub>K</sub>),  $\mathbf{K}(\varphi \wedge \psi) \in \Delta$ . So  $|\varphi \wedge \psi|_L = |\varphi|_L \cap |\psi|_L \in \Delta$ .

**$f^C$  is closed under superset.**  $L$  contains all instances of (DISTR<sub>K</sub>) and the necessitation rule, from which we can derive:

$$(\text{RIGHT-MONOTONE}_K) \quad \text{If } \varphi \rightarrow \psi \in L \text{ then } \mathbf{K}\varphi \rightarrow \mathbf{K}\psi \in L$$

Suppose  $|\varphi|_L \in f^C(\Delta)$ , and  $|\varphi|_L \subseteq |\psi|_L$ . The former fact gives us  $\mathbf{K}\varphi \in \Delta$ . The latter gives us, for all maximally consistent  $\Delta$ , if  $\varphi \in \Delta$  then  $\psi \in \Delta$ , i.e.  $\varphi \rightarrow \psi \in L$  [Is this correct? Probably not; we need to close the canonical model under superset]. By (RIGHT-MONOTONE<sub>K</sub>), we have  $\mathbf{K}\psi \in \Delta$ , i.e.  $|\psi|_L \in f^C(\Delta)$ .

**$f^C$  contains the unit.**  $L$  is closed under (NEC) for  $\mathbf{K}$ , from which we can derive:

$$(\text{TOP}_K) \quad \mathbf{K}\top$$

That is,  $\mathbf{K}\top \in \Delta$  for all maximally consistent  $\Delta$ . So  $|\top|_L \in f^C(\Delta)$ , i.e.  $W^C \in f^C(\Delta)$ .

**$f^C$  is reflexive.** First, let  $\Delta \in W^C$ , and suppose  $|\varphi|_L \in f^C(\Delta)$ . By definition of  $f^C$ ,  $\mathbf{K}\varphi \in \Delta$ . By (REFL<sub>K</sub>),  $\varphi \in \Delta$ . Since  $\varphi$  was chosen arbitrarily, we have for all  $\varphi$ , if  $|\varphi|_L \in f^C(\Delta)$  then  $\varphi \in \Delta$ . In other words,  $\Delta \in \bigcap_{|\varphi|_L \in f^C(\Delta)} |\varphi|_L = \cap f^C(\Delta)$ .

**$f^C$  is transitive.** Suppose  $|\varphi|_L \in f^C(\Delta)$ . By definition of  $f^C$ ,  $\mathbf{K}\varphi \in \Delta$ . By the (TRANS<sub>K</sub>) axiom,  $\mathbf{K}\mathbf{K}\varphi \in \Delta$ . But this means that  $|\mathbf{K}\varphi|_L \in f^C(\Delta)$ . By definition of proof set, we have  $\{\Sigma \mid \mathbf{K}\varphi \in \Sigma\} \in f^C(\Delta)$ . That is,  $\{\Sigma \mid |\varphi|_L \in f^C(\Sigma)\} \in f^C(\Delta)$ , and we are done.

**$f^C$  is antisymmetric.** Suppose  $\Delta_1 \in \cap f^C(\Delta_2)$  and  $\Delta_2 \in \cap f^C(\Delta_1)$ . By definition of core,  $\Delta_1 \in \bigcap_{|\varphi|_L \in f^C(\Delta_2)} |\varphi|_L$  and  $\Delta_2 \in \bigcap_{|\varphi|_L \in f^C(\Delta_1)} |\varphi|_L$ , i.e. we have both of the following:

1.  $\forall \varphi$ , if  $\mathbf{K}\varphi \in \Delta_2$  then  $\varphi \in \Delta_1$
2.  $\forall \varphi$ , if  $\mathbf{K}\varphi \in \Delta_1$  then  $\varphi \in \Delta_2$

We want to show that  $\Delta_1 = \Delta_2$ . For contradiction, suppose not; without loss of generality, say  $\varphi \in \Delta_1$ , but  $\varphi \notin \Delta_2$ .

From  $\varphi \notin \Delta_2$  and (2) we get  $\mathbf{K}\varphi \notin \Delta_1$ . Since  $\Delta_1$  is maximal,  $\neg \mathbf{K}\varphi \in \Delta_1$ . Since  $\varphi \in \Delta_1$  and  $\neg \mathbf{K}\varphi \in \Delta_1$ ,  $\varphi \wedge \neg \mathbf{K}\varphi \in \Delta_1$ . Since  $\Delta_1$  is consistent,  $\neg(\varphi \wedge \neg \mathbf{K}\varphi) \equiv \neg\varphi \vee \mathbf{K}\varphi \equiv \varphi \rightarrow \mathbf{K}\varphi \notin \Delta_1$ . From (1) we have  $\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \notin \Delta_2$ .

