

Neural Net Semantics with Modal Operators

Step 0. Find a paper I enjoy, and read it. Try to understand its ideas, with an eye towards extending it/altering it.

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](#).

Step 1. Look for an extension/open problem that makes me think “What the fuck? That's still open? No way, this shit is low-hanging fruit, free paper here I come.” i.e. something *easy* and *straightforward*, without complications.

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. ψ is in the propagation of the signal φ), 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.

Step 2. Follow-up question (only answer after Step 1): Is the extension *interesting* or *surprising*? What do we learn by extending the result?

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 Γ , we can *build* a neural network \mathcal{N} satisfying Γ .

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”

Step 3. Two things to do at this point:

- **Make a new Texmacs file named “PAPERNAME-master-notes.tm”. Transcribe the key definitions, examples, lemmas, and results from the paper. This makes it easier to later copy-paste parts of proofs, and also ensures that I don't reinvent the wheel later (it's tempting to redefine everything yourself!)**

- Go to <https://www.connectedpapers.com/> and download any major nearby papers. Upload the papers to paperless-ngx and make a point to read them (understanding context helps a lot!).

Related Papers:

Neural Network Semantics / Semantic Encodings.

Classic Papers. [17]

Conditional Logic (Feedforward Net). [2], [14], [15], [8] (soundness), [9] (model-building)
 [Any other relevant work by the Garcez lab?]

Description Logic w. Typicality. [10], [11] [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. [5], [6]

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. [7]

Nominals (Hybrid Logic). [4]

Step 4. Write up my new definitions & proof in the Texmacs file. Again, should be a very straightforward extension, and the proof (proofs are just unit-tests for definitions) shouldn't take up too much room at all (1-2 pages, including defs)

1 Interpreted Neural Nets

1.1 Basic Definitions

DEFINITION 1.1. An **interpreted ANN** (Artificial Neural Network) is a pointed directed graph $\mathcal{N} = \langle N, E, W, A, O, V \rangle$, 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)
- O is a function which maps each $n \in N$ to $O^{(n)}: \mathbb{R} \rightarrow \{0, 1\}$ (the **output function** for n)
- $V: \text{propositions} \cup \text{nominals} \rightarrow \mathcal{P}(N)$ is an assignment of nominals to individual neurons (the **valuation function**). If i is a nominal, we require $|V(i)| = 1$, i.e. a singleton.

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

- **Feed-forward**, i.e. E does not contain any cycles
- **Binary**, i.e. the output of each neuron is in $\{0, 1\}$
- $O^{(n)} \circ A^{(n)}$ is **zero at zero** in the first parameter, i.e.

$$O^{(n)}(A^{(n)}(\vec{0}, \vec{w})) = 0$$

- $O^{(n)} \circ A^{(n)}$ is **strictly monotonically increasing** in the second parameter, i.e. for all $\vec{x}, \vec{w}_1, \vec{w}_2 \in \mathbb{R}^k$, if $\vec{w}_1 < \vec{w}_2$ then $O^{(n)}(A^{(n)}(\vec{x}, \vec{w}_1)) < O^{(n)}(A^{(n)}(\vec{x}, \vec{w}_2))$. We will more often refer to the equivalent condition:

$$\vec{w}_1 \leq \vec{w}_2 \quad \text{iff} \quad O^{(n)}(A^{(n)}(\vec{x}, \vec{w}_1)) \leq O^{(n)}(A^{(n)}(\vec{x}, \vec{w}_2))$$

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 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:

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

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

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

Proof. We have:

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

□

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 $(m_i, n) \in E$, $\text{Activates}_{\text{Prop}(S)}(\vec{m}, n) = 1$.

PROPOSITION 1.8. (Adapted from [14, Remark 4]) Let $\mathcal{N} \in \text{Net}$. For all $S, S_1, S_2 \in \text{Set}$, Prop satisfies

(Inclusion). $S \subseteq \text{Prop}(S)$

(Idempotence). $\text{Prop}(S) = \text{Prop}(\text{Prop}(S))$

(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 $i, j \in \{0, \dots, n\}$

Proof. We check each in turn:

(Inclusion). If $n \in S$, then $n \in \text{Prop}(S)$ by the base case of Prop .

