NEURAL NETWORK SEMANTICS

Caleb Schultz Kisby

Submitted to the faculty of the University Graduate School

in partial fulfillment of the requirements

for the degree

Doctor of Philosophy

in the School of Informatics, Computing, and Engineering,

Indiana University

May 23, 2025

Accepted by the Graduate Faculty, Indiana University, in partial fulfillment of the requirements for the degree of Doctor of Philosophy.

Doctoral Committee

Lawrence S. Moss, PhD (Co-advisor)

Saúl A. Blanco, PhD (Co-advisor)

David B. Leake, PhD

Chung-chieh Shan, PhD

May 23, 2025

Copyright © 2025 Caleb Schultz Kisby [Dedication page]

# Acknowledgments

[Write the acknowledgments page]

# Preface

[Write the preface, which includes the different conferences + workshops this work was presented at.]

# Caleb Schultz Kisby

# **Neural Network Semantics**

sdfafsfl;asjdfa;lkjf l;sdaj fl;asjf kl;asjfd ;lasjf ;lkasjf ;klsajf l;ksaj f;lksajf;klasj f;klsaj f;klsaj f;klasj f;klsaj f;klsaj f;klsaj f;lkjs a;flkj sa;lkfj s;klafj ;slakjf ;lasfj

**References List:** 

History of neural network semantics. [44] [7] [38] [39] [17] [40] [34] [35] [28] [?] [29]

Social network semantics. [8]

**Dynamic logics for learning.** [11] [9] [10]

**Dynamic epistemic logic; Belief revision.** [69] [66] [71] [12] [13] [56] [36] [70] [68] [14]

General neuro-symbolic AI. [5] [60] [15] [31] [20] [26] [6] [27] [43] [19]

Neural network verification. [3]

AI/neural network Alignment.

**Neural networks as automata.** [46] [64] [48] [47]

Neural network descriptive complexity. [33] [42] [21]

Neural networks & Category theory.

**General conditional logic.** [37]

**General modal logic.** [49] [55]

**Classic papers in AI.** [32] [63] [30] [54] [61] [45] [59] [72] [65] [62]

General TSC/mathematics. [2] [57]

**General cognitive science.** [51]

Systems and frameworks. [50] [23] [1]

[Note to self so I don't forget—use new-dpage and new-dpage\* commands! (follow them with noindent!)]

Lawrence S. Moss, PhD (Co-advisor)

Saúl A. Blanco, PhD (Co-advisor)

David B. Leake, PhD

Chung-chieh Shan, PhD

May 23, 2025

# Contents

Dedication		
Acknowledgments		
Preface		
Abstract vii		
Introduction		
Background: Defeasible Reasoning in Artificial Intelligence		
1 Defeasible Reasoning in Conditional Logic 5		
2 Defeasible Reasoning in Modal Logic		
3 Dynamic Epistemic Logic and Belief Revision		
4 Defeasible Reasoning in Neural Networks		
Neural Network Semantics		
1 Introduction		
2 Neural Network Models		
3 Properties of Clos, Reach, and Reach <sup><math>\downarrow</math></sup>		
4 Neural Network Semantics for Conditional Logic		
5 Neural Network Semantics for the Modal Logic of C		
6 Proof System and Soundness 27		
7 Reflections on Neural Network Verification		
Interlude: A History and Survey of Neural Network Semantics		
Neural Network Model Constructions		
1 Introduction		

2	The Canonical Neural Network Construction		
3	Filtration: Building a Finite Neural Network		
4	Model Building and Completeness Theorems		
5	The Modelling Power of Neural Networks		
6	Reflections on Model Building and Interpretability		
D	<b>Dynamic Update in Neural Network Semantics</b>		
1	Introduction		
2	Hebbian Learning: A Simple Neural Network Update Policy		
3	Properties of Hebb and Hebb*		
4	Neural Network Semantics for Hebbian Update		
5	A Complete Logic of Iterated Hebbian Update		
6	Reflections on Neural Network Alignment		
N	eural Network Semantics for First-Order Logic		
1	Introduction		
2	Lifting a Modal Logic to First-Order Logic		
3	Neural Network Semantics for First-Order Logic		
4	Axioms, Soundness, and Frame Conditions		
5	Reflections on First-Order Reasoning using Neural Networks		
С	onclusions		
1	Open Questions		
B	ibliography		
A	ppendix A		
A	.1 Appendix for Section 3		
A	.2 Appendix A2		
А	.3 Appendix A3		

### **Chapter 1**

#### Introduction

In the last 15 years, modern machine learning systems have achieved unprecedented success in learning from data with little human guidance. Consider, for example, large language models such as Llama and GPT [1; 23; 72], which have taken the world by storm with their ability to learn to converse in English solely from unstructured text data scraped from the Internet. Or, consider AlphaGo [62], which learned to play Go at a human expert level by repeatedly playing against itself. These breakthroughs in machine learning are in large part thanks to the widespread use of neural networks—brain-inspired computational models that are flexible and excel at learning from unstructured data.

But the danger of neural networks is that they come with no guarantees of safety, reliability, or correctness. Neural systems can carry all sorts of misconceptions, make silly logical mistakes, and are quite happy to spread disinformation [25; 41; 52; 65]. Diagnosing, correcting, and understanding the source of these errors is not feasible, so in practice neural networks are treated as 'black-boxes' [cite]. In this sense, neural networks lack transparency and interpretability.

How can we "open the black box" to better understand, reason about, and guide the behavior of neural networks? The field of *neuro-symbolic AI* has emerged to answer this question [5; 15; 60]. Consider symbolic systems—logics—which have been used in theoretical computer science to reason about and guide the behavior of computational systems. Historically, logics were used to model intelligent behavior prior to the rapid growth of neural network systems [cite]. Whereas neural networks are not easily interpreted, logics provide transparent access to their reasoning via a human-interpretable language and compositional semantics. But also unlike neural networks, traditional logics fail to model flexible learning or update. (This is known as the *frame problem* in AI [22; 45; 61]).

The primary goal of neuro-symbolic AI is to combine the complementary strengths of neural networks and logic, retaining the advantages of both. There are now many distinct proposals for

neuro-symbolic systems—to many to count! To give a short but representative list, these systems include Logic Tensor Networks [6], Distributed Alignment Search [27], DeepProbLog [43], Logic Explained Networks [19], neural networks as automata [46; 48; 73], and neural network fibring [26]. These systems and algorithms are disparate, and there is little agreement among them on how to interface neural and symbolic systems. Together these proposals form a scattered picture, and some unifying perspective or theory is needed. In the preface to a recent neuro-symbolic survey book [15], Frank van Harmelen writes:

What are the possible interactions between knowledge and learning? Can reasoning be used as a symbolic prior for learning ... Can symbolic constraints be enforced on data-driven systems to make them safer? Or less biased? Or can, vice versa, learning be used to yield symbolic knowledge? ... And how to avoid the inevitable bias seeping from the data into the resulting knowledge base? ... **neuro-symbolic systems currently lack a theory that even begins to ask these questions, let alone answer them.** 

In this dissertation, I will explore a class of neuro-symbolic systems that tightly couple logic with neural network dynamics. This may seem like adding yet another neuro-symbolic proposal to the pile. But in fact, the point of this work is to develop and promote a broader perspective that deepens our understanding of well-known neuro-symbolic systems and informs the development of neuro-symbolic interfaces in the future.

Here is the perspective I'm promoting: Neural networks can be thought of as models for formal logic, i.e., logics can be interpreted directly on neural networks. Speaking more technically, logical operators such as conditionals  $\Rightarrow$ , modalities  $\Box$ , dynamic updates [*P*], and quantifiers  $\forall x$  can be interpreted by semantics that refer to, e.g., the input-output behavior, signal propagation, or update of weights in a neural network. I will soon show that many neuro-symbolic systems can be considered from this point of view; I will call any specification of formal semantics that falls under this umbrella "*neural network semantics*."

In previous work, the input-output behavior of (binary) neural networks was shown to be an

appropriate semantics for defeasible conditionals  $\varphi \Rightarrow \psi$  (read "typically, if  $\varphi$  then  $\psi$ ") [7; 17; 20; 38; 39]. In this dissertation, I develop this same idea in more expressive logical contexts, including modal logic [16; 67], Dynamic Epistemic Logic (DEL) [68; 70; 71], and first-order logic (FOL) [24]. I will focus on three, highly nonstandard, logics that are interpreted on binary neural networks:

- 1. A defeasible modal logic whose central modality  $C\varphi$  is interpreted via the closure (or forward propagation) of the signal  $\varphi$  in a neural network. (Chapters 3 and 4)
- A DEL whose update operator [P] is interpreted as Hebbian update (on input P) [32], a simple neural network learning policy that changes weights according to the rule "neurons that fire together wire together." I will consider a logic of single-step Hebbian update, as well as a logic of iterated Hebbian update. (Chapter 5)
- 3. A defeasible FOL whose central quantifier  $\forall x \varphi$  is interpreted via the closure (or forward propagation) of the signal  $\varphi$  in a specially-designed neural network that uses variable assignments in place of neurons. (Chapter 6)

As a logician, I will take each of these semantics seriously. I will develop them by mathematically proving formal logical results, including soundness, completeness, expressivity, and frame correspondence theorems. These may seem like results that merely satisfy a logician's curiosity. But my point is that, in the context of neural network semantics, these theorems enable us to directly answer fundamental questions about neural network inference, learning, and neuro-symbolic interfaces. Here is an explicit list of questions this perspective helps us answer, along with the chapter in which I answer them:

- What algebraic properties hold for neural network inference? Can we design an expressive logic that can be used to formally verify these properties? (**Chapter 3**)
- Can we build a neural network that obeys a given set of constraints on its inference (stated as formulas in an expressive logic)? (Chapter 4)
- How powerful are neural network models, when compared to other classes of models? What kinds of models can neural nets simulate? And vice-versa, what kind of models can

be simulated by neural nets? (Chapter 4)

- What algebraic properties hold for neural network learning policies? Can we design a *dynamic* logic that can be used to formally verify properties of neural network learning?
  (Chapter 5)
- In general, can we build a neural network that obeys a set of constraints on its inference *before and after learning takes place*? Note: This is one of the keys to the AI Alignment problem. (Chapter 5)
- Can we design a neuro-symbolic system that supports reasoning in full first-order logic?
  What kinds of neural network properties guarantee the soundness of FOL axioms? (Chapter 6)

In summary, I will defend the following thesis statement:

**Thesis Statement.** Neural networks can be treated as a class of models in formal logic, where logics are interpreted directly on the neural net ("*neural network semantics*"). By developing this idea in richer logical contexts, we are able to answer fundamental questions about neural network inference, learning, and neuro-symbolic interfaces in general.

### Chapter 2

#### **Background: Defeasible Reasoning in Artificial Intelligence**

The connection between neural networks and formal logic begins with defeasible reasoning (aka nonmonotonic reasoning, or reasoning by default). In standard treatments of logic, the facts you infer are non-revocable, i.e., they cannot be withdrawn in light of new information. But we live in a world of change, partial information, and exceptions—in order to effectively reason, an agent must jump to conclusions about what is "normally" or "plausibly" the case, and be ready to withdraw these inferences. For these reasons, defeasible reasoning is a central to a theoretical understanding of artificially intelligent agents.

Here's a classic example: If you know Tweety is a bird, you should conclude (assuming we're in a "normal" situation) that Tweety flies. But if you then discover that Tweety is a penguin, you must retract that conclusion. The standard material implication fails to model this: If Tweety $\rightarrow$ penguin, penguin $\rightarrow$ bird, and bird $\rightarrow$ flies we must conclude that Tweety flies.

In this chapter, I will give a tour of many different ways to model defeasible reasoning in formal logic. I will focus on the "preferential" or "plausibility" approach to defeasible reasoning, which branches from the classic papers [37] and [cite Shoham 1988]. First, I will present the standard plausibility semantics for conditional logics (where  $\varphi \Rightarrow \psi$  expresses "typically,  $\varphi$  are  $\psi$ "). Then I will discuss many different ways to transfer these semantics to more expressive modal logics. I will present the logic of  $\mathbf{C}\varphi$  ("the current state is the best one where  $\varphi$  holds"), which forms the backbone of my work connecting neural networks and logic. Finally, I will introduce neural networks, and discuss how they may be seen as models of defeasible reasoning as well. This will set us up for the central plot of my thesis: Developing a neural network semantics for the logic of  $\mathbf{C}\varphi$ .

#### **1** Defeasible Reasoning in Conditional Logic

I will now present the standard way to model nonmonotonic inference in conditional logic, in the KLM tradition [37]. The language is stratified—sentences are conditionals  $\varphi \Rightarrow \psi$ , where  $\varphi, \psi \in \mathcal{L}_{prop}$  are propositional formulas connected by  $\neg, \land, \rightarrow$  in the usual way. Sentences  $\varphi \Rightarrow \psi$ cannot be nested within each other, nor within propositional formulas. This odd feature is due to the original conception in [37] that  $\varphi \Rightarrow \psi$  specify inference rules, but are not themselves propositions. The intended meaning of  $\varphi \Rightarrow \psi$  is "typically (normally),  $\varphi$  are  $\psi$ ", e.g., bird $\Rightarrow$ flies reads "typically, birds fly."

Kraus, Lehmann, and Magidor use the following models to interpret these conditional sentences. I will be moving on pretty quickly to modal logic syntax and semantics, so I won't dwell on these models too long. Let W be an underlying set of worlds (propositional valuations) for  $\mathcal{L}_{prop}$ (not necessarily the set of all worlds for  $\mathcal{L}_{prop}$ ).

**Definition 1.1.** A cumulative-ordered model is  $\mathcal{M} = \langle S, l, \prec \rangle$ , where

- S is a nonempty set of states
- $l: S \to \mathcal{P}(W) \{\emptyset\}$  (a *labelling* of states)
- $\prec: S \times S$  (the *plausibility order*, or *preference relation*)

The plausibility order  $\prec$  is required to be a strict order relation (irreflexive and transitive) satisfying the Smoothness Condition (which I will state shortly).  $S_1 \prec S_2$  intuitively means that the agent considers the state  $S_1 \in S$  to be more plausible, or more normal, than  $S_2 \in S$ . In order to reason about the most plausible (normal) states, we can look at the  $\prec$ -minimal states. Formally, each cumulative-ordered model determines a function best $_{\prec}: S \rightarrow S$ 

$$best_{\prec}(S) = \{w \in l(S) \mid For all \ u \in l(S), \neg u \prec w\}$$

For propositional formulas  $\varphi \in \mathcal{L}_{\text{prop}}$ ,  $[\![\varphi]\!] = \{S \in S \mid w \models \varphi \text{ for all } w \in l(S)\}$ , i.e., the set of states where  $\varphi$  is true everywhere. I said before that  $\prec$  must satisfy the Smoothness Condition [37]—this condition says that for any propositional formula  $\varphi \in \mathcal{L}_{\text{prop}}$ ,  $[\![\varphi]\!]$  has no infinitely descending  $\prec$ chains, i.e., every non-empty  $[\![\varphi]\!]$  has at least one minimal element.

**Postulate 1.2.** For all cumulative-ordered models  $\mathcal{M}$ , states  $S \in S$ , and all  $\varphi \in \mathcal{L}_{prop}$ , if  $S \in [\![\varphi]\!]$  then either  $S \in \text{best}_{<}([\![\varphi]\!])$ , or there is some  $S' \prec S$  better than S that *is* the best, i.e.  $S' \in \text{best}_{<}([\![\varphi]\!])$ .

I can now give the KLM intepretation of conditional sentences:

$$\mathcal{M} \models \varphi \Rightarrow \psi$$
 iff  $\mathsf{best}_{<}(\llbracket \varphi \rrbracket) \subseteq \llbracket \psi \rrbracket$ 

That is, in the most plausible (normal) states where  $\varphi$  holds,  $\psi$  holds, which was our intended reading. There is a lot more to say about conditional logics like these (expressivity, proof systems, soundness and completeness, their rich history), but I must move on. I will conclude with an example demonstrating that these semantics do in fact resolve our earlier issue with Tweety the penguin.

Example 1.3. [Give an example of these semantics successfully modeling defeasible reasoning. Maybe the Tweety example?]

#### 2 Defeasible Reasoning in Modal Logic

The inability to nest conditionals  $\varphi \Rightarrow \psi$  makes conditional logics somewhat flat and inexpressive. Additionally,  $\varphi \Rightarrow \psi$  only allows us to refer to the plausibility of the premise  $\varphi$ , and not the antecedent  $\psi$ . For example, the following sentences are not expressible in the conditional language above:

- If birds typically fly, then Tweety does.
- The car normally drives, but the check engine light is always on.
- This was not your typical criminal.
- If this isn't normal, I don't know what is.

We can overcome this by transferring the main ideas of the semantics to a more expressive language—in particular, to modal logics. [is there anything more I need to say here to motivate the reader?]

### 2.1 A Brief Crash Course in Modal Logic

Let's briefly introduce the basics of modal logic. [cite a standard modal logic text or two!] A modal logic extends propositional logic with "modal formulas"  $\Box \varphi$  and  $\diamond \varphi$  ( $\Box \varphi$  is read "it is necessary that  $\varphi$ ";  $\diamond \varphi$  is read "it is possible that  $\varphi$ ." Standard (normal) modal logics are interpreted using a relational (Kripke) model, which is just an ordinary graph equipped with a valuation of propositions.

**Definition 2.1.** A relational model is  $\mathcal{M} = \langle W, R, V \rangle$ , where *W* is a set of nodes (*worlds*, aka *states*), *R*: *W* × *W* an edge relation (the *accessibility relation*), and *V*: propositions  $\rightarrow \mathcal{P}(W)$  (the *valuation function*).

**Definition 2.2.** Let **Rel** be the class of all relational models, and let  $\text{Rel}_{S4}$  be the class of all whose accessibility relation *R* is reflexive and transitive.

Unlike conditional logic, in modal logics we evaluate a formula *locally*. That is, instead of  $\varphi$  being true or false, we consider the set of worlds where  $\varphi$  is true. We write  $\mathcal{M}, w \Vdash \varphi$  to indicate that  $\varphi$  holds at world w. The semantics of propositions and boolean connectives  $\neg$ ,  $\land$  are what you might expect.  $\diamond \varphi$  is defined as  $\neg \Box \neg \varphi$ . The key case is for  $\Box \varphi$ :

 $\mathcal{M}, w \Vdash \Box \varphi$  iff for all *u* such that *wRu*, we have  $\mathcal{M}, u \Vdash \varphi$ 

That is,  $\Box \varphi$  holds if  $\varphi$  holds everywhere accessible from the current state ( $\varphi$  is *necessarily* true).

The accessibility relation can have many different interpretations depending on what phenomenon we are trying to model. For example, if *R* indicates which states are *possibly known* (i.e., *epistemically accessible*), then  $\Box \varphi$  takes on the reading " $\varphi$  is known (by some agent)," written box $\varphi$ . Similarly,  $\Box$  can be cast as belief **B**, obligation **O**, provability **P**, etc. There may be one, or many, modal operators in a modal logic. We may also index modalities  $\Box_i$ , indicating a modal attitude for each agent *i* (in a multi-agent setting), or for different relations  $R_i$  within the same agent.

#### 2.2 Defeasible Modal Logics

Unfortunately, this usual treatment of modal logics cannot model defeasible reasoning. [cite Chellas or something] This is because all normal modal logics satisfy the axiom

$$\Box(\varphi \to \psi) \to \Box \varphi \to \Box \psi$$

For concreteness, let's read  $\Box$  as belief, and suppose  $\Box$ (bird $\rightarrow$ flies), i.e., for all things we could possibly believe we see from the current state, if we see a bird then it flies. Now say we believe we see a bird ( $\Box$ bird). Then the axiom says we *necessarily* believe it flies ( $\Box$ flies), leaving no room for revoking our initial conclusion.

There is a wide variety of different ways to resolve this, to rework the idea of defeasibility into modal logic. This is not the right time or place for a thorough literature review, but I will tour a representative sample to give you a sense of what can be done. Certain approaches use the cumulative-ordered models defined for conditional logic above; others use relational models, but interpret the relation *R* to be a plausibility ordering; others still use both. For my purposes, I will define plausibility models as Krikpe models, but with an irreflexive plausibility relation  $\prec$  (there is no distinction between epistemic accessibility and plausibility). [cite plausibility models, the word is used by Baltag & Smets]

**Definition 2.3.** A plausibility model is  $\mathcal{M} = \langle W, \langle, V \rangle$ , where

- *W* is a set of *worlds* or *states*
- *R*: *W* × *W* (the *epistemic accessibility relation*)
- $\prec: W \times W$  (the *plausibility order*)
- *V*: [todo] (the *propositional valuation*)

[(Basically a cumulative-ordered model along with an epistemic frame)] I require *R* to be reflexive and transitive [Connected? Symmetric?]. As with cumulative-ordered models, I require < to be irreflexive, transitive, antisymmetric. In cases where we want to refer to the reflexive extension of <, I write  $u \le v$  to mean u < v or u = v. As before, each plausibility model determines a best<sub><</sub> function, whose definition now simplifies to

$$best_{\prec}(S) = \{w \in S \mid For all \ u \in S, \neg u \prec w\}$$

Similarly, we require the best<sub><</sub> operator to satisfy a similar Smoothness condition. For *any* formula  $\varphi$  (not just propositional formulas), say we have defined some semantics  $\parallel$ , and let  $\llbracket \varphi \rrbracket = \{w \in W \mid \mathcal{M}, w \Vdash \varphi\}$ . Smoothness says:

**Postulate 2.4.** For all plausibility models  $\mathcal{M}$ ,  $w \in W$ , and formulas  $\varphi$ , if  $w \in \llbracket \varphi \rrbracket$  then either  $w \in best_{<}(\llbracket \varphi \rrbracket)$ , or there is some  $v \prec w$  better than w that *is* the best, i.e.  $v \in best_{<}(\llbracket \varphi \rrbracket)$ .

[Should I move this postulate down past the semantics? It would be nice if I defined the language first, so I could say  $\varphi \in \mathcal{L}_{C}$ .]

**Definition 2.5.** Let **Plaus** be the class of all such plausibility models.

Here are some of the ways we can transfer defeasibility into a modal logic setting:

# **Boutilier's Modal Treatment.**

#### Baltag & Smets' Safe Belief.

• There are many ways to rework the idea of typicality in modal logic: conditional belief  $B^{\varphi}\psi$ , regular belief  $B\varphi$ , typicality  $\cdot, T$ , "defeasible modalities"

#### **3** Dynamic Epistemic Logic and Belief Revision

[This is really one of the best upshots of modeling something using modal logic! If formulas are nestable, then we can nest things inside and within update operators.] [Introduce Dynamic Epistemic Logic, and various operators Cond, Lex, Consr for belief revision.]

[Include some parts of this in this section]

**Other Dynamic Logics for Learning.** My approach to modeling learning in neural networks takes inspiration from Dynamic Epistemic Logic (DEL) [68; 71]. Perhaps the closest logics to mine are the logics for plausibility upgrade, in particular conditionalization (Cond), lexicographic (Lex), and conservative (Consr) upgrade [66; 70]. In the Expressivity portion of my dissertation, I will explore whether Hebbian update Hebb\* can be simulated by plausibility upgrade, and vice-versa whether Cond, Lex, Consr, and vice-versa whether these plausibility upgrades can be simula.

In my thesis, I mainly focus on the effects of single-step updates. But recent literature on learning in DEL goes beyond this by considering iterated update and convergence to the truth ("learning in the limit") [9; 10; 11; 14]. The key questions here are: How can we compare the

learning power of different iterated update policies? How can we axiomatize important properties of learning? These questions are answered in terms of updates on more classical graphs and plausibility structures. Although in this thesis I don't consider iterated update, I do lay down the groundwork to import neural network learning into this setting.

#### **4** Defeasible Reasoning in Neural Networks

Introduce neural networks, for a broad audience (including logicians that know nothing about them). Explain how inference works in a neural network, and how neural networks can be thought of as performing defeasible reasoning. (Hannes explains it pretty well, maybe I should borrow his example.)

Explain learning in a neural network at a high level, both unsupervised (Hebbian learning is representative of unsupervised learning algorithms, mention the relationship with Principal Components Analysis) and supervised (backpropagation is representative here, and the efficient computation of backpropagation is what has made neural networks so successful).

Inference in a neural network is "like" a conditional inference—but this analogy goes further. Many authors have already studied a formal correspondence between the input-output behavior of a neural network and defeasible conditionals. [cite d'Avila Garces, Hannes, Giordano, really all the people here who have made the observation before me. This is a good time to break down the history of how it happened.] [Talk about both "soundness" and "completeness" here]

In the rest of this chapter, I will extend this analogy by giving a neural network interpretation for the more general logic of  $\mathbf{C}$ . My main point in considering a modal language is the same as before: It buys us expressive power over conditionals, and in particular sets us up to express *neural network update* using Dynamic Epistemic Logic. Towards the end of this work, I will extend these neural network semantics for  $\mathbf{C}$  with a dynamic operator for a simple Hebbian update policy over neural networks.

# **Chapter 3**

# **Neural Network Semantics**

### **1** Introduction

[I wrote the following short paragraph for my thesis proposal. But now I have more space to say more, slowly, about where neural network semantics comes from, what the underlying idea is. It would be nice to introduce it using an example of a neural network in practice—justify why binary, interpreted neural networks are the "right" thing to look at!]

I will now give an overview of the particular neural network semantics [rephrase] I've developed during my PhD. First, I will discuss the simplifying assumptions that make it possible to use neural networks as models, and introduce the *closure* (or forward propagation) of a signal in the net. This closure operator allows us to express the inference, or input-output behavior, of the net. I will give a modal logic whose key operator is given by this closure operator. I will then turn to dynamic update in neural networks and introduce iterated Hebbian learning, one of the simplest learning policies over nets. Finally, I will give a dynamic logic whose formulas express the behavior of a neural network before and after Hebbian update.

#### **2** Neural Network Models

A model of neural network semantics is an artificial neural network (ANN), along with a valuation function which interprets propositions as sets of neurons. I will make a few more simplifying assumptions soon, but this is the basic idea.

**Definition 2.1.** A neural network model is  $\mathcal{N} = \langle N, \text{bias}, E, W, A, \eta, V \rangle$ , where

- *N* is a finite nonempty set (the set of *neurons*)
- bias  $\in N$  is a fixed node (the *bias* node)
- Each  $E \subseteq N \times N$  (the *edge relation*)
- $W: E \to \mathbb{Q}$  (the *edge weights*)
- $A: \mathbb{Q} \to \mathbb{Q}$  (the activation function)
- $\eta \in \mathbb{Q}, \eta \ge 0$  (the *learning rate*)
- $V: \mathcal{L}_{\text{prop}} \rightarrow \mathcal{P}(N)$  (the valuation function)

In general, a *state* is just a possible activation pattern or configuration of the net. In practice, a neural network's nodes take on fuzzy activation values in [0, 1]. But we would like to associate each state with a binary set of neurons—either a neuron is active (1) or it is not (0). To do this, I assume that the activation function A is a (nonzero) binary step function ( $A: \mathbb{Q} \rightarrow \{0, 1\}$ ). [Definition:  $\mathcal{N}$  has a threshold,  $\exists t \in \mathbb{Q}$  with A(t) = 1;  $\mathcal{N}$  is nondecreasing. These things amount to A being a binary step function.] It turns out this binary constraint is also a common theoretical assumption in work that analyzes neural networks as automata [46; 48; 73]. In their terminology, we say our nets are *saturated*.