From here, we can apply the (T<sub>K</sub>) axiom to get  $\mathbf{K}\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \notin \Delta_2$ . Since  $\Delta_2$  is maximal,  $\neg \mathbf{K}\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \in \Delta_2$ . And so  $\neg \mathbf{K}\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \vee \mathbf{K}\varphi \in \Delta_2$  (we can disjunct with anything; I happened to choose  $\mathbf{K}\varphi$ ). But this is just  $\mathbf{K}\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \mathbf{K}\varphi \in \Delta_2$ .

We can apply the (K) axiom to re-distribute the **K**:  $\mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \in \Delta_2$ . Since  $\varphi \notin \Delta_2$ ,  $\neg\varphi \in \Delta_2$ . So in particular,  $\mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \wedge \neg\varphi \in \Delta_2$ . Factoring the negation out, we have  $\neg[\neg \mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \vee \varphi] \equiv \neg[\mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \rightarrow \varphi] \in \Delta_2$ . But this contradicts the (GRZ<sub>K</sub>) axiom  $\mathbf{K}(\mathbf{K}(\varphi \rightarrow \mathbf{K}\varphi) \rightarrow \varphi) \rightarrow \varphi \in \Delta_2$ !

**$g^C$  contains the unit.** Similar to the proof for  $f^C$ , but apply necessitation for **T** instead of **K**.

**$g^C$  is reflexive.** Similar to the proof for  $f^C$ , but apply (REFL<sub>T</sub>) instead of (REFL<sub>K</sub>).

**$g^C$  is transitive.** Similar to the proof for  $f^C$ , but apply (TRANS<sub>T</sub>) instead of (TRANS<sub>K</sub>).

**$g^C$  contains  $f^C$ .** Suppose  $|\varphi|_L \in f^C(\Delta)$ . By definition of  $f^C$ ,  $\mathbf{K}\varphi \in \Delta$ . By the (INCL) axiom,  $\mathbf{T}\varphi \in \Delta$ . And so  $|\varphi|_L \in f^C(\Delta)$ .

**$f^C$  is the skeleton of  $g^C$ .** Let  $\varphi, \psi$  be formulas, and suppose  $\Delta \in \cap f^C(|\psi|_L)$ . We would like to show:

$$|\varphi|_L \in g^C(\Delta) \quad \text{iff} \quad \{\Sigma \mid |\varphi|_L \in g^C(\Sigma)\} \cup (\cap f^C(|\psi|_L))^c \in g^C(\Delta)$$

First, let's unpack our hypothesis. By definition,  $\Delta \in \bigcup_{\Theta \in |\psi|_L} \cap f^C(\Theta)$ . In other words,  $\Delta \in \{\Sigma \mid \exists \Theta \text{ such that } \psi \in \Theta \text{ and } \Sigma \in \cap f^C(\Theta)\}$ . Well,  $\psi \in \Theta$  holds iff  $\Theta \in |\psi|_L$  iff  $\mathcal{M}^C, \Theta \models \psi$  (by the Truth Lemma). So we have  $\Delta \in \{\Sigma \mid \exists \Theta \text{ such that } \mathcal{M}^C, \Theta \models \psi \text{ and } \Sigma \in \cap f^C(\Theta)\}$ , i.e.  $\Delta \in \{\Sigma \mid \mathcal{M}^C, \Sigma \models \langle \mathbf{K}^\downarrow \rangle \psi\}$ . Applying the Truth Lemma again, this gives us  $\Delta \in |\langle \mathbf{K}^\downarrow \rangle \psi|$ , i.e.  $\langle \mathbf{K}^\downarrow \rangle \psi \in \Delta$ .

Now let's prove the original claim. We will show the ( $\rightarrow$ ) direction (the ( $\leftarrow$ ) direction is similar). Suppose  $|\varphi|_L \in g^C(\Delta)$ . By definition of  $g^C$ ,  $\mathbf{T}\varphi \in \Delta$ . Since  $\langle \mathbf{K}^\downarrow \rangle \psi \in \Delta$ , the (SKEL) axiom gives us  $\mathbf{T}(\langle \mathbf{K}^\downarrow \rangle \psi \rightarrow \mathbf{T}\varphi) \in \Delta$ . By definition of  $g^C$  again,  $|\langle \mathbf{K}^\downarrow \rangle \psi \rightarrow \mathbf{T}\varphi|_L \in g^C(\Delta)$ . Well,