(Idempotence). 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 $(m_i, n) \in E$,

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

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

(Cumulative). For the (\subseteq) direction, let $n \in \text{Prop}(S_1)$. We proceed by induction on $\text{Prop}(S_1)$.

Base Step. Suppose $n \in S_1$. Well, $S_1 \subseteq S_2 \subseteq \text{Prop}(S_2)$, so $n \in \text{Prop}(S_2)$.

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

$$\text{Activates}_{\text{Prop}(S_1)}(\vec{m}, n) = 1$$

By inductive hypothesis, $m_i \in \text{Prop}(S_1)$ iff $m_i \in \text{Prop}(S_2)$. By Proposition 1.6, $\text{Activates}_{\text{Prop}(S_2)}(\vec{m}, n)$, and so $n \in \text{Prop}(S_2)$.

Now consider the (\supseteq) direction. The Inductive Step holds similarly (just swap S_1 and S_2). As for the Base Step, if $n \in S_2$ then since $S_2 \subseteq \text{Prop}(S_1)$, $n \in S_1$.

(Loop). Let $n \geq 0$ and suppose the hypothesis. Our goal is to show that for each i , $\text{Prop}(S_i) \subseteq \text{Prop}(S_{i-1})$, and additionally $\text{Prop}(S_0) \subseteq \text{Prop}(S_n)$. This will show that all $\text{Prop}(S_i)$ contain each other, and so are equal. Let $i \in \{0, \dots, n\}$ (if $i = 0$ then $i - 1$ refers to n), and let $e \in \text{Prop}(S_i)$. We proceed by induction on $\text{Prop}(S_i)$.

Base Step. $e \in S_i$, and since $S_i \subseteq \text{Prop}(S_{i-1})$ by assumption, $e \in \text{Prop}(S_{i-1})$.

Inductive Step. For those m_1, \dots, m_k such that $(m_i, e) \in E$,

$$\text{Activates}_{\text{Prop}(S_i)}(\vec{m}, e) = 1$$

By inductive hypothesis, $m_j \in \text{Prop}(S_i)$ iff $m_j \in \text{Prop}(S_{i-1})$. By Proposition 1.6, $\text{Activates}_{\text{Prop}(S_{i-1})}(\vec{m}, e) = 1$, and so $e \in \text{Prop}(S_{i-1})$. \square

DEFINITION 1.9. Let $\text{Reach}: \text{Set} \rightarrow \text{Set}$ be defined recursively as follows: $n \in \text{Reach}(S)$ iff either

Base Case. $n \in S$, or

Constructor. There is an $m \in \text{Reach}(S)$ such that $(m, n) \in E$.

PROPOSITION 1.10. Let $\mathcal{N} \in \text{Net}$. For all $S, S_1, S_2 \in \text{Set}$, $n, m \in N$, Reach satisfies

(Inclusion). $S \subseteq \text{Reach}(S)$

(Idempotence). $\text{Reach}(S) = \text{Reach}(\text{Reach}(S))$

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

(Contains Prop). $\text{Prop}(S) \subseteq \text{Reach}(S)$

Proof. We check each in turn:

(Inclusion). Similar to the proof of Inclusion for Prop .

(Idempotence). Similar to the proof of Idempotence for Prop .

(Monotonicity). Let $n \in \text{Reach}(S_1)$. We proceed by induction on $\text{Reach}(S_1)$.

Base Step. $n \in S_1$. So $n \in S_2 \subseteq \text{Reach}(S_2)$.

Inductive Step. There is an $m \in \text{Reach}(S_1)$ such that $(m, n) \in E$. By inductive hypothesis, $m \in \text{Reach}(S_2)$. And so by definition, $n \in \text{Reach}(S_2)$.