- [Todo, put somewhere in this section, optional property]
- *N* is Fully connected: ∀m, n ∈ N, either (m, n) ∈ E, (n, m) ∈ E, or m and n have exactly the same predecessors and successors.
- In machine learning practice, "fully connected" means that there is an edge from every node in layer *l* to every node in the *following* layer *l* + 1. But here we mean something much stronger: the graph is fully connected, including "highway edges" that cut between layers, as shown in [DIAGRAM]. (This intuition comes from work on highway networks [63].)

This assumption is crucial for our results about iterated Hebbian learning, and we expect that letting it go will not be easy.

Since I'm only considering binary neural networks, either a neuron is active (1) or it is not (0). A *state* (activation pattern) of the net is just a set of neurons that happen to be active. Additionally, I assume there is a special bias node that is active in every state. This bias node provides provides persistent input to the neural network; its purpose is to avoid the edge case where the input  $S_0$  is empty, and so the net doesn't have any "fuel" to activate any new nodes. (Since bias is always active, we do not need any edges going into it. So I assume bias has no predecessors.) Putting all this together, the states of the net are defined as follows.

State 
$$\mathcal{N} = \{S \mid S \subseteq N \text{ and } bias \in S\}$$

Usually  $\mathcal{N}$  is understood from context, and I'll just write State without the subscript.

**Example 2.2.** [Give an example of a neural network, with weights and bias, and valuation function, and give examples of states (including the bias node!) Does the valuation function necessarily have to be a state? i.e. does it necessarily have to include bias when mapping propositions? I believe it shouldn't, and this example could illustrate that (e.g. we map p to  $\{a, b, c\}$ , but despite the mapping note that this set is not a *state* of the net, since bias is not in it. We can extend it to a state by adding bias.)]

[**TODO:** Add a note here where I discuss the difference between neural network models and standard graph models for logic. The key thing here is that neural network models are *connectionist* / represent their inferences in a distributed way. But I have to make that somewhat precise.]

#### 2.1 Clos: The Network Closure Operator

We can describe the input-output behavior of neural networks in terms of their state. Say we are given an input state A consisting of input-layer nodes, and a possible classification state B consisting of output-layer nodes. Active neurons in A subsequently activate new neurons, which

activate yet more neurons, until eventually the state of the net stabilizes. If this final state includes the output *B*, we say "the net *classifies A* as *B*".

Consider for example the neural network in [EXAMPLE where we classify Tweety as flying. Make the connection clearer—for binary nets, this is exactly what (binary) classification means!]

The state at the fixed point of this process is called the *closure* or *forward propagation* of the signal *A* through the net, Clos(A). This closure operator is central to the semantics, since it captures the underlying dynamics involved in neural network inference. Formally, each choice of *E*, *W*, *A* specifies a transition function from state  $S \in State$  to the next state. Given an initial state  $S_0$ , this transition function  $F_{S_0}$ : State<sub>N</sub>  $\rightarrow$  State<sub>N</sub> is given by

$$F_{S_0}(S) = S_0 \cup \left\{ n \mid A\left(\sum_{m \in \operatorname{preds}(n)} W(m, n) \cdot \chi_S(m)\right) = 1 \right\}$$

where  $\chi_S(m) = 1$  iff  $m \in S$  is the indicator function. In other words,  $F_{S_0}(S)$  consists of the initial state  $S_0$ , along with all those nodes that are activated by their predecessors in *S*. Notice that  $F_{S_0}(S)$  is extensive: all nodes in the initial state will stay activated in successive states.

These neural network models have one final constraint: This transition function  $F_{S_0}$  eventually gives a unique *fixed point* (or stable state) under the input  $S_0$ , i.e. a unique state *S* such that  $F_{S_0}(S) = S$ . This guarantees that the closure Clos(S) exists.

**Postulate 2.3.** I assume that for all states  $S_0$ ,  $F_{S_0}$  applied repeatedly to  $S_0$ , i.e.

$$S_0, F_{S_0}(S_0), F_{S_0}(F_{S_0}(S_0)), \dots, F_{S_0}^k(S_0), \dots$$

eventually reaches a finite fixed point, and moreover this state is the *only* fixed point under  $S_0$ . Formally, this means that for all  $S_0 \in \text{State}_N$  there is some  $k \in \mathbb{N}$  such that:

- 1.  $F_{S_0}(F_{S_0}^k(S_0)) = F_{S_0}^k(S_0)$ . That is, the activation pattern under  $F_{S_0}$  will eventually stabilize.
- 2.  $F_{S_0}^k(S_0)$  is the only state  $S \in \text{State}_N$  such that  $F_{S_0}(S) = S$ . In other words, the final state is unique for each initial state  $S_0$ .

Let the closure  $\text{Clos}: \text{State}_{\mathcal{N}} \to \text{State}_{\mathcal{N}}$  be the function that produces that least fixed point:  $\text{Clos}(S) = F_{S_0}^k(S_0)$  for that  $k \in \mathbb{N}$  above. Finally, let **Net** be the class of all binary neural network models that



Figure 2.1. [caption!]

satisfy this postulate.

**Example 2.4.** At this point, let's walk through an example run of the closure of a set Clos(S). Consider the neural network shown in Figure 2.1, and take the activation function to be A(x) = 1 iff  $x \ge 1$ . The first cell shows the neural network at initial state *S*, where only the bias node and



Figure 2.2. [Separate into (a) left, and (b) right. TODO—caption!]

one other node are active. For each neuron *n*, we activate that neuron if the weighted sum of its predecessors is  $\geq 1$ . (Note that currently inactive neurons do not contribute to this sum.)

Now repeat this process for multiple rounds k, until no new neurons activate in the current round. The first two rows of Figure 2.1 display the k = 1 round, and the final row shows an abbreviated run of the k = 2 round. The final cell gives exactly those neurons activated in the closure Clos(S). Take a moment to check that this is in fact a stable state, i.e., the remaining nodes in the final cell do not subsequently activate. (Hint: You can look at the previous cells for the weights of those nodes.)

**Note.** Not every neural network has a unique fixed point Clos(S) for every  $S \in State_N$ . Consider the recurrent neural network N and the state S shown in Figure 2.2. Take the activation function to be A(x) = 1 iff  $x \ge 0$ . The active bias neuron in S subsequently activates a, then b, then c. But once c

is active, *a* is no longer active (the weighted sum of its predecessors is now -2 + 1 = -1). And so the net continues to oscillate between these states, never stabilizing.

Characterizing those nets that have a unique fixed point for every state (i.e., characterizing **Net**) is a big open problem. Clearly, feed-forward networks do, since without recurrent edges the subsequent activations will eventually terminate. But the difference between Clos being well-defined and not is often subtle. In the net  $\mathcal{N}$  from before, if we change the weight W(bias, a) to 3 instead of 1, the resulting net will no longer oscillate!

**Open Question 1.** Is there an algebraic or topological characterization of the class of recurrent nets whose closure Clos(S) reaches a unique fixed point (i.e., does not diverge or oscillate)?

# 2.2 Reach and Reach<sup>+</sup>: The Graph Reachability Operators

[As I mentioned before] [TODO I no longer mention this! I need a different segue], a key feature of Clos is that it is not monotonic. Let's consider two closure operators over neural networks that *are* monotonic—graph reachability Reach, and "reverse" graph reachability Reach<sup>↓</sup>. Reach(*S*) just returns the set of all neurons graph-reachable from *S*, i.e. Reach: State<sub>N</sub>  $\rightarrow$  State<sub>N</sub> is given by  $n \in$  Reach(*S*) iff there exists  $m \in S$  with an *E*-path from *m* to *n*. Similarly, Reach<sup>↓</sup>(*S*) returns the set of all neurons that graph-reach some node in *S*, i.e. Reach<sup>↓</sup>: State<sub>N</sub>  $\rightarrow$  State<sub>N</sub> is given by  $n \in$  Reach<sup>↓</sup>(*S*) iff there exists  $m \in S$  with an *E*-path from *m* to *n*.

Reach and Reach<sup> $\downarrow$ </sup> are not very interesting operators in their own right, but I consider them in this discussion for two reasons. First, graph reachability is necessary for reasoning about Hebbian learning (see Chapter ?). Second, in Chapter ? I would like to compare the expressive power of neural networks against many classes of models, including relational (graph) models.

# powerp

# 3 Properties of Clos, Reach, and Reach<sup>+</sup>

The following theorem, due to [38], says that we can neatly characterize the algebraic structure of Clos a closure operator. Note that Leitgeb proves this for *inhibition nets*, i.e. weightless neural

networks with both excitatory and inhibitory connections. But inhibition nets and our nets  $\mathcal{N} \in \mathbf{Net}$  are equivalent with respect to their propagation structure—I prove this result again for **Net** as a kind of "sanity check" that my definitions are correct.

**Proposition 3.1.** (Leitgeb, [38; 39]) For all  $S, S_1, \ldots, S_k \in \text{State}_N$ ,

**Inclusion.**  $S \subseteq Clos(S)$ 

**Idempotence.** Clos(Clos(S)) = Clos(S)

**Cumulative.** If  $S_1 \subseteq S_2 \subseteq \text{Clos}(S_1)$ , then  $\text{Clos}(S_1) = \text{Clos}(S_2)$ 

**Proof.** I'll prove each in turn:

**Inclusion.** By definition,  $Clos(S) = F_S^k(S)$  for some  $k \in \mathbb{N}$ , where  $F_S^k$  is the transition function from Definition [which?]. By induction on this k (the number of iterations needed to construct the closure):

**Base Step.** k = 0, and so Clos(S) = S. So if  $n \in S$ , then  $n \in Clos(S)$ .

**Inductive Step.** Let  $k \ge 0$ , and suppose  $n \in S$ . We have  $Clos(S) = F_S(K) = F_S(F_S^{k-1}(S))$ . Expanding this term out, we have

$$F_{S}(F_{S}^{k-1}(S)) = S \cup \left\{ n \mid A\left(\sum_{m \in \operatorname{preds}(n)} W(m,n) \cdot \chi_{F_{S}^{k-1}(S)}(m)\right) = 1 \right\}$$

Since  $n \in S$ , *n* is in the left-hand side of this union. And so  $n \in Clos(S)$ .

**Idempotence.** I will prove a stronger claim: For all  $k \in \mathbb{N}$ ,

$$F_S^k(\operatorname{Clos}(S)) = \operatorname{Clos}(S)$$

In other words, applying the transition function  $F_S$  any number of times to Clos(S) has no effect. Since  $Clos(Clos(S)) = F_S^k(Clos(S))$  for some  $k \in \mathbb{N}$ , the Idempotence property follows from this. Let's proceed by induction on k.

**Base Step.** k = 0, and so the goal simplifies to Clos(S) = Clos(S), which is true.

**Inductive Step.** Let  $k \ge 0$ . We have  $F_S^k(\operatorname{Clos}(S)) = F_S(F_S^{k-1}(\operatorname{Clos}(S)))$ . Our inductive hypothesis here is that, for k-1,  $F_S^{k-1}(\operatorname{Clos}(S)) = \operatorname{Clos}(S)$ . So we can subsitute that to get  $F_S^k(\operatorname{Clos}(S)) = F_S(\operatorname{Clos}(S)) = \operatorname{Clos}(S)$ , which is what we wanted to show.



Figure 3.1. [Separate into (a) left, and (b) right. TODO—caption!]

**Cumulative.** Suppose  $S_1 \subseteq S_2 \subseteq \text{Clos}(S_1)$ , as illustrated in Figure 3.1 (a). First, since  $\text{Clos}(S_2)$  is a fixed point under  $S_2$ , we have  $F_{S_2}(\text{Clos}(S_2)) = \text{Clos}(S_2)$ . But I will show that  $\text{Clos}(S_1)$  is *also* a fixed point under  $S_2$ , i.e.  $F_{S_2}(\text{Clos}(S_1)) = \text{Clos}(S_1)$  (notice the difference in the subscripts). Since we assumed that there is a *unique* fixed point under  $S_2$ , it will follow that these two states must be the same. In other words,  $\text{Clos}(S_1) = \text{Clos}(S_2)$ .

Let's expand  $F_{S_2}(Clos(S_1))$ . By definition of F,

$$F_{S_2}(\operatorname{Clos}(S_1)) = S_2 \cup \left\{ n \mid A\left(\sum_{m \in \operatorname{preds}(n)} W(m, n) \cdot \chi_{\operatorname{Clos}(S_2)}(m)\right) = 1 \right\}$$

Compare this to  $F_{S_1}(Clos(S_1))$ :

$$F_{S_1}(\operatorname{Clos}(S_1)) = S_1 \cup \left\{ n \mid A\left(\sum_{m \in \operatorname{preds}(n)} W(m, n) \cdot \chi_{\operatorname{Clos}(S_2)}(m)\right) = 1 \right\}$$

Putting the two together, we can see that

$$F_{S_2}(\operatorname{Clos}(S_1)) = (F_{S_1}(\operatorname{Clos}(S_1)) - S_1) \cup S_2$$
  
= (Clos(S\_1) - S\_1) \cup S\_2 (since Clos(S\_1) is fixed under S\_1)  
= Clos(S\_1) (since S\_1 \subseteq S\_2 \subseteq \operatorname{Clos}(S\_1))

In the terminology of [37], Prop is a *cumulative* closure operator (it satisfies the cumulative property). An important feature of Clos is that it is nonmonotonic: it is not the case that for all

 $A, B \in \text{State}$ , if  $A \subseteq B$  then  $\text{Clos}(A) \subseteq \text{Clos}(B)$ . Intuitively, this is because our net's weights can be negative, and so Clos(B) can inhibit the activation of new neurons that would otherwise be activated by Clos(A).

[**TODO:** Add a note here where I talk about this proof strategy used in the proof of the Cumulative property. This proof trick is going to show up everywhere I want to prove something important about neural network closure patterns. I borrow the trick from Hannes, but my applications of the trick (from completeness to Hebbian update to first-order properties) is the biggest technical contribution of this dissertation, and is important for the reader to internalize.]

**Proposition 3.2.** (Leitgeb, [38; 39]) It is not the case that for all  $S_1, S_2 \in \text{State}_N$ , if  $S_1 \subseteq S_2$ , then  $\text{Clos}(S_1) \subseteq \text{Clos}(S_2)$ .

**Proof.** Consider the BFNN  $\mathcal{N}$  in Figure 3.1 (b). Let the activation function A be A(x) = 1 iff x > 0. We have  $S_1 = \{a\} \subseteq \{a, b\} = S_2$ , and so the hypothesis holds. But  $Clos(S_1) = \{a, c\} \nsubseteq \{a, b\} = Clos(S_2)$ . Observe that c does not get activated in  $Clos(S_2)$  because the weights cancel each other out.  $\Box$ 

Next, let's check that Reach is in fact a monotonic closure operator.

**Proposition 3.3.** For all  $S, A, B \in \text{State}_{\mathcal{N}}$ ,

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

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

**Monotonic.** If  $A \subseteq B$  then  $\operatorname{Reach}(A) \subseteq \operatorname{Reach}(B)$ 

**Closed under**  $\cup$ . Reach $(A \cup B) = \text{Reach}(A) \cup \text{Reach}(B)$ 

**Proof.** I'll prove each in turn:

**Inclusion.** Suppose  $n \in S$ . We have the trivial path from  $n \in S$  to itself.

**Idempotent.** The  $(\leftarrow)$  direction follows from inclusion. As for  $(\rightarrow)$ , suppose  $n \in \text{Reach}(\text{Reach}(S))$ , i.e. there is a path from some  $m \in \text{Reach}(S)$  to n. By definition of Reach again, there is a path from some  $x \in S$  to m. But we can put these together to obtain a path from  $x \in S$  to n.

- **Monotonic.** Suppose  $A \subseteq B$ , and let  $n \in \text{Reach}(A)$ . By definition of Reach, we have a path from some  $m \in A$  to n. But since  $A \subseteq B$ ,  $m \in B$ . So we have a path from  $m \in B$  to n, i.e.  $n \in \text{Reach}(B)$ .
- **Closed under**  $\cup$ . For the  $(\rightarrow)$  direction, suppose  $n \in \text{Reach}(A \cup B)$ . So there is a path from some  $m \in A \cup B$  to n. We have two cases:  $m \in A$ , in which case we have a path from  $m \in A$  to n; or  $m \in B$ , in which case we have a path from  $m \in B$  to n.

As for the  $(\leftarrow)$  direction, suppose  $n \in \text{Reach}(A) \cup \text{Reach}(B)$ . Similarly, we have two cases, and in either case we have a path from  $n \in A \cup B$ . So  $n \in \text{Reach}(A \cup B)$ .

Reach<sup> $\downarrow$ </sup> is as well, and the proof for this is the same, mutatis mutandis the direction of the path.

**Proposition 3.4.** For all  $S, A, B \in \text{State}_{\mathcal{N}}$ ,

**Inclusion.**  $S \subseteq \operatorname{Reach}^{\downarrow}(S)$ 

**Idempotent.** Reach<sup> $\downarrow$ </sup>(Reach<sup> $\downarrow$ </sup>(S)) = Reach<sup> $\downarrow$ </sup>(S)

**Monotonic.** If  $A \subseteq B$  then  $\operatorname{Reach}^{\downarrow}(A) \subseteq \operatorname{Reach}^{\downarrow}(B)$ 

**Closed under**  $\cup$ . Reach<sup> $\downarrow$ </sup> $(A \cup B) = \text{Reach}<sup><math>\downarrow$ </sup> $(A) \cup \text{Reach}<sup><math>\downarrow$ </sup>(B)

[What interaction property do Clos, Reach and Reach<sup> $\downarrow$ </sup> share?]

# **4** Neural Network Semantics for Conditional Logic

[Rewrite this section so that I'm basically mentioning Hannes' semantics and results, and not proving anything new myself! Give an illustrative example, which I will be borrowing later for the logic of  $\mathbb{C}$ ]

[Introduce this all slowly. I will now explain how we can use neural networks as models to interpret formulas in logic. First, I will give Hannes' semantics for conditional logic. Then I will introduce my own semantics for modal logic based on his.]

**Definition 4.1.** Formulas in our conditional language  $\mathcal{L}_{\Rightarrow}$  are given by [todo]

**Definition 4.2.** The semantics for  $\mathcal{L}_{\Rightarrow}$  is given as follows. [todo]

**Definition 4.3.** We write  $\models \alpha \Rightarrow \beta$  to mean all nets  $\mathcal{N} \models \alpha \Rightarrow \beta$ , and  $\Gamma \models \alpha \Rightarrow \beta$  to mean every model  $\mathcal{N}$  of  $\Gamma$ , i.e.  $\mathcal{N} \models \gamma \Rightarrow \delta$  for all  $\gamma \Rightarrow \delta \in \Gamma$  also models  $\alpha \Rightarrow \beta$ .

[Prove a few key properties for forward propagation, we can read the axioms directly off of these, then the proofs for the axioms' soundness follows]

**Definition 4.4.** The proof system for the conditional logic over  $\mathcal{L}_{\Rightarrow}$  is given as follows:  $\vdash \varphi$  iff [todo]

**Definition 4.5.** [Definition of  $\Gamma \vdash \alpha \Rightarrow \beta$ ]

**Theorem 4.6.** (Hannes Leitgeb, [38; 39]) This proof system is sound; for all  $\Gamma \subseteq \mathcal{L}_{\Rightarrow}$  and  $\varphi \in \mathcal{L}_{\Rightarrow}$ , if  $\Gamma \vdash \varphi$  then  $\Gamma \models \varphi$ . [Rephrase for conditionals  $\alpha \Rightarrow \beta$ !]

Proof. [todo]

**Lemma 4.7.** Let  $\mathcal{M} = \langle W, \prec, V \rangle$  be a plausibility model, and  $\mathcal{N}$  be given by the NAND construction above. For all conditional *terms*  $\alpha \in [I \text{ need a symbol for this...}], [[\alpha]]_{\mathcal{N}} = [[\alpha]]_{\mathcal{M}}^{\mathbb{C}}$  (the *complement* of  $[[\alpha]]_{\mathcal{M}}!$ )

**Proof.** We proceed by induction on  $\alpha$ .

**Lemma 4.8.** Let  $\mathcal{M} = \langle W, \prec, V \rangle$  be a plausibility model, and  $\mathcal{N}$  be given by the NAND construction above. For all conditional formulas  $\alpha \Rightarrow \beta \in \mathcal{L}^{\Rightarrow}$ , where  $\alpha, \beta \in [todo]$ ,

$$\mathcal{N} \models \alpha \Rightarrow \beta \text{ iff } \mathcal{M} \models \alpha \Rightarrow \beta$$

**Proof.** Combining the previous two lemmas, we have

$$\begin{split} \mathcal{N} &\models \alpha \Rightarrow \beta \quad \text{iff} \quad [\![\beta]\!]_{\mathcal{N}} \subseteq \text{Clos}([\![\alpha]\!]_{\mathcal{N}}) & \text{(by definition)} \\ & \text{iff} \quad [\![\beta]\!]_{\mathcal{M}}^{\mathbb{C}} \subseteq \text{Clos}([\![\alpha]\!]_{\mathcal{M}}^{\mathbb{C}}) & \text{(by Lemma 5.17)} \\ & \text{iff} \quad [\![\beta]\!]_{\mathcal{M}}^{\mathbb{C}} \subseteq \text{best}_{<}([\![\alpha]\!]_{\mathcal{M}})^{\mathbb{C}} & \text{(by Lemma 5.17)} \\ & \text{iff} \quad \text{best}_{<}([\![\alpha]\!]_{\mathcal{M}}) \subseteq [\![\beta]\!]_{\mathcal{M}} & \text{(flipping \subseteq and complementing both sides)} \\ & \text{iff} \quad \mathcal{M} \models \alpha \Rightarrow \beta & \text{(by definition)} \\ \end{split}$$

**Theorem 4.9.** (Model Building for  $\mathcal{L}^{\Rightarrow}$ ) For all consistent  $\Gamma \subseteq \mathcal{L}^{\Rightarrow}$ , there is finite  $\mathcal{N}$  such that  $\mathcal{N} \models \Gamma$ .

# Proof. []

**Corollary 4.10.** (Completeness for  $\mathcal{L}^{\Rightarrow}$ ) For all consistent  $\Gamma \subseteq \mathcal{L}^{\Rightarrow}$  and all conditionals  $\alpha \Rightarrow \beta \in \mathcal{L}^{\Rightarrow}$ ,

if 
$$\Gamma \models \alpha \Rightarrow \beta$$
 then  $\Gamma \vdash \alpha \Rightarrow \beta$ 

# Proof. []

- Need to also have modal logic semantics for plausibility models. The big thing here is that I have to prove we can still build a finite plausibility model in this setting!!!—completeness on the plausibility model end is going to be the hard part, and will involve temporal logic tricks.
- We will get a lot more mileage out of Lemma [todo]!

### 5 Neural Network Semantics for the Modal Logic of C

I can now state the specific logic and neural network semantics that I will consider. Let p,q,... be finitely many atomic propositions. These represent fixed states, corresponding to features in the external world that we know ahead of time. Usually these are input and output states, although they could be intermediate "hidden" states if we know these features ahead of time. For example, p might be the set of neurons that represent the color pink. For more complex formulas,

**Definition 5.1.** Formulas in the base modal language  $\mathcal{L}_{C}$  are given by

$$\varphi, \psi \coloneqq p \mid \neg \varphi \mid \varphi \land \psi \mid \mathbf{A}\varphi \mid \Box \varphi \mid \mathbf{C}\varphi$$

 $\top, \bot, \lor, \rightarrow, \leftrightarrow$  and the dual modal operators  $\mathbf{E}, \diamond, \langle \mathbf{C} \rangle$  are defined in the usual way.

The operator **A** is just the modal universal quantifier;  $\mathbf{A}\varphi$  reads " $\varphi$  holds everywhere." The  $\Box$  operator is the standard relational modal one, with intended reading [todo]. Recall that we read conditionals  $\varphi \Rightarrow \psi$  classically as "the best  $\varphi$ -states are  $\psi$ -states." [UPDATE READING The modal operator **C** is intended to be the "best" modality that occurs in this reading. In other words,  $\mathbf{C}\varphi$  is meant to read "the current state (neuron) *n* is a best (most normal, plausible, typical) one."]

[Move to plausibility model discussion: Note that the semantics for  $C\varphi$  here are not the most general possible; [cite Johan, Sonja + Alexandru] consider best<sub><,*i*,*w*</sub> relations that depend on the agent *i*, as well as the current state *w* where  $C\varphi$  is being evaluated. [Define  $\models_{Plaus}$ !!!]]

It is not immediately clear how these readings are justified; in my dissertation, I will justify these readings by connecting the neural network semantics I give here to more traditional semantics for  $\mathbf{A}, \Box, \Box^{\downarrow}$ , and  $\mathbf{C}$ .

At last, here are the neural network semantics for  $\mathcal{L}_{C}$ . I will define the semantics for the dual operators  $\mathbf{E}, \diamond, \langle \mathbf{C} \rangle$  instead of  $\mathbf{A}, \Box, \mathbf{C}$ ; either choice is equivalent, but this choice is somewhat more intuitive, since  $\diamond$  and  $\langle \mathbf{C} \rangle$  directly map to the neural network operators Reach and Clos.

**Definition 5.2.** For all  $\mathcal{N} \in \mathbf{Net}$ ,  $n \in N$ :

$$\begin{array}{lll} \mathcal{N},n \Vdash p & \text{iff} & n \in V(p) \\ \mathcal{N},n \Vdash \neg \varphi & \text{iff} & \mathcal{N},n \not\models \varphi \\ \mathcal{N},n \Vdash \varphi \land \psi & \text{iff} & \mathcal{N},n \Vdash \varphi \text{ and } \mathcal{N},n \Vdash \psi \\ \mathcal{N},n \Vdash \mathbf{E}\varphi & \text{iff} & \llbracket \varphi \rrbracket \neq \emptyset \\ \mathcal{N},n \Vdash \Diamond \varphi & \text{iff} & n \in \text{Reach}(\llbracket \varphi \rrbracket \cup \{\text{bias}\}) \\ \mathcal{N},n \Vdash \langle \mathbf{C} \rangle \varphi & \text{iff} & n \in \text{Clos}(\llbracket \varphi \rrbracket \cup \{\text{bias}\}) \end{array}$$

where  $\llbracket \varphi \rrbracket = \{n \in N \mid \mathcal{N}, n \Vdash \varphi\}.$ 

**Proposition 5.3.** By the definition of the duals  $\mathbf{E}, \diamond, \langle \mathbf{C} \rangle$ , we can instead define the semantics for  $\mathbf{A}, \Box, \mathbf{C}$  as follows. For all  $\mathcal{N} \in \mathbf{Net}$ ,  $n \in N$ :

$$\mathcal{N}, n \Vdash \mathbf{A}\varphi \quad \text{iff} \quad \llbracket \varphi \rrbracket = N$$
$$\mathcal{N}, n \Vdash \Box \varphi \quad \text{iff} \quad n \in (\text{Reach}(\llbracket \varphi \rrbracket^{\mathbb{C}} \cup \{\text{bias}\}))^{\mathbb{C}}$$
$$\mathcal{N}, n \Vdash \mathbf{C}\varphi \quad \text{iff} \quad n \in (\text{Clos}(\llbracket \varphi \rrbracket^{\mathbb{C}} \cup \{\text{bias}\}))^{\mathbb{C}}$$

**Definition 5.4.** We write  $\mathcal{N} \models_{\text{Net}} \varphi$  ("the net *models*  $\varphi$ ") to mean  $\mathcal{N}, n \models \varphi$  for all  $n \in N$ ;  $\models_{\text{Net}} \varphi$  to mean all nets  $\mathcal{N} \models \varphi$ ; and finally,  $\Gamma \models_{\text{Net}} \varphi$  to mean every model  $\mathcal{N}$  of  $\Gamma$ , i.e.  $\mathcal{N} \models_{\text{Net}} \psi$  for all  $\psi \in \Gamma$  also models  $\varphi$ .