$$\begin{aligned} |\langle \mathbf{K}^\downarrow \rangle \psi \rightarrow \mathbf{T}\varphi|_L &= |\neg \langle \mathbf{K}^\downarrow \rangle \psi \vee \mathbf{T}\varphi|_L \\ &= |\mathbf{T}\varphi|_L \cup |\neg \langle \mathbf{K}^\downarrow \rangle \psi|_L \\ &= |\mathbf{T}\varphi|_L \cup \{\Sigma \mid \mathcal{M}^C, \Sigma \models \langle \mathbf{K}^\downarrow \rangle \psi\}^c && \text{(Truth Lemma)} \\ &= |\mathbf{T}\varphi|_L \cup \{\Sigma \mid \exists \Theta \text{ such that } \mathcal{M}^C, \Theta \models \psi \text{ and } \Sigma \in \cap f^C(\Theta)\}^c && \text{(by our semantics)} \\ &= |\mathbf{T}\varphi|_L \cup \{\Sigma \mid \exists \Theta \text{ such that } \Theta \in |\psi|_L \text{ and } \Sigma \in \cap f^C(\Theta)\}^c && \text{(Truth Lemma)} \\ &= |\mathbf{T}\varphi|_L \cup \{\Sigma \mid \exists \Theta \text{ such that } \psi \in \Theta \text{ and } \Sigma \in \cap f^C(\Theta)\}^c && \text{(definition of } |\psi|_L) \\ &= |\mathbf{T}\varphi|_L \cup \left( \bigcup_{\Theta \in |\psi|_L} \cap f^C(\Theta) \right)^c \\ &= |\mathbf{T}\varphi|_L \cup (\cap f^C(|\psi|_L))^c \\ &= \{\Sigma \mid \mathbf{T}\varphi \in \Sigma\} \cup (\cap f^C(|\psi|_L))^c && \text{(definition of } |\mathbf{T}\varphi|_L) \\ &= \{\Sigma \mid |\varphi|_L \in g^C(\Sigma)\} \cup (\cap f^C(|\psi|_L))^c && \text{(definition of } g^C) \end{aligned}$$



So this latter set is in  $g^C(\Delta)$ . But that is exactly what we wanted to show!  $\square$

**THEOREM 5.13. (Model Building)** Given any consistent  $\Gamma \subseteq \mathcal{L}$ , we can construct a BFNN  $\mathcal{N}$  and neuron  $n \in N$  such that  $\mathcal{N}, n \models \Gamma$ .

**Proof.** Extend  $\Gamma$  to maximally consistent  $\Delta$  using Lemma 5.4. Let  $\mathcal{M}^C$  be a canonical model for  $\mathbf{L}$  guaranteed by Lemma 5.9. By the Truth Lemma (Lemma 5.10),  $\mathcal{M}^C, \Delta \models \Delta$ . So in particular,  $\mathcal{M}^C, \Delta \models \Gamma$ .

By the Finite Model Property (Lemma 5.11), we can construct a finite model  $\mathcal{M}'$  satisfying exactly the same formulas at all worlds. By Proposition 5.12,  $\mathcal{M}'$  is a preferential filter.

From here, we can build our net  $\mathcal{N}^\bullet$  as before, satisfying exactly the same formulas as  $\mathcal{M}$  at all neurons (by Theorem 4.7). And so  $\mathcal{N}^\bullet, \Delta \models \Gamma$ .  $\square$

**THEOREM 5.14. (Completeness)** For all consistent  $\Gamma \subseteq \mathcal{L}$ , if  $\Gamma \models_{\text{BFNN}} \varphi$  then  $\Gamma \vdash \varphi$

**Proof.** Suppose contrapositively that  $\Gamma \not\models \varphi$ . This means that  $\Gamma \cup \{\neg\varphi\}$  is consistent, i.e. by Theorem 5.13 we can build a BFNN  $\mathcal{N}$  and neuron  $n$  such that  $\mathcal{N}, n \models \Gamma \cup \{\neg\varphi\}$ . In particular,  $\mathcal{N}, n \not\models \varphi$ . But then we must have  $\Gamma \not\models \varphi$ .  $\square$

## TODO:

- Change activation functions so that they aren't dependent on weights
- Make diagrams in Tikz
- Think about unsupervised Hebbian update  $[\varphi]^*$
- Figure out if we need activation function to be monotonically increasing
- Prove Loop & Cumulative from our properties
- Do filtration/finite model property
- Get bound on the size of the finite model (if it seems easy)
- Close canonical models under  $\subseteq$
- Write careful paragraph about leaving fuzzy to future work, and why we make this choice despite it being an over-simplified case.
- Make corrections Saul gave earlier (whoops, I forgot!)