(Contains Prop). 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 $(m_i, n) \in E$,

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

Since $O \circ A$ 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

[Todo: Check that Cumulative and Loop actually follow from what we have! (Will have to use Acyclic property of Reach .)]

1.3 Neural Network Semantics

DEFINITION 1.11. 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 1.12. Let $\mathcal{N} \in \text{Net}$. The semantics $\llbracket \cdot \rrbracket : \mathcal{L} \rightarrow \text{Set}$ for \mathcal{L} are defined recursively as follows:

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

DEFINITION 1.13. (**Truth at a neuron**) $\mathcal{N}, n \Vdash \varphi$ iff $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}$.

DEFINITION 1.14. (**Truth in a net**) $\mathcal{N} \models \varphi$ iff $\mathcal{N}, n \Vdash \varphi$ for all $n \in N$.

DEFINITION 1.15. (**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 $\cap f(w)$.

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.

Proof. [Todo]

\square

PROPOSITION 2.5. [19, [Which?]] If $\mathcal{F} = \langle W, f \rangle$ is a proper filter, then for all $w \in W$, $Y^c \in f(w)$ iff $Y \notin f(w)$.

Proof. (\rightarrow) Suppose for contradiction that $Y^c \in f(w)$ and $Y \in f(w)$. Since \mathcal{F} is closed under intersection, $Y^c \cap Y = \emptyset \in f(w)$, which contradicts the fact that \mathcal{F} is proper.

(\leftarrow) Suppose for contradiction that $Y \notin f(w)$, yet $Y^c \notin f(w)$. Since \mathcal{F} is closed under intersection, $\cap f(w) \in f(w)$. Moreover, since \mathcal{F} is closed under superset we must have $\cap f(w) \not\subseteq Y$ and $\cap f(w) \not\subseteq Y^c$. But this means $\cap f(w) \not\subseteq Y \cap Y^c = \emptyset$, i.e. there is some $x \in \cap f(w)$ such that $x \in \emptyset$. This contradicts the definition of the empty set. \square

DEFINITION 2.6. Let $\mathcal{F} = \langle W, f, g \rangle$ be a frame. \mathcal{F} is a **preferential filter** iff:

- W is finite
- $\langle W, f \rangle$ forms a proper filter, and g contains the unit
- f is **acyclic**: for all $u_1, \dots, u_n \in W$, if $u_1 \in \cap f(u_2), \dots, u_{n-1} \in \cap f(u_n), u_n \in \cap f(u_1)$ then all $u_i = u_j$.
- 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 $X \in g(w)$.

PROPOSITION 2.7. Let $\mathcal{F} = \langle W, f \rangle$ be a frame. Suppose f is reflexive, transitive, and **asymmetric**, i.e. $u_1 \in \cap f(u_2)$ and $u_2 \in \cap f(u_1)$ implies $u_1 = u_2$. Then f is acyclic.

Proof. Let $u_1, \dots, u_n \in W$, and suppose $u_1 \in \cap f(u_2), \dots, u_{n-1} \in \cap f(u_n), u_n \in \cap f(u_1)$. WLOG we will show that $u_1 = u_n$. [Todo] \square

2.2 Neighborhood Semantics

DEFINITION 2.8. [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 2.9. [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	$\{u \mid \mathcal{M}, u \Vdash \varphi\} \in f(w)$
$\mathcal{M}, w \Vdash \mathbf{T}\varphi$	iff	$\{u \mid \mathcal{M}, u \Vdash \varphi\} \in g(w)$

In neighborhood semantics, the operators **K**, and **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:

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{K} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin f(w) \\ \mathcal{M}, w \Vdash \langle \mathbf{T} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin g(w) \end{aligned}$$

DEFINITION 2.10. [19, Definition 1.13] (**Truth in a model**) $\mathcal{M} \models \varphi$ iff $\mathcal{M}, w \Vdash \varphi$ for all $w \in W$.

