# Tree-Structured Algorithms as Causal Abstractions of Neural Networks
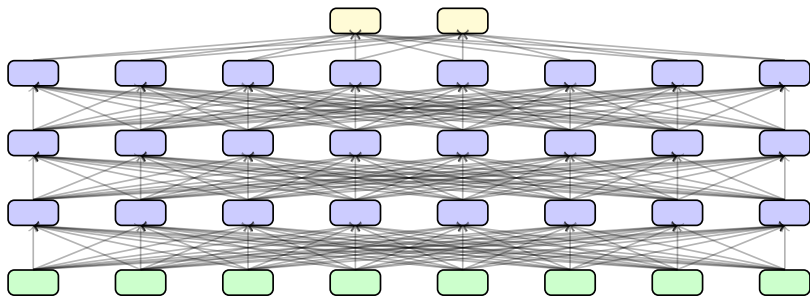
## Atticus Geiger

Stanford Linguistics and the Stanford NLP Group
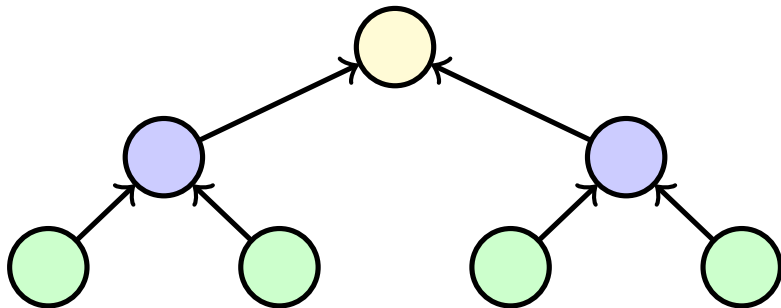
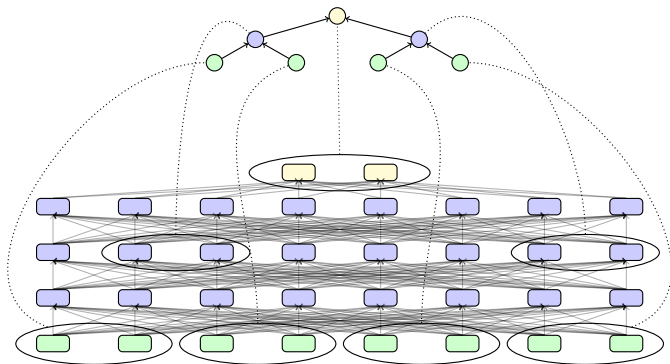# Introduction

# Deep Learning Model

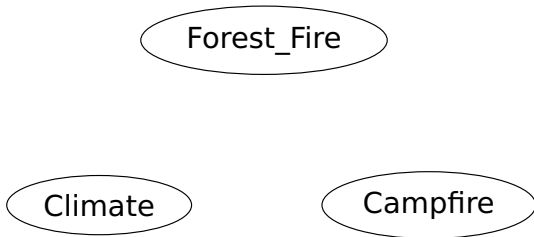# Tree Structured Algorithm

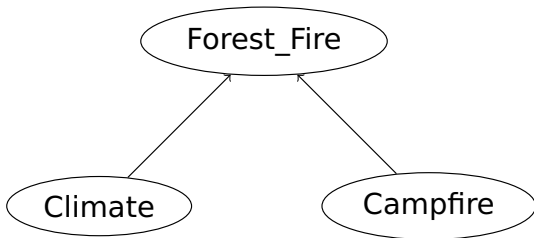# Constructive Causal Abstraction

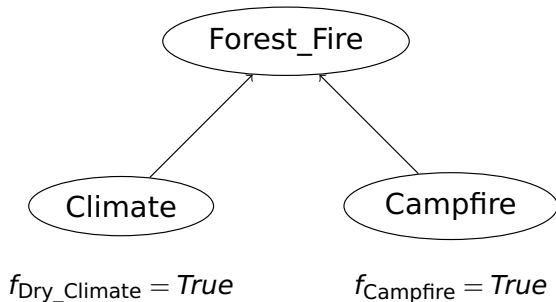# Acyclic Causal Models

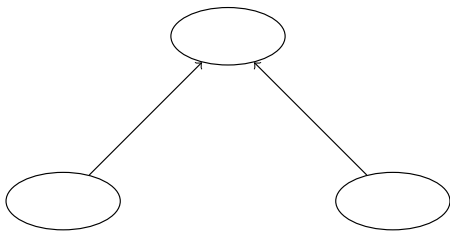# Dynamics with an Intervention

# Dynamics with an Intervention

# Dynamics with an Intervention

$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \land \text{camp fire}$



$f_{\text{Dry\_Climate}} = \textit{True}$          $f_{\text{Campfire}} = \textit{True}$

# Dynamics with an Intervention

$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \wedge \text{camp fire}$
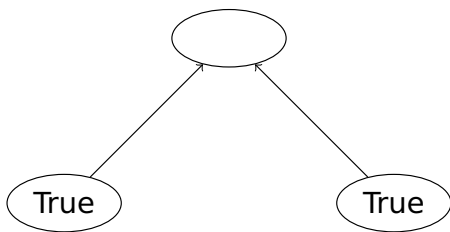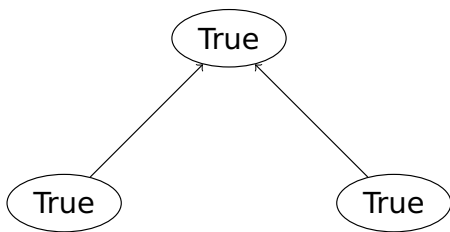


$f_{\text{Dry\_Climate}} = \textit{True}$          $f_{\text{Campfire}} = \textit{True}$

# Dynamics with an Intervention

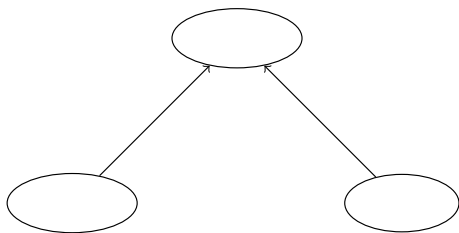$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \wedge \text{camp fire}$
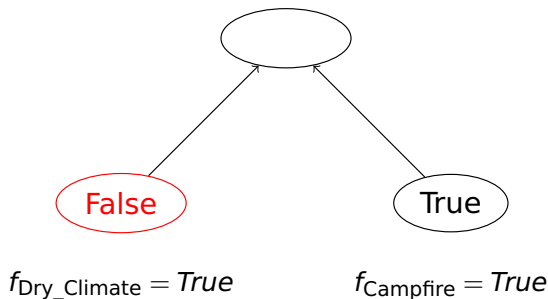


$f_{\text{Dry\_Climate}} = \textit{True}$          $f_{\text{Campfire}} = \textit{True}$
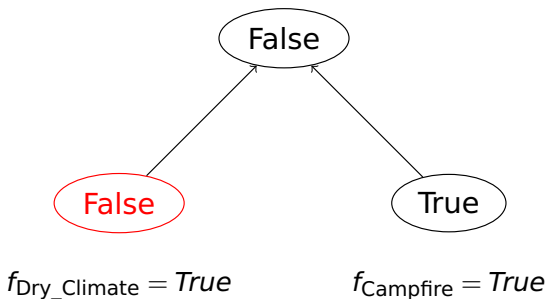
# Dynamics with an Intervention

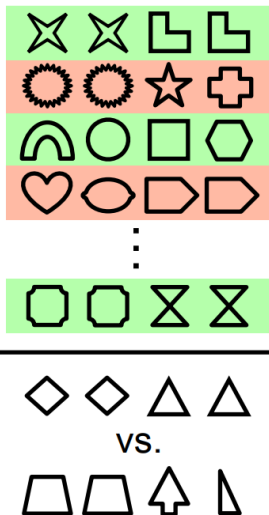$$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \land \text{camp fire}$$



$f_{\text{Dry\_Climate}} = \textit{True}$        $f_{\text{Campfire}} = \textit{True}$

# Dynamics with an Intervention

$$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \wedge \text{camp fire}$$



$f_{\text{Dry\_Climate}} = \textit{True}$    $f_{\text{Campfire}} = \textit{True}$

# Dynamics with an Intervention

$$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \wedge \text{camp fire}$$



$f_{\text{Dry\_Climate}} = \textit{True}$        $f_{\text{Campfire}} = \textit{True}$

# Dynamics with an Intervention

$$f_{\text{Forest\_Fire}}(\text{dry\_climate}, \text{campfire}) = \text{dry\_climate} \wedge \text{camp fire}$$



$f_{\text{Dry\_Climate}} = \textit{True}$          $f_{\text{Campfire}} = \textit{True}$

# Hierarchical Equality Task

# Hierarchical Equality Task

# Algorithms as Acyclic Causal Models

Introduction    Acyclic Causal Models    Hierarchical Equality Task    **Algorithms**    Deep Learning Models    Constructive Causal Abstraction

○○○○     ○○     ○○     ○●○     ○○     ○○○○○○

# Tree-Structured Algorithm

$$\mathbf{Id_1}(\,\cdot\,,\,\cdot\,) =$$

|   | △ | ■ | ⬠ |
|---|---|---|---|
| △ | T | F | F |
| ■ | F | T | F |
| ⬠ | F | F | T |

$$\mathbf{Id_2}(\,\cdot\,,\,\cdot\,) =$$

|   | T | F |
|---|---|---|
| T | T | F |
| F | F | T |

**function** EQUALITYTASK(*shape1*, *shape2*, *shape3*, *shape4*)

     *same1* ← $\mathbf{Id_1}$(*shape1*, *shape2*)

     *same2* ← $\mathbf{Id_1}$(*shape3*, *shape4*)

     *same3* ← $\mathbf{Id_2}$(*same1*, *same2*)

     **return** *same3*

# Causal Model of Algorithm



$$f_O(v_1, v_2) = \textbf{Id}_2(v_1, v_2)$$

$$f_{V_1}(i_1, i_2) = \textbf{Id}_1(i_1, i_2) \qquad f_{V_2}(i_3, i_4) = \textbf{Id}_1(i_3, i_4)$$

# Deep Learning Models as Acyclic Causal Models
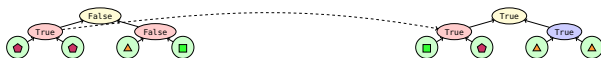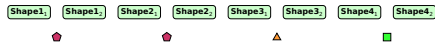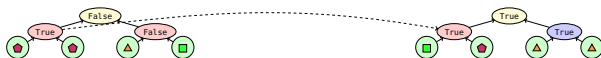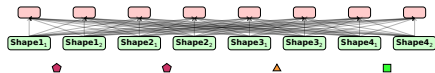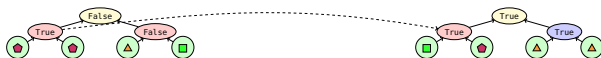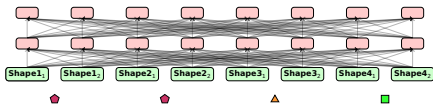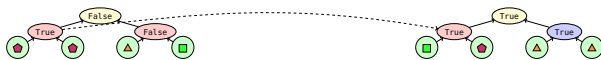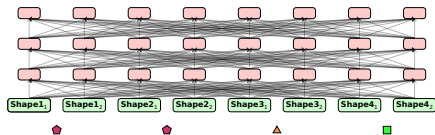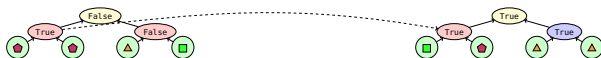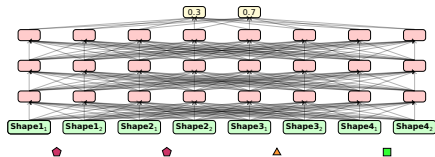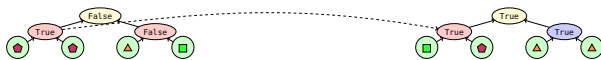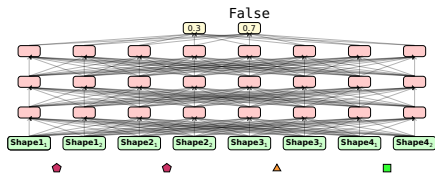
# Deep Learning Models