## References

- [1] Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration-a structured survey. *ArXiv preprint cs/0511042*, 2005.
- [2] Christian Balkenius and Peter Gärdenfors. Nonmonotonic Inferences in Neural Networks. In *KR*, pages 32–39. 1991.
- [3] Vaishak Belle. Logic Meets Learning: From Aristotle to Neural Networks. In *Neuro-Symbolic Artificial Intelligence: The State of the Art*, pages 78–102. IOS Press, 2021.
- [4] Reinhard Blutner. Nonmonotonic inferences and neural networks. In *Information, Interaction and Agency*, pages 203–234. Springer, 2004.
- [5] Antony Browne and Ron Sun. Connectionist inference models. *Neural Networks*, 14(10):1331–1355, 2001.
- [6] Dov M Gabbay, Ian Hodkinson, and Mark A Reynolds. Temporal logic: mathematical foundations and computational aspects. 1994.
- [7] Artur S d'Avila Garcez, Krysia Broda, and Dov M Gabbay. Symbolic knowledge extraction from trained neural networks: a sound approach. *Artificial Intelligence*, 125(1-2):155–207, 2001.
- [8] Artur S d'Avila Garcez, Luis C Lamb, and Dov M Gabbay. *Neural-symbolic cognitive reasoning*. Springer Science Business Media, 2008.
- [9] Laura Giordano, Valentina Gliozzi, and Daniele Theseider Dupré. From common sense reasoning to neural network models through multiple preferences: An overview. *CoRR*, abs/2107.04870, 2021.
- [10] Laura Giordano, Valentina Gliozzi, and Daniele Theseider Dupré. A conditional, a fuzzy and a probabilistic interpretation of self-organizing maps. *Journal of Logic and Computation*, 32(2):178–205, 2022.
- [11] Donald Hebb. *The organization of behavior: A neuropsychological theory*. Wiley, New York, 1949.
- [12] The Third AI Summer, AAAI Robert S. Engelmore Memorial Award Lecture. AAAI, 2020.

- [13] Caleb Kisby, Saúl Blanco, and Lawrence Moss. The logic of hebbian learning. In *The International FLAIRS Conference Proceedings*, volume 35. 2022.
- [14] Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets. *Artificial Intelligence*, 128(1-2):161–201, 2001.
- [15] Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets II. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 11(supp02):105–135, 2003.
- [16] Hannes Leitgeb. Neural Network Models of Conditionals. In *Introduction to Formal Philosophy*, pages 147–176. Springer, 2018.
- [17] Warren S McCulloch and Walter Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, 1943.
- [18] Simon Odense and Artur d'Avila Garcez. A semantic framework for neural-symbolic computing. *ArXiv preprint arXiv:2212.12050*, 2022.
- [19] Eric Pacuit. *Neighborhood semantics for modal logic*. Springer, 2017.
- [20] Md Kamruzzaman Sarker, Lu Zhou, Aaron Eberhart, and Pascal Hitzler. Neuro-symbolic artificial intelligence: current trends. *ArXiv preprint arXiv:2105.05330*, 2021.
- [21] Dongran Yu, Bo Yang, Dayou Liu, and Hui Wang. A survey on neural-symbolic systems. *ArXiv preprint arXiv:2111.08164*, 2021.

## Talk Abstract

Artificial Intelligence is in the midst of a crisis. Its two main paradigms — connectionist neural networks and classical (logic) models — seem diametrically opposed, with no clear way to reconcile the two. Neural networks learn flexibly from unstructured data, and are more cognitively plausible. But it is difficult to interpret what a neural network *knows* or *has learned*. On the other hand, classical models *do* represent knowledge explicitly and transparently, but are notoriously rigid and are often criticized for being completely cognitively implausible. This talk is an introduction to a developing theory of *Neuro-Symbolic Artificial Intelligence* that aims to bridge the gap. The key insight of this theory is that neural networks and classical models are really two different representations of the same information; it is in principle possible to translate between the two. I will illustrate how this is possible by offering (1) a translation for feed-forward nets that is provably correct, and (2) a logical account for “naive” Hebbian learning. Although this talk is primarily from the point-of-view of AI, I hope to convince you that this is a deep result for cognitive science at large.

### Bio:

Caleb Schultz Kisby is a fifth-year Computer Science PhD student at Indiana University, co-advised by Saúl Blanco and Larry Moss. He is interested in combining neural and symbolic accounts of reasoning, especially for systems that learn (which has long been neglected by Neuro-Symbolic AI). Before joining IU, he received his B.S. in Computer Science and Mathematics at the University of South Carolina.