



# Reasoning about Neural Network Learning

Caleb Kisby, Saúl Blanco, and Lawrence Moss  
Luddy School of Informatics, Computing, and Engineering

## Reasoning about Static Nets

### Monotonicity Axioms

$\text{know } (A \rightarrow B) \rightarrow (\text{know } A \rightarrow \text{know } B)$   
 $(\text{typ } A_1 \rightarrow A_2) \dots (\text{typ } A_n \rightarrow A_1) \rightarrow$   
 $(\text{typ } A_i \leftrightarrow A_j)$

### Basic Modal Axioms

$\text{know } A \rightarrow A$   
 $\text{know } A \rightarrow \text{know know } A$   
 $\text{typ } A \rightarrow A$   
 $\text{typ } A \rightarrow \text{typ typ } A$   
 $\text{know } A \rightarrow \text{typ } A$

### Syntax

**p**  
 $A$  **and**  $B$   
 $A \rightarrow B$   
**know**  $A$   
**typ**  $A$   
 $A \Rightarrow B$   
 $[\text{hebb } A] B$   
 $[\text{hebb}^* A] B$

### Classical Meaning

proposition  
 $A$  and  $B$   
 $A$  implies  $B$   
the agent **knows**  $A$   
**typically**  $A$   
 $\text{typ } A \rightarrow B$   
incremental pref upgrade on  $A$   
preference upgrade on  $A$

### Neural Network

a (fuzzy) set of neurons  
 $A \cup B$   
 $A \supseteq B$   
the set of neurons **reachable** from  $A$   
the set of neurons **activated** by  $A$   
on input  $A$  the net **predicts**  $B$   
**learn**  $A$  (Hebbian)  
**repeatedly learn**  $A$  (Hebbian)

## Reasoning about Learning

### Induction Axioms

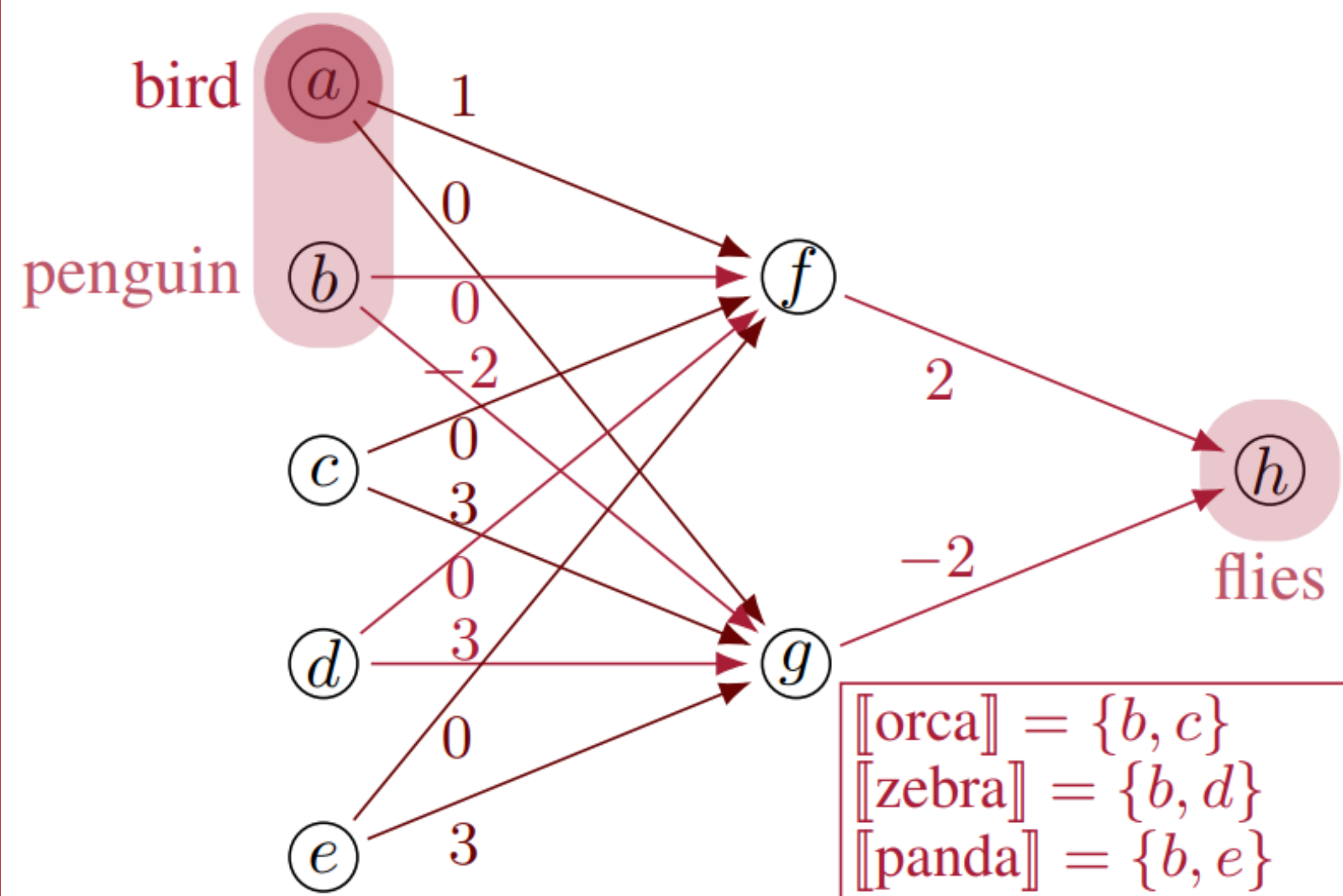
$[\text{hebb}^* A] B \rightarrow B$  and  $[\text{hebb } A] [\text{hebb}^* A] B$   
 $[\text{hebb}^* A] (B \rightarrow [\text{hebb } A] B) \rightarrow [\text{hebb}^* A] B$

### What The Net Learns

$[\text{hebb}^* A] \text{typ } B \leftrightarrow$   
 $\begin{cases} \text{typ } [\text{hebb}^* A] B & \text{if typ } A \text{ or typ } B \text{ is } \emptyset \\ \text{typ } [\text{hebb}^* A] B \text{ and } & \\ (\text{typ } A \text{ or know } B) & \text{otherwise} \end{cases}$

## Model Checking

### Task: Does the net satisfy P?



$\mathcal{N} \models \text{typ penguin} \rightarrow \text{flies}$ , but  
 $\mathcal{N} \not\models [\text{hebb orca}] [\text{hebb zebra}] [\text{hebb panda}]$   
 $\text{typ penguin} \rightarrow \text{flies}$

```
>>> print(model.is_model("typ penguin → flies"))
True
```

```
>>> print(model.is_model("[hebb orca] [hebb zebra] [hebb panda] \
typ penguin → flies))
False
```

## Model Building

### Task: Build a net that satisfies P.

GOAL. (Binary, feedforward) nets are equivalent to a certain class of classical modal frames.

COROLLARY. Given a knowledge base  $\Gamma$ , we can construct a net  $\mathcal{N}$  such that  $\mathcal{N} \models \Gamma$

COROLLARY. The axioms for reasoning about **know**, **typ**, and **[hebb\* A]** are complete.

## Work in Progress

- Use Lean to verify model checking code
- Finish proof for model building
- Extend system to reason about fuzzy sets
- Extend with **[backprop A]** (backpropagation)

**Acknowledgments:** This work was funded by the US Department of Defense (Contract W52P1J2093009). Thank you for your support!

[github.com/ais-climber/neural-semantics](https://github.com/ais-climber/neural-semantics)