The following fact says that these neural semantics for **C** match up with Leitgeb's neural semantics for conditionals  $\varphi \Rightarrow \psi$ . Formally, conditionals  $\varphi \Rightarrow \psi$  are expressible in  $\mathcal{L}_{\mathbf{C}}$  by  $\mathbf{A}(\mathbf{C}\varphi \rightarrow \psi)$ . This justifies our reading of  $\mathbf{C}\varphi$  as "the current state (neuron) *n* is a best (most normal, plausible, typical) one."

**Proposition 5.5.** Let  $\varphi \Rightarrow \psi \in \mathcal{L}_{\Rightarrow}$  (that is,  $\varphi, \psi \in \mathcal{L}_{prop}$ ). Then for all  $\mathcal{N} \in \mathbf{Net}$ ,

$$\mathcal{N} \models_{\mathbf{Net}} \varphi \Rightarrow \psi \text{ iff } \mathcal{N} \models_{\mathbf{Net}} \mathbf{A}(\mathbf{C}\varphi \rightarrow \psi)$$

Proof. [Todo]

**Example 5.6.** [Give an example of a neural network model from before, but this time with our semantics in action. Give a number of constraints involving bird, penguin, and fly, and evaluate the truth of these at particular states]

[Discussion of semantics]

[Note that the syntax is backwards/dualed from the syntax for neural network semantics. But because they're duals of each other, it doesn't matter. (It's a known modal logic trick that you can do induction on either version/form.)]

[Talk about the possible worlds interpretation of these, & the local evaluation at particular neurons!]

[talk about  $\mathcal{N} \models \varphi$ , which in this context means " $\varphi$  activates *all* of the neurons in  $\mathcal{N}$ ", or " $\varphi$  holds across the entire net." Also define  $\Gamma \models \varphi$ .]

[A reader might be confused about the "swaps" from  $\mathsf{Reach}^{\downarrow}$  to  $\diamond$ , as well as the choice to use diamond operators instead of the normal box ones.] [How can we justify the readings we had above? (at least intuitively)]

[I now have space to say these points in more detail (there's lots to say!)] Although these semantics are based on Leitgeb's [40], they differ in a few key ways. First, his semantics uses conditionals  $\varphi \Rightarrow \psi$  to capture neural network inference, whereas mine instead centers on the modal operator  $\langle \mathbf{C} \rangle$ . Second, I include these additional operators  $\Box$  and  $\Box^{\downarrow}$  that are not mentioned in his work. Finally, Leitgeb battles with the issue of how to correctly interpret negation; I sidestep this issue by using neural networks for interpreting  $\langle \mathbf{C} \rangle \varphi$  (where the "action" happens), but not for  $\neg$  and  $\land$ . The bottom line is this: proving completeness for this logic is not necessarily just a matter of importing the proof from [40].
As with possible-worlds semantics for modal logic, we evaluate the truth of formulas *locally* at particular neurons.

I have tried to keep these semantics as close as possible to the ordinary presentation of semantics for modal logic. Basically,  $\neg$ ,  $\land$ ,  $\lor$ ,  $\rightarrow$ , **A**, **E** are all treated classically. But we defer to the neural network when faced with evaluating  $\langle \mathbf{C} \rangle$ ,  $\diamond$ , and  $\diamond^{\downarrow}$ . Since  $\llbracket \varphi \rrbracket$  can be any arbitrary set of neurons, it may not include the bias node. This means  $\llbracket \varphi \rrbracket$  is not necessarily be a valid state of the neural network, since the bias might not be in it. So in order to apply the neural network functions Reach<sup> $\downarrow$ </sup>, Reach, and Clos, we provide the bias on input by giving  $\llbracket \varphi \rrbracket \cup \{\text{bias}\} \in \text{State}_{\mathcal{N}}$ . **[TODO: I will have to do this for the conditional logic semantics too!]** 

### 5.1 Why Consider this *Modal* Logic?

[The conditional logic came first, and I could have just re-used it in order to save myself a lot of work. But I think it's important to explain my aesthetic choices in moving to the modal logic, rather than just using Hannes' conditional logic. But this is such a technical detail that I should at least advise the reader to skip this section if they're not interested.]

## 5.2 What Makes these Semantics "Neural"

What makes these semantics "neural" or "connectionistic"? Easy answer: the key operators for inference are implemented in a neural network. Better answer: no single neuron/node holds the information for  $\langle \mathbf{C} \rangle \varphi$  How is this different from relational or neighborhood semantics? Is this a meaningful difference? What does Hannes have to say about it in his dissertation?

The maybe it just means "loosely inspired by real neuron and synapse behavior", but even then there are probably properties we can write down and check.

#### 6 Proof System and Soundness

For C alone, [38] proves that the properties in Proposition [which?] are complete for Clos over

Axioms for D: Axioms for A: (**Dual**)  $\Diamond \varphi \leftrightarrow \neg \Box \neg \varphi$ (Dual)  $\mathbf{E}\varphi \leftrightarrow \neg \mathbf{A} \neg \varphi$ **(Distr)**  $\Box(\varphi \land \psi) \leftrightarrow (\Box \varphi \land \Box \psi)$ (Distr)  $\mathbf{A}(\varphi \rightarrow \psi) \rightarrow (\mathbf{A}\varphi \rightarrow \mathbf{A}\psi)$ (Refl)  $\Box \varphi \to \varphi$ (Refl)  $A\phi \rightarrow \phi$ (**Trans**)  $\Box \varphi \rightarrow \Box \Box \varphi$ (5)  $\mathbf{E}\boldsymbol{\varphi} \rightarrow \mathbf{A}(\mathbf{E}\boldsymbol{\varphi})$ (Interact)  $A\varphi \rightarrow \Box \varphi$ **Axioms for C: Rules of Inference:**  $\langle \mathbf{C} \rangle \varphi \leftrightarrow \neg \mathbf{C} \neg \varphi$ (Dual) **(MP)** From  $\vdash \varphi \rightarrow \psi$  and  $\vdash \varphi$ (Refl)  $\mathbf{C}\varphi \rightarrow \varphi$ we can infer  $\vdash \psi$  $\mathbf{C}\varphi \rightarrow \mathbf{C}\mathbf{C}\varphi$ (A-Nec) From  $\vdash \varphi$ , we can infer  $\vdash A\varphi$ (Trans)  $A(C\varphi \rightarrow \psi) \rightarrow$ (CM) ( $\Box$ -**Rep**) From  $\vdash \varphi \leftrightarrow \psi$ , infer  $(\mathbf{C}(\varphi \land \psi) \rightarrow \mathbf{C}\varphi)$  $\vdash \Box \varphi \leftrightarrow \Box \psi$ (Interact)  $\Box \varphi \rightarrow C \varphi$ (**C-Rep**) From  $\vdash \varphi \leftrightarrow \psi$ , infer  $\vdash \mathbf{C} \varphi \leftrightarrow \mathbf{C} \psi$ 

Figure 6.1. Axioms and rules of inference for [todo]

binary, feed-forward nets. We transcribe these into our modal language.

**Definition 6.1.** The proof system  $\vdash$  for neural network semantics  $\models_{\text{Net}}$  over  $\mathcal{L}_{C}$  is given as follows:  $\vdash \varphi$  iff either  $\varphi$  is valid in propositional logic, or  $\varphi$  is one of the axioms listed in Figure 6.1, or  $\varphi$  follows from some previously obtained formulas by one of the inference rules.

The axioms for  $\Box$  and **A** form "*MLU*" (modal logic with the universal quantifier), and this is *almost* the standard complete axiomatization of this logic [Cite! Johan van Benthem mentions it in Modal Logic for Open Minds];  $\Box$  does not satisfy a (Nec) rule, so instead we have the rule of replacement (**Rep**) from classical modal logic. Aside from this difference, you can think of  $\Box$  as the modal logic *S4*. **A** is just the modal logic *S5* with an additional interaction axiom  $\mathbf{A}\varphi \rightarrow \Box \varphi$  stating that if  $\varphi$  holds everywhere, then it holds everywhere above the current world.

The axioms and rules for **C** come from the rules for cumulative conditionals  $\varphi \Rightarrow \psi$ . Let me explain how I derived them. First, from [cite Giordano's paper], if in conditional logic we allow for *propositional* nesting of formulas (we still cannot nest  $\varphi \Rightarrow \psi$ ), we can express the axioms for cumulative logic as follows. Let  $\vdash_{\text{Prop}}\varphi$  mean that  $\varphi$  is provable in the underlying propositional

system.

(Reflexivity).  $\varphi \Rightarrow \varphi$ (Left Equivalence). If  $\vdash_{\text{Prop}} \varphi \leftrightarrow \psi$  then  $(\varphi \Rightarrow \rho) \rightarrow (\psi \rightarrow \rho)$ (Right Weakening). If  $\vdash_{\text{Prop}} \varphi \rightarrow \psi$  then  $(\rho \Rightarrow \varphi) \rightarrow (\rho \rightarrow \psi)$ (Cautious Monotonicity).  $(\varphi \Rightarrow \psi) \land (\varphi \Rightarrow \rho) \rightarrow (\varphi \land \psi \rightarrow \rho)$ (Cautious Cut).  $(\varphi \land \psi \Rightarrow \rho) \land (\varphi \Rightarrow \psi) \rightarrow (\varphi \Rightarrow \rho)$ 

We can straightforwardly translate these into the language of  $\mathcal{L}_{C}$  by replacing each  $\varphi \Rightarrow \psi$  with  $A(C\varphi \rightarrow \psi)$ . From here, each of the above cumulative logic axioms follows from modal ones: (**Reflexivity**) follows from (**Refl**); (**Left Equivalence**) follows from propositional axioms plus (**Rep**); (**Right Weakening**) follows from propositional axioms plus modus ponens (**MP**); (**Cautious Monotonicity**) follows from the modal axiom (**CM**). [TODO: What about (**Cautious Cut**) now? Is it sound for **Net**? Does it need to be replaced by something else?]

Additionally, I have added the axioms (**Trans**) and (**Interact**) to reflect the fact that Clos is idempotent, and for all  $S \in \text{State}_{\mathcal{N}}$ ,  $\text{Clos}(S) \subseteq \text{Reach}(S)$ .

As a modal operator, **C** is non-normal and classical. Instead of a (**Distr**) axiom, **C** satisfies the weaker (**CM**). The **C** operator also satisfies the replacement rule (**Rep**) from classical modal logic in lieu of (**Nec**) [cite (**Rep**) rule].

**Definition 6.2.** If  $\Gamma \subseteq \mathcal{L}_{\mathbb{C}}$  is a set of formulas and  $\varphi \in \mathcal{L}_{\mathbb{C}}$  a formula, then  $\Gamma \vdash \varphi$  whenever there are finitely many  $\psi_1, \ldots, \psi_k \in \Gamma$  such that  $\vdash \psi_1 \land \ldots \land \psi_k \rightarrow \varphi$ .

**Theorem 6.3.** ( Soundness for  $\mathcal{L}_{C}$  over  $\models_{Net}$ ) These rules and axioms are sound for neural network models; for all consistent  $\Gamma \subseteq \mathcal{L}_{C}$  and  $\varphi \in \mathcal{L}_{C}$ , if  $\Gamma \vdash \varphi$  then  $\Gamma \models_{Net} \varphi$ .

**Proof.** Suppose  $\Gamma \vdash \varphi$ . That is, there are finitely many  $\psi_1, \ldots, \psi_k \in \Gamma$  such that  $\vdash \psi_1 \land \ldots \land \psi_k \rightarrow \varphi$ , which in turn means (by (**MP**)) that if  $\vdash \psi_1, \ldots, \vdash \psi_k$ , then  $\vdash \varphi$ . Now let  $\mathcal{N} \models \Gamma$ . In particular, this means  $\mathcal{N} \models \psi_1, \ldots, \psi_k$ . I now need to show that  $\mathcal{N} \models \varphi$ . Since  $\vdash \varphi, \varphi$  is itself an axiom or follows from previously obtained formulas by the inference rules. In order to prove  $\mathcal{N} \models \varphi$ , it's enough to show that the axioms and rules of inference are valid (hold for all  $\mathcal{N} \in \mathbf{Net}$  at all neurons  $n \in N$  whatsoever).

To check:	Use:	To check:	Use:
( <b>Distr</b> ) for $\Box$	Monotonicity of Reach	(Refl) for C	Inclusion of Clos
( <b>Refl</b> ) for $\Box$	Inclusion of Reach	(Trans) for C	Idempotence of Clos
(Trans) for $\Box$	Idempotence of Reach	(CM) for C	Cumulativity of Clos
(Interact), A+□	$Reach(S) \subseteq W$	(Interact), $\Box + C$	Reach contains Clos
( <b>Rep</b> ) for $\square$	Reach is a function	(Rep) for C	Clos is a function

Figure 6.2. [caption!]

The propositional axioms and (MP) are totally classical, and are known to be sound. The (Dual) axioms for A,  $\Box$ , and C are also sound (by definition of their duals). The semantics for A do not make use of the neural network at all, so the axioms (Distr), (Refl), (5) and the (Nec) rule for A are classically sound as well. Let  $\mathcal{N} \in \text{Net}, n \in N$ . I will now check the remaining rules and axioms in detail. This is a bit tedious, so see Figure 6.2 for a summary.

#### (Distr) for □.

- (**Refl**) for  $\Box$ . Suppose contrapositively  $\mathcal{N}, n \not\models \varphi$ . So  $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}^{c}$ . Note that in particular  $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}^{c} \cup \{\text{bias}\}$  (we can simply add the bias node). By Inclusion of Reach,  $n \in \text{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{c} \cup \{\text{bias}\})$ . By the semantics,  $\mathcal{N}, n \models \diamond \neg \varphi$ . By (**Dual**) for  $\Box, \mathcal{N}, n \models \neg \Box \varphi$ , i.e.,  $\mathcal{N}, n \not\models \Box \varphi$ , and we are done.
- (**Trans for**  $\Box$ ). Suppose contrapositively  $\mathcal{N}, n \not\models \Box \Box \varphi$ . This gives us  $\mathcal{N}, n \not\models \neg \Box \Box \varphi$ . Applying (**Dual**) twice, we have  $\mathcal{N}, n \not\models \Diamond \Diamond \neg \varphi$ . By the semantics,  $n \in \text{Reach}(\text{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\}) \cup \{\text{bias}\})$ . Note that bias is always in the Reach of a set already containing bias, so we can simplify this to  $n \in \text{Reach}(\text{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\}))$ . By Idempotence of Reach,  $n \in \text{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\})$ . But this means  $\mathcal{N}, n \not\models \Diamond \neg \varphi$ , and by (**Dual**) we conclude that  $\mathcal{N}, n \not\models \Box \varphi$ .
- (Interact) for A and  $\Box$ . Suppose contrapositively  $\mathcal{N}, n \not\models \Box \varphi$ . By (Dual) for  $\Box, \mathcal{N}, n \models \diamond \neg \varphi$ . By the semantics,  $n \in \text{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\}) \subseteq N$ . [Todo]
- (**Rep**) for  $\Box$ . Suppose  $\models_{Net} \varphi \leftrightarrow \psi$ , i.e.,  $\varphi \leftrightarrow \psi$  holds for all nets, for all neurons. So for all  $\mathcal{N} \in \mathbf{Net}$ ,  $\|\varphi\|_{\mathcal{N}} = \|\psi\|_{\mathcal{N}}$ . Since Reach is a function,  $\operatorname{Reach}(\|\varphi\|_{\mathcal{N}} \cup \{\operatorname{bias}\}) = \operatorname{Reach}(\|\psi\|_{\mathcal{N}} \cup \{\operatorname{bias}\}) = \operatorname{Reach}(\|\psi\|_{\mathcal{N}} \cup \{\operatorname{bias}\})$  for all  $\mathcal{N}$ . So in particular,  $\operatorname{Reach}(\|\varphi\|_{\mathcal{N}}^{\mathbb{C}} \cup \{\operatorname{bias}\}) = (\operatorname{Reach}(\|\psi\|_{\mathcal{N}}^{\mathbb{C}}) \cup \{\operatorname{bias}\})^{\mathbb{C}}$  for all  $\mathcal{N}$ , which means  $\models_{\operatorname{Net}} \Box \varphi \leftrightarrow \Box \psi$ .

- (**Refl**) for **C**. Suppose contrapositively  $\mathcal{N}, n \not\models \varphi$ . So  $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}}$ . Note that in particular  $n \in \llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\}$  (we can simply add the bias node). By Inclusion of Clos,  $n \in \text{Clos}(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{\text{bias}\})$ . By the semantics,  $\mathcal{N}, n \models \langle \mathbf{C} \rangle \neg \varphi$ . By (**Dual**) for  $\mathbf{C}, \mathcal{N}, n \models \neg \mathbf{C} \varphi$ , i.e.,  $\mathcal{N}, n \not\models \mathbf{C} \varphi$ , and we are done.
- (Trans) for C. Suppose contrapositively  $\mathcal{N}, n \not\models CC\varphi$ . This gives us  $\mathcal{N}, n \not\models \neg CC\varphi$ . Applying (Dual) twice, we have  $\mathcal{N}, n \not\models \langle C \rangle \langle C \rangle \neg \varphi$ . By the semantics,  $n \in Clos(Clos(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{bias\}) \cup \{bias\})$ . Note that bias is always in the closure of a set, so we can simplify this to  $n \in Clos(Clos(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{bias\}))$ . By Idempotence of  $Clos, n \in Clos(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{bias\})$ . But this means  $\mathcal{N}, n \not\models \langle C \rangle \neg \varphi$ , and by (Dual) we conclude that  $\mathcal{N}, n \not\models C\varphi$ .
- (CM) for C. [Todo]

### (Interact) for $\Box$ and C.

(**Rep**) for C. Suppose  $\models_{Net} \varphi \leftrightarrow \psi$ , i.e.,  $\varphi \leftrightarrow \psi$  holds for all nets, for all neurons. So for all  $\mathcal{N} \in \mathbf{Net}$ ,  $\llbracket \varphi \rrbracket_{\mathcal{N}} = \llbracket \psi \rrbracket_{\mathcal{N}}$ . Since Clos is a function,  $Clos(\llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{bias\}) = Clos(\llbracket \psi \rrbracket_{\mathcal{N}} \cup \{bias\})$ for all  $\mathcal{N}$ . So in particular,  $Clos(\llbracket \varphi \rrbracket_{\mathcal{N}}^{\mathbb{C}} \cup \{bias\}) = (Clos(\llbracket \psi \rrbracket_{\mathcal{N}}^{\mathbb{C}}) \cup \{bias\})^{\mathbb{C}}$  for all  $\mathcal{N}$ , which means  $\models_{Net} C\varphi \leftrightarrow C\psi$ .

## 7 Reflections on Neural Network Verification

[Here's where I can discuss things like property verification vs verifying a particular neural network vs constraints + model building (alignment), extraction, "valuation search", and Thomas Icards' method]

[say something about neuro-symbolic systems more generally in this framework] The basis for many neuro-symbolic systems is that they encode logical information into neural networks, or conversely, encode neural networks as models in logic (see [53]). Given these translations, certain neural networks and logic models are able to represent the same information.

#### Interlude: A History and Survey of Neural Network Semantics

[Rewrite this to flow better, and to relate neuro-symbolic systems to the work in this section. I'm trying to get across a broader perspective, not simply a single new system!] My thesis work builds on existing logics that use neural network semantics, and shares similarities with logics for modeling social networks. Additionally, my approach to modeling learning takes inspiration from work on learning in Dynamic Epistemic Logic (DEL). Here I'll take a moment to situate my thesis in this broader context and clarify what my contribution is in each case.

[Make it clear that this point of view is not my own original idea; but applying it broadly and developing it as far as I do *is* novel. Give a history of this perspective, and also make an effort to cast other neuro-symbolic systems within this framework (including the FLaNN work, etc.) This is the place for a proper literature review.]

The core idea behind neural network semantics—that neural networks can be treated as models for logic—actually dates back to the very first paper on neural networks. In McCulloch and Pitts [44], logical formulas are mapped directly to individual neurons in the net. This approach suffers from the well-known "grandmother cell" problem [30]: it is cognitively implausible that an individual neuron could represent a complex concept such as "grandmother". Instead, concepts in brain networks are distributed across multiple neurons at once.

Neural network semantics is based on a recent reimagining of this approach [7; 40], where logical formulas are mapped to sets of neurons rather than to individual neurons. Early work established the correspondence between inference in a neural network and nonmonotonic conditionals [7; 17; 38; 39]. More recently, [28; 29] proved soundness for inference over fuzzy neural networks. In my thesis work so far [34; 35], I applied ideas from Dynamic Epistemic Logic to model a simple update policy, Hebbian learning, in neural network semantics. The key results of this work are the first ever soundness and completeness theorems for any learning policy on neural networks.

[Go through popular neuro-symbolic proposals, and explain their relationship with neural network semantics. What order???] Semantic Encodings. Logics with Social Network Semantics. Logic Tensor Networks. Distributed Alignment Search. DeepProbLog.

Logic Explained Networks.

Neural Networks as Automata.

Neural Network Fibring.

Logics with Social Network Semantics. [Make this a single point] It has recently come to my attention that a similar approach is being used to model group behavior in social networks. In these social network logics [4; 8; 18], nodes in the graph represent individual agents, and each formula is mapped to the set of agents that adopt a certain social attitude. Agents influence each other, and the spread of their attitudes is modeled much in the same way as forward propagation of a signal in a neural network.

This work shares essentially the same premise and techniques as neural network semantics; I personally view this as a case of parallel discovery. But the two approaches still differ in interesting ways. First, in some sense the two semantics are operating on different "levels": social networks model interactions between multiple agents, whereas neural networks model interactions between components of the same (single) agent. Second, the two differ in subject matter. Social network semantics focuses on different social links between agents, and how these links change [4]. Neural network semantics, my own work included, instead focuses on inferences and updates inspired by artificial and natural neural networks.

## **Chapter 4**

## **Neural Network Model Constructions**

## **1** Introduction

- model-building techniques
- applications: completeness + expressivity for neural networks viewed as logic models
- expressivity, including how to build neural network models from other types of models, and vice-versa

I will now show that these are all the axioms for neural network semantics, i.e., these axioms are complete. This is philosophically significant for the theory of neuro-symbolic AI—at a high level, completeness says that this logic exactly matches the closure behavior of binary neural networks. At the crux of this proof is the task of neural network model building: Given any set of constraints  $\Gamma$  consistent in this logic, I will need to construct a neural network model  $\mathcal{N} \in \mathbf{Net}$  satisfying  $\Gamma$ . This is just as interesting, but is more practical: it means we can custom-design neural networks that model exactly the constraints over  $\mathcal{L}_{\mathbf{C}}$  we would like them to have.

In previous work by Leitgeb [38; 39], completeness for the cumulative conditional logic over  $\mathcal{L}_{\Rightarrow}$  was shown by way of plausibility models **Plaus**. Leitgeb showed that any plausibility model  $\mathcal{M} \in$  **Plaus** can be transformed into an equivalent neural network  $\mathcal{N}' \in$  **Net**: For all  $\varphi \Rightarrow \psi \in \mathcal{L}_{\Rightarrow}$ ,

$$\mathcal{N}' \models \varphi$$
 iff  $\mathcal{M} \models \varphi$ 

It would be nice if this technique lifted to the more expressive language  $\mathcal{L}_{\mathbf{C}}$ . But we have no such luck. In Section [todo], I will show that the language  $\mathcal{L}_{\mathbf{C}}$  can express formulas that hold for all  $\mathcal{M} \in \mathbf{Plaus}$ , yet do not hold in every  $\mathcal{N} \in \mathbf{Net}$ . In a precise sense, plausibility models are *not* the correct analogy for explaining the closure behavior of neural networks.

Instead, I will directly build a neural network model using a modified version of the standard canonical model construction in modal logic. At first, I will build a net with infinitely many neu-

rons. Then, in section [todo] I will use the standard modal logic technique of filtration to transform this network into a finite one.

The proof I give here is non-constructive, since it first requires building an infinite neural network. Unfortunately, I do not know of a constructive algorithm for building a neural network model from  $\Gamma$ . For now, I will leave it as an open question.

**Open Question 2.** For a given finite set of constraints  $\Gamma \subseteq \mathcal{L}_{\mathbf{C}}$ , is there a constructive algorithm that produces a model  $\mathcal{M} \in \mathbf{Plaus}$  and state  $w \in W$  such that  $\mathcal{M}, w \models \Gamma$ ?

#### 2 The Canonical Neural Network Construction

**Lemma 2.1.** (Lindenbaum's Lemma [cite!]) We can extend any consistent set  $\Gamma$  to a maximally consistent set  $\Delta \supseteq \Gamma$ .

First, lets start by constructing a relational model inspired by the standard Kripke-style canonical model. I will use this model as a base for the canonical neural network model.

**Definition 2.2.** Let the canonical neural network model be  $\mathcal{N}^c = \langle N^c, \text{bias}^c, E^c, W^c, A^c, \eta^c, V^c \rangle$ , where

- N<sup>c</sup> = {Δ | Δ is maximally consistent over L<sub>C</sub>}. There are countably many; let's fix an order on these nodes Δ<sub>0</sub>, Δ<sub>1</sub>,...
- $\Delta' E^c \Delta$  iff for all  $\varphi \in \mathcal{L}_{\mathbf{C}}$ , if  $\Box \varphi \in \Delta$  then  $\varphi \in \Delta'$ .
- bias<sup>c</sup> is [todo]
- Suppose  $\Delta_i \in N$  has predecessors  $\Delta'_{i1}, \Delta'_{i2}, \dots, \Delta'_{ik}$ . For each  $\Delta'_{ij}E^c\Delta_i$ , let

$$W^c(\Delta'_{ii}, \Delta_i) = (p_i)^j$$

where  $p_i$  is the *i*<sup>th</sup> prime number. The idea is that each prime  $p_i$  uniquely codes for the node  $\Delta_i$ , and the weight between  $\Delta_i$  and its predecessor  $\Delta'_{ij}$  is a power of  $p_i$  that uniquely codes for  $\Delta'_{ij}$ . Consequently, any activation value *x* we care about is going to be a sum of powers of  $p_i$ , from which we can reconstruct precisely the  $\Delta_i$  being activated and the predecessors  $\Delta'_{ij}$  that were involved in activating it.