DEFINITION 2.11. [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:

- $V_{\mathcal{M}^\bullet}(p) = V_{\mathcal{N}}(p)$;
- $V_{\mathcal{M}^\bullet}(i) = V_{\mathcal{N}}(i)$

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 φ . The propositional, $\neg\varphi$, and $\varphi \wedge \psi$ cases are trivial.

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

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

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

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

□

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 . Since Reach is monotonic, **[Make this a lemma!]** we have $\bigcup_i \text{Reach}(X_i^c) = \text{Reach}(\bigcup_i X_i^c) = \text{Reach}((\bigcap_i X_i)^c)$. 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 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 acyclic. Suppose $u_1, \dots, u_n \in W$, with $u_1 \in \cap f(u_2), \dots, u_{n-1} \in \cap f(u_n), u_n \in \cap f(u_1)$. That is, each $u_i \in \cap_{X \in f(u_{i+1})} X$. By choice of f , each $u_i \in \cap_{u_{i+1} \notin \text{Reach}(X^c)} X$. Substituting X^c for X we get $u_i \in \cap_{u_{i+1} \notin \text{Reach}(X)} X^c$. In other words, $u_1 \in \text{Reach}^{-1}(u_2), \dots, u_{n-1} \in \text{Reach}^{-1}(u_n), u_n \in \text{Reach}^{-1}(u_1)$. **[Update!]** By Proposition ?, each $u_i = u_j$.

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 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.

□

4 From Frames to Nets

This is the harder (“completeness”) direction!

DEFINITION 4.1. Suppose we have net \mathcal{N} and node $n \in N$ with incoming nodes m_1, \dots, m_k , $(m_i, n) \in E$. Let $\text{hash}: \mathcal{P}(\{m_1, \dots, m_k\}) \times \mathbb{N}^k \rightarrow \mathbb{N}$ be defined by

$$\text{hash}(S, \vec{w}) = \prod_{m_i \in S} w_i$$

PROPOSITION 4.2. $\text{hash}(S, \vec{W}(m_i, n)): \mathcal{P}(\{m_1, \dots, m_k\}) \rightarrow P_k$, where

$$P_k = \{n \in \mathbb{N} \mid n \text{ is the product of some subset of primes } \{p_1, \dots, p_k\}\}$$

is bijective (and so has a well-defined inverse hash^{-1}).

DEFINITION 4.3. Let \mathcal{M} be a model based on preferential filter $\mathcal{F} = \langle W, f, g \rangle$. Its **simulation net** $\mathcal{N}^* = \langle N, E, W, A, O, V \rangle$ 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 $(m_i, n) \in E$.

- $W(m_i, n) = p_i$, the i th prime number.
- $A^{(n)}(\vec{x}, \vec{w}) = \text{hash}(\{m_i \mid x_i = 1\}, \vec{w})$
- $O^{(w)}(x) = 1$ iff $(\text{hash}^{-1}(x)[0])^c \notin g(n)$
- $V_{\mathcal{N}^*}(p) = V_{\mathcal{M}}(p)$

LEMMA 4.4. Let $\vec{m} = m_1, \dots, m_k$ be those nodes such that $(m_i, n) \in E$. Then

$$\text{Activates}_S(\vec{m}, n) = 1 \quad \text{iff} \quad \{m_i \mid m_i \in S\}^c \notin g(n)$$

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

$$\begin{aligned} & O^{(n)}(A^{(n)}((\chi_S(m_1), \dots, \chi_S(m_k)); (W(m_1, n), \dots, W(m_k, n)))) = 1 \\ \text{iff} & \text{hash}^{-1}(\text{hash}(\{m_i | m_i \in S\}; (W(m_1, n), \dots, W(m_k, n))))[0]^\complement \notin g(n) \\ \text{iff} & \{m_i | m_i \in S\}^\complement \notin g(n) \end{aligned}$$

□

CLAIM 4.5. \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 for contradiction that E contains a cycle, i.e. distinct $u_1, \dots, u_n \in N$ such that $u_1 E u_2, \dots, u_{n-1} E u_n, u_n E u_1$. Then we have $u_1 \in \cap f(u_2), \dots, u_{n-1} \in \cap f(u_n), u_n \in \cap f(u_1)$, which contradicts the fact that f is acyclic.