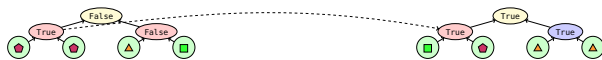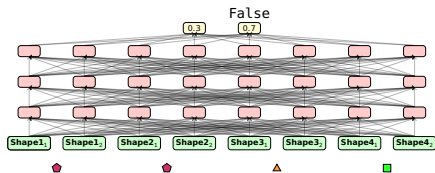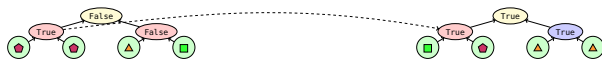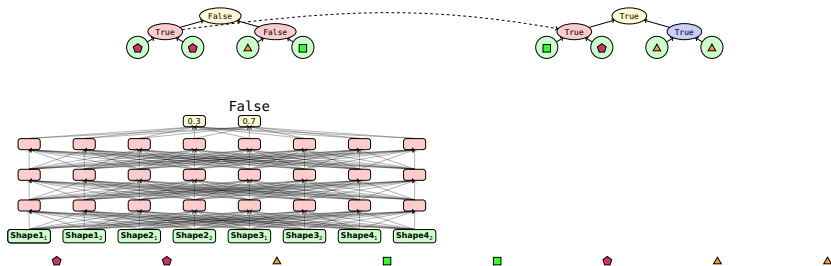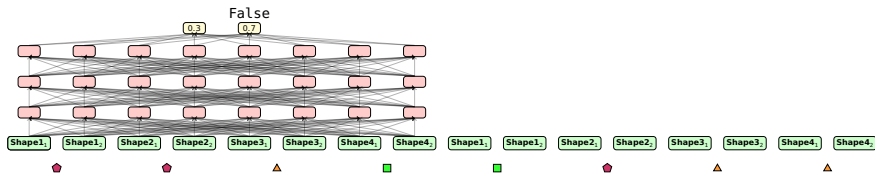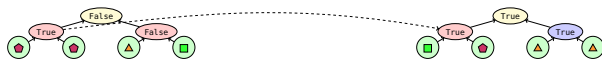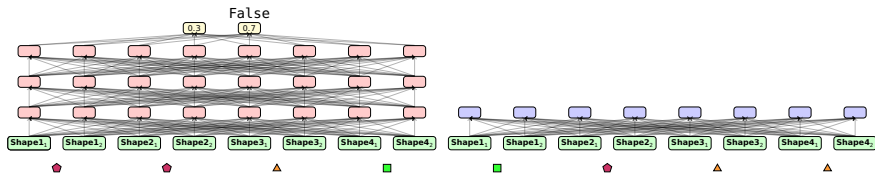# Constructive Causal Abstraction

# Alignment

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions
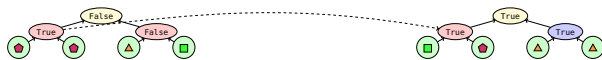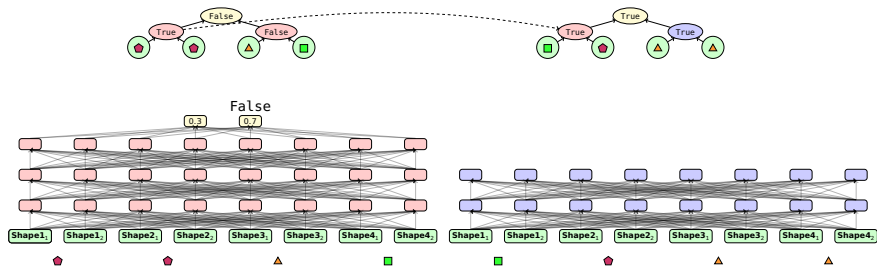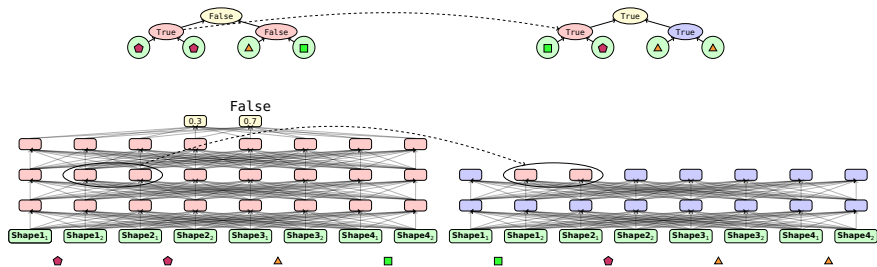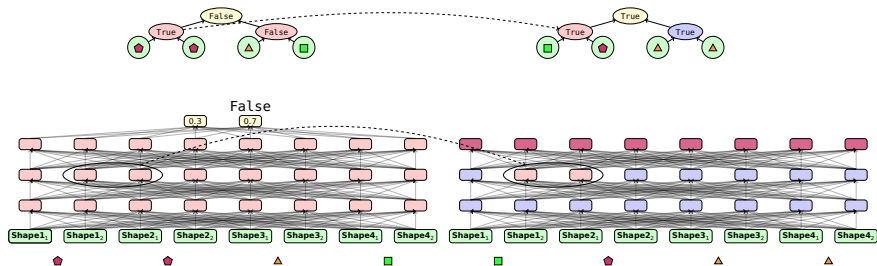
# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions
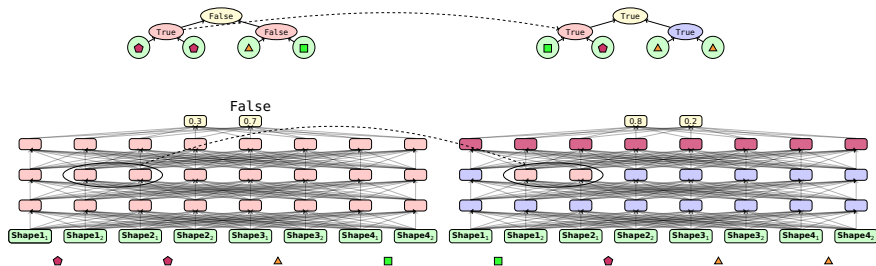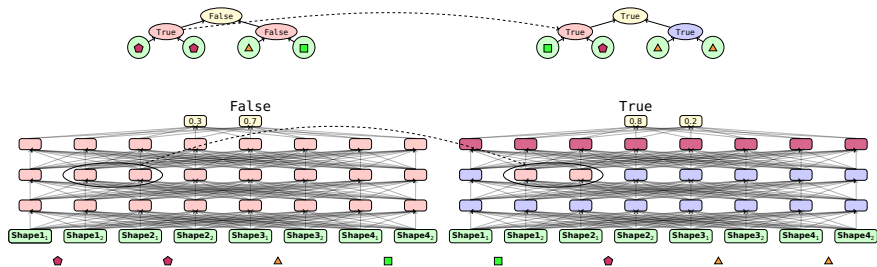
# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions

# Interchange Interventions
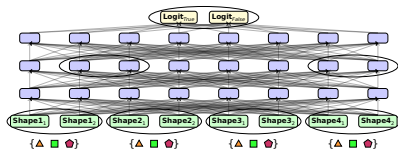
# Interchange Interventions
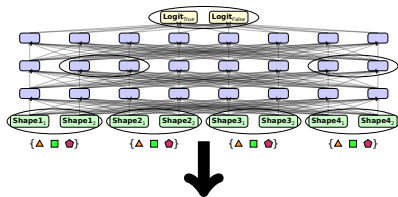
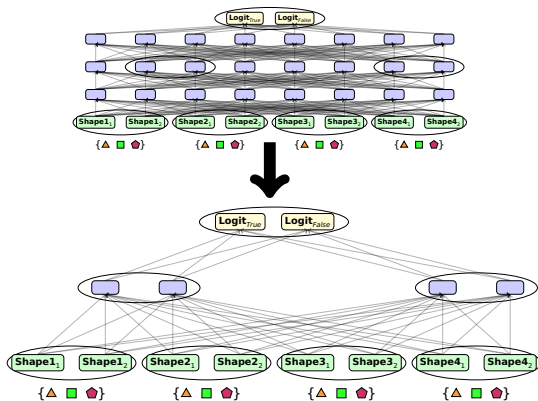# Interchange Interventions

# Interchange Interventions
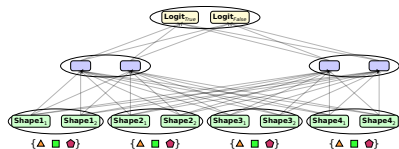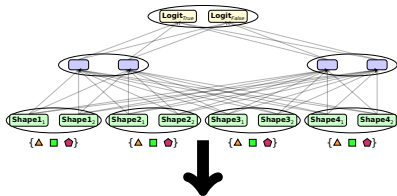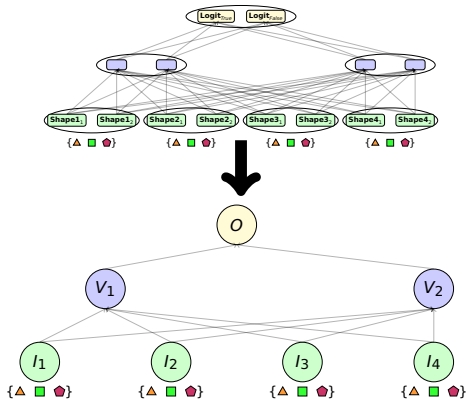
# Interchange Interventions

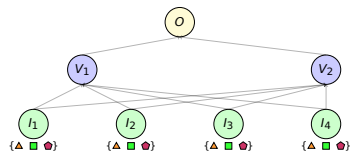# Marginalize

# Marginalize
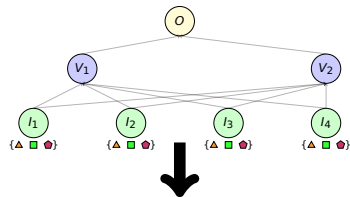
# Marginalize

# Variable Merge

# Variable Merge

# Variable Merge

# Value Change

# Value Change

# Value Change