Let x∈Q, and suppose that x = ∑<sub>j∈X</sub> (p<sub>i</sub>)<sup>j</sup> for some i∈ {1,..., |N<sup>c</sup>|} and some subset X ⊆ {1,..., k}, where k is the number of predecessors of node Δ<sub>i</sub>. In other words, x uniquely codes for a subset {Δ'<sub>ij</sub> | Δ'<sub>ij</sub>E<sup>c</sup>Δ<sub>i</sub> and j∈X} of predecessors of Δ<sub>i</sub>, as explained just above. Let the activation function A<sup>c</sup>(x) be defined as follows:

$$A^{c}(x) = 1 \text{ iff for all } \psi, \langle \mathbb{C} \rangle \psi \in \bigcap_{\Delta'_{ij} E^{c} \Delta_{i} \text{ and } j \in X} \Delta'_{ij} \quad \text{implies } \langle \mathbb{C} \rangle \psi \in \Delta_{i}$$

(If *x* does not code for any valid subset *X*, then simply set  $A^{c}(x) = 0$ .)

[This is my current best guess for how to define  $A^{c}$ ! It's tricky!]

•  $\Delta \in V^c(p)$  iff  $p \in \Delta$ 

This construction is still a bit opaque. First, the choices for  $N^c$  and  $V^c$  are totally standard fare when building a canonical model. The basic idea is that the nodes are maximally consistent sets  $\Delta$ , and we can control what must be true at the current world by moving to the right  $\Delta$ . The choice for  $E^c$  is the reverse of the canonical Kripke relation  $R^c$ , since neural networks look "downwards" rather than "upwards" for a set that caused the activation of n.

All of the action is in the choice for the activation function  $A^c$ . Let  $\Delta_i$  have predecessors  $\Delta'_{i1}$ ,  $\Delta'_{i2}, \ldots, \Delta'_{ik}$  as before, and let X denote the subset of them that are currently active. Intuitively, we want  $A^c(x) = 1$  precisely whenever  $\Delta_i$  "says" that the subset X "should" activate  $\Delta_i$ . But since we don't yet have access to the semantics—we are currently *defining* the semantics!—we have to say this syntactically. [This is difficult to say syntactically, but what's written here is my current best guess!]

**Proposition 2.3.** (CP) The canonical neural network model  $\mathcal{N}^c$  is a well-defined (infinite) neural network model.

**Proof.** [TODO] I need to show is that, for all  $S \in \text{State}_{\mathcal{N}^c}$ ,  $\text{Clos}_{\mathcal{N}^c}(S)$  is in fact the unique fixed point under *S*. [how the hell do I prove that...?] [What constraints to I need on  $E^c$ ??]

**Lemma 2.4.** (C) Truth Lemma for  $\mathcal{L}_{\mathbf{C}}$ ) We have, for all  $\Delta \in N^c$ ,  $\varphi \in \mathcal{L}_{\mathbf{C}}$ ,

$$\mathcal{N}^c, \Delta \Vdash \varphi \text{ iff } \varphi \in \Delta$$

**Proof.** [WARNING: I do not have this result yet! It's a pain in the ass to get right.]

By induction on  $\varphi$ . The propositional and boolean cases are straightforward. The  $A\varphi$  case is totally standard, and follows from the usual lemmas about maximally consistent sets [cite]. I'll skip to the most relevant cases. It's enough to consider the  $\diamond$  and  $\langle C \rangle$  cases:

♦ Case.

**(C)** Case. Observe that the Truth Lemma in this case,  $\mathcal{N}^c$ ,  $\Delta \Vdash \langle \mathbf{C} \rangle \varphi$  iff  $\langle \mathbf{C} \rangle \varphi \in \Delta$ , is equivalent to the claim:

$$\mathsf{Clos}_{\mathcal{N}^c}(\llbracket \varphi \rrbracket) = \{ \Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta \}$$

First, by construction  $\operatorname{Clos}_{\mathcal{N}^c}(\llbracket \varphi \rrbracket)$  is the unique fixed point of the transition function  $F_{\llbracket \varphi \rrbracket}$ under  $\llbracket \varphi \rrbracket$ . I will show here that the set  $\{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\}$  is *also* such a fixed point, i.e.,  $F_{\llbracket \varphi \rrbracket}(\{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\}) = \{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\}$ . Since  $\operatorname{Clos}_{\mathcal{N}^c}(\llbracket \varphi \rrbracket)$  is the *unique* fixed point under  $\llbracket \varphi \rrbracket$ , it will follow that these two states are the same. In other words,  $\operatorname{Clos}_{\mathcal{N}^c}(\llbracket \varphi \rrbracket) = \{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\}$ .

For the  $(\rightarrow)$  direction, suppose  $\Delta \in F_{\llbracket \varphi \rrbracket}(\{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\})$ . By definition of  $F_{\llbracket \varphi \rrbracket}$ , we have two cases:

- **Case 1.**  $\Delta \in [\![\varphi]\!]$ . So  $\mathcal{N}^c, \Delta \Vdash \varphi$ , and so by the inductive hypothesis  $\varphi \in \Delta$ . By (**Refl**) and (**Dual**),  $\langle \mathbf{C} \rangle \varphi \in \Delta$ .
- **Case 2.**  $\Delta$  is activated by its predecessors. For concreteness, let's say  $\Delta = \Delta_i$  for some  $i \in \mathbb{N}$ , and let  $\Delta'_{i1}, \Delta'_{i2}, \dots, \Delta'_{ik}$  be its  $E^c$ -predecessors. Formally, we have

$$A^{c}\left(\sum_{\Delta'_{ij}E^{c}\Delta_{i}}W^{c}(\Delta'_{ij},\Delta_{i})\cdot\chi_{\{\Delta\mid \langle \mathbf{C}\rangle\varphi\in\Delta\}}(\Delta'_{ij})\right)=1$$

Let this inner sum be denoted by x (that is,  $A^c(x) = 1$ ). Since  $A^c(x) = 1$ , we must have  $x = \sum_{j \in X} (p_i)^j$  for some subset  $X \subseteq \{1, ..., k\}$ . As explained above, from this sum of prime powers we can reconstruct the subset of predecessors that are active:  $\{\Delta'_{ij} \mid \Delta'_{ij}E^c\Delta_i \text{ and } j \in X\}$ . By construction of  $A^c$ , we have:

For all 
$$\psi$$
,  $\langle \mathbf{C} \rangle \psi \in \bigcap_{\Delta'_{ij} E^c \Delta_i \text{ and } j \in X} \Delta'_{ij}$  implies  $\langle \mathbf{C} \rangle \psi \in \Delta_i$ 

Notice that in the sum above,  $\Delta'_{ij}$  is active  $(j \in X)$  exactly when  $\langle \mathbf{C} \rangle \varphi \in \Delta'_{ij}$ . It follows that  $\langle \mathbf{C} \rangle \varphi \in \bigcap_{\Delta'_{ij}E^c \Delta_i \text{ and } j \in X} \Delta'_{ij}$  holds. And so we conclude that  $\langle \mathbf{C} \rangle \varphi \in \Delta_i = \Delta$ , which is what I wanted to show.

As for the  $(\leftarrow)$  direction, suppose  $\langle \mathbf{C} \rangle \varphi \in \Delta$ . Let's say  $\Delta = \Delta_i$  for some  $i \in \mathbb{N}$ , and let  $X = \{j \mid \Delta'_{ij} E^c \Delta_i$  and  $\langle \mathbf{C} \rangle \varphi \in \Delta'_{ij}\}$  be a set that selects active predecessors of  $\Delta$ . I want to show that  $\Delta \in F_{\llbracket \varphi \rrbracket}(\{\Delta \mid \langle \mathbf{C} \rangle \varphi \in \Delta\})$ , i.e., these predecessors activate  $\Delta$ :

$$A^{c}\left(\sum_{\Delta'_{ij}E^{c}\Delta_{i}}W^{c}(\Delta'_{ij},\Delta_{i})\cdot\chi_{\{\Delta+\langle \mathbf{C}\rangle\varphi\in\Delta\}}(\Delta'_{ij})\right)=1$$

By construction of  $A^c$  it's enough to show:

For all 
$$\psi$$
,  $\langle \mathbf{C} \rangle \psi \in \bigcap_{\Delta'_{ij} E^c \Delta_i \text{ and } \langle \mathbf{C} \rangle \varphi \in X} \Delta'_{ij}$  implies  $\langle \mathbf{C} \rangle \psi \in \Delta_i$ 

Let  $\psi \in \mathcal{L}_{\mathbf{C}}$  be a formula, and suppose  $\langle \mathbf{C} \rangle \psi \in \bigcap_{\Delta'_{ij}E^c \Delta_i \text{ and } j \in X} \Delta'_{ij}$ . I will now drop the subscripts and re-phrase this to be as general as possible: For all maximally consistent sets  $\Delta'$ , if  $\Delta' E^c \Delta$  and  $\langle \mathbf{C} \rangle \varphi \in \Delta'$ , then  $\langle \mathbf{C} \rangle \psi \in \Delta'$ .

#### **3** Filtration: Building a Finite Neural Network

The fact that we can build this canonical neural network is not enough. The canonical net  $\mathcal{N}^c$  is huge—it has infinitely many nodes, each of which are infinite sets of formulas! In this section, I will transform  $\mathcal{N}^c$  into an equivalent finite neural network model  $\mathcal{N} \in \mathbf{Net}$ . [Relate this to the finite model property: every formula  $\varphi$  that has a model has a finite countermodel.] (Polynomial size would be nice, but I will leave this to future work. We probably won't be so lucky.)

The canonical neural network construction is not totally bizarre from a modal logic standpoint, so we can apply the standard technique of *filtration* to construct the finite model.

## **Definition 3.1.** [Define filtration!]

[For computer science reader: this is similar to the Myhill-Nerode theorem about finding a

minimal DFA equivalent to a given DFA,]

**Theorem 3.2.** [A filtration of  $\mathcal{M}$  preserves the formulas true in  $\mathcal{M}$ , at every world.]

**Proposition 3.3.** [A filtration of a (possibly infinite)  $\mathcal{N}$  is still a well-defined neural network model.]

**Theorem 3.4.** [If  $\Gamma$  is finite, then any filtration of a (possibly infinite) neural network model  $\mathcal{N}$  through  $\Gamma$  is also finite.]

**Definition 3.5.** [Define the fine filtration (I need to pick a specific filtration to see that one exists. In order to pick one that may prove to be more useful as a neural network, I pick the fine filtration (the most densely connected filtration) over the coarse filtration (the least densely connected filtration)]

**Proposition 3.6.** [The fine filtration is in fact a filtration]

**Corollary 3.7.** (Finite Model Property) [State the finite model property, putting it all together. We can build a finite ]

**Corollary 3.8.** (Finite Model Building) [For any finite set of constraints  $\Gamma$ , we can construct a *finite*  $\mathcal{N} \in \mathbf{Net}$  satisfying  $\Gamma$ ]

[Note that this doesn't actually give us a constructive *algorithm* for building the finite model, since we have to build the infinite canonical model first.] [After talking to Saul and Larry about this, I should just leave it as an open question and point to Larry's work on constructively building finite models from formulas. It is likely this trick will work here too.]

**Open Question 3.** For a given finite set of constraints  $\Gamma \subseteq \mathcal{L}_{\mathbf{C}}$ , is there a constructive algorithm that produces a model  $\mathcal{M} \in \mathbf{Plaus}$  and state  $w \in W$  such that  $\mathcal{M}, w \models \Gamma$ ?

## **4** Model Building and Completeness Theorems

[I can only state + prove these after doing filtration!!]

**Theorem 4.1.** (C) Model Building for  $\models_{\text{Net}}$ ) For all consistent  $\Gamma \subseteq \mathcal{L}_{C}$ , there is some finite  $\mathcal{N} \in \text{Net}$ and neuron  $n \in N$  such that  $\mathcal{N}, n \Vdash_{\text{Net}} \Gamma$ .

**Proof.** Let  $\Gamma$  be consistent, and let  $\mathcal{M} \in \mathbf{Plaus}$ ,  $w \in W$  be some finite plausibility model and world such that  $\mathcal{M}, w \Vdash_{\mathbf{Plaus}} \varphi$  (which we know exists by Theorems [model building] and [filtration], see Appendix [] for details). Now let  $\mathcal{N} \in \mathbf{Net}$  be given by the NAND construction above. We have  $w \in W \subset N$ , and so we can evaluate  $\Vdash_{\mathbf{Net}}$  at  $w \in N$ . By Lemma 5.19,  $\mathcal{N}, w \Vdash_{\mathbf{Net}} \varphi$ , and we are done.  $\Box$ 

**Corollary 4.2.** (Completeness for  $\models_{Net}$ ) For all consistent  $\Gamma \subseteq \mathcal{L}$  and all formulas  $\varphi \in \mathcal{L}$ ,

if 
$$\Gamma \models_{\mathbf{Net}} \varphi$$
 then  $\Gamma \vdash \varphi$ 

**Proof.** Since the language  $\mathcal{L}_{\mathbf{C}}$  has negation, completeness follows from model building in the usual way; this proof is entirely standard. Suppose contrapositively that  $\Gamma \not\models \varphi$ . It follows that  $\Gamma \vdash \neg \varphi$ . So  $\Gamma \cup \{\neg \varphi\}$  is consistent, and by Theorem 4.1 we have  $\mathcal{N} \in \mathbf{Net}$  and  $n \in N$  such that  $\mathcal{N}, n \models_{\mathbf{Net}} \Gamma \cup \{\neg \varphi\}$ . But then  $\mathcal{N}, n \models_{\mathbf{Net}} \Gamma$  yet  $\mathcal{N}, n \not\models \varphi$  [fix notation here], which is what we wanted to show.  $\Box$ 

**Example 4.3.** (Building a Neural Network Model I) The discussion above is all very abstract. Let's see how the neural network model building is done by example. Consider our birds example from before, where  $\mathcal{L}_{prop} = \{$ bird, penguin, fly $\}$ . Say we want to build a model satisfying the non-monotonic constraints from before, i.e.,

$$\Gamma = \{ bird \Rightarrow fly, penguin \rightarrow bird, penguin \rightarrow \neg fly \}$$

That is, birds typically fly, penguins are birds, but penguins do not fly.

First, Theorem 4.1 says that we can construct a neural network *and a neuron* at which these constraints are true. But say we want these constraints to be true everywhere. This is simple enough: just add the universal quantifier **A** in front of each constraint. Additionally, bird $\Rightarrow$ fly is not in the language of  $\mathcal{L}_{\mathbf{C}}$ , so we need to translate it. Recall that the conditional bird $\Rightarrow$ fly is equivalent to **A**(**C**bird $\rightarrow$ fly). Our set of constraints is now:

$$\Gamma = \{A(Cbird \rightarrow fly), A(penguin \rightarrow bird), A(penguin \rightarrow \neg fly)\}$$

First, we need to construct a finite plausibility model  $\mathcal{M} \in \mathbf{Plaus}$ . It is always possible to construct



**Figure 4.1.** [TODO—caption! Separate into (a) left, and (b) right] T, P, S, N denote Tweety, Piper, Skipper, and Nils, respectively. For (b), mention that the activation function is A(x) = 1 iff x > -1.

a plausibility model with constraints in  $\mathcal{L}_{C}$ —see Appendix [TODO] for the details. However, the proof in Appendix [TODO] is not constructive, so I do not currently have an algorithm for doing this. For small examples like this one, it's not too difficult to hand-craft a model.

Take the following model, illustrated by Figure 4.1 (a). Let  $W = \{\text{Tweety}, \text{Piper}, \text{Skipper}, \text{Nils}\}$ , and let the propositional assignment V be given by  $[[bird]] = \{\text{Tweety}, \text{Piper}, \text{Skipper}, \text{Nils}\}$ , [[pen $guin]] = \{\text{Skipper}, \text{Nils}\}$ , and  $[[fly]] = \{\text{Tweety}, \text{Piper}\}$ . That is, Skipper and Nils are the only penguins, and Tweety and Piper are the only birds that fly. Finally, let the plausibility order ( $\prec$ ) = {(Tweety, Skipper), (Tweety, Nils), (Piper, Skipper), (Piper, Nils)}.

In order to transform this model into a neural network, we apply the NAND construction. Let N = W, V = V,  $E = \prec$ , and add a new node bias. Connect the bias node to Skipper and Nils (since they are not  $\prec$ -minimal). Then set all weights to  $W(m,n) = -\frac{1}{3}$ . Finally, pick A(x) = 1 iff x > -1, and take an arbitrary learning rate  $\eta$ . The resulting neural network  $\mathcal{N}$  is shown in Figure 4.1 (b).

Let's check that  $\mathcal{N}$  does in fact satisfy the constraints  $\Gamma$ , i.e.,  $\mathcal{N} \models \Gamma$ . Well,

- $\mathcal{N} \models \mathbf{A}(\text{penguin} \rightarrow \text{bird})$ , since  $[[\text{penguin}]] = \{\text{Skipper, Nils}\} \subseteq [[\text{bird}]]$  (for all  $w \in N, w \Vdash \text{penguin implies } w \Vdash \text{bird})$ .
- Similarly,  $\mathcal{N} \models \mathbf{A}(\text{penguin} \rightarrow \neg \text{fly})$ , since  $[[\text{penguin}]] = \{\text{Skipper}, \text{Nils}\} \subseteq [[\text{flies}]]^{\mathbb{C}}$ .
- Finally, let's check  $\mathcal{N} \models \mathbf{A}(\mathbf{Cbird} \rightarrow \mathbf{fly})$ . In terms of neural network semantics, this says that

for all  $w \in N$ , if  $w \in (\text{Clos}(\llbracket \text{bird} \rrbracket^{C} \cup \{\text{bias}\}))^{C}$  then  $w \in \llbracket \text{fly} \rrbracket$ . Observe that

 $(Clos(\llbracketbird\rrbracket^{C} \cup \{bias\}))^{C} = (Clos(\{bias\}))^{C}$  $= \{bias, Skipper, Nils\}^{C}$  $= \{Tweety, Piper\}$ 

From here we see that all *w* in this set {Tweety, Piper} fly.

Example 4.4. (Building a Neural Network Model II) [TODO]

## **5** The Modelling Power of Neural Networks

So far, I have presented one important neural network model construction: the canonical construction. The canonical construction shows us how a neural network model  $\mathcal{N}$  can simulate any finite choice of formulas  $\Gamma$  over the logic  $\mathcal{L}_{C}$ . This gives us a sense of how powerful **Net** is—of what sorts of constraints neural networks (as logical models) are capable of representing. In this section, I will give a richer sense of how powerful neural networks are by comparing **Net** against other model classes. Along the way, I will catalog specialized neural network constructions that we can use to have **Net** simulate (or be simulated by) these other model classes. I aim to answer the questions:

- What kinds of models can **Net** simulate?
- What kind of models can be simulated by Net?

I will introduce a number of models for modal logics in the literature: relational models **Rel**, neighborhood models **Nbhd**, plausibility models **Plaus**, and a special case of social network models **SocialNet**<sup>MAJ</sup>. I will include in the mix:

- *Pointed* neural network models **Net**\*, briefly discussed in Section []
- "Universal" models **Univ**, which is degenerate case of **Rel** where all modalities are interpreted as the universal modality **A**.

To make the comparison fair, all models will share a generalized multi-modal language  $\mathcal{L}_{\Box}$  given by

$$\varphi, \psi \coloneqq p \mid \neg \varphi \mid \varphi \land \psi \mid \{\Box_i\}_{i \in I} \varphi$$

where I is some fixed set of indices. Notice that  $\mathcal{L}_{C}$  is a special case of  $\mathcal{L}_{\Box}$ , where the operators  $A, \Box, C$  are different choices for  $\Box_{i}$  that come with additional interaction axioms. More generally, you can think of each  $\Box_{i}$  as representing a different modality (e.g., belief vs knowledge). Or, you could instead think of each  $i \in I$  as an agent in a multi-agent setting, and  $\Box_{i}$  as agent *i*'s belief (or knowledge, etc.).

**Generalized Neural Network Models.** This change in language requires a somewhat generalized neural network semantics. In this work so far, I have only defined the neural network closure operators Clos and Reach. But in principle we could define other closure operators, each reflecting a different kind of modality or conditional. I want to characterize what  $\Box_i$ -formulas a neural network can *in principle model*, and for that I need a more general definition.

**Definition 5.1.** A neural network model is  $\mathcal{N} = \langle N, \{E_i\}_{i \in I}, \{W_i\}_{i \in I}, \{A_i\}_{i \in I}, \eta, V \rangle$ , where we now have a choice of  $E_i$ ,  $W_i$ ,  $A_i$  for each  $i \in I$ .

**Definition 5.2.** A *pointed* neural network model is  $\mathcal{N} = \langle N, \text{bias}, \{E_i\}_{i \in I}, \{W_i\}_{i \in I}, \{A_i\}_{i \in I}, \eta, V \rangle$ , whose states are given by  $\text{State}_{\mathcal{N}} = \{S \mid S \subseteq N \text{ and } \text{bias} \in S\}.$ 

Everything from Section  $\square$  generalizes: Each choice  $E_i$ ,  $W_i$ ,  $A_i$  specifies a state transition function  $F_{i,S_0}$ : State<sub>N</sub>  $\rightarrow$  State<sub>N</sub>, which in turn defines a closure function  $Clos_i$ : State<sub>N</sub>  $\rightarrow$  State<sub>N</sub>. For the purposes of this section, let **Net** refer to the class of this generalized neural network model, and **Net**\* refer to the pointed variation.

The semantics for  $\diamond_i$  generalizes as follows.

$$\mathcal{N}, n \Vdash \diamond_i \varphi \text{ iff } n \in \operatorname{Clos}_i(\llbracket \varphi \rrbracket)$$

Similarly, for pointed models I define

$$\mathcal{N}, n \Vdash \diamond_i \varphi \text{ iff } n \in \operatorname{Clos}_i(\llbracket \varphi \rrbracket \cup \{ \text{bias} \})$$

As expected,  $\Box_i$  is defined to be the dual of  $\diamond_i$ .

**Note.** I claimed above that the operators  $\mathbf{A}, \Box, \mathbf{C}$  can be viewed as instances of  $\Box_i$ . For  $\mathbf{C}$ , this is fairly obvious, since the semantics for  $\langle \mathbf{C} \rangle$  are exactly the same as  $\diamond_i$ . To see how  $\mathbf{A}$  can be

simulated via  $\Box_i$ , see the subsection on simulating **Univ** (where  $\Box_i$  is defined to be exactly **A**). Finally, the  $\Box$  operator was originally given semantics in terms of graph-reachability Reach. To see how it can be simulated via  $\Box_i$ , see the subsection on simulating **Rel**, and then reverse the edges in that construction. (The reason for this is that relational models "look upward" for some world *u* where  $\varphi$  holds, whereas Reach "looks downward.")

An Outline and Summary. I will now get in the weeds and formally compare these model classes. First, I will formally define simulations that will allow us to compare the modeling power of Net against the rest. Then, for each model class, I will prove simulation results that position them around Net. Figure [] illustrates a summary of these results. The figure displays a lattice where each class can simulate the models in the classes above it.

#### 5.1 Measuring Modeling Power via Simulations

To compare the modeling power of neural networks with these other classes, I need to measure their ability to simulate each other. Formally,

**Definition 5.3.** Let  $\mathscr{C}_1, \mathscr{C}_2$  be two model classes, with semantics defined over  $\mathcal{L}_{\Box}$ . A *simulation f*:  $\mathscr{C}_1 \to \mathscr{C}_2$  is an injective function such that for all models  $\mathcal{M} \in \mathscr{C}_1$  and formulas  $\varphi \in \mathcal{L}_{\Box}$ ,

$$\mathcal{M} \models \varphi \text{ iff } f(\mathcal{M}) \models \varphi$$

If a simulation exists, we say that  $C_2$  simulates  $C_1$ . If  $C_2$  simulates  $C_1$ , but not conversely, we say that the simulation is *strict*.

In model theory, the *theory* of a class of models  $\operatorname{Th}(\mathscr{C})$  is the set of validities, i.e., formulas that *must* be true for every  $\mathcal{M} \in \mathscr{C}$ . This is often used as a measure of the modeling power of  $\mathscr{C}$ : if  $\operatorname{Th}(\mathscr{C}_1) \subseteq \operatorname{Th}(\mathscr{C}_2)$ , then  $\mathscr{C}_1$  is more general than  $\mathscr{C}_2$  in the sense that it requires fewer axioms. (In other words,  $\mathscr{C}_1$  can satisfy formulas that  $\mathscr{C}_2$  cannot.)

**Definition 5.4.** Let  $\mathscr{C}$  be a class of models. The *theory* of  $\mathscr{C}$  (over  $\mathcal{L}_{\Box}$ ) is

Th(
$$\mathscr{C}$$
) = { $\varphi \in \mathcal{L}_{\Box}$  | for all  $\mathcal{M} \in \mathscr{C}$  we have  $\mathcal{M} \models \varphi$ }

The following proposition states the relationship between simulations and  $\operatorname{Th}(\mathscr{C})$ . In words, if  $\mathscr{C}_2$  simulates  $\mathscr{C}_1$ , then  $\mathscr{C}_2$  is the more general one. This proposition also gives us an easy test for showing that  $\mathscr{C}_2$  does not simulate  $\mathscr{C}_1$ : just try to show  $\operatorname{Th}(\mathscr{C}_2) \not\subseteq \operatorname{Th}(\mathscr{C}_1)$  by finding a formula  $\varphi \in \operatorname{Th}(\mathscr{C}_1)$  such that  $\varphi \notin \operatorname{Th}(\mathscr{C}_2)$ .

**Proposition 5.5.** Let  $\mathscr{C}_1, \mathscr{C}_2$  be two model classes. If  $\mathscr{C}_2$  simulates  $\mathscr{C}_1$ , then  $\operatorname{Th}(\mathscr{C}_2) \subseteq \operatorname{Th}(\mathscr{C}_1)$ .

**Proof.** Suppose  $\mathscr{C}_2$  simulates  $\mathscr{C}_1$ . This means there is an injection  $f: \mathscr{C}_1 \to \mathscr{C}_2$  such that  $\mathcal{M} \models \varphi$  iff  $f(\mathcal{M}) \models \varphi$ . Now let  $\varphi \in \text{Th}(\mathscr{C}_2)$ , and let  $\mathcal{M} \in \mathscr{C}_1$ . Well,  $f(\mathcal{M}) \in \mathscr{C}_2$ , and so by the definition of  $\text{Th}(\mathscr{C}_2), f(\mathcal{M}) \models \varphi$ . By our simulation, this gives us  $\mathcal{M} \models \varphi$ , and we are done.  $\Box$ 

**Note.** The notion of simulation that I've defined here only tells the semantic side of the story. If we want to compare different *languages* in addition to model classes, we need to define *infomorphisms* between two logics  $(\mathcal{L}, \mathscr{C})$ . (Pairs of languages, along with their model class, are also known as *institutions* in Institution Theory [cite institution theory, "Information Flow: The Logic of Distributed Systems", "Categories, Allegories", "Institution Theory"]). Infomorphisms are defined as follows.