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

$O^{(n)} \circ A^{(n)}$ is monotonically increasing. Let \vec{w}_1, \vec{w}_2 be such that hash is well-defined (i.e. are vectors of prime numbers), and suppose $\vec{w}_1 < \vec{w}_2$. If $O^{(v)}(A^{(v)}(\vec{x}, \vec{w}_1)) = 1$, then $(\text{hash}^{-1}(\text{hash}(\vec{x}, \vec{w}_1)))[0]^\complement \notin g(v)$. But this just means $\{m_i | x_i = 1\}^\complement \notin g(v)$. And so $(\text{hash}^{-1}(\text{hash}(\vec{x}, \vec{w}_2)))[0]^\complement \notin g(v)$. But then $O^{(n)}(A^{(n)}(\vec{x}, \vec{w}_2)) = 1$.

The main point here is just that \vec{w}_1 and \vec{w}_2 are just indexing the set in question, and their actual values don't affect the final output (we don't need the $\vec{w}_1 < \vec{w}_2$ hypothesis!). The real work happens within $g(v)$. □

LEMMA 4.6. $\text{Reach}_{\mathcal{N}^\bullet}(S) = \{n | S^\complement \notin f(n)\}$

Proof. For the (\supseteq) direction, let $n \in N$ be such that $S^\complement \notin f(n)$. By Proposition 2.5 and the fact that $\langle W, f \rangle$ forms a proper filter, $S \in f(n)$. By the definition of core, $\cap f(n) \subseteq S$. f is reflexive, which means in particular that $n \in \cap f(n) \subseteq S$. By the base case of Reach , we have $n \in \text{Reach}_{\mathcal{N}^\bullet}(S)$.

Now for the (\subseteq) direction. Suppose $n \in \text{Reach}(S)$, and proceed by induction on Reach .

Base step. $n \in S$. Suppose for contradiction that $S^\complement \in f(n)$. By definition of core, $\cap f(n) \subseteq S^\complement$. But since \mathcal{F} is reflexive, $n \in \cap f(n)$. So $n \in S^\complement$, which contradicts $n \in S$.

Inductive step. There is $m \in \text{Reach}_{\mathcal{N}^\bullet}(S)$ such that $(m, n) \in E$ (and so $m \in \cap f(n)$). By inductive hypothesis, $S^\complement \notin f(m)$. Now suppose for contradiction that $S^\complement \in f(n)$. Since f is transitive, $\{t | S^\complement \in f(t)\} \in f(n)$. By definition of core, $\cap f(n) \subseteq \{t | S^\complement \in f(t)\}$. Since $m \in \cap f(n)$, $S^\complement \in f(m)$. But this contradicts $S^\complement \notin f(m)$! □

LEMMA 4.7. $\text{Prop}_{\mathcal{N}^\bullet}(S) = \{n | S^\complement \notin g(n)\}$

Proof. For the (\supseteq) direction, let $n \in N$, suppose $S^\complement \notin g(n)$. Since g contains f , $S^\complement \notin f(n)$. By Proposition 2.5 and the fact that $\langle W, f \rangle$ forms a proper filter, $S \in f(n)$. By the definition of core, $\cap f(n) \subseteq S$. f is reflexive, which means in particular that $n \in \cap f(n) \subseteq S$. By the base case of Prop , $n \in \text{Prop}_{\mathcal{N}^\bullet}(S)$.

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

Base step. $n \in S$. Suppose for contradiction that $S^\complement \in g(n)$. Since \mathcal{G} is reflexive, $n \in \cap g(n)$. By definition of core, we have $\cap g(n) \subseteq S^\complement$. But then $n \in \cap g(n) \subseteq S^\complement$, i.e. $n \in S^\complement$, 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

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

By Lemma 4.4, this means that $\{m_i | m_i \in \text{Prop}_{\mathcal{N}^\bullet}(S)\}^\complement \notin g(n)$. But by our inductive hypothesis, $\{m_i | m_i \in \text{Prop}_{\mathcal{N}^\bullet}(S)\} = \{m_i | S^\complement \notin g(n)\}$. For convenience, let T be this latter set, i.e. $T = \{m_i | S^\complement \notin g(n)\}$. So we have $T^\complement \notin g(n)$.

