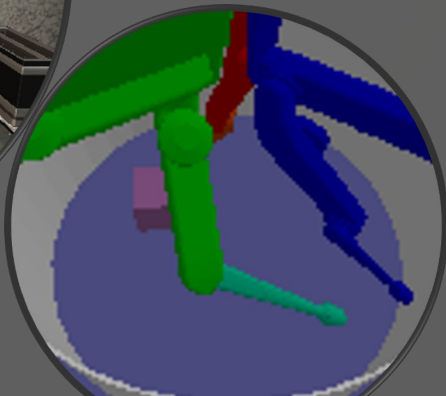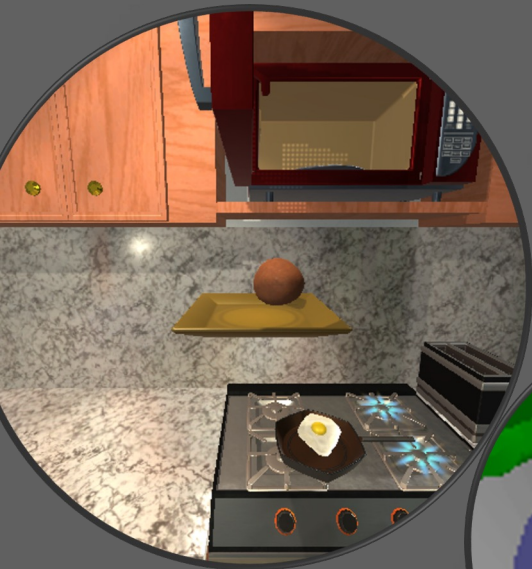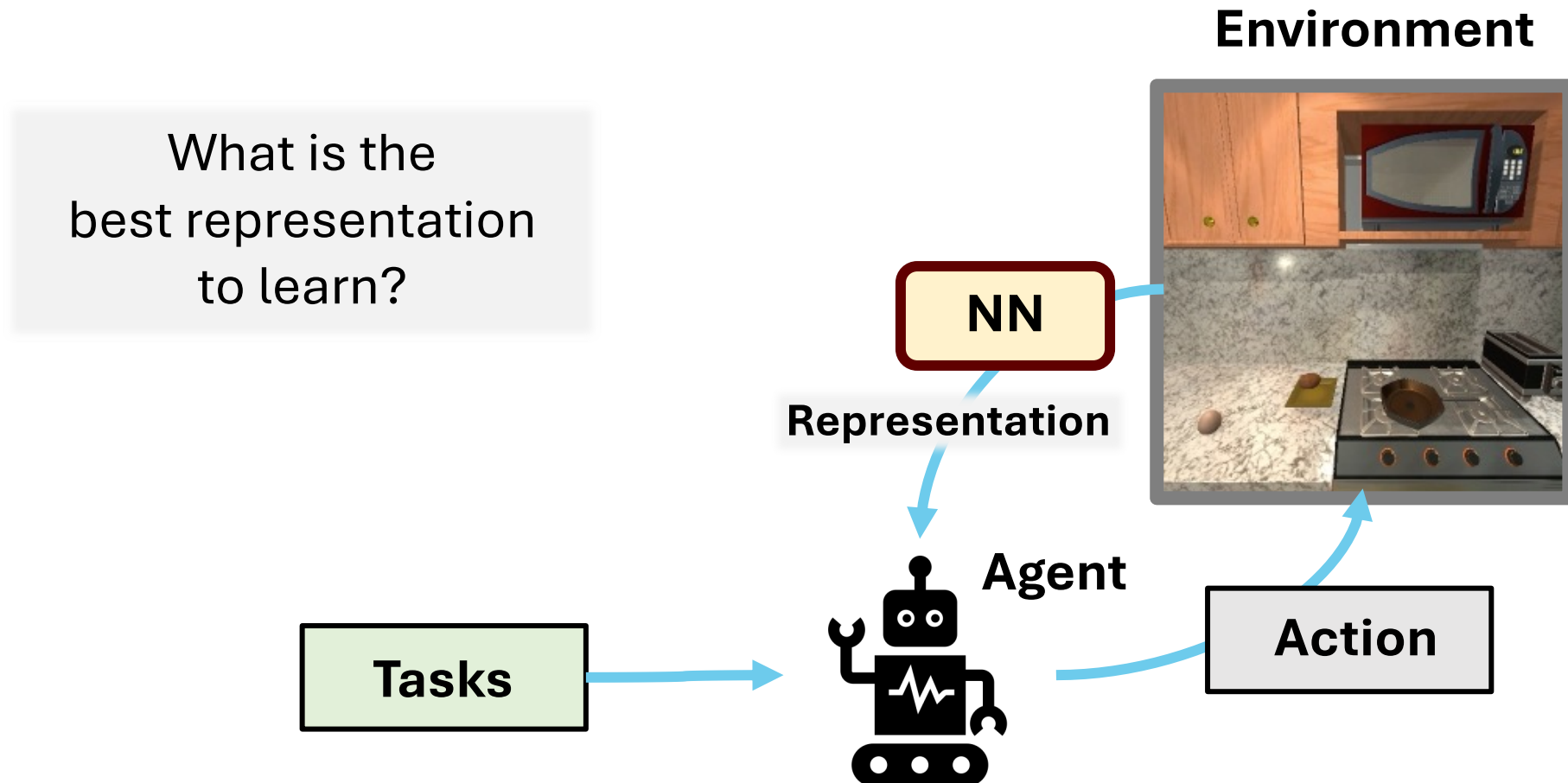# BISCUIT: Causal Representation Learning from Binary Interactions

*Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, Efstratios Gavves*

*July 18, 2023*

# Problem Setup

**Environment**

What is the
best representation
to learn?

**NN**

**Representation**

**Agent**

**Action**
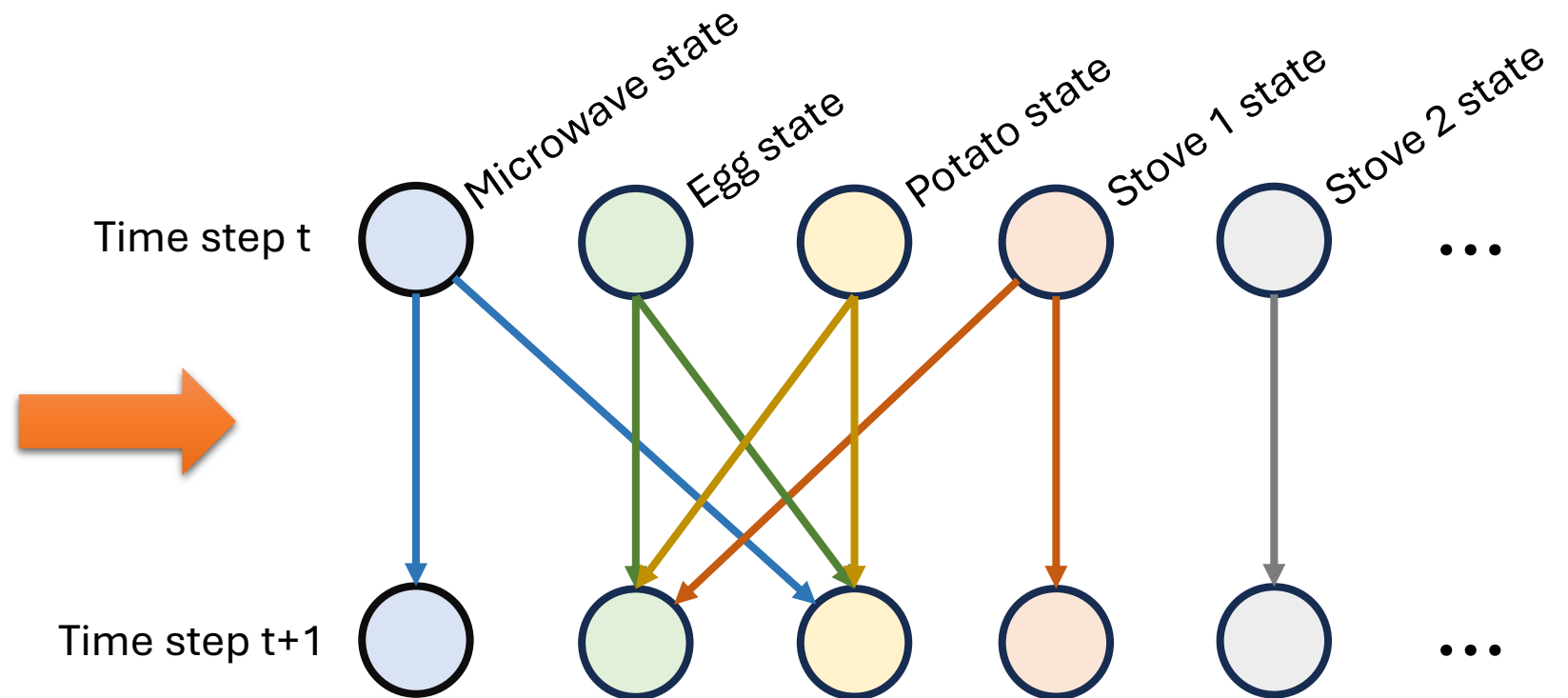
**Tasks**

# Causal Representation Learning
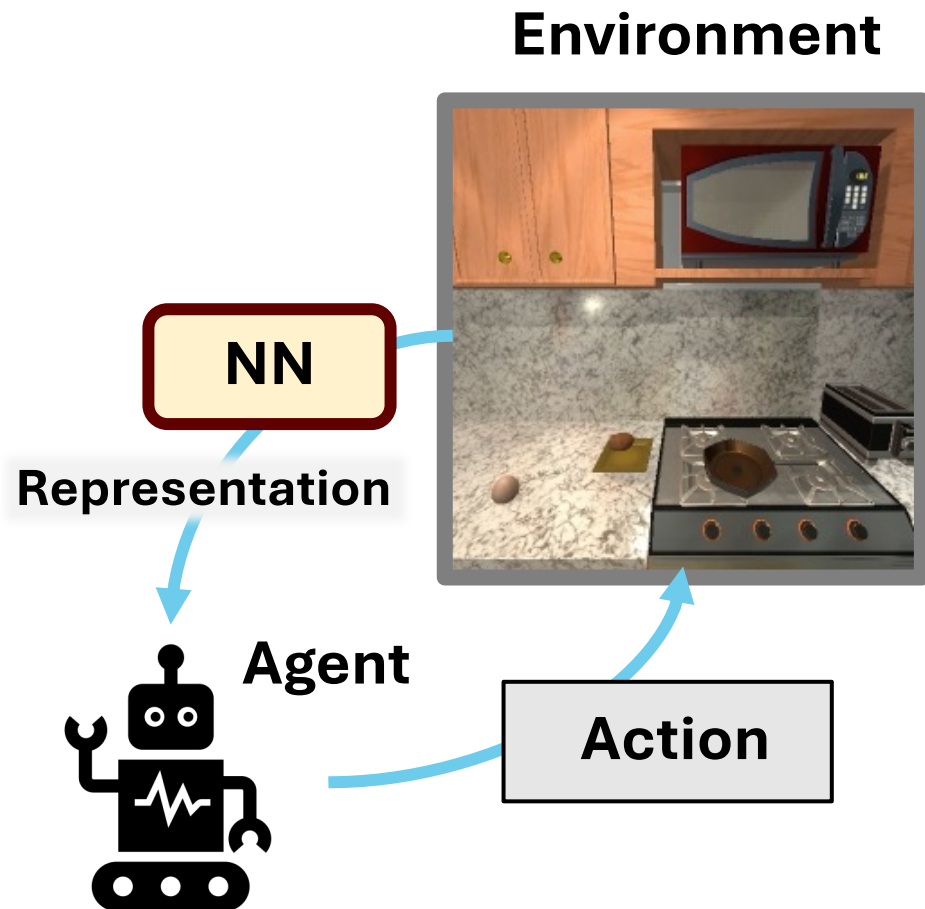


*Dense Representation*

**Interpretable?**

**Generalizable?**

**Reasoning-oriented?**

# Causal Representation Learning



**Goal of Causal Representation Learning**

# Causal Representation Learning

## Environment



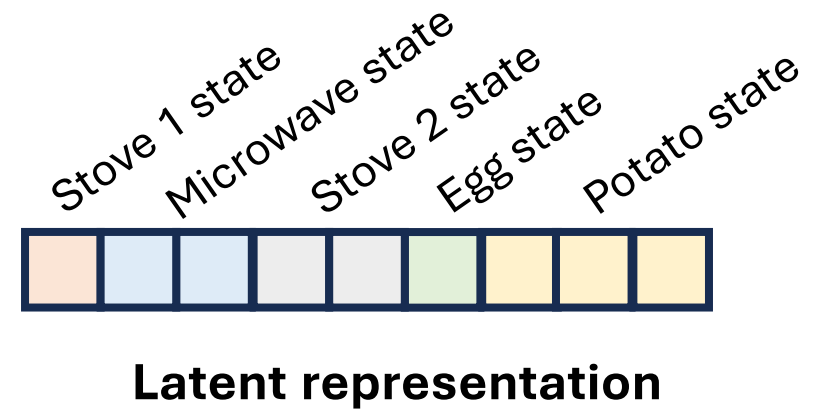**NN**

Representation

**Agent**

**Action**

## Representation Learning Tasks
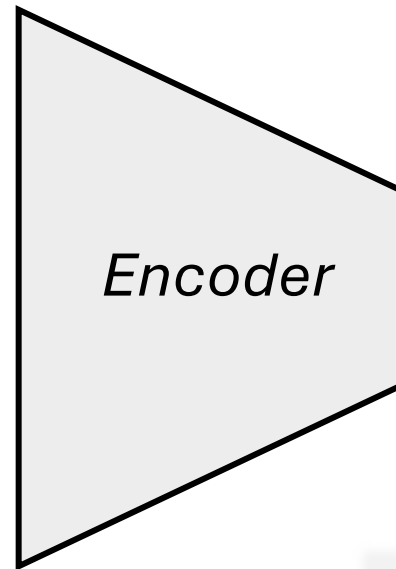
What are the causal variables of the environment?

How do they interact with each other?

How can the agent intervene on causal variables?

# Challenges in Causal Representation Learning



*Encoder*

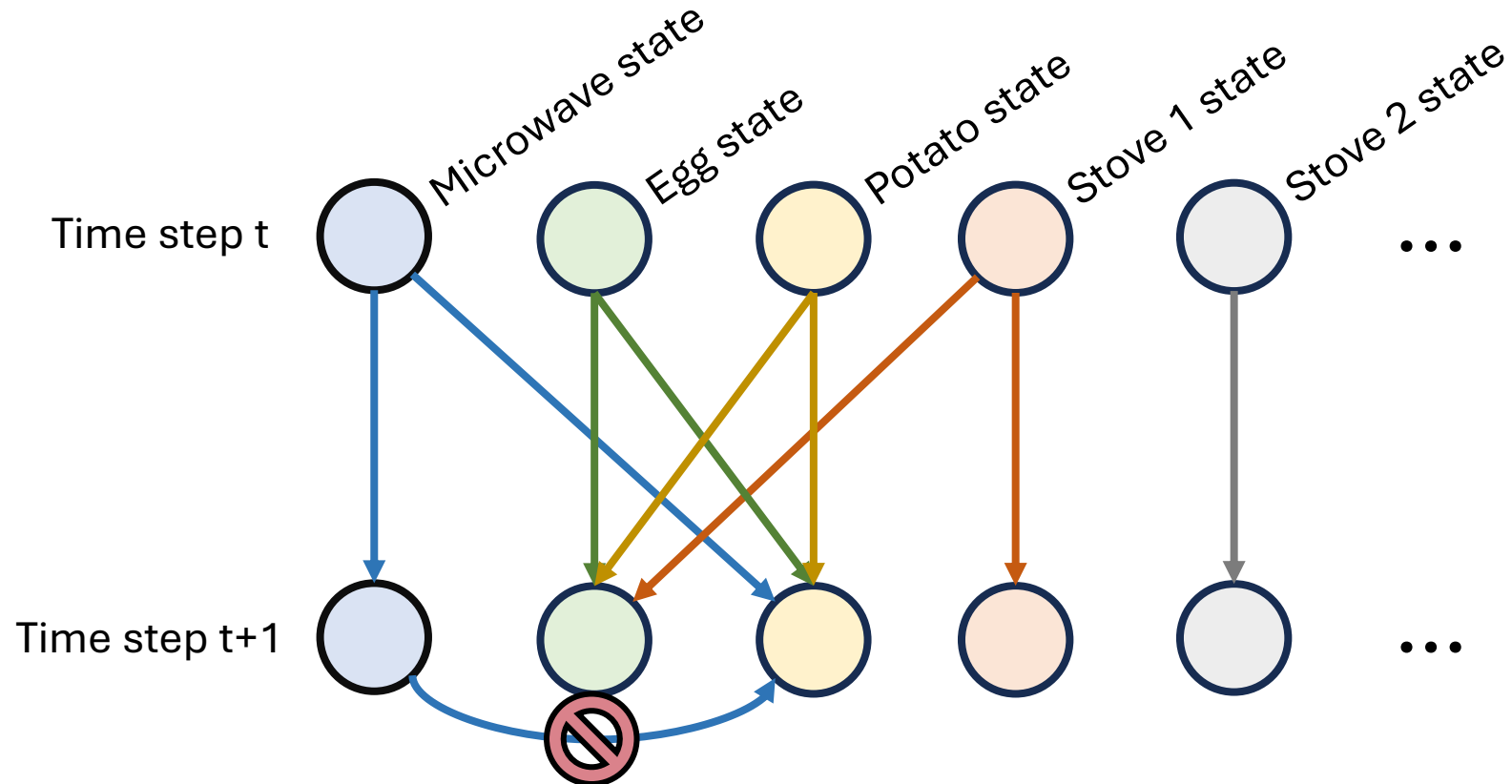Stove 1 state   Microwave state   Stove 2 state   Egg state   Potato state

**Latent representation**

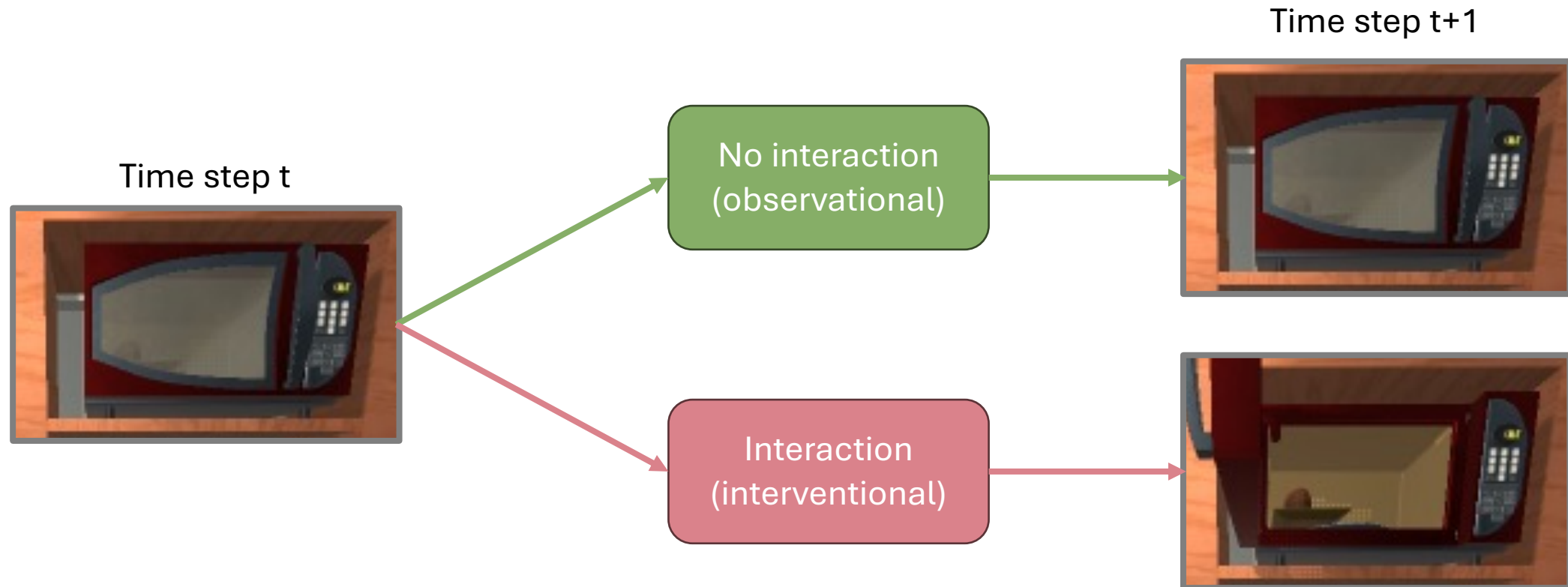How can we ensure that the causal variables are identified in the latent space?

# 🍪 BISCUIT: Learning Causal Representations

Assumption 1: Causal Relations can be resolved over time

# 🍪 BISCUIT: Learning Causal Representations

Assumption 2: Interactions between the agent and causal variables can be described
by **binary variables**



Time step t

No interaction
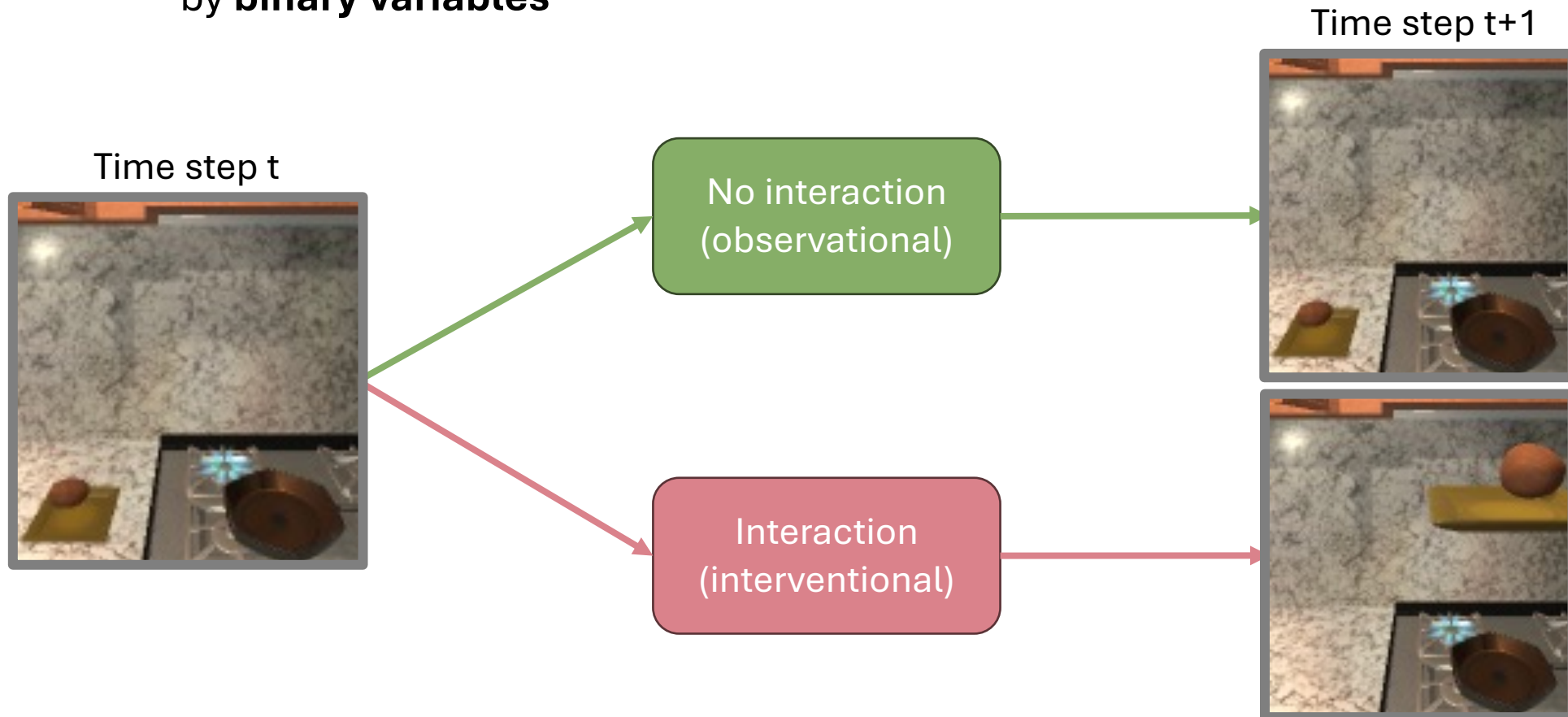(observational)

Interaction
(interventional)

Time step t+1

# 🍪 BISCUIT: Learning Causal Representations

Assumption 2: Interactions between agent and causal variables can be described
by **binary variables**



Time step t

No interaction
(observational)

Interaction
(interventional)

Time step t+1

# 🍪 BISCUIT: Theoretical Results

Assumption 1: Causal Relations can be resolved over time

Assumption 2: Interactions between agent and causal variables can be described
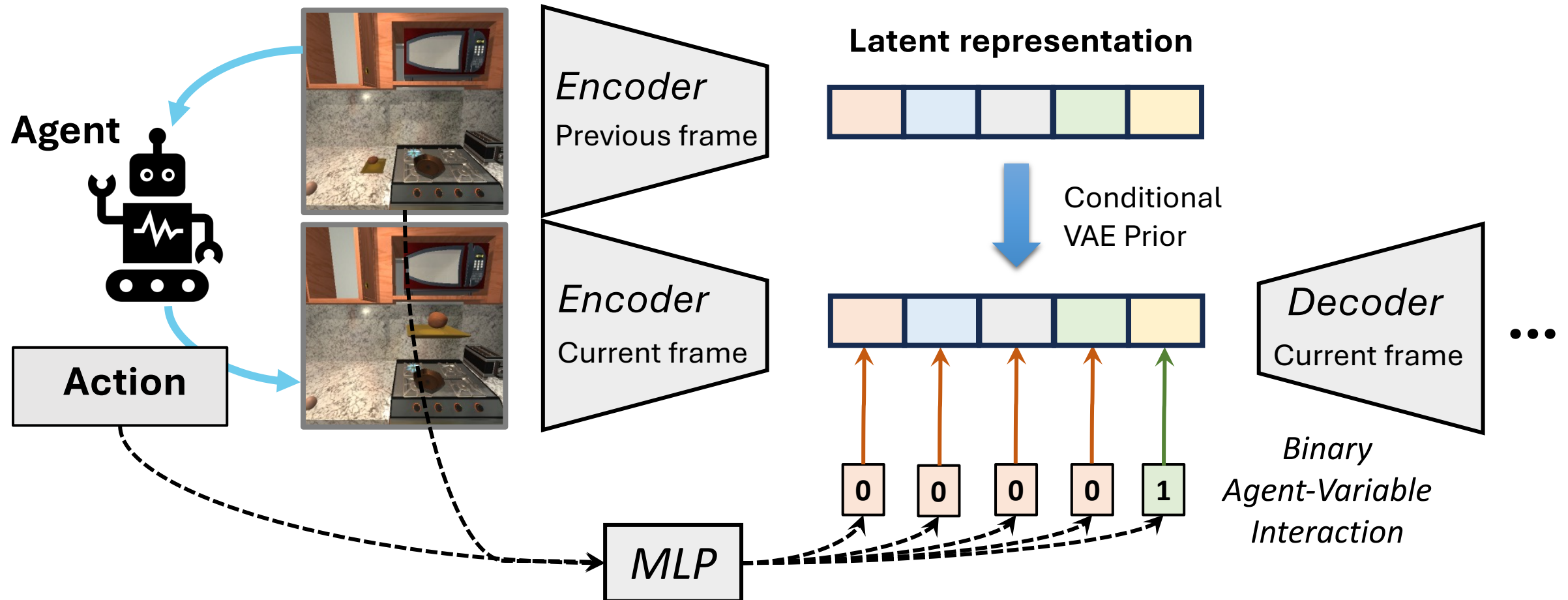by **binary variables**

Assumption 3: All causal variables have different interaction patterns

Assumption 4: The causal mechanisms need to sufficiently vary on *interventions* or *over time*
(allows for additive Gaussian noise models)

**BISCUIT Theoretical Result**
Under these assumptions, causal variables can be identified from videos with low-level actions.
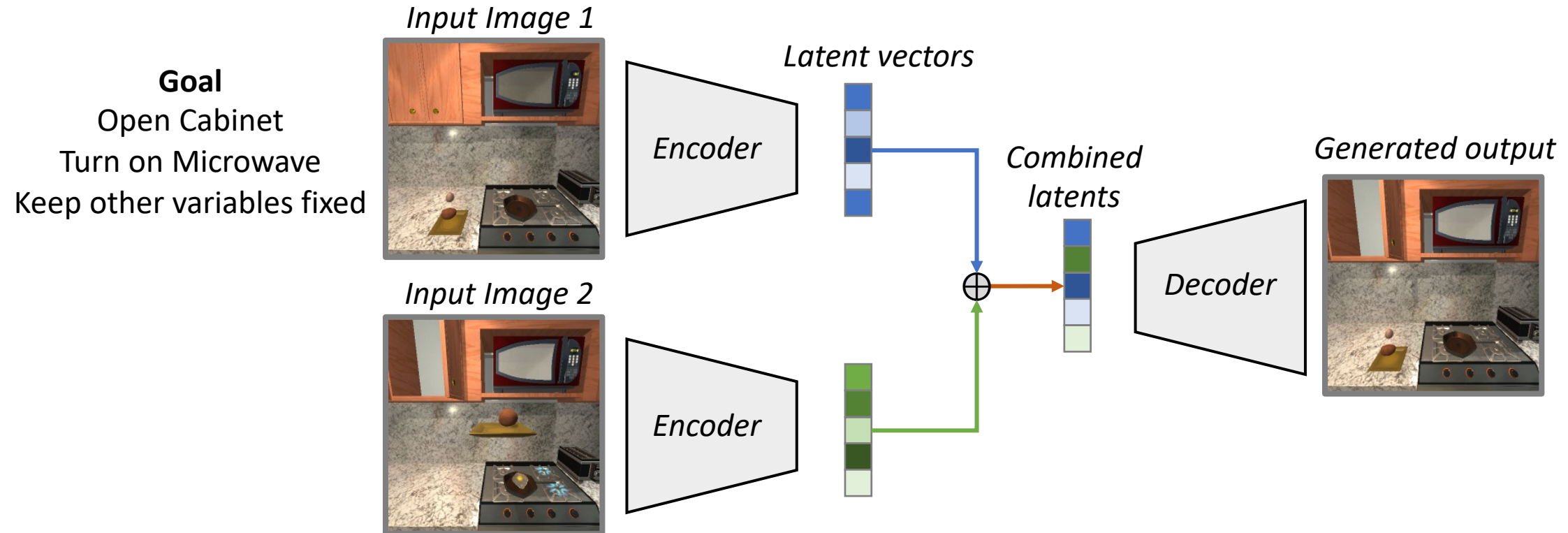
# 🍪 BISCUIT: Architecture

# **Experiments – iTHOR**

- Kitchen environment with 10 causal variables
  - Cabinet (open/closed)
  - Microwave (open/closed)
  - Microwave (on/off)
  - Egg (position, broken, cooked)
  - Potato (position)
  - 4x Stove burner (on/off, burning)
  - Toaster (on/off)

- Close-to random policy

- Actions represented as x-y coordinate of a randomly sampled object pixel

# iTHOR – Simulate Latent Interventions

- BISCUIT accurately identifies causal variables

- Validated by performing interventions in latent space



**Goal**
Open Cabinet
Turn on Microwave
Keep other variables fixed

# iTHOR – Simulate Latent Interventions



Input image 1     Input image 2     Generated Output     Latents from image 2

Microwave Open

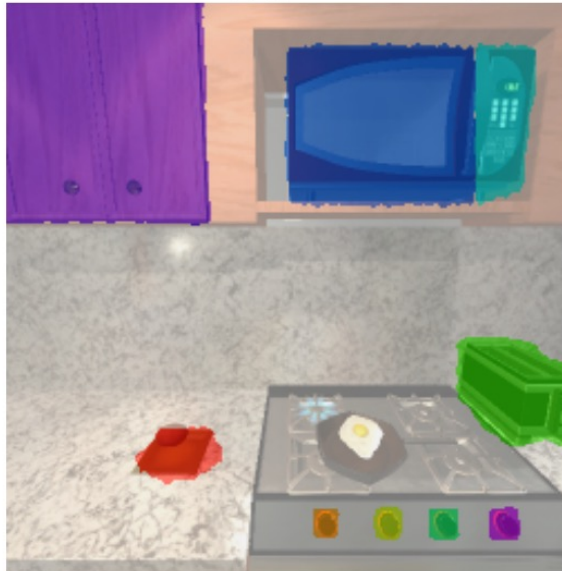# iTHOR – Simulate Latent Interventions

# iTHOR – Interaction Maps

- Visualize learned interaction variables by the x-y locations they are active

- Each causal variable shown in different color

# CausalWorld – Robotic Trifinger

- Tri-finger robot interacting with its environment and objects

  - Causal variables include object position, frictions, colors, etc.

- Action: 9-dimensional motor angles (3 per finger)

- BISCUIT identifies causal variables accurately

Accuracy of learned causal variables
(higher is better / lower is better)

| Models | CausalWorld |
|---|---|
| iVAE (Khemakhem et al., 2020a) | 0.28 / 0.00 |
| LEAP (Yao et al., 2022b) | 0.30 / 0.00 |
| DMS (Lachapelle et al., 2022b) | 0.32 / 0.00 |
| BISCUIT-NF (Ours) | **0.97** / 0.01 |

# CausalWorld – Learned Interactions



**F1 scores for learned interaction variables**

| Learned Interaction Variables | Finger 1 - Color | Finger 2 - Color | Finger 3 - Color | Floor Friction | Stage Friction | Cube Friction | Cube State |
|---|---|---|---|---|---|---|---|
| Finger 1 - Color | 45.1 | 7.1 | 8.9 | 5.2 | 4.8 | 3.5 | 16.6 |
| Finger 2 - Color | 6.2 | 47.2 | 8.6 | 4.8 | 5.1 | 3.1 | 24.7 |
| Finger 3 - Color | 8.5 | 6.6 | 50.1 | 3.5 | 3.6 | 3.9 | 20.2 |
| Floor Friction | 4.3 | 3.9 | 4.8 | 94.8 | 3.4 | 3.9 | 4.1 |
| Stage Friction | 4.4 | 5.4 | 3.6 | 4.5 | 96.8 | 4.8 | 3.1 |
| Cube Friction | 4.8 | 3.5 | 3.2 | 5.8 | 5.9 | 93.2 | 5.4 |
| Cube State | 18.0 | 16.0 | 21.8 | 4.3 | 3.4 | 4.5 | 72.1 |

Ground Truth Interaction Variables

# **Conclusion**

- BISCUIT identifies causal variables from interactive environments

- Key assumption: binary interaction variables describe agent-causal variable interactions

- Applicable to a variety of robotic and embodied AI environments

- Ability to 'imagine' by performing latent interventions

- Identifies actions to perform interventions

Project website and demo: [phlippe.github.io/BISCUIT/](phlippe.github.io/BISCUIT/)