**Definition 5.6.** Let  $\mathscr{C}_1, \mathscr{C}_2$  be model classes, and let  $\mathcal{L}_1, \mathcal{L}_2$  be languages. There is an *infomorphism* (aka *translation*) from  $(\mathcal{L}_1, \mathscr{C}_1)$  into  $(\mathcal{L}_2, \mathscr{C}_2)$  if there exist  $f: \mathscr{C}_2 \to \mathscr{C}_1, \tau: \mathcal{L}_1 \to \mathcal{L}_2$  such that for all  $\varphi \in \mathcal{L}_1, \mathcal{M} \in \mathscr{C}_2$ 

$$f(\mathcal{M}) \models \varphi \text{ iff } \mathcal{M} \models \tau(\varphi)$$

Since I have fixed the language  $\mathcal{L}_{\Box}$ , I won't necessarily need to use this idea. But I mention it here in order to set readers who are curious about the expressivity of ( $\mathcal{L}_{\Box}$ , **Net**) on the right foot.

#### 5.2 Relational Models (Rel)

Relational models are just Kripke models, which I briefly introduced in Section [].

A relational model is  $\mathcal{M} = \langle W, \{R\}_{i \in I}, V \rangle$ , where

- *W* is some finite set of worlds (or states)
- Each  $R_i \subseteq W \times W$  (the accessibility relations)
- V: Proposition  $\rightarrow \mathcal{P}(W)$  (the valuation function)

Define **Rel** to be the class of all such models, and define  $\operatorname{Rel}_{S4}$  to be the class of all such models where *R* is additionally reflexive and transitive. The semantics for both classes is given by: [Todo, extend to the full language of  $\mathcal{L}_{C}$ !]

 $\begin{array}{lll} \mathcal{M}, w \Vdash p & \text{iff} & w \in V(p) \\ \mathcal{M}, w \Vdash \neg \varphi & \text{iff} & \mathcal{M}, w \not \models \varphi \\ \mathcal{M}, w \Vdash \varphi \land \psi & \text{iff} & \mathcal{M}, w \Vdash \varphi \text{ and } \mathcal{M}, w \Vdash \psi \\ \mathcal{M}, w \Vdash \Box_i \varphi & \text{iff} & \text{for all } u \text{ with } wR_i u, \mathcal{M}, u \Vdash \varphi \end{array}$ 

[Mention axioms, soundness, completeness (refer to the appendix!)]

## 5.3 Universal Models (Univ)

#### 5.4 Neighborhood Models (Nbhd)

## 5.5 Plausibility Models (Plaus)

[(this is also where I can talk about **Net**\*)]

# 5.6 Majority-Vote Social Networks (SocialNet<sup>MAJ</sup>)

• Let's now do model translations to get at the expressivity of neural networks (over Modal and Conditional logic). Here's a hierarchy of the models above, to start:

- To make the comparison with neural networks fair, I will only consider the reflexive and transitive variants **Rel**<sub>*S*4</sub>, **Nbhd**<sub>*S*4</sub> of relational and neighborhood models.
- [DIAGRAM]
- The translations from Nbhd<sub>S4</sub> to Plaus and from Plaus to Rel<sub>S4</sub> are folklore. Instead of giving these translations, I will instead translate from Net to Rel<sub>S4</sub> and from Nbhd<sub>S4</sub> to Net in order to explicitly show how to translations involving neural networks (neural network model constructions). The equivalence between Plaus and Net is already known, but the backwards direction has never been proven with an explicit model construction. So although the results in this section are already known, the proofs I give here are totally new.
- First, show translation from **Net** to **Rel**<sub>*S*4</sub>, and show there is no translation the other way around (axiom in **Rel**<sub>*S*4</sub> that is not an axiom in **Net**)
- Next, show translation from **Nbhd**<sub>*S*4</sub> to **Net**, and show there is no translation the other way around (axiom in **Net** that is not an axiom in **Nbhd**<sub>*S*4</sub>)
- Next, from the model-building construction in the completeness chapter, we have a translation from **Net** to **Plaus**.
- Explain that completeness implies that *in principle* we have a translation the other way (**Plaus** to **Net**), but it doesn't actually give the explicit model building procedure! Here is where I will give my own.
- Make a note here about the social majority logic above: There is a strict translation from
  Net to the social majority logic (Net is more general, in the sense that it requires fewer
  axioms).
- Give a brief note here on the expressive power of the conditional logic **Net** vs the modal logic **Net**.

**Plausibility Models.** A plausibility model, first introduced in [37], is  $\mathcal{M} = \langle W, \{R\}_{i \in \mathbf{I}}, V \rangle$ , i.e. the models themselves are just relational models. As before, I assume that *W* is finite, and as with **Rel**<sub>*S*4</sub>, each *R*<sub>*i*</sub> is reflexive and transitive. The key difference is that we interpret  $\Box_i \varphi$  to hold in the best (or most plausible) states satisfying  $\varphi$ . Formally, let  $\text{best}_{R_i}(S) = \{w \in S \mid \text{for all } u \in S, \neg uR_iw\}$ 

(the  $R_i$ -minimal states over *S*). We additionally impose the following "smoothness condition" [37] on best<sub> $R_i$ </sub>:

**Postulate 5.7.** For all models  $\mathcal{M}$ ,  $i \in I$ , sets S, and all  $w \in W$ , if  $w \in S$  then either  $w \in \text{best}_{R_i}(S)$ , or there is some  $vR_iw$  better than w that *is* the best, i.e.  $v \in \text{best}_{R_i}(S)$ .

The new semantics for  $\Box_i$  is

$$\mathcal{M}, w \Vdash \Box_i \varphi$$
 iff  $w \in \mathsf{best}_{R_i}(\llbracket \varphi \rrbracket)$ 

where  $\llbracket \varphi \rrbracket = \{u \mid \mathcal{M}, u \models \varphi\}$ . In practice, plausibility semantics coexist alongside relational semantics, so I allow some  $\Box_i \varphi$  to be given relational semantics instead. Let **Plaus** be the class of all such models. Since we include relational operators, note that **Rel**<sub>S4</sub>  $\subseteq$  **Plaus**.

Any plausibility operator  $\Box_i$  picks out a corresponding conditional:  $\Box_i \varphi \rightarrow \psi$  reads "the best  $\varphi$  are  $\psi$ ," which in the KLM tradition is the semantics for the conditional  $\varphi \Rightarrow \psi$ .

[Mention axioms, soundness, completeness (refer to the appendix!)]

#### Social Network Models.

- Introduce social network models, an example with a [DIAGRAM] would be nice!
- In these social network logics [4; 8; 18], nodes in the graph represent individual agents, and each formula is mapped to the set of agents that adopt a certain social attitude. Agents influence each other, and the spread of their attitudes is modeled much in the same way as forward propagation of a signal in a neural network.
- Give a concrete social network logic: Social majority. Make sure to emphasize that social majority is one of many (!) choices, and is a relatively simple choice to model.
- Both kinds of models use fundamentally the same approach ("This work shares essentially the same premise and techniques as neural network semantics"): distributed information over several connected nodes, modal operator interpreted as the fixed-point of some diffusion
- But the two approaches still differ in interesting ways. First, in some sense the two semantics are operating on different "levels": social networks model interactions between multiple

agents, whereas neural networks model interactions between components of the same (single) agent. Second, the two differ in subject matter. Social network semantics focuses on different social links between agents, and how these links change [4]. Neural network semantics, my own work included, instead focuses on inferences and updates inspired by artificial and natural neural networks.

#### Neighborhood Models.

A neighborhood model is  $\mathcal{M} = \langle W, \{N\}_{i \in \mathbf{I}}, V \rangle$ , where *W* and *V* are as before and each  $N_i: W \to \mathcal{P}(\mathcal{P}(W))$  is an accessibility *function*. The intuition is that  $N_i$  maps each state *w* to the "formulas" (sets of states) that hold at *w*. Define **Nbhd** to be the class of all neighborhood models.

Moreover, the *core* of *N* is  $\cap N(x) = \bigcap_{X \in N(w)} X$ . As with **Rel**, let **Nbhd**<sub>S4</sub> be the class of all neighborhood models that are additionally reflexive  $(\forall w, w \in \cap N(w))$  and transitive  $(\forall w, if X \in N(w))$  then  $\{v \mid X \in N(v)\} \in N(w)$ ).

The semantics for both classes is the same as the previous classes, except the  $\Box_i$  case is now:

$$\mathcal{M}, w \Vdash \Box_i \varphi \text{ iff } [\![\varphi]\!] \in N_i(w)$$

where again  $\llbracket \varphi \rrbracket = \{ u \mid \mathcal{M}, u \Vdash \varphi \}.$ 

**Proposition 5.8.** (C) There is a strict translation from (Modal, Net) to (Modal,  $\text{Rel}_{S4}$ ).

Proof. []

**Proposition 5.9.** (C) There is a strict translation from (Modal, Nbhd<sub>S4</sub>) to (Modal, Net).

**Proposition 5.10.** There is a translation from ( $\mathcal{L}_{C}$ , Net) to ( $\mathcal{L}_{C}$ , Plaus).

**Proof.** [This is just a corollary of completeness: Mention that we just use the NAND construction from the Completeness proof.] □

An easy corollary of completeness is that exactly the same axioms hold over **Plaus** and **Net**, i.e. [This is no longer true!!! Rewrite this whole story here!]

**Corollary 5.11.** Th(Plaus) = Th(Net).

#### **Proof.** We have:

$\varphi \in \mathrm{Th}(\mathbf{Plaus})$	iff	$\models_{\text{Plaus}} \varphi$	(By definition)
	iff	$\vdash_{\mathbf{C}} \varphi$	(By weak completeness for Plaus)
	iff	$\models_{\text{Net}} \varphi$	(By weak completeness for Net)
	iff	$\varphi \in \mathrm{Th}(\mathbf{Plaus})$	(By definition)

Based on this, and the completeness result, we might naturally expect that there is a translation from ( $\mathcal{L}_{C}$ , **Plaus**) to ( $\mathcal{L}_{C}$ , **Net**). But surprisingly, there is *no* such translation!

**Theorem 5.12.** (C) There is no translation from  $(\mathcal{L}_{C}, \text{Plaus})$  to  $(\mathcal{L}_{C}, \text{Net})$ .

**Proof.** [First, I can prove it if  $\tau(\varphi) = \varphi$ ! I should try to generalize this for any  $\tau$ !] For  $\tau(\varphi) = \varphi$ , let  $f: \mathbf{Net} \to \mathbf{Plaus}$  be any arbitrary way to transform a neural net into a plausibility model. I need to show that there is some net  $\mathcal{N} \in \mathbf{Net}$  and some formula  $\varphi \in \mathcal{L}_{\mathbf{C}}$  such that one of  $f(\mathcal{N}) \models \varphi$  or  $\mathcal{N} \models \varphi$ holds and the other does not. Let  $\varphi = \mathbf{C}p$  for some proposition p, and let  $\mathcal{N}$  be the net in Figure [DIAGRAM] (a). Since  $[\![\mathbf{C}p]\!]_{\mathcal{M}} = \mathsf{best}_{\prec_{f(\mathcal{N})}}([\![p]\!])$  and  $[\![\mathbf{C}p]\!]_{\mathcal{N}} = (\mathsf{Clos}_{\mathcal{N}}([\![p]\!]^{\mathsf{C}}))^{\mathsf{C}}$ , it's enough to show that

There is some 
$$\llbracket p \rrbracket \in \text{State}_{\mathcal{N}}$$
 such that  $\text{best}_{\leq_{f(\mathcal{N})}}(\llbracket p \rrbracket) \neq (\text{Clos}_{\mathcal{N}}(\llbracket p \rrbracket^{\mathbb{C}}))^{\mathbb{C}}$ 

Well, for this particular net we have  $\text{State}_{\mathcal{N}} = \{\{\text{bias}\}, \{\text{bias}, a\}, \{\text{bias}, b\}, \{\text{bias}, a, b\}\}$  (all subsets of the nodes that contain the bias node). Any  $[\![p]\!] \in \text{State}_{\mathcal{N}}$  must be one of these four sets.

Suppose for contradiction that for all such  $\llbracket p \rrbracket$ ,  $best_{\prec_{f(\mathcal{N})}}(\llbracket p \rrbracket) = (Clos_{\mathcal{N}}(\llbracket p \rrbracket^{C}))^{C}$ . In particular,  $best_{\prec_{f(\mathcal{N})}}(\{bias, a, b\}) = (Clos_{\mathcal{N}}(\{bias, a, b\}^{C}))^{C} = \{bias, a, b\}$ . (See the calculation in Figure [DIA-GRAM] (b).) But this implies that bias, *a*, and *b* must all be mutually incomparable via  $\prec_{f(\mathcal{N})}$ . But then

$$\text{best}_{\leq_{\ell(\mathcal{N})}}(\{\text{bias}, b\}) = \{\text{bias}, b\} \neq \{\text{bias}\} = (\text{Clos}_{\mathcal{N}}(\{\text{bias}, b\}^{C}))^{C}$$

(The first equality comes from bias, a, b incomparable; the second follows from the calculation in Figure [DIAGRAM] (b).)

In particular, the neural network used in the proof [refer to it] cannot be transformed into an equivalent plausibility model. Moreover, this net isn't a particularly strange one—it's just an ordinary feed-forward weighted net! So we have an intuition that neural networks are more general in

some sense. But Corollary 5.11 says that no axiom  $\varphi \in \mathcal{L}_{\mathbf{C}}$  witnesses this difference. This means that there *is* some difference between them, but  $\mathcal{L}_{\mathbf{C}}$  is not expressive enough to point to it. [the interesting conclusion of this reasoning—what axiom can give us the difference??]

**Theorem 5.13.** (C) There is a strict translation from (Modal, Net) to the social majority logic [give it a name]

Proof. []

**Theorem 5.14.** (C) There is a strict translation from (Modal, Net) to (Conditional, Net). [does it go this way, or the other way??]

Proof. []

[Integrate the following into this chapter!]

**The Graph-Reachability Construction.** I will show here how the more general Clos operator can be used to simulate Reach and Reach<sup> $\downarrow$ </sup>. Suppose we are given a graph  $\langle N, E \rangle$  and a valuation function *V*. How can we build a neural network whose closure Clos is graph-reachability Reach? We want to build a net:

$$\mathcal{N} = \langle N, \text{bias}, E, W, A, \eta, V \rangle$$

[what to do about bias and  $\eta$ ?] For the weights, pick

$$W(m,n) = \begin{cases} 1 & \text{if } mEn \\ 0 & \text{otherwise} \end{cases}$$

Then pick the activation function A(x) = 1 iff x > 0. Recall that  $n \in F_{S_0}(S)$  iff  $n \in S_0$  or is activated by its predecessors in *S*. In this case,  $n \in F_{S_0}(S)$  whenever  $n \in S_0$  or at least one *E*-predecessor *m* of *n* is in *S*. I call this the *graph-reachability construction* because the closure Clos(S) produces exactly those nodes graph-reachable from *S*:

**Proposition 5.15.** For all states  $S \in \text{State}$ , Clos(S) = Reach(S).

**Proof.** First, the  $(\subseteq)$  direction. Let  $n \in \text{Clos}(S) = F_S^k(S)$  for some  $k \in \mathbb{N}$ . By induction on k.

**Base Step.**  $n \in F_S^0(S) = S$ . So there is a trivial  $E_i$ -path (length=0) from  $n \in S$  to itself.

**Inductive Step.** Let  $k \ge 0$ . We have  $n \in F_S^k(S) = F_S(F_S^{k-1}(S))$ . By construction of  $F_S$ , we have two cases: Either  $n \in F_S^{k-1}(S)$  or at least one *E*-predecessor *x* of *n* is in  $F_S^{k-1}(S)$ . In the first case, our inductive hypothesis gives a path from some  $m \in S$  to *n*. In the second case, our inductive hypothesis gives a path from some  $m \in S$  to *x*. But since *xEn*, we can extend this path to be from *m* to *n*.

As for the  $(\supseteq)$  direction, suppose there is an *E*-path from some  $m \in S$  to *n*. We proceed by induction on the length of that path.

- **Base Step.** The path is trivial, i.e. has length 0. So  $n \in S$ . But  $S = F_i^0(S) \subseteq \text{Clos}_i(S)$ , and so  $n \in \text{Clos}_i(S)$ .
- **Inductive Step.** Say the path is of length *l*≥0. Let *x* be some immediate *E<sub>i</sub>*-predecessor of *n*. By the inductive hypothesis,  $x \in Clos_i(S)$ , and so  $x \in F_i^k(S)$  for some natural *k*. But since *x* is an *E<sub>i</sub>*-predecessor of *n*, by construction of *F<sub>i</sub>*,  $n \in F_i(F_i^k(S)) = F_i^{k+1}(S)$ . Since  $Clos_i(S)$  is a closure, it includes  $F_i^{k+1}(S)$ . So  $n \in Clos_i(S)$ .

As for Reach<sup> $\downarrow$ </sup>, we first *reverse* the edges *E*, and then do the graph-reachability construction. In other words, let *mE'n* iff *nEm*,

$$W(m,n) = \begin{cases} 1 & \text{if } mE'n \\ 0 & \text{otherwise} \end{cases}$$

and pick A,  $\eta$ , V the same as above. For this construction, the closure Clos(S) produces exactly those nodes which reach some node in S. The proof is similar to the proof for Reach. [Check that that's actually true!!]

**Proposition 5.16.** For all states  $S \in \text{State}$ ,  $\text{Clos}(S) = \text{Reach}^{\downarrow}(S)$ .

**The Social Majority Construction.** [Introduce the idea of social networks here if I haven't already, and mention the "social majority" propagation/diffusion (tell it slowly, like a story). I will show that our neural networks can simulate this simple social majority operator (make the social majority operator a bit more formal).]

As before, we want to build a neural network  $\mathcal{N}$  where the graph  $\langle N, E_i \rangle$ , bias, and evaluation V are given. This time, pick  $W_i(m,n) = \frac{1}{|\text{preds}(n)|}$ , and then pick  $A_i(x) = 1$  iff  $x \ge \frac{1}{2}$ . Visually, for each



Figure 5.1. [TODO—caption! Separate into (a) left, and (b) right]

node *n* and its predecessors  $m_1, \ldots, m_r$  we have

#### [DIAGRAM!]

This gives us  $n \in F_{S_0}(S)$  if  $n \in S_0$  or if the majority (more than half) of *E*-predecessors are in *S*. In this case, the closure Clos can be interpreted as the diffusion of an opinion or attitude through a social network. This is one of the choices that [8] consider for modelling influence in social networks. [this paragraph is a bit terse now that I'm writing a longer version of this.]

**The NAND Construction.** Suppose we have plausibility model  $\mathcal{M} = \langle W, R, \prec, V \rangle$  and we want to construct an equivalent neural network  $\mathcal{N}$ . In [38], Hannes first does this for *inhibition nets*, i.e., nets with inhibitory edges that block excitatory edges. (He handles weighted nets later.) I will first consider his construction, and then modify it for weighted nets.

Here's the inhibition net construction: First, create a fresh node bias. Take  $N = W \cup \{\text{bias}\}$  (so  $\mathcal{N}$  is still finite), V = V, let the excitatory edges be exactly uEv iff vRu (E is the reverse of R). Next, create an edge from bias to every n that is not E-minimal (in other words, if n has any predecessors at all, then bias is one of them). Then for each node n and its predecessors bias  $= m_0, m_1, \ldots, m_r$ , connect inhibition edges as illustrated in Figure 5.1 (a).

That is, each node  $m_i$  is inhibited by  $m_{i-1}$  (bias =  $m_0$  inhibited by  $m_r$ ). This has the following effect: if all  $m_i$  activate, they each inhibit each other, and so n does not activate. If only *some*  $m_i$ activate, then there is some  $m_i$  that is uninhibited, and so n activates. And finally, since bias is always active we cannot have *no*  $m_i$  active. In other words,  $n \in F_{S_0}(S)$  iff  $n \in S_0$ , or *not all non-bias predecessors* m *are in* S. (Since bias is always active, this results in a NAND-like output.)

We can simulate this effect with weighted neural networks. Let  $W(m, n) = -\left(\frac{1}{r+1}\right)$  (the extra +1 accounts for the bias), and let pick A(x) = 1 iff x > -1. For now, the choice of learning rate  $\eta$  is arbitrary (see Section [todo]). This construction is illustrated in Figure 5.1 (b). Take a moment to check that  $n \in F_{S_0}(S)$  iff  $n \in S_0$ , or at least one non-bias predecessor  $m \notin S$ . [also, the negative values ensure *A* is nondecreasing!]

What is the relationship between this neural network's fixed points Clos(S) and the plausibility model's minimal states best<sub><</sub>(S)? It turns out that for this NAND construction, Clos is precisely the *dual* of best<sub><</sub>:

**Lemma 5.17.** Let  $\mathcal{M} = \langle W, R, \prec, V \rangle$  be a plausibility model such that  $R = (\prec_{\text{reverse}})$ . Let  $\mathcal{N}$  be given by the NAND construction above. For all  $S \subseteq \text{State}_{\mathcal{N}}$ ,  $\text{Clos}_{\mathcal{N}}(S) = (\text{best}_{\prec}(S^{\mathbb{C}}))^{\mathbb{C}}$ . In a slogan:

$$Clos_{\mathcal{N}}(S)$$
 is the dual of  $best_{\prec}(S)$ .

**Note.** There is some ambiguity over what universe the complements in  $(\text{best}_{<}(S^{C}))^{C}$  are taken. The bias node of course occurs in *S*, since *S* is a state of the net. But we should also allow for bias to occur in the final complement of  $(\text{best}_{<}(S^{C}))^{C}$ , in order for this term to be exactly  $\text{Clos}_{\mathcal{N}}(S)$ . So I interpret all of these complements over universe *N* (rather than *W*). Explicitly, the claim is:

For all 
$$S \subseteq \text{State}_{\mathcal{N}}, \text{Clos}_{\mathcal{N}}(S) = N \setminus (\text{best}_{\prec}(N \setminus S))$$

**Proof.** Once again, I will take advantage of the fact that fixed points of the transition function  $F_S$  are unique. First, since Clos(S) is a fixed point under *S*, we have  $F_S(Clos(S)) = Clos(S)$ . But I will show that  $(best_{<}(S^{C}))^{C}$  is *also* a fixed point under *S*, i.e.  $F_S((best_{<}(S^{C}))^{C}) = (best_{<}(S^{C}))^{C}$ . Since we assumed that there is a *unique* fixed point under *S*, it will follow that these two sets must be the same. In other words,  $Clos(S) = (best_{<}(S^{C}))^{C}$ .

Recall that the type of  $F_S$  is  $F_S$ : State<sub>N</sub>  $\rightarrow$  State<sub>N</sub>. We need to make sure that  $(\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}} \in$ State<sub>N</sub>, i.e., it is in fact a state of the neural network and we can apply  $F_S$  to it. In particular, we need to check that  $\text{bias} \in (\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ . Well,  $S \subseteq \text{State}_N$ , and so  $\text{bias} \in S$ . This means  $\text{bias} \notin S^{\mathbb{C}}$ , and by the contrapositive of best-inclusion,  $\text{bias} \notin (\text{best}_{<}(S^{\mathbb{C}}))$ , which is what I wanted to show. So the term  $F_S((\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}})$  is well-defined.

We are now ready to show that  $F_S((\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}) = (\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ . For the  $(\rightarrow)$  direction, suppose  $n \in F_S((\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}})$ . By construction of  $F_S$ , we have two cases:

- **Case 1.**  $n \in S$ . In this case, we trivially have  $n \notin \text{best}_{\prec}(S^{\mathbb{C}})$ , since *n* is not even in  $S^{\mathbb{C}}$ . And so  $n \in (\text{best}_{\prec}(S^{\mathbb{C}}))^{\mathbb{C}}$ .
- **Case 2.** At least one non-bias predecessor mEn is  $m \notin (\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ . So  $m \in \text{best}_{<}(S^{\mathbb{C}})$ . Since m is a predecessor of n (mEn), nRm (reverse it). But then since  $R = (<_{\text{reverse}})$ , m < n. So m is <-better than n (that is, m < n), which implies that n cannot be a best  $S^{\mathbb{C}}$ -element:  $n \notin \text{best}_{<}(S^{\mathbb{C}})$ . So  $n \in (\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ .

As for the  $(\leftarrow)$  direction, suppose contrapositively that  $n \notin F_S((\text{best}_{\prec}(S^{\mathbb{C}}))^{\mathbb{C}})$ . By construction, this means that  $n \notin S$  and for all *m* with  $m \prec n$  (including bias of course),  $m \in (\text{best}_{\prec}(S^{\mathbb{C}}))^{\mathbb{C}}$ .

[all predecessors m with mEn m of n (including bias of course) are in  $(\text{best}_{\leq}(S^{\mathbb{C}}))^{\mathbb{C}}$ .]

We already have  $n \in S^{\mathbb{C}}$ ; from here I'd like to show that *n* is the *best*  $S^{\mathbb{C}}$ -element. Suppose for contradiction that  $n \notin \text{best}_{<}(S^{\mathbb{C}})$ . By the smoothness condition (see Appendix [todo]) there must be some  $m \in S^{\mathbb{C}}$ , m < n that is the best, i.e.,  $m \in \text{best}_{<}(S^{\mathbb{C}})$ . Since m < n and  $R = (<_{\text{reverse}})$ , we have nRm, and so *m* is an *E*-predecessor of *n* (*mEn*). Note that we always have bias  $\in S$ , and in particular this means bias  $\notin \text{best}_{<}(S^{\mathbb{C}})$  (by best-inclusion, since  $\text{best}_{<}(S^{\mathbb{C}}) \subseteq S^{\mathbb{C}}$ ). So *m* cannot be the bias node. Complementing, we see that  $m \notin (\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ . In other words, some non-bias *E*-predecessor of *n* is not in  $(\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}}$ . By construction of  $F_S$ , this means  $n \in F_S((\text{best}_{<}(S^{\mathbb{C}}))^{\mathbb{C}})$ , which contradicts our initial hypothesis.

**Lemma 5.18.** Let  $\mathcal{M} = \langle W, R, \prec, V \rangle$  be a plausibility model such that  $R = (\prec_{\text{reverse}})$ . Let  $\mathcal{N}$  be given by the NAND construction above. For all  $S \subseteq \text{State}_{\mathcal{N}}$ :

• Reach<sub> $\mathcal{N}$ </sub>(*S*) = {*w* | there is some *u*  $\in$  *S* with *uRw*}

• Reach  $\mathcal{N}(S) = \{w \mid \text{there is some } u \in S \text{ with } wRu\}$ 

**Proof.** I will prove it for Reach; the proof for Reach<sup> $\downarrow$ </sup> is similar. Let  $w \in W \subseteq N$ .

 $w \in \operatorname{Reach}_{\mathcal{N}}(S) \quad \text{iff there exists } u \in S \text{ with an } E \text{-path from } w \text{ to } u$ iff there exists  $u \in S$  with an R-path from u to w (By constr of  $\mathcal{N}$ ) iff  $u \in S$  and uRw (Since R is transitive)

The following lemma says that the constructed net  $\mathcal{N}$  is in fact equivalent to  $\mathcal{M}$ , i.e., the two satisfy exactly the same formulas at exactly the same worlds over the modal language of  $\mathcal{L}_{\mathbf{C}}$ . Notice that the claim is only made for worlds  $w \in W$  and not for the bias node, since we cannot evaluate  $\parallel_{\mathbf{Plaus}}$  at the bias node.

**Lemma 5.19.** (C) Let  $\mathcal{M} = \langle W, \prec, V \rangle$  be a plausibility model such that  $R = (\prec_{\text{reverse}})$ , and let  $\mathcal{N}$  be given by the NAND construction above. For all modal formulas  $\varphi \in \mathcal{L}_{\mathbf{C}}$  and all  $w \in W$ ,

$$\mathcal{N}, w \Vdash_{\operatorname{Net}} \varphi \text{ iff } \mathcal{M}, w \Vdash_{\operatorname{Plaus}} \varphi$$

**Proof.** By induction on  $\varphi$ . I will show the key inductive cases  $\diamond^{\downarrow}\varphi$ , and  $\langle \mathbf{C} \rangle \varphi$ . The inductive case for  $\mathbf{E}\varphi$  is straightforward, and the  $\diamond\varphi$  case is similar to the case for  $\diamond^{\downarrow}\varphi$  (substitute Reach<sup> $\downarrow$ </sup> for Reach).

Case  $\Box \varphi$ .

$$\mathcal{N}, w \Vdash_{\mathbf{Net}} \Box \varphi \quad \text{iff} \quad w \in \mathsf{Reach}^{\downarrow}(\llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\mathsf{bias}\}) \qquad (By \text{ neural net semantics})$$

$$\text{iff} \quad \exists u \in \llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\mathsf{bias}\} \text{ with } wRu \qquad (By \text{ Lemma 5.18})$$

$$\text{iff} \quad [\mathsf{TODO}]$$

$$\text{iff} \quad \exists u \in W \text{ with } wRu \text{ and } \mathcal{M}, u \Vdash \varphi$$

$$\text{iff} \quad \mathcal{M}, w \Vdash_{\mathbf{Plaus}} \Box \varphi \qquad (By \text{ plausibility semantics})$$

Case  $\diamond^{\downarrow} \varphi$ .

$\mathcal{N}, w \Vdash_{\mathbf{Net}} \diamond^{\downarrow} \varphi$	iff	$w \in \operatorname{Reach}(\llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\operatorname{bias}\})$	(By neural net semantics)
		$\exists u \in \llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\text{bias}\} \text{ with } uRw$	(By Lemma 5.18)
	iff	$\exists u \in W$ with $uRw$ and $\mathcal{M}, u \Vdash \varphi$	
	iff	$\mathcal{M}, w \Vdash_{Plaus} \diamond^{\downarrow} \varphi$	(By plausibility semantics

**Case**  $(C)\varphi$ . In this case, we need to take care to be very explicit about what universe our complements are taken over. So I will write  $N \setminus S$  and  $W \setminus S$  in place of  $S^{\mathbb{C}}$ .

$$\begin{split} &\mathcal{N}, w \Vdash_{\mathbf{Net}} \langle \mathbf{C} \rangle \varphi \quad \text{iff} \quad w \in \mathsf{Clos}(\llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\mathsf{bias}\}) & (\mathsf{By neural net semantics}) \\ &\text{iff} \quad w \in N \setminus (\mathsf{best}_{<}(N \setminus (\llbracket \varphi \rrbracket_{\mathcal{N}} \cup \{\mathsf{bias}\}))) & (\mathsf{By Lemma 5.17}) \\ &\text{iff} \quad w \in N \setminus (\mathsf{best}_{<}((N \setminus \llbracket \varphi \rrbracket_{\mathcal{N}}) \cap (N \setminus \{\mathsf{bias}\}))) & (\mathsf{Grouping terms}) \\ &\text{iff} \quad w \in N \setminus (\mathsf{best}_{<}((N \setminus \llbracket \varphi \rrbracket_{\mathcal{N}}) \cap W)) & (\mathsf{Since } W = N \setminus \{\mathsf{bias}\}) \\ &\text{iff} \quad w \in N \setminus (\mathsf{best}_{<}(W \setminus \llbracket \varphi \rrbracket_{\mathcal{N}})) & (\mathsf{Since } w \in W \text{ and } w \neq \mathsf{bias}) \\ &\text{iff} \quad w \in W \setminus (\mathsf{best}_{<}(W \setminus \llbracket \varphi \rrbracket_{\mathcal{N}})) & (\mathsf{By inductive hypothesis}) \\ &\text{iff} \quad \mathcal{M}, w \Vdash_{\mathbf{Plaus}} \neg \mathbf{C} \neg \varphi & (\mathsf{By plausibility semantics}) \\ &\text{iff} \quad \mathcal{M}, w \Vdash_{\mathbf{Plaus}} \langle \mathbf{C} \rangle \varphi \end{array}$$

# 6 Reflections on Model Building and Interpretability

[Todo]

## Chapter 5

## **Dynamic Update in Neural Network Semantics**

## **1** Introduction

[I now have space to express these points more slowly, in detail.] The neural network semantics presented so far shows us how we can use neural networks as models for modal logic. Neural network inference can be expressed in this logic using  $\langle \mathbf{C} \rangle \varphi$ , which denotes the forward propagation of the signal  $[\![\varphi]\!]$  through the net. However, as discussed in the introduction, the mystery about neural networks is how their inference interacts with their *learning*. In this section, I will show how to extend these semantics to model learning and update in a neural net.

As previously mentioned, I formalize neural network update using the methodology of Dynamic Epistemic Logic. Our static operators  $\diamond$ ,  $\diamond^{\downarrow}$ , and  $\langle \mathbf{C} \rangle$  are interpreted by examining the state of the neural net. The DEL trick is to introduce a new "dynamic" operator [*P*] which *changes* the net in response to some observed formula *P*. First, we extend the language  $\mathcal{L}_{\mathbf{C}}$  to  $\mathcal{L}_{Update}$ , which includes these dynamic operators:

$$\varphi, \psi \coloneqq p \mid \neg \varphi \mid \varphi \land \psi \mid \mathbf{E}\varphi \mid \Diamond \varphi \mid \ \Diamond^{\downarrow}\varphi \mid \langle \mathbf{C} \rangle \varphi \mid [P]\varphi$$

Here,  $[P]\varphi$  reads "after the agent observes  $P, \varphi$  is true".

Let Update: Net × State  $\rightarrow$  Net be any function which takes a neural network, some state *S*, and "updates" the net somehow in response to *S*. We can interpret [*P*] as performing this update by adding the following line to the semantics:

$$\mathcal{N}, n \Vdash [P] \varphi$$
 iff Update $(\mathcal{N}, [P]), n \Vdash \varphi$ 

In other words, in order to evaluate  $[P]\varphi$ , we simply evaluate  $\varphi$  in the updated net Update $(\mathcal{N}, [P])$ .

From a DEL perspective, this is a standard move to make. But from a machine learning perspective, there are a couple caveats that I should mention. First, [P] does not model learning in the sense of "iterated update until convergence", but rather only models a single step of update. Second, we should think of [P] as modeling *unsupervised learning*—the model updates in response to an input P, but no "expected answer" y is given alongside P. It is an open problem to formalize supervised learning (in this machine learning sense) in DEL in a non-trivial way.

## 2 Hebbian Learning: A Simple Neural Network Update Policy

[I now have space to express these points more slowly, in detail.] So far, I've discussed learning and update in very general terms. For my thesis, I will model a simple update policy over neural networks: Hebbian learning. The point in starting with Hebbian learning is to get the details right on a simpler example before lifting these ideas to, say, gradient descent through backpropagation [58].

Hebb's classic learning rule [32] states that when two adjacent neurons are simultaneously and persistently active, the connection between them strengthens ("neurons that fire together wire together"). In contrast with backpropagation, Hebbian learning is errorless and unsupervised. Another key feature is that Hebbian update is local — the change in a weight  $\Delta W(m,n)$  depends only on the activation of the immediately adjacent neurons. For this reason, the Hebbian family of learning policies is often considered more biologically plausible than backpropagation.

There are many variations of Hebbian learning, but I will only consider the most basic form of Hebb's rule:  $\Delta W(m,n) = \eta x_m x_n$ , where  $\eta$  is the learning rate and  $x_m, x_n$  are the outputs of adjacent neurons *m* and *n*. This is the *unstable* variation of Hebb's rule; repeatedly applying the rule will make the weights arbitrarily large. I will not consider stabilizing variants such as Oja's rule [54].

Single-Step Hebbian Update. First, consider what happens in a single step of Hebbian update. Given a net  $\mathcal{N}$  and a state *S*, we first propagate *S* forward through  $\mathcal{N}$ . Any edges that are involved in this propagated activation pattern Clos(*S*) simply have their weights strengthened. Formally,

**Definition 2.1.** Let Hebb: Net  $\times$  State  $\rightarrow$  Net be given by

 $\mathsf{Hebb}(\langle N, \mathsf{bias}, E, W, A, \eta, V \rangle, S) = \langle N, \mathsf{bias}, E, W', A, \eta, V \rangle$ 

where  $W'(m,n) = W(m,n) + \eta \cdot \chi_{\operatorname{Clos}(S)}(m) \cdot \chi_{\operatorname{Clos}(S)}(n)$ .

Note that Hebb does not affect the edges, activation function, or evaluation of propositions. This means the resulting net is still binary, and closures Clos(S) still exist and are unique. Therefore Hebb is well-defined. This also means that Hebb does not affect the Reach or Reach<sup> $\downarrow$ </sup> operators.

**Proposition 2.2.** Reach<sub>Hebb</sub> $(\mathcal{N},A)(B) = \operatorname{Reach}_{\mathcal{N}}(B)$ 

**Proof.** A single step of Hebbian update  $\text{Hebb}(\mathcal{N}, A)$  doesn't change the edge relation *E* of the graph. So if  $n \in N$ , any path from  $m \in B$  to *n* in  $\text{Hebb}(\mathcal{N}, A)$  is the same path in  $\mathcal{N}$ .  $\Box$ 

And similarly:

**Proposition 2.3.** Reach  $\downarrow_{\mathsf{Hebb}(\mathcal{N},A)}(B) = \mathsf{Reach}_{\mathcal{N}}^{\downarrow}(B)$ 

The following is easy to see [I now have space to explain] (since  $\eta \ge 0$ ).

**Proposition 2.4.** Let  $m, n \in N$ . We have:

- $W_{\mathcal{N}}(m,n) \leq W_{\mathsf{Hebb}(\mathcal{N},S)}(m,n)$
- If either  $m \notin \text{Clos}(S)$  or  $n \notin \text{Clos}(S)$ , then  $W_{\text{Hebb}(\mathcal{N},S)}(m,n) = W_{\mathcal{N}}(m,n)$ .

**Proof.** For the first part, observe:

$$\begin{split} W_{\mathcal{N}}(m,n) &\leq W_{\mathcal{N}}(m,n) + \eta \qquad (\text{since } \eta \geq 0) \\ &\leq W_{\mathcal{N}}(m,n) + \eta \cdot \chi_{\mathsf{Clos}(S)}(m) \cdot \chi_{\mathsf{Clos}(S)}(n) \quad (\text{since for all } S,n, \chi_{S}(n) \geq 0) \\ &= W_{\mathsf{Hebb}(\mathcal{N},S)}(m,n) \end{split}$$

As for the second part, if either  $m \notin Clos(S)$  or  $n \notin Clos(S)$ , then by definition of Hebb,

$$W_{\mathsf{Hebb}(\mathcal{N},S)}(m,n)$$

$$= W_{\mathcal{N}}(m,n) + \eta \cdot \chi_{\mathsf{Clos}(S)}(m) \cdot \chi_{\mathsf{Clos}(S)}(n)$$

$$= W_{\mathcal{N}}(m,n) + \eta \cdot 0$$

$$= W_{\mathcal{N}}(m,n) + \eta$$

$$= W_{\mathcal{N}}(m,n)$$

**Iterated Hebbian Update.** In addition to the single-step Hebb operator, in my thesis work I have also modelled *iterated* Hebbian update Hebb\*. The idea is this: what happens when we propagate a signal *S* through the net, and then *repeatedly* strengthen the weights of the edges that are involved? Recall that our single-step Hebb is unstable; if we repeat Hebb on a single input state *S*, the net's

weights within Clos(S) will be so high that *any* activation pattern that makes contact with Clos(S) will "rip through" it entirely. Repeating Hebb on *S* further will not change the Clos(S)-structure, i.e., the update has reached a fixed point. Hebb\* returns the net at this fixed point.

Instead of reasoning abstractly about this fixed point, I formalize it by explicitly defining the number of iterations iter needed to reach it. The idea is to set iter to be so high, all updated weights W'(m,n) overpower any negative weights that would otherwise cancel their effect. The following definitions might look like black magic, but they are set up to capture this intuition (I verified in Lean that this is the right choice for iter, see [35]).

**Definition 2.5.** Let  $\mathcal{N}$  be a net,  $n \in N$ , and let  $m_1, \ldots, m_k$  list the predecessors of n. The *negative weight score* of n is the sum of all the negative weights of n's predecessors, i.e.,

$$\mathsf{nws}(n) = \sum_{m \in \mathsf{preds}(n)} \begin{cases} W(m,n) & \text{if } W(m,n) < 0\\ 0 & \text{otherwise} \end{cases}$$

**Definition 2.6.** The *minimum negative weight score* is simply

$$\mathsf{mnws} = \min_{n \in N} \mathsf{nws}(n)$$

**Proposition 2.7.** For all  $S \in \text{State}$ ,  $m, n \in N$ , we have mnws  $\leq W(m, n) \cdot \chi_S(m)$ .

**Proof.** Let *m*, *n* be any nodes in *N*. We have:

$$\begin{aligned} & \text{mnws} \leq \text{nws}(n) \\ & = \sum_{m \in \text{preds}(n)} \begin{cases} W(m,n) & \text{if } W(m,n) < 0 \\ 0 & \text{otherwise} \end{cases} & \text{(by definition)} \\ & = \sum_{m \in \text{preds}(n)} \begin{cases} W(m,n) \cdot \chi_{S}(m) & \text{if } W(m,n) < 0 \\ 0 & \text{otherwise} \end{cases} & \text{(since each } W(m,n) < 0 \\ & \text{and } \chi_{S}(m) \in \{0,1\} ) \\ & \leq W(m,n) \cdot \chi_{S}(m) \end{cases} & \text{(the sum of negative terms is } \leq \\ & \text{any particular term} ) \end{aligned}$$

**Definition 2.8.** Recall that the activation function *A* is nonzero, i.e. there is some  $t \in \mathbb{Q}$  such that A(t) = 1. We set the number of iterations iter to be exactly

$$iter = \begin{cases} \left\lceil \frac{t - |N| \cdot mnws}{\eta} \right\rceil & \text{if } \ge 1\\ 1 & \text{otherwise} \end{cases}$$

Note that iter will always be a positive integer, and so iterating iter times is well-defined. This choice for iter may seem opaque, but we will see in Lemma [which] why it guarantees that the updated weights overpower competing edge weights.

**Definition 2.9.** Let  $Hebb^*$ : Net  $\times$  State  $\rightarrow$  Net be given by

 $\mathsf{Hebb}^*(\langle N, \mathsf{bias}, E, W, A, \eta, V \rangle, S) = \langle N, \mathsf{bias}, E, W', A, \eta, V \rangle$ 

where  $W'(m,n) = W(m,n) + \text{iter} \cdot \eta \cdot \chi_{\text{Clos}(S)}(m) \cdot \chi_{\text{Clos}(S)}(n)$ .

As with Hebb, Hebb<sup>\*</sup> does not affect the edges, activation function, or evaluation of propositions. Therefore Hebb<sup>\*</sup> is well-defined. This also means that Hebb<sup>\*</sup> does not affect the Reach or Reach<sup> $\downarrow$ </sup> operators.

**Proposition 2.10.** Reach<sub>Hebb\*( $\mathcal{N}$ ,A)(B) = Reach<sub> $\mathcal{N}$ </sub>(B)</sub>

**Proposition 2.11.** Reach  $\downarrow_{\text{Hebb}^*(\mathcal{N},A)}(B) = \text{Reach}_{\mathcal{N}}^{\downarrow}(B)$ 

Similar to Proposition [todo], we have the following:

**Proposition 2.12.** Let  $m, n \in N$ . We have:

- $W_{\mathcal{N}}(m,n) \leq W_{\mathsf{Hebb}^*(\mathcal{N},S)}(m,n)$
- If either  $m \notin \text{Clos}(S)$  or  $n \notin \text{Clos}(S)$ , then  $W_{\text{Hebb}^*(\mathcal{N},S)}(m,n) = W_{\mathcal{N}}(m,n)$

Proof. []

The following fact about Hebb\* is the most important. It is a formal expression of our statement before: Updated weights  $W_{\text{Hebb}^*(\mathcal{N},A)}(B)$  are so high that if *m* is active in Hebb\* then *n* must be as well.

**Lemma 2.13.** (C) Let  $A, B \in \text{State}, m, n \in N$ . If  $m \in \text{preds}, m, n \in \text{Clos}(A)$ , and  $m \in \text{Clos}(B)$ , then

$$A(\sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m_i,n) \cdot \chi_{\operatorname{Clos}(B)}(m)) = 1$$

(Take care to notice the different subscripts for *W* and  $\chi$ !)

**Proof.** *A* is a binary step function, which in particular means it is binary, has a threshold, some  $t \in \mathbb{Q}$  with A(t) = 1, and is nondecreasing. Since *A* is nondecreasing, it's enough for us to show

$$t \leq \sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m_i,n) \cdot \chi_{\operatorname{Clos}(B)}(m)$$
Well, we have

$$\begin{split} \sum_{m_i \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m_i,n) \cdot \chi_{\operatorname{Clos}(B)}(m_i) \\ &= \sum_{m_i \in \operatorname{preds}(n), \text{ and } m_i \neq m} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m_i,n) \cdot \chi_{\operatorname{Clos}(B)}(m_i) \\ &+ W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m,n) \cdot \chi_{\operatorname{Clos}_{\operatorname{Hebb}^*(\mathcal{N},A)}(B)}(m) \\ &\geq (|N|-1) \cdot \operatorname{mnws} + W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m,n) \cdot \chi_{\operatorname{Clos}(B)}(m) \\ & \text{(by Proposition [todo], since we are adding |N|-1 terms)} \\ &= (|N|-1) \cdot \operatorname{mnws} + W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m,n) \cdot 1 \\ & \text{(since } m \in \operatorname{Clos}(B)) \\ &= (|N|-1) \cdot \operatorname{mnws} + W_{\mathcal{N}}(m,n) + \operatorname{iter} \cdot \chi_{\operatorname{Clos}(A)}(m) \cdot \chi_{\operatorname{Clos}(A)}(n) \\ & \text{(by definition of Hebb}^*) \\ &= (|N|-1) \cdot \operatorname{mnws} + W_{\mathcal{N}}(m,n) + \operatorname{iter} \\ & \text{(since } m, n \in \operatorname{Clos}(A)) \\ &\geq (|N|-1) \cdot \operatorname{mnws} + \operatorname{mnws} + \operatorname{iter} \cdot \eta \cdot \chi_{\operatorname{Clos}(A)}(m) \cdot \chi_{\operatorname{Clos}(A)}(n) \\ & \text{(the sum of negative weights is } \leq any particular weight)} \\ &= |N| \cdot \operatorname{mnws} + \operatorname{iter} \cdot \eta \\ & \text{(grouping like terms)} \end{split}$$

So at this point we need to show:

$$t \leq |N| \cdot \text{mnws} + \text{iter} \cdot \eta$$

Rearranging this to solve for iter, it suffices to show:

$$\frac{t - |N| \cdot \text{mnws}}{\eta} \le \text{iter}$$

But we defined iter to be exactly the integer ceiling of this expression on the left (and 1 if the expression on the left is negative)!  $\Box$ 

## **3** Properties of Hebb and Hebb\*

We have the following algebraic properties for Hebb.

One worry we might have is that, in each iteration, we always update by Clos(S) in the *original* net. But it turns out that this Clos(S) doesn't change with each iteration, i.e.

**Proposition 3.2.** (C)  $Clos_{Hebb(\mathcal{N},S)}(S) = Clos_{\mathcal{N}}(S)$ 

**Proof.** Let  $F_S$  be the state transition function for  $\mathcal{N}$  under S, and  $F'_S$  be the state transition function for Hebb $(\mathcal{N}, S)$  under S. First, since  $Clos_{Hebb}(\mathcal{N},S)(S)$  is a fixed point under S in Hebb $(\mathcal{N}, S)$ , we have  $F'_S(Clos_{Hebb}(\mathcal{N},S)(S)) = Clos_{Hebb}(\mathcal{N},S)(S)$ . But I will show that  $Clos_{\mathcal{N}}(S)$  is *also* a fixed point under S in Hebb $(\mathcal{N}, S)$ , i.e.  $F'_S(Clos_{\mathcal{N}}(S)) = Clos_{\mathcal{N}}(S)$ . Since we assumed there is a *unique* fixed point under S in Hebb $(\mathcal{N}, S)$ , it will follow that these two states must be the same. In other words,  $Clos_{Hebb}(\mathcal{N},S)(S) = Clos_{\mathcal{N}}(S)$ .

For the  $(\leftarrow)$  direction, suppose  $n \in \operatorname{Clos}_{\mathcal{N}}(S)$ . Since  $\operatorname{Clos}_{\mathcal{N}}(S)$  is a fixed point under S in  $\mathcal{N}$ ,  $\operatorname{Clos}_{\mathcal{N}}(S) = F_S(\operatorname{Clos}_{\mathcal{N}}(S))$ . By definition of F, either  $n \in S$  (in which case we are done), or n is activated by its predecessors m in  $\operatorname{Clos}_{\mathcal{N}}(S)$  over  $\mathcal{N}$ , i.e.

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\mathcal{N}}(m,n) \cdot \chi_{\operatorname{Clos}_{\mathcal{N}}(S)}(m)\right) = 1$$

By the first part of Proposition [todo], each  $W_{\mathcal{N}}(m, n) \leq W_{\mathsf{Hebb}(\mathcal{N}, S)}(m, n)$ . So the inner sum using the former is  $\leq$  the inner sum using the latter. Since *A* is nondecreasing, we have

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}(\mathcal{N},S)}(m,n) \cdot \chi_{\operatorname{Clos}_{\mathcal{N}}(S)}(m)\right) = 1$$

But this implies that  $n \in F'_S(\text{Clos}_{\mathcal{N}}(S))$ .

As for the  $(\rightarrow)$  direction, suppose  $n \in F'_S(\text{Clos}_{\mathcal{N}}(S))$ . By definition of F', either  $n \in S$  (in which case we are done), or *n* is activated by its predecessors *m* in  $\text{Clos}_{\mathcal{N}}(S)$  over  $\text{Hebb}(\mathcal{N}, S)$ , i.e.

$$A\left(\sum_{m \in \text{preds}(n)} W_{\text{Hebb}(\mathcal{N},S)}(m,n) \cdot \chi_{\text{Clos}_{\mathcal{N}}(S)}(m)\right) = 1$$

Suppose for contradiction that  $n \notin Clos_{\mathcal{N}}(S)$ . By the second part of Proposition [todo], each  $W_{\text{Hebb}(\mathcal{N},S)}(m,n) = W_{\mathcal{N}}(m,n)$ , and so we have

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\mathcal{N}}(m,n) \cdot \chi_{\operatorname{Clos}_{\mathcal{N}}(S)}(m)\right) = 1$$

but this implies that  $n \in F_S(\text{Clos}_{\mathcal{N}}(S)) = \text{Clos}_{\mathcal{N}}(S)$ , which contradicts  $n \notin \text{Clos}_{\mathcal{N}}(S)$ . and so Hebb\* is equivalent to repeatedly applying Hebb until we reach a fixed point [35]. [Elaborate on this point, it's said a little too quickly for the reader to internalize it! (maybe a picture would help?)] We have the following algebraic properties for Hebb<sup>\*</sup>. Before proving these, I'll give some intuition for what these properties say about Hebb<sup>\*</sup>. [(1) is just used to show (2)] Part (2) expresses a lower bound for  $Clos_{Hebb^*(\mathcal{N},A)}(B)$ , whereas (3) gives an upper bound within Clos(A). [There's got to be a better way to explain what these mean... this isn't clear or intuitive at all.]

**Proposition 3.3.** (C)  $Clos_{Hebb^*(\mathcal{N},S)}(S) = Clos_{\mathcal{N}}(S)$ 

**Proof.** The proof is the same as the proof for Hebb (see Proposition [todo]). In place of Propositions [todo] and [todo], we instead apply Propositions [todo] and [todo], respectively.

**Theorem 3.4.** (C) Let  $\mathcal{N} \in \mathbf{Net}$ , and suppose  $\mathcal{N}$  is fully connected. For all  $A, B \in \mathbf{State}$ ,

$$\mathsf{Clos}_{\mathsf{Hebb}^*(\mathcal{N},A)}(B) = \mathsf{Clos}(B \cup (\mathsf{Clos}(A) \cap \mathsf{Reach}(\mathsf{Clos}(A) \cap \mathsf{Clos}(B))))$$

**Proof.** Let  $F_B$  be the state transition function for  $\mathcal{N}$  under B, and  $F_B^*$  be the state transition function for Hebb<sup>\*</sup>( $\mathcal{N}$ , A) under B. For notational convenience, let T be the set inside Clos on the right-hand side, i.e.

$$T = B \cup (Clos(A) \cap Reach(Clos(A) \cap Clos(B)))$$

This proof follows the major plot beats of the proof for Theorem **[TODO]**. First, since  $Clos_{Hebb^*(\mathcal{N},A)}(B)$  is a fixed point under *B* in Hebb<sup>\*</sup>( $\mathcal{N}, A$ ), we have  $F_B^*(Clos_{Hebb^*(\mathcal{N},A)}(B)) = Clos_{Hebb^*(\mathcal{N},A)}(B)$ . But I will show that Clos(T) is *also* a fixed point under *B* in Hebb( $\mathcal{N}, A$ ), i.e.  $F_B^*(Clos(T)) = Clos(T)$ . Since we postulated that there is a *unique* fixed point under *B* in Hebb<sup>\*</sup>( $\mathcal{N}, A$ ), it will follow that these two states must be the same:  $Clos_{Hebb^*(\mathcal{N},A)}(B) = Clos(T)$ .

Let's show that  $F_B^*(\operatorname{Clos}(T)) = \operatorname{Clos}(T)$ . For the  $(\leftarrow)$  direction, suppose  $n \in \operatorname{Clos}(T)$ . Since  $\operatorname{Clos}(T)$  is a fixed point under T in  $\mathcal{N}$ ,  $\operatorname{Clos}(T) = F_T(\operatorname{Clos}(T))$ . By definition of F, we have two cases:

**Case 1.**  $n \in T$ , i.e.  $n \in B \cup (\operatorname{Clos}(A) \cap \operatorname{Reach}(\operatorname{Clos}(A) \cap \operatorname{Clos}(B)))$ . If  $n \in B$ , then we're done by the definition of  $F_B^*$  (the next state includes all nodes in *B*). Otherwise, we have  $n \in \operatorname{Clos}(A)$ and a path from some  $m \in \operatorname{Clos}(A) \cap \operatorname{Clos}(B)$  to n in  $\mathcal{N}$ . Since  $\mathcal{N}$  is fully connected, mis in fact a predecessor of n. Moreover,  $m \in \operatorname{Clos}(T)$ , since  $m \in \operatorname{Clos}(A) \cap \operatorname{Clos}(B)$  and by inclusion of Reach and Clos. So we have  $m \in \operatorname{preds}(n), m, n \in \operatorname{Clos}(A)$ , and  $m \in \operatorname{Clos}(T)$ . But these are exactly the conditions of Lemma [todo]! This means we have

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^{*}(\mathcal{N},A)}(m,n) \cdot \chi_{\operatorname{Clos}(T)}(m)\right) = 1$$

which implies that  $n \in F_B^*(Clos(T))$ .

**Case 2.** *n* is activated by its predecessors *m* in Clos(T) over  $\mathcal{N}$ , i.e.

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\mathcal{N}}(m,n) \cdot \chi_{\operatorname{Clos}(T)}(m)\right) = 1$$

By the first part of Proposition [todo], each  $W_{\mathcal{N}}(m,n) \leq W_{\mathsf{Hebb}^*(\mathcal{N},S)}(m,n)$ . So the inner sum using the former is  $\leq$  the inner sum using the latter. Since *A* is nondecreasing, we have

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m,n) \cdot \chi_{\operatorname{Clos}(T)}(m)\right) = 1$$

But this immediately implies that  $n \in F_B^*(Clos(T))$ .

As for the  $(\rightarrow)$  direction, suppose  $n \in F_B^*(Clos(T))$ . By definition of  $F^*$ , we have two cases:

**Case 1.**  $n \in B$ . So  $n \in B \cup (Clos(A) \cap Reach(Clos(A) \cap Clos(B))) = T$ . By inclusion of Clos, we have  $n \in Clos(T)$ .

**Case 2.** *n* is activated by its predecessors *m* in Clos(T) over  $Hebb^*(\mathcal{N}, A)$ , i.e.

$$A\left(\sum_{m \in \operatorname{preds}(n)} W_{\operatorname{Hebb}^*(\mathcal{N},A)}(m,n) \cdot \chi_{\operatorname{Clos}(T)}(m)\right) = 1$$

Suppose for contradiction that  $n \notin Clos(T)$ . By the second part of Proposition [todo], each  $W_{Hebb^*(\mathcal{N},A)}(m,n) = W_{\mathcal{N}}(m,n)$ , and so we have

$$A\left(\sum_{m \in \text{preds}(n)} W_{\mathcal{N}}(m, n) \cdot \chi_{\text{Clos}(T)}(m)\right) = 1$$

but this implies that  $n \in F_T(Clos(T)) = Clos(T)$ , which contradicts  $n \notin Clos(T)$ . So we must have  $n \in Clos(T)$ .

**Corollary 3.5.** If  $Clos(A) \cap Clos(B) = \emptyset$ , then  $Clos_{Hebb^*(\mathcal{N},A)}(B) = Clos(B)$ .

**Proof.** Suppose  $Clos(A) \cap Clos(B) = \emptyset$ . We have

$$\begin{aligned} \mathsf{Clos}_{\mathsf{Hebb}^*(\mathcal{N},A)}(B) &= \mathsf{Clos}(B \cup (\mathsf{Clos}(A) \cap \mathsf{Reach}(\mathsf{Clos}(A) \cap \mathsf{Clos}(B)))) \\ &= \mathsf{Clos}(B \cup (\mathsf{Clos}(A) \cap \mathsf{Reach}(\emptyset))) \\ &= \mathsf{Clos}(B \cup \emptyset) \\ &= \mathsf{Clos}(B) \end{aligned}$$

### 4 Neural Network Semantics for Hebbian Update

**[TODO]** Give official languages  $\mathcal{L}_{\text{Hebb}}$  and  $\mathcal{L}_{\text{Hebb}*}$  for the logics of Hebb and Hebb\*, i.e. using operators  $[P]_{\text{Hebb}}\varphi$  and  $[P]_{\text{Hebb}*}\varphi$  and their semantics!!! And also just say that the definitions  $\mathcal{N}\models_{\text{Hebb}}\varphi$  (same for Hebb\*),  $\Gamma\models_{\text{Hebb}}\varphi$  (same for Hebb\*) are what you'd expect.

#### **5** A Complete Logic of Iterated Hebbian Update

- In the previous section, I gave sound axioms for Hebb\*. It turns out those axioms are nearly complete! I'll show this here, and give the complete set of "reduction axioms"
- I already proved soundness, but I have to do it again for these new axioms!

#### **Definition 5.1.** [Reduction axioms for Hebb\*]

**Theorem 5.2.** (C) These reduction axioms are sound; for all  $\Gamma \subseteq \mathcal{L}^*$  and  $\varphi \in \mathcal{L}^*$ , if  $\Gamma \vdash \varphi$  then  $\Gamma \models \varphi$ .

Proof. [todo]			

Proof. [todo]

#### 5.1 Example: Verifying a Neural Networks Behavior After Learning

do example using single-step Hebbian learning, since iterated is a bit more abstract...

**Definition 5.3.** [Term rewriting translation system]

**Proposition 5.4.** (C) Each tr( $\varphi$ ) is update-operator free [todo, fix this statement!]

Proof. [todo]

**Proposition 5.5.** (C) Each tr( $\varphi$ ) actually terminates [todo, fix this statement!]

**Proposition 5.6.** (C)  $\vdash \varphi \leftrightarrow tr(\varphi)$  [todo, fix this statement, which  $\vdash$  is this?]

**Proof.** [todo]

**Theorem 5.7.** ( Model Building for  $\mathcal{L}^*$ ) For all consistent  $\Gamma \subseteq \mathcal{L}^*$ , there is finite  $\mathcal{N}$  such that  $\mathcal{N} \models \Gamma$ .

**Proof.** [todo]

**Theorem 5.8.** (C) Completeness for  $\mathcal{L}^*$ ) For all consistent  $\Gamma \subseteq \mathcal{L}^*$  and all formulas  $\varphi \in \mathcal{L}^*$ ,

if 
$$\Gamma \models \varphi$$
 then  $\Gamma \vdash \varphi$ 

**Theorem 5.9.** (C) For all  $\Gamma^* \subseteq \mathcal{L}^*$ , there is  $\mathcal{N}$  such that  $\mathcal{N} \models \Gamma^*$ .

**Proof.** Let  $\Gamma^* \subseteq \mathcal{L}^*$ . As outlined in the paper, our plan is to define rewrite rules based on our reduction axioms that "translate away" all of the dynamic formulas  $\langle \varphi \rangle_{\text{Hebb}} \psi$  in  $\Gamma^*$ , resulting in  $\Gamma^{tr} \subseteq \mathcal{L}$ . By our assumption, we have a net  $\mathcal{N} \models \Gamma^{tr}$ , and we show that this very same net  $\mathcal{N} \models \Gamma^*$ .

It's easy to see intuitively how this translation should go. For example, given the formula

$$\langle p \rangle_{\mathsf{Hebb}} (\langle p \rangle_{\mathsf{Hebb}} \langle \mathbf{B} \rangle \land \diamondsuit) \in \Gamma^*$$

we would recursively apply our reduction axioms, pushing  $\langle p \rangle_{\text{Hebb}}$  further into the expression until we can eliminate the propositional cases  $\langle p \rangle_{\text{Hebb}} q$ .

We define the term-rewriting system that does the translation  $\tau(\varphi)$  for all  $\varphi$  as follows.

- $\tau(p) = p$
- $\tau(\neg \varphi) = \neg \tau(\varphi)$
- $\tau(\varphi \land \psi) = \tau(\varphi) \land \tau(\psi)$
- $\tau(\mathbf{K}\varphi) = box(\tau(\varphi))$
- $tr(\langle \varphi \rangle_{\text{Hebb}}p) = tr(p)$
- $tr(\langle \varphi \rangle_{\text{Hebb}} \neg \psi) = tr(\neg \langle \varphi \rangle_{\text{Hebb}} \psi)$
- $tr(\langle \varphi \rangle_{\text{Hebb}}(\psi \land \rho)) = tr(\langle \varphi \rangle_{\text{Hebb}}\psi \land \langle \varphi \rangle_{\text{Hebb}}\rho)$
- $tr(\langle \varphi \rangle_{\text{Hebb}} \diamond) = tr(\diamond)$

- $tr(\langle \varphi \rangle_{\text{Hebb}} \langle \mathbf{B} \rangle) = tr(\langle \mathbf{B} \rangle)$
- $tr(\langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho) = tr(\langle \varphi \rangle_{\text{Hebb}} (tr(\langle \psi \rangle_{\text{Hebb}} \rho)))$

Formally, the term-rewriting system takes a formula  $\varphi$  and recursively applies these equational rules to  $\varphi$  (from left-to-right). We just need to check that

- 1. For all  $\psi$ ,  $tr(\psi)$  is update-operator-free
- 2. This term rewriting actually terminates

The work involved in showing termination is long and tedious. The usual approach is to define a measure on formulas  $c(\varphi)$  that *decreases* with each application of our reduction axioms (from left-to-right). In particular, we need *c* to satisfy

- If  $\psi$  is a subexpression of  $\varphi$ ,  $c(\varphi) > c(\psi)$
- $\bullet \quad c(\langle \varphi \rangle_{\mathsf{Hebb}} p) \! > \! c(p)$
- $c(\langle \varphi \rangle_{\text{Hebb}} \neg \psi) > c(\neg \langle \varphi \rangle_{\text{Hebb}} \psi)$
- $c(\langle \varphi \rangle_{\mathsf{Hebb}}(\psi \land \rho)) > c(\langle \varphi \rangle_{\mathsf{Hebb}}\psi \land \langle \varphi \rangle_{\mathsf{Hebb}}\rho)$
- $c(\langle \varphi \rangle_{\text{Hebb}} \diamondsuit) > c(\diamondsuit)$
- $c(\langle \varphi \rangle_{\text{Hebb}} \langle \mathbf{B} \rangle) > c(\langle \mathbf{B} \rangle)$
- $c(\langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho) > c(\langle \varphi \rangle_{\text{Hebb}} (tr(\langle \psi \rangle_{\text{Hebb}} \rho)))$

But coming up with a measure c that works is tricky, and is dependent on the specific reduction axioms. For the gritty details involved in coming up with this measure, as well as proving termination for the term rewriting system, see [12].

From here, we assume we have this measure c. We now have two things left to show:

**Proposition.** (C) For all  $\varphi \in \Gamma^*$ , we have  $\vdash \varphi \leftrightarrow \tau(\varphi)$ 

**Proof.** By induction on  $c(\varphi)$ .

**Base Step.** If  $\varphi$  is a proposition *p*, then we (trivially) have  $\vdash p \leftrightarrow p$ .

**Inductive Step.** We consider each possible inductive case, and suppose the claim holds for formulas  $\psi$  with smaller  $c(\psi)$ . The  $\neg \phi$ ,  $\phi \land \psi$ , **K**, and **B** cases all follow from applying the translation, and then applying inductive hypothesis on the subexpression that results from this.

Here are the rest of the cases. Notice that we apply the inductive hypothesis to terms whose c-cost is smaller (this is why we needed the decreasing properties of c before).

 $\langle \varphi \rangle_{\text{Hebb}} p$  case. We have

$$tr(\langle \varphi \rangle_{\mathsf{Hebb}} p) = tr(p) = p$$

and so we need to show that

$$\vdash \langle \varphi \rangle_{\mathsf{Hebb}} p \leftrightarrow p$$

but this holds by our propositional reduction axiom.

 $\langle \varphi \rangle_{\text{Hebb}} \neg \psi$  case. We have:

$$\begin{array}{l} \neg \langle \varphi \rangle_{\mathsf{Hebb}} \neg \psi \\ \leftrightarrow \neg \langle \varphi \rangle_{\mathsf{Hebb}} \psi & \text{(by the reduction axiom)} \\ \leftrightarrow tr(\neg \langle \varphi \rangle_{\mathsf{Hebb}} \psi) & \text{(inductive hypothesis)} \\ = tr(\langle \varphi \rangle_{\mathsf{Hebb}} \neg \psi) & \text{(by our translation)} \end{array}$$

 $\langle \varphi \rangle_{\text{Hebb}} \psi \wedge \rho$  case.. We have:

$$\begin{split} \vdash \langle \varphi \rangle_{\text{Hebb}}(\psi \land \rho) \\ \leftrightarrow \langle \varphi \rangle_{\text{Hebb}} \psi \land \langle \varphi \rangle_{\text{Hebb}} \rho & \text{(by the reduction axiom)} \\ \leftrightarrow tr(\langle \varphi \rangle_{\text{Hebb}} \psi \land \langle \varphi \rangle_{\text{Hebb}} \rho) & \text{(inductive hypothesis)} \\ = tr(\langle \varphi \rangle_{\text{Hebb}}(\psi \land \rho)) & \text{(by our translation)} \end{split}$$

 $\langle \varphi \rangle_{\text{Hebb}}$ K case.. We have:

$$\begin{split} \vdash \langle \varphi \rangle_{\mathsf{Hebb}} \mathbf{K} & \\ \leftrightarrow \mathbf{K} & \text{(by the reduction axiom)} \\ \leftrightarrow tr(\mathbf{K}) & \text{(inductive hypothesis)} \\ = tr(\langle \varphi \rangle_{\mathsf{Hebb}} \mathbf{K}) & \text{(by our translation)} \end{split}$$

 $\langle \varphi \rangle_{\text{Hebb}}$ B case.. We have:

$$\begin{array}{l} -\langle \varphi \rangle_{\mathsf{Hebb}} \mathbf{B} \\ \leftrightarrow \langle \mathbf{B} \rangle \\ (by \text{ the reduction axiom}) \\ \leftrightarrow tr(\langle \mathbf{B} \rangle) \\ (inductive hypothesis) \\ = tr(\langle \varphi \rangle_{\mathsf{Hebb}} \mathbf{B}) \\ (by \text{ our translation}) \end{array}$$

 $\{\varphi\}_{\text{Hebb}}\{\psi\}_{\text{Hebb}}\rho$  case.. This case is more interesting. First, notice our translation for this case:

$$tr(\langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho) = tr(\langle \varphi \rangle_{\text{Hebb}} tr(\langle \psi \rangle_{\text{Hebb}} \rho))$$

That is, we translate the inner expression first, then translate the outer expression. This inner  $tr(\langle \psi \rangle_{\text{Hebb}} \rho)$  is equivalent to some update-operator-free formula  $\chi$ :

$$\vdash \chi \leftrightarrow tr(\langle \psi \rangle_{\text{Hebb}} \rho) \leftrightarrow \langle \psi \rangle_{\text{Hebb}} \rho \tag{5.1}$$

(This last equivalence follows from our inductive hypothesis, which we can apply because  $\langle \psi \rangle_{\text{Hebb}} \rho$  is a subexpression of  $\langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho$ .)

What about  $tr(\langle \varphi \rangle_{\text{Hebb}} \chi)$ ? Well, since  $\chi$  is update-operator-free, this reduces to our previous inductive cases. So we have

$$\vdash tr(\langle \varphi \rangle_{\text{Hebb}} \chi) \leftrightarrow \langle \varphi \rangle_{\text{Hebb}} \chi \tag{5.2}$$

Putting this all together, we have:

$$\begin{array}{l} \langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho \\ \leftrightarrow \langle \varphi \rangle_{\text{Hebb}} \chi & (by (5.1)) \\ \leftrightarrow tr(\langle \varphi \rangle_{\text{Hebb}} \chi) & (by (5.2)) \\ \leftrightarrow tr(\langle \varphi \rangle_{\text{Hebb}} (tr(\langle \psi \rangle_{\text{Hebb}} \rho))) & (by (5.1)) \\ \leftrightarrow tr(\langle \varphi \rangle_{\text{Hebb}} \langle \psi \rangle_{\text{Hebb}} \rho) & (by our translation) \end{array}$$

**Theorem 5.10.** (Completeness) The logic of Hebbian Learning is completely axiomatized by the base axioms [ref!], along with the above reduction axioms. That is, for all consistent  $\Gamma^* \subseteq \mathcal{L}^*$ , if  $\Gamma^* \models \varphi$  then  $\Gamma^* \vdash \varphi$ .

**Proof.** Since our language  $\mathcal{L}^*$  has negation, completeness follows from model building in the usual way; this proof is entirely standard. Suppose contrapositively that  $\Gamma^* \not\models \varphi$ . It follows that  $\Gamma^* \vdash \neg \varphi$ . So  $\Gamma^* \cup \{\neg \varphi\}$  is consistent, and by Theorem [todo—need to modify the construction to make sure the net is fully connected. Quote: "But remember that our nets are also fully connected! So we need to modify the model construction from [38] by introducing a zero weight edge between every pair of previously unconnected nodes. Note that this change does not affect the \$\Prop\$-structure of the net."], we have  $\mathcal{N} \in \mathbb{N}$  such that  $\mathcal{N} \models \Gamma^* \cup \{\neg \varphi\}$ . But then  $\mathcal{N} \models \Gamma^*$  yet  $\mathcal{N} \not\models \varphi$ , which is what we wanted to show.

#### 5.2 Example: Building a Neural Network with Learning Constraints

## 6 Reflections on Neural Network Alignment

[Prove soundness of axioms for both single-step and iterated Hebbian update!]

## **Chapter 6**

## **Neural Network Semantics for First-Order Logic**

## **1** Introduction

## 2 Lifting a Modal Logic to First-Order Logic

- A relatively new technique in modal logic allows us to "lift" the modal semantics to firstorder logic. In this section, I describe this technique, and explain how it is used to give semantics for classical first-order logic.
- Explain where this technique comes from: A way to explain how first-order logic can be derived as a special kind of modal logic. This is widely known within certain logic circles, but is not widely known among those who use logics in AI. I am using it in a totally non-controversial way: I follow it as instructions for how to get a predicate logic from a modal one.
- Basic idea: Possible worlds are variable assignments
- This is done in two stages: Stage 1 is to treat each quantifier per variable ∀x like a dynamic modal operator [x]. Stage 2 is to consider substitution [y/x] like a 2-variable dynamic modal operator as well.
- This is first-order logic *without* equality. We can optionally introduce equality as a predicate =(*x*, *y*) [I don't know the consequences of that choice...].
- In Section 3, I will follow Stage 1 and interpret  $\exists x$  as the closure of a signal in a neural network (following closely to the modal logic interpretation). In Section 4, I will consider *classical* substitution and explore how it interacts with the neural network closure variant

#### Figure 2.1.

of  $\exists x$ . (I need to define free variables and "y is free for x in  $\varphi$ ")

• Here I can make a table of first-order logic axioms and the frame properties they correspond to. Figure 2.1.

### **3** Neural Network Semantics for First-Order Logic

## 3.1 Variable-Assignment Networks

I should first say that instead of an underlying set of propositions, there is an underlying set
of predicates p(x<sub>1</sub>,...,x<sub>k</sub>)

**Definition 3.1.** A variable-assignment network is  $\mathcal{N} = \langle N, \{E_x\}_{x \in VAR}, \{W_x\}_{x \in VAR}, \{A_x\}_{x \in VAR}, V \rangle$ , where

- N [type it correctly] is a finite nonempty set of *variable assignments*, which we treat as nodes of the network. I will use lowercase greek symbols α, β,... to denote these special variable assignment neurons, in order to distinguish them from ordinary neurons n, m,...
- The edges, weights, and activation function are defined exactly as before, but now I define one for each variable:  $E_x$ ,  $W_x$ ,  $A_x$  for  $x \in VAR$ .
- V: [] → P(N) is a *valuation* that maps each atomic predicate p(x<sub>1</sub>,...,x<sub>k</sub>) to the set of variable assignments that make p true (equivalently, the set of variable assignments that are active in state p).

#### [define State<sub>x</sub> for variable-assignment models ( $\mathcal{N}$ subscript will be hidden)]

Think of a variable-assignment network like an ensemble of neural networks  $\mathcal{N}_x = \langle N, E_x, W_x, A_x, V \rangle$  for each variable *x*. Each of these nets  $\mathcal{N}_x$  in the ensemble specifies a particular accessibility

structure for assignments  $\alpha \in N$  under variable *x*. Unlike classical assignment models however, this accessibility structure is determined by a neural network transition function  $F_{x,S_0}$  from state  $S \in \text{State}_x$  to the next state (for a particular variable *x*). This function  $F_{x,S_0}$  is defined exactly as before; the only difference is that there is now one per  $x \in \text{VAR}$ .

$$F_{x,S_0}(S) = S_0 \cup \left\{ \eta \mid A_x \left( \sum_{m \in \operatorname{preds}_x(n)} W_x(m,n) \cdot \chi_S(m) \right) = 1 \right\}$$

where  $\chi_S(m) = 1$  iff  $m \in S$  is the indicator function.

The closure for each  $x \in VAR$  is defined similarly. I postulate that for all states  $S_0$ ,  $F_{x,S_0}$  applied repeatedly to  $S_0$  has a unique fixed point under  $S_0$  (it is the only state S such that  $F_{x,S_0}(S) = S$ ). Let  $Clos_x$ : State<sub>x</sub>  $\rightarrow$  State<sub>x</sub> be the function that produces that least fixed point, for network  $\mathcal{N}_x$ . [define class of VANs]

### 3.2 Semantics for Quantifiers and Predicates

**Definition 3.2.** Formulas in the [extended first-order language] are given by

$$\varphi, \psi \coloneqq p(x_1, \dots, x_k) \mid \neg \varphi \mid \varphi \land \psi \mid U\varphi \mid \forall x\varphi \mid \forall x\varphi$$

 $\top, \bot, \lor, \rightarrow, \leftrightarrow$  and the dual quantifiers  $\exists x, \exists x \text{ are defined in the usual way.}$ 

 $\forall x \varphi$  is intended to be the minimal universal quantifier, which says " $\varphi$  is true in all variable assignments that are accessible from the current one." The special operator  $\forall x \varphi$  is intended to be the universal quantifier we get by interpreting  $\exists x \varphi$  as the closure of signal  $\varphi$  in the neural network  $\mathcal{N}_x$ .

**Definition 3.3.** For all  $\mathcal{N} \in \mathbf{Net}$ , and variable assignments  $\alpha \in N$ :

$$\begin{array}{lll} \mathcal{N}, \alpha \models p(x_1, \dots, x_k) & \text{iff} & \alpha \in V(p(x_1, \dots, x_k)) \\ \mathcal{N}, \alpha \models \neg \varphi & \text{iff} & \mathcal{N}, \alpha \not\models \varphi \\ \mathcal{N}, \alpha \models \varphi \land \psi & \text{iff} & \mathcal{N}, \alpha \models \varphi \text{ and } \mathcal{N}, \alpha \models \psi \\ \mathcal{N}, \alpha \models U\varphi & \text{iff} & [\text{todo}] \\ \mathcal{N}, \alpha \models \exists x \varphi & \text{iff} & \text{there is an } E_x\text{-path from } \llbracket \varphi \rrbracket \text{ to } \alpha \\ \mathcal{N}, \alpha \models \exists x \varphi & \text{iff} & \alpha \in \text{Clos}_x(\llbracket \varphi \rrbracket) \end{array}$$

where  $\llbracket \varphi \rrbracket = \{ \alpha \in N \mid \mathcal{N}, \alpha \models \varphi \}.$ 

**Proposition 3.4.** By the definition of the duals  $\exists x, \exists x$ , we can instead define the semantics for  $\forall x$ ,  $\forall x$  as follows. For all  $\mathcal{N} \in [], \alpha \in N$ :

$$\mathcal{N}, \alpha \models \forall x \varphi$$
 iff there is *no*  $E_x$ -path from  $\llbracket \varphi \rrbracket^{\mathbb{C}}$  to  $\alpha$   
 $\mathcal{N}, \alpha \models \forall x \varphi$  iff  $\alpha \in (\operatorname{Clos}(\llbracket \varphi \rrbracket^{\mathbb{C}}))^{\mathbb{C}}$ 

#### 3.3 Semantics for Variable Substitution

The story up to here is incomplete; variable substitution is at the heart of FOL, and I have not yet said anything about substitution in variable-assignment networks. Following [cite van Ben-them] closely, I consider a substitution relation  $\sim_{x,y}$ . The intended meaning of  $\alpha \sim_{x,y} \beta$  is " $\beta$  is the result of assigning  $x \coloneqq y$  in  $\alpha$ ." Formally,

$$\alpha \sim_{x,y} \beta$$
 iff  $\beta(x) = \alpha(y)$  and  $\alpha(z) = \beta(z)$  for all  $z \neq x$ 

For convenience, I assume that  $\neg_{x,y}$  is *functional*, i.e. for each  $\alpha \in N$  there is exactly one  $\beta \in N$  with  $\alpha \neg_{x,y} \beta$ .

Van Benthem also considers a generalized notion of substitution, where  $\sim_{x,y}$  is just an abstract relation on assignments. For my purposes I will opt for the concrete definition, since substitution is a syntactic notion and we don't gain much by interpreting it in a neural network structure.

From here, I will extend the language [] to [], which has substitution operators [y/x] for *x*,  $y \in VAR$ . Formulas are given by

$$\varphi, \psi \coloneqq p(x_1, \dots, x_k) \mid \neg \varphi \mid \varphi \land \psi \mid U\varphi \mid \forall x\varphi \mid \forall x\varphi \mid [y/x]\varphi$$

The semantics for [y/x] is given by

 $\mathcal{N}, \alpha \models [y/x]\varphi$  iff for that unique  $\beta$  such that  $\alpha \sim_{x,y} \beta$ , we have  $\mathcal{N}, \beta \models \varphi$ 

#### 4 Axioms, Soundness, and Frame Conditions

#### 4.1 Axioms for Quantifiers and Predicates

- In this section, I will explore which first-order logic axioms are sound for ∀*x* interpreted in a variable-assignment neural network.
- Sound ones: (Dual), (Rep), (CM), (Refl), and (Trans) get inherited by the underlying modal logic, and are sound for the same reasons. (e.g., (CM) holds because each Clos<sub>x</sub> is cumulative.)
- Are there any sound variable interactions (for two variables *x*, *y*)? All the ones listed appear to not necessarily be sound, but is there something we can say that's true?
- The unsound ones. For each of them, I should (1) prove that it isn't sound with a counterexample, and (2) think about *what property of* Clos<sub>x</sub> *would make it true*? (Think of each of these axioms as somewhat negotiable; it's perfectly fine for neural networks to model dependent quantifiers ∀x rather than the classical *independent* quantifiers.)

(Nec).

(Distr).

(Eucl).

(???).

(Exch). [Give the neural network that is a countermodel here] [Also, talk about what this means—neural networks are capable of modeling *dependent* quantifiers!]

**Proposition 4.1.** The exchange axiom (Exch) is sound if and only if for all nets  $\mathcal{N} \in [], \varphi \in []$  and  $x, y \in VAR$ ,  $Clos_x(Clos_y(\llbracket \varphi \rrbracket)) = Clos_y(Clos_x(\llbracket \psi \rrbracket))$ .

Proof. [todo]

(**PR**).

(Confl).

		Neural	Generalized	Classical
(Dual)	$\exists x \varphi \leftrightarrow \neg \forall x \neg \varphi$	$\bigotimes$	$\odot$	$\odot$
(Nec)	If $\vdash \varphi$ then $\vdash \forall x \varphi$	$\oslash$	$\oslash$	$\oslash$
(CM)	$U(\forall x \varphi \rightarrow \psi) \rightarrow (\forall x (\varphi \land \psi) \rightarrow \forall x \varphi)$	$\oslash$	$\oslash$	$\oslash$
(Distr)	$\forall x(\varphi \to \psi) \to (\forall x\varphi \to \forall x\psi)$	$\bigotimes$	$\oslash$	$\oslash$
(Refl)	$\forall x \varphi \to \varphi$	$\oslash$	$\bigotimes$	0
(Trans)	$\forall x \varphi \leftrightarrow \forall x \forall x \varphi$	$\oslash$	$\bigotimes$	$\odot$
(Eucl)	$\exists x \varphi \to \forall x \exists x \varphi$	$\bigotimes$	$\bigotimes$	$\odot$
(???)	$\varphi \rightarrow \forall x \varphi$ for x not free in $\varphi$	$\bigotimes$	$\bigotimes$	$\odot$
(Exch)	$\forall x \forall y \varphi \leftrightarrow \forall y \forall x \varphi$	$\bigotimes$	$\bigotimes$	$\odot$
(Confl)	$\exists x \forall y \varphi \rightarrow \forall y \exists x \varphi$	$\bigotimes$	$\bigotimes$	$\bigcirc$

#### **Semantics for First-Order Logic**

**Figure 4.1.** Soundness for various FOL axioms, under three different semantics: neural semantics using variable-assignment networks [], the generalized semantics [], and classical semantics for FOL []. For the generalized and classical semantics,  $\forall x$  and  $\exists x$  should be interpreted as the usual  $\forall x$  and  $\exists x$ .  $\bigcirc$  indicates that the axiom is sound, i.e., holds for all models in that class without introducing any new frame properties.  $\bigotimes$  indicates that the axiom does not hold for all models in that class, and may require additional frame properties to make it hold.

Abbreviations: (**CM**) is Cautious Monotonicity (from conditional logic); (**Eucl**) is Euclidean; (**Exch**) is the standard FOL exchange principle; (**Conff**) is Confluence.

## 4.2 Axioms for Variable Substitution

Let's now consider soundness for this extended language with [y/x]. Since the operator [y/x] is defined classically, it satisfies all the usual FOL substitution axioms over predicates,  $\land$ , and  $\neg$ :

- $[y/x]p(x) \leftrightarrow p(y)$ , and  $[y/x]p(z) \leftrightarrow p(z)$  for  $z \neq x$ .
- $[y/x](\varphi \land \psi) \leftrightarrow ([y/x]\varphi \land [y/x]\psi)$
- $[y/x] \neg \varphi \leftrightarrow \neg [y/x] \varphi$  (since substitution is *functional*)

[It satisfies the usual interaction axioms with the classical/minimal quantifier  $\forall x$  so long as I add in the right frame conditions, which I should do above.]

I will now explore the different possible interactions between classical substitution [y/x] and

neural quantifiers  $\forall x$ , i.e., interactions between the substitution relation  $\neg_{x,y}$  and  $\text{Clos}_x$ . In contrast to the axioms with quantifiers alone, I consider the [y/x] interactions to be relatively non-negotiable; in order to interpret a neural network as a FOL reasoner, it's important for the net's activation function  $\text{Clos}_x$  to get along with substitution (the core mechanic of FOL). Here is a list of the standard interactions for classical FOL (these versions with the [y/x] operator come from [cite van Benthem], who drew the standard axioms from [cite Enderton]):

- 1.  $[y/x] \forall x \varphi \leftrightarrow \forall x \varphi$
- 2.  $[y/x] \forall z \varphi \leftrightarrow \forall z [y/x] \varphi$  for  $z \neq x$
- 3.  $\forall x \varphi \rightarrow [y/x] \varphi$ , if y free for x in  $\varphi$

[The second axiom includes the case  $[y/x] \forall y \varphi \leftrightarrow \forall y [y/x] \varphi$ . I'm not listing existential interactions such as  $[y/x] \exists x \varphi \leftrightarrow \exists x \varphi$ ,  $[y/x] \exists y \varphi \leftrightarrow \exists y [y/x] \varphi$ , and  $[y/x] \exists z \varphi \leftrightarrow \exists z [y/x] \varphi$  for  $z \neq x, y$ , since for functional [y/x] they are each equivalent to their duals.]

[I'm also only interested in the  $(\rightarrow)$  direction of (1) and (2) above, ]

Note that these are *all* sound in classical FOL, and in generalized FOL they are each only sound alongside frame conditions. For variable-assignment network semantics, it's not necessary for  $Clos_x$  to interact with  $\prec_{x,y}$  at all—by default, none of these interaction axioms are sound. But we can now ask: what neural network frame conditions are sufficient to make these axioms hold? In other words, what class of variable-assignment nets  $\mathcal{N}$  interact in the expected way with substitution?

Consider the following three properties of variable-assignment nets  $\mathcal{N}$ . The first, that  $\mathcal{N}$  respects substitutions, says that if  $\beta$  is the result of assigning  $x \coloneqq y$  in  $\alpha$ , and  $\alpha$  is activated by signal  $[\![\varphi]\!]$ , then  $\beta$  is activated as well.

**Definition 4.2.** A variable-assignment net  $\mathcal{N} \in []$  respects substitutions iff for all  $\alpha, \beta \in N, x$ ,  $y \in VAR, \varphi \in []$ , if  $\alpha \sim_{x,y} \beta$  and  $\alpha \in Clos_x(\llbracket \varphi \rrbracket)$ , then  $\beta \in Clos_x(\llbracket \varphi \rrbracket)$ .

**Example 4.3.** [Give an example of a neural network that respects substitutions (but does not necessarily reflect them)]

The second property,  $\mathcal{N}$  reflects substitutions, says that if  $\beta$  is the result of assigning  $x \coloneqq y$  in  $\alpha$ , and substituting y for x in  $\varphi$  does not result in any binding issues, if  $\beta$  is activated by signal  $[\![\varphi]\!]$ , then  $\alpha$  is activated as well.

**Definition 4.4.** A variable-assignment net  $\mathcal{N} \in []$  *reflects substitutions* iff for all  $\alpha, \beta \in N, x, y \in VAR, \varphi \in []$ , if  $\alpha \sim_{x,y} \beta, \beta \in Clos_x(\llbracket \varphi \rrbracket)$ , and and *y* is free for *x* in  $\varphi$ , we have  $\alpha \in Clos_x(\llbracket \varphi \rrbracket)$ .

**Example 4.5.** [Give an example of a neural network that reflects substitutions (but does not necessarily respect them)]

A neural network both *respects and reflects substitutions* iff whenever  $\alpha \sim_{x,y} \beta$ ,  $\alpha$  and  $\beta$  are activated by exactly the same signals  $[\![\varphi]\!]$  (provided *y* is free for *x* in  $\varphi$ ).

**Definition 4.6.** [A third condition for the other axiom]

Example 4.7. [Give an example of a neural network with this third condition, but not the others]

The main result of this section is that these three properties are sufficient for the neural network interpreted  $\forall x$  to interact normally with [y/x]. [Answer: what kinds of neural network ensembles are these ones like?]

Example 4.8. [Give an example of a neural network with all three!]

I will now prove that these three properties are sufficient:

**Proposition 4.9.** Suppose a variable-assignment net  $\mathcal{N}$  *respects* substitutions. Then the  $(\rightarrow)$  direction of the axiom  $[y/x] \forall x \varphi \leftrightarrow \forall x \varphi$  is sound.

**Proof.** Let  $\mathcal{N} \in []$  and let  $\alpha \in N$  be any assignment. Let  $\beta$  be that unique assignment such that  $\alpha \sim_{x,y} \beta$ . Suppose  $\mathcal{N}, \alpha \models [y/x] \forall x \varphi$ . By the semantics for substitution,  $\mathcal{N}, \beta \models \forall x \varphi$ , and then by the semantics for  $\forall x$  we have  $\beta \notin \text{Clos}_x([\![\varphi]\!]^{\mathbb{C}})$ . Since  $\mathcal{N}$  respects substitutions,  $\alpha \notin \text{Clos}_x([\![\varphi]\!]^{\mathbb{C}})$ . But this means  $\mathcal{N}, \alpha \models \forall x \varphi$ , and we are done.

**Proposition 4.10.** Suppose a variable-assignment net  $\mathcal{N}$  *reflects* substitutions. Then the axiom " $\forall x \varphi \rightarrow [y/x]\varphi$ , if y free for x in  $\varphi$ " is sound.

**Proof.** Let  $\mathcal{N} \in []$  and let  $\alpha \in N$  be any assignment. Let  $\beta$  be that unique assignment such that  $\alpha \sim_{x,y} \beta$ . Suppose  $\mathcal{N}, \alpha \models \forall x \varphi$ , and suppose *y* is free for *x* in  $\varphi$ . By the semantics for  $\forall x$  we have

 $\alpha \notin \operatorname{Clos}_{x}(\llbracket \varphi \rrbracket^{\mathbb{C}}) = \operatorname{Clos}_{x}(\llbracket \neg \varphi \rrbracket).$ 

Now, suppose for contradiction that  $\mathcal{N}, \alpha \not\models [y/x]\varphi$ , i.e.,  $\beta \in \llbracket \varphi \rrbracket^{\mathbb{C}}$ . By Inclusion of  $\operatorname{Clos}_x, \beta \in \operatorname{Clos}_x(\llbracket \varphi \rrbracket^{\mathbb{C}}) = \operatorname{Clos}_x(\llbracket \neg \varphi \rrbracket)$ . Since *y* is free for *x* in  $\varphi$ , *y* is free for *x* in  $\neg \varphi$ . So we can apply the fact that  $\mathcal{N}$  reflects substitutions, which gives us  $\alpha \in \operatorname{Clos}_x(\llbracket \varphi \rrbracket^{\mathbb{C}})$ . But this directly contradicts our hypothesis. So we must conclude that  $\mathcal{N}, \alpha \models [y/x]\varphi$ .

**Corollary 4.11.** Suppose a variable-assignment net  $\mathcal{N}$  respects substitutions. Then the  $(\leftarrow)$  direction of the axiom  $[y/x] \forall x \varphi \leftrightarrow \forall x \varphi$  is sound.

**Proof.** If  $\mathcal{N}$  respects substitutions, then " $\forall x \varphi \rightarrow [y/x] \varphi$ , if *y* free for *x* in  $\varphi$ " is sound. I will show that the formula  $\forall x \varphi \rightarrow [y/x] \forall x \varphi$  follows. Let  $\alpha \in N$  and suppose  $\mathcal{N}, \alpha \models \forall x \varphi$ . By the soundness of (**Trans**),  $\mathcal{N}, \alpha \models \forall x \forall x \varphi$ . Now, notice that *y* is free for *x* in the expression  $\forall x \varphi$  (since *x* is not free at all, it is safe to substitute *y* for *x*). So by hypothesis,  $\mathcal{N}, \alpha \models [y/x] \forall x \varphi$ , which was the goal.  $\Box$ 

**Proposition 4.12.** Suppose a variable-assignment net  $\mathcal{N}$  is [todo]. Then the following substitution interaction axioms are sound:

• [todo]

Proof. [todo]

#### **5** Reflections on First-Order Reasoning using Neural Networks

#### 5.1 How to Interpret Variable-Assignment Networks

- Part of my conclusion is this: In order for a variable-assignment net to behave like full FOL reasoner, it must have [list of properties]
- In addition to whatever else I do here, I should probably sanity check and make sure that there *is* actually a variable-assignment net that satisfies all of the properties necessary to make it act like a full FOL reasoner. (If there isn't one that has all the properties, then one of the properties is saying too much / together they are trivial.)

## 5.2 Are variable-assignment networks cognitively plausible?

Working list of papers that might bear on this question:

- Corey, Neural basis for generalized quantifier comprehension, 2005
- Kiela, Variable Binding in Biologically Plausible Neural Networks, 2011

The questions here are (1) whether there really are distinct **types** of edges in neural networks, and (2) whether neurons can be thought of as carrying variable-assignment information.

## 5.3 Related Work: Other Neuro-Symbolic Systems for First-Order Reasoning

Working list of neuro-symbolic systems that interpret quantifiers or variable substitution in any formal way:

- Logic Tensor Networks (over "real logic")—see Badreddine, Logic Tensor Networks, 2022
- Logical Neural Networks—see Riegel, Gray, Logical neural networks, 2020
- Kiela, Variable Binding in Biologically Plausible Neural Networks, 2011
- What about DeepProbLog?

#### Chapter 7

#### Conclusions

In the Introduction, I said that this dissertation serves to promote the point of view that **neural networks can be thought of as models for formal logic**. Let's now take a step back and see how far this perpective has taken us.

[I like the way Levin Hornischer wrote his: A summary, followed by a list of results, followed by a list of open questions]

I could also give a fun sampler of some future applications of neural network semantics. [For these 4 directions, I already have some basic ideas for how to proceed. So I will finish out by outlining what I think the next steps are. I don't necessarily have to prove anything difficult, but I think it's important to show that I've been thinking about these deeply, and I have a start for how to proceed (and it's nice to be able to claim the ideas!)] [Maybe it's better in this section to pivot to an "open questions" format and keep it brief...]

**Neural Network Updates from Epistemic Ones.** My original goal was to explore what "classical" updates correspond to neural network updates such as Hebb\*, and conversely what neural updates correspond to plausibility updates such as Cond, Lex, and Consr. This theorem doesn't seem to clarify that; it's constructive (we *do* in fact build the corresponding updates), but by translating back and forth at the static level we don't define the update in the original updates "native environment." [Elaborate/clarify this point]. For the remainder of this section, I will put in the extra work to see what Hebb\*, Cond, Lex, and Consr "look like" on the other side.

**Definition 0.1.** [Define the plausibility update that simulates Hebb\*]

**Definition 0.2.** [Define the neural update that simulates Cond]

**Definition 0.3.** [Define the plausibility update that simulates Lex]

**Definition 0.4.** [Define the plausibility update that simulates Consr]

**The Learning Power of Hebbian Learning.** [Differentiate between (single-step) update and (in the limit) learning. Cite papers on learning power of different epistemic updates] [Go ahead and give the mathematical setup needed to ask the question formally.] [An interesting question, but broad: what can we say about neural network learning in the context of computational learning theory?]

A Dynamic Logic of Backpropagation. [Of course I don't have any results about this (I'm focusing on giving a complete story about the unsupervised upgrade), but in principle we can do the same thing for gradient descent (implemented as backpropagation over a neural network). Formulas  $[P;Q]\varphi$  need to be used. I can go ahead and give what the semantics would look like over neural network models. And it would be nice if I could give at least one sound property of backpropagation here!]

#### **1** Open Questions

[First, list all the open questions! Then talk in detail about a few!]

- Give a characterization of those (recurrent) nets whose closure Clos(S) reaches a unique fixed point (i.e. does not oscillate)
- 2. Extend the semantics to account for recurrent neural networks whose Clos(S) does oscillate!
- 3. Extend the semantics for fuzzy (and possibly probabilistic) activation patterns Clos(S)
- 4. Give semantics and sound axioms for other neural network architectures that use new ideas,e.g. transformers & attention mechanisms. (Brand new ideas are needed, this isn't just a simple extension!)
- 5. It's an open problem to prove formal soundness results for neuro-symbolic systems like LTNs and DeepProbLog? Do the techniques from this work shed light on how to do this, or what axioms are needed?
- 6. Extend the semantics to give sound (and possibly complete) axioms for other neural network updates, especially gradient descent implemented as backpropagation! Can we use

this to give a "classical" belief revision operator equivalent to backpropagation?

- 7. Extend the semantics to give sound (and possibly complete) neural network semantics for first-order logic! Proving soundness & completeness for *any* neural network interpretation of first-order logic is a huge open problem in neuro-symbolic AI.
- 8. Extend the semantics to stabilized Hebbian update, using Oja's rule
- 9. What concrete neural network learning policy corresponds to Lex? To Consr? What about Rott's 21 different belief revision operators? What classical belief revision policies correspond to Hebbian update? Oja's rule? Backpropagation?
- 10. Compare the learning power of neural network updates (Hebbian update, backprop) in the limit against other kinds of updates, e.g. epistemic updates, or against each other.
- 11. It's clear that every unsupervised belief revision operator is a supervised one. But is every supervised belief revision operator equivalent to an unsupervised one? Does this result carry over to neural networks?
- 12. Characterize different belief revision operators & learning policies using *frame correspondences*. Modal logic is really the perfect environment for asking this question in this way.

#### Bibliography

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat et al. GPT-4 technical report. *ArXiv preprint arXiv:2303.08774*, 2023.
- [2] Alfred V. Aho, Michael R Garey, and Jeffrey D. Ullman. The transitive reduction of a directed graph. SIAM Journal on Computing, 1(2):131–137, 1972.
- [3] Aws Albarghouthi et al. Introduction to neural network verification. *Foundations and Trends*® *in Programming Languages*, 7(1–2):1–157, 2021.
- [4] Edoardo Baccini, Zoé Christoff, and Rineke Verbrugge. Dynamic logics of diffusion and link changes on social networks. *Studia Logica*, pages 1–71, 2024.
- [5] Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration A structured survey. In *We Will Show Them! Essays in Honour of Dov Gabbay, Volume 1*, pages

167–194. College Publications, 2005.

- [6] Samy Badreddine, Artur d'Avila Garcez, Luciano Serafini, and Michael Spranger. Logic Tensor Networks. *Artificial Intelligence*, 303:103649, 2022.
- [7] Christian Balkenius and Peter G\u00e4rdenfors. Nonmonotonic inferences in neural networks. In *KR*, pages 32–39. Morgan Kaufmann, 1991.
- [8] Alexandru Baltag, Zoé Christoff, Rasmus K Rendsvig, and Sonja Smets. Dynamic epistemic logics of diffusion and prediction in social networks. *Studia Logica*, 107:489–531, 2019.
- [9] Alexandru Baltag, Nina Gierasimczuk, Aybüke Özgün, Ana Lucia Vargas Sandoval, and Sonja Smets. A dynamic logic for learning theory. *Journal of Logical and Algebraic Methods in Programming*, 109:100485, 2019.
- [10] Alexandru Baltag, Nina Gierasimczuk, and Sonja Smets. Truth-tracking by belief revision. *Studia Logica*, 107:917–947, 2019.
- [11] Alexandru Baltag, Dazhu Li, and Mina Young Pedersen. On the right path: A modal logic for supervised learning. In *International Workshop on Logic, Rationality and Interaction*, pages 1–14. Springer, 2019.
- [12] Alexandru Baltag, Lawrence S Moss, and Sławomir Solecki. Logics for epistemic actions: completeness, decidability, expressivity. *Logics*, 1(2):97–147, 2023.
- [13] Alexandru Baltag, Lawrence S Moss, and Slawomir Solecki. The logic of public announcements, common knowledge, and private suspicions. In *Proceedings of the 7th conference on Theoretical aspects of rationality and knowledge*, pages 43–56. 1998.
- [14] Alexandru Baltag and Sonja Smets. Group belief dynamics under iterated revision: Fixed points and cycles of joint upgrades. In *Proceedings of the 12th Conference on Theoretical Aspects of Rationality and Knowledge*, pages 41–50. 2009.
- [15] Tarek R Besold, Artur d'Avila Garcez, Sebastian Bader, Howard Bowman, Pedro Domingos, Pascal Hitzler, Kai-Uwe Kühnberger, Luis C Lamb, Priscila Machado Vieira Lima, Leo de Penning et al. Neural-symbolic learning and reasoning: A survey and interpretation. In *Neuro-Symbolic Artificial Intelligence: The State of the Art*, pages 1–51. IOS press, 2021.

- [16] Patrick Blackburn, Maarten de Rijke, and Yde Venema. *Modal Logic*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press, 2001.
- [17] Reinhard Blutner. Nonmonotonic inferences and neural networks. *Synthese*, 142:143–174, 2004.
- [18] Zoé Christoff and Jens Ulrik Hansen. A logic for diffusion in social networks. *Journal of Applied Logic*, 13(1):48–77, 2015.
- [19] Gabriele Ciravegna, Pietro Barbiero, Francesco Giannini, Marco Gori, Pietro Lió, Marco Maggini, and Stefano Melacci. Logic Explained Networks. *Artificial Intelligence*, 314:103822, 2023.
- [20] Artur 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.
- [21] Walter Dean. Computational complexity theory. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, 2021.
- [22] Daniel C Dennett. Cognitive wheels: the frame problem of ai. *The philosophy of artificial intelligence*, 147:170, 1990.
- [23] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan et al. The Llama 3 herd of models. *ArXiv preprint arXiv:2407.21783*, 2024.
- [24] Herbert B Enderton. A mathematical introduction to logic. Elsevier, 2001.
- [25] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in Large Language Models: A survey. *Computational Linguistics*, pages 1–79, 2024.
- [26] Artur SD'Avila Garcez, Luis C Lamb, and Dov M Gabbay. *Neural-Symbolic Cognitive Rea-soning*. Springer Science & Business Media, 2008.
- [27] Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. Finding alignments between interpretable causal variables and distributed neural represen-

tations. In Causal Learning and Reasoning, pages 160-187. PMLR, 2024.

- [28] 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.
- [29] Laura Giordano and Daniele Theseider Dupré. Weighted defeasible knowledge bases and a multipreference semantics for a deep neural network model. In *Logics in Artificial Intelligence: 17th European Conference, JELIA 2021, Virtual Event, May 17–20, 2021, Proceedings 17*, pages 225–242. Springer, 2021.
- [30] Charles G Gross. Genealogy of the "grandmother cell". *The Neuroscientist*, 8(5):512–518, 2002.
- [31] Frankvan Harmelen. Preface: The 3rd AI wave is coming, and it needs a theory. In *Neuro-Symbolic Artificial Intelligence: The State of the Art*, page 0. IOS Press BV, 2022.
- [32] Donald Hebb. *The Organization of Behavior*. Psychology Press, apr 1949.
- [33] Neil Immerman. *Descriptive Complexity*. Springer Science & Business Media, 1998.
- [34] Caleb Kisby, Saúl Blanco, and Lawrence Moss. The logic of Hebbian learning. In *The International FLAIRS Conference Proceedings*, volume 35. 2022.
- [35] Caleb Schultz Kisby, Saúl A Blanco, and Lawrence S Moss. What do Hebbian learners learn? Reduction axioms for iterated Hebbian learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 14894–14901. 2024.
- [36] Dexter Kozen and Rohit Parikh. An elementary proof of the completeness of PDL. *Theoretical Computer Science*, 14(1):113–118, 1981.
- [**37**] Sarit Kraus, Daniel Lehmann, and Menachem Magidor. Nonmonotonic reasoning, preferential models and cumulative logics. *Artificial intelligence*, 44(1-2):167–207, 1990.
- [38] Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets. *Artificial Intelligence*, 128(1-2):161–201, 2001.
- [39] Hannes Leitgeb. Nonmonotonic reasoning by inhibition nets II. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 11(supp02):105–135, 2003.

- [40] Hannes Leitgeb. Neural network models of conditionals. In *Introduction to Formal Philos-ophy*, pages 147–176. Springer, 2018.
- [**41**] Hector J Levesque. *Common sense, the Turing test, and the quest for real AI*. MIT press, 2017.
- [42] Leonid Libkin. *Elements of Finite Model Theory*, volume 41. Springer, 2004.
- [43] Robin Manhaeve, Sebastijan Dumančić, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. Neural probabilistic logic programming in DeepProbLog. *Artificial Intelligence*, 298:103504, 2021.
- [44] 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, dec 1943.
- [45] Drew McDermott. A critique of pure reason. *Computational intelligence*, 3(3):151–160, 1987.
- [46] William Merrill. Sequential neural networks as automata. *ArXiv preprint arXiv:1906.01615*, 2019.
- [47] William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. *ArXiv preprint arXiv:2310.07923*, 2023.
- [48] William Merrill, Gail Weiss, Yoav Goldberg, Roy Schwartz, Noah A Smith, and Eran Yahav.A formal hierarchy of RNN architectures. *ArXiv preprint arXiv:2004.08500*, 2020.
- [49] Lawrence S Moss. Finite models constructed from canonical formulas. *Journal of Philosophical Logic*, 36:605–640, 2007.
- [50] Leonardo de Moura and Sebastian Ullrich. The Lean 4 theorem prover and programming language. In *Automated Deduction–CADE 28: 28th International Conference on Automated Deduction, Virtual Event, July 12–15, 2021, Proceedings 28*, pages 625–635. Springer, 2021.
- [51] Gregory Murphy. *The Big Book of Concepts*. MIT press, 2004.
- [52] Cathy O'neil. Weapons of math destruction: How big data increases inequality and threatens democracy. Crown, 2017.
- [53] Simon Odense and Artur d'Avila Garcez. A semantic framework for neurosymbolic compu-

tation. Artificial Intelligence, 340:104273, 2025.

- [54] Erkki Oja. Simplified neuron model as a principal component analyzer. *Journal of mathematical biology*, 15:267–273, 1982.
- [55] Eric Pacuit. *Neighborhood Semantics for Modal Logic*. Springer, 2017.
- [56] Jan A. Plaza. Logics of public communications. *Synthese*, 158:165–179, 2007.
- [57] George Polya. Mathematics and Plausible Reasoning: Induction and Analogy in Mathematics, volume 2. Princeton University Press, 1954.
- **[58]** David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. *Biometrika*, 71(599-607):6, 1986.
- [59] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- [60] Md Kamruzzaman Sarker, Lu Zhou, Aaron Eberhart, and Pascal Hitzler. Neuro-Symbolic Artificial Intelligence: Current Trends. *AI Communications*, 34, 2022 2022.
- [61] Murray Shanahan. The frame problem. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, 2016.
- [62] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- [63] Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28, page 0. Curran Associates, Inc., 2015.
- [64] Lena Strobl, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. What formal languages can transformers express? A survey. *Transactions of the Association for Computational Linguistics*, 12:543–561, 2024.
- [65] Alex Tamkin, Miles Brundage, Jack Clark, and Deep Ganguli. Understanding the capabilities, limitations, and societal impact of Large Language Models. *ArXiv preprint arXiv:2102.02503*, 2021.

- [66] Johan van Benthem. Dynamic logic for belief revision. *Journal of applied non-classical logics*, 17(2):129–155, 2007.
- [67] Johan van Benthem. Modal logic for open minds. Number 199. Center for the Study of Language and Information Stanford, 2010.
- [68] Johan van Benthem. *Logical Dynamics of Information and Interaction*. Cambridge University Press, 2011.
- [69] Johan van Benthem and Fenrong Liu. Dynamic logic of preference upgrade. *Journal of Applied Non-Classical Logics*, 17(2):157–182, 2007.
- [70] Johan van Benthem and Sonja Smets. Dynamic logics of belief change. In H. Van Ditmarsch, J. Halpern, W. van der Hoek, and B. Kooi, editors, *Handbook of Epistemic Logic*, pages 313–393. College Publications, London, UK, 2015.
- [71] Hans Van Ditmarsch, Wiebe van Der Hoek, and Barteld Kooi. *Dynamic Epistemic Logic*, volume 337. Springer, 2007.
- [72] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [73] Gail Weiss, Yoav Goldberg, and Eran Yahav. On the practical computational power of finite precision RNNs for language recognition. *ArXiv preprint arXiv:1805.04908*, 2018.

# Appendix A

# A.1 Appendix for Section 3.

Theorem 4.1. blah blah	
Proof.	
Proof Sketch. dfadsl;fjal;sdfj	

# A.2 Appendix A2

## A.3 Appendix A3