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 $T^c \notin g(v)$. \square

THEOREM 4.8. 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 φ . Again, the propositional, $\neg\varphi$, and $\varphi \wedge \psi$ cases are trivial.

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

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{K} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin f(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}\} \notin f(w) \text{ (Inductive Hypothesis)} \\ & \text{ iff } \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}^c \notin g(w) \\ & \text{ iff } w \in \text{Reach}_{\mathcal{N}^\bullet}(\llbracket \varphi \rrbracket) \text{ (by Lemma 4.6)} \\ & \text{ iff } w \in \llbracket \langle \mathbf{K} \rangle \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ (by definition)} \\ & \text{ iff } \mathcal{N}^\bullet, w \Vdash \langle \mathbf{K} \rangle \varphi \text{ (by definition)} \end{aligned}$$

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

$$\begin{aligned} \mathcal{M}, w \Vdash \langle \mathbf{T} \rangle \varphi & \text{ iff } \{u \mid \mathcal{M}, u \not\Vdash \varphi\} \notin g(w) \text{ (by definition)} \\ & \text{ iff } \{u \mid u \notin \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}\} \notin g(w) \text{ (Inductive Hypothesis)} \\ & \text{ iff } \llbracket \varphi \rrbracket_{\mathcal{N}^\bullet}^c \notin g(w) \\ & \text{ iff } w \in \text{Prop}_{\mathcal{N}^\bullet}(\llbracket \varphi \rrbracket) \text{ (by Lemma 4.7)} \\ & \text{ iff } w \in \llbracket \langle \mathbf{T} \rangle \varphi \rrbracket_{\mathcal{N}^\bullet} \text{ (by definition)} \\ & \text{ iff } \mathcal{N}^\bullet, w \Vdash \langle \mathbf{T} \rangle \varphi \text{ (by definition)} \end{aligned}$$

\square

COROLLARY 4.9. $\mathcal{M} \models \varphi$ 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

(K). $\mathbf{K}(\varphi \rightarrow \psi) \rightarrow (\mathbf{K}\varphi \rightarrow \mathbf{K}\psi)$

(T_K). $\mathbf{K}\varphi \rightarrow \varphi$

(4_K). $\mathbf{K}\varphi \rightarrow \mathbf{K}\mathbf{K}\varphi$

(Grz). [Update! – try to use Grz with Refl + Trans to get acyclic!]

(T_T). $\mathbf{T}\varphi \rightarrow \varphi$

(4_T). $\mathbf{T}\varphi \rightarrow \mathbf{T}\mathbf{T}\varphi$

(K-T). $\mathbf{K}\varphi \rightarrow \mathbf{T}\varphi$

that is closed under:

(Necessitation). If $\varphi \in \mathbf{L}$ then $\Box\varphi \in \mathbf{L}$ for $\Box \in \{\mathbf{K}, \mathbf{T}\}$

DEFINITION 5.2. [19, Definition 2.30] **(Deduction for L)** $\vdash \varphi$ iff either φ is an axiom, or φ 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] Γ is **consistent** iff $\Gamma \not\vdash \perp$. Γ is **maximally consistent** if Γ 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 set Γ 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 Δ be maximally consistent, and let $\Box \in \{\mathbf{K}, \mathbf{T}\}$. 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 $\Sigma = \{\psi \mid \Box\psi \in \Delta\}$ [why is Σ maximally consistent?]. Then by construction, for all ψ $\Box\psi \in \Delta$ implies $\psi \in \Sigma$, but $\varphi \notin \Sigma$ (since $\Box\varphi \notin \Delta$). \square

5.2 Soundness

THEOREM 5.7. (**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:

(K)	Monotonicity of Reach
(T _K)	Inclusion of Reach
(4 _K)	Idempotence of Reach
(Grz)	Proposition ?[Check! – and update, since the def changed]
(T _T)	Inclusion of Prop
(4 _T)	Idempotence of Prop
(T–K)	Reach contains Prop
(Necessitation)	$\forall w, w \notin \text{Reach}(\emptyset), \text{Prop}(\emptyset)$

Alternative:

$\langle W, f \rangle$ forms a filter
Reflexivity of f
Transitivity of f
f is acyclic [Check!]
Reflexivity of g
Transitivity of g
g contains f
f, g contain the unit

\square

5.3 Model Building

Given a set $\Gamma \subseteq \mathcal{L}$, I will show that we can build a net \mathcal{N} that models Γ . 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 5.8. [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 5.9. [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 φ . The propositional, and boolean cases are straightforward.

K case.

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

T case.

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

□

THEOREM 5.10. [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 5.11. If \mathcal{M} is finite and satisfies the Truth Lemma, then \mathcal{M} is a preferential filter.

Proof. W^C is finite by assumption. Since \mathbf{L} contains all instances of **(K)**, **(T)**, **(4)**, **(T)**, **(4)** it follows that f^C is a reflexive, transitive, proper filter, and g^C is reflexive and transitive (this is another classical result, see Pacuit's book). The only things left to show are that f^C is acyclic and f^C is the skeleton of g^C .

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. \mathbf{L} contains all instances of **(K)**, from which we can derive all instances of:

$$\text{(C)} \quad \mathbf{K}\varphi \wedge \mathbf{K}\psi \rightarrow \mathbf{K}(\varphi \wedge \psi)$$

Suppose $|\varphi|_{\mathbf{L}}, |\psi|_{\mathbf{L}} \in f^C(\Delta)$. 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 **(C)**, $\mathbf{K}(\varphi \wedge \psi) \in \Delta$. So $|\varphi \wedge \psi|_{\mathbf{L}} = |\varphi|_{\mathbf{L}} \cap |\psi|_{\mathbf{L}} \in \Delta$.

f^C is closed under superset. \mathbf{L} contains all instances of **(K)** and the necessitation rule, from which we can derive:

$$\text{(RM)} \quad \text{If } \varphi \rightarrow \psi \in \mathbf{L} \text{ then } \mathbf{K}\varphi \rightarrow \mathbf{K}\psi \in \mathbf{L}$$

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

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

$$\text{(N)} \quad \mathbf{K}\mathbf{T}$$

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

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

f^C is transitive. Suppose $|\varphi|_L \in f^C(\Delta)$. By definition of f^C , $\mathbf{K}\varphi \in \Delta$. By the **(4_K)** 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 $\{\Sigma | \mathbf{K}\varphi \in \Sigma\} \in f^C(\Delta)$. That is, $\{\Sigma | |\varphi|_L \in f^C(\Sigma)\} \in f^C(\Delta)$, and we are done.

f^C is acyclic. [Update! – no more nominals, acyclic rule has changed to Grz!] Since f^C is reflexive and transitive, by Proposition 2.7 it's enough to show that f^C is asymmetric. 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 $\Delta_1 = \Delta_2$. I'll show the (\subseteq) direction (the other direction is similar). Suppose for contradiction that $\varphi \in \Delta_1$, but $\varphi \notin \Delta_2$ (i.e. $\neg\varphi \in \Delta_2$).

Since Δ_1 is named, some $i \in \Delta_1$. By **(Antisym)**, $\mathbf{K}(\langle \mathbf{K} \rangle i \rightarrow i) \in \Delta_1$. By (2), $\langle \mathbf{K} \rangle i \rightarrow i \in \Delta_2$. Rewriting, we get $\neg i \rightarrow \mathbf{K}\neg i \in \Delta_2$. [What next?]

