



BISCUIT: Causal Representation Learning from Binary Interactions

Phillip Lippe¹, Sara Magliacane^{2,3}, Sindy Löwe², Yuki M. Asano¹, Taco Cohen⁴, Efstratios Gavves¹

¹University of Amsterdam, QUVA lab, ²University of Amsterdam, AMLab, ³MIT-IBM Watson AI Lab, ⁴Qualcomm AI Research, The Netherlands

What is BISCUIT?