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 **(T_T)** instead of **(T_K)**.

g^C is transitive. Similar to the proof for f^C , but apply **(4_T)** instead of **(4_K)**.

g^C contains f^C . [Check]

□

THEOREM 5.12. (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 Γ to maximally consistent Δ using Lemma 5.4. Let \mathcal{M}^C be a canonical model for **L** guaranteed by Lemma 5.8. By the Truth Lemma (Lemma 5.9), $\mathcal{M}^C, \Delta \models \Delta$. So in particular, $\mathcal{M}^C, \Delta \models \Gamma$.

By the Finite Model Property (Lemma 5.10), we can construct a finite model \mathcal{M}' satisfying exactly the same formulas at all worlds. By Proposition 5.11, \mathcal{M}' is a preferential filter.

From here, we can build our net \mathcal{N}^* as before, satisfying exactly the same formulas as \mathcal{M} at all neurons (by Theorem 4.8). And so $\mathcal{N}^*, \Delta \models \Gamma$. □

THEOREM 5.13. (Completeness) For all consistent $\Gamma \subseteq \mathcal{L}$, if $\Gamma \models_{\text{BFNN}} \varphi$ then $\Gamma \vdash \varphi$

Proof. Suppose contrapositively that $\Gamma \not\vdash \varphi$. This means that $\Gamma \cup \{\neg\varphi\}$ is consistent, i.e. by Theorem 5.12 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 $\Gamma \not\models \varphi$. □

TODO:

- Double-check properties for canonical model & completeness
- Do filtration/finite model property
- Get bound on the size of the finite model.
- Think about complexity of decidability of the logic (but only if it seems easy)
- Copy-paste flipping \wedge, \vee, \neg considerations
- Write up fuzzy network considerations (in a crisp (non-fuzzy) language) — fuzzy nets satisfy *exactly* the same crisp formulas as binary nets
- Make drawings in Tikz
- Make corrections Saul gave
- Close the canonical model under superset
- Put the page number/theorem number for each result

- Rename the axioms to something more readable ((T_T) is confusing as hell)

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] Patrick Blackburn, Maarten De Rijke, and Yde Venema. *Modal logic: graph. Darst*, volume 53. Cambridge University Press, 2001.
- [5] Reinhard Blutner. Nonmonotonic inferences and neural networks. In *Information, Interaction and Agency*, pages 203–234. Springer, 2004.
- [6] Antony Browne and Ron Sun. Connectionist inference models. *Neural Networks*, 14(10):1331–1355, 2001.
- [7] Dov M Gabbay, Ian Hodkinson, and Mark A Reynolds. Temporal logic: mathematical foundations and computational aspects. 1994.
- [8] 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.
- [9] Artur S d'Avila Garcez, Luis C Lamb, and Dov M Gabbay. *Neural-symbolic cognitive reasoning*. Springer Science Business Media, 2008.
- [10] 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.
- [11] 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.
- [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.

Appendix A Helper Proofs

Proof. (of Proposition 4.2) To show that hash is injective, suppose $\text{hash}(S_1) = \text{hash}(S_2)$. So $\prod_{m_i \in S_1} p_i = \prod_{m_i \in S_2} p_i$, and since products of primes are unique, $\{p_i | m_i \in S_1\} = \{p_i | m_i \in S_2\}$. And so $S_1 = S_2$.

To show that hash is surjective, let $x \in P_k$. Now let $S = \{m_i | p_i \text{ divides } x\}$. Then $\text{hash}(S) = \prod_{m_i \in S} p_i = \prod_{(p_i \text{ divides } x)} p_i = x$. \square

Step 5. Step away (for a few days). Come back and check the proof *slowly* to make sure there aren't any missing edge cases or conditions.

- If it's all good — congratulations, you got a free paper!
- Usually there will be some idiotic mistake in the proof. It may seem like *you're* the idiot for trying it — but in fact, it's now your job to figure out *what conditions will make this naive proof work!*

Step 10. Move on to the write-up stage. But otherwise, step away from the problem — there are too many other interesting things to spend all of your time on this one. Trust that one day a different solution will come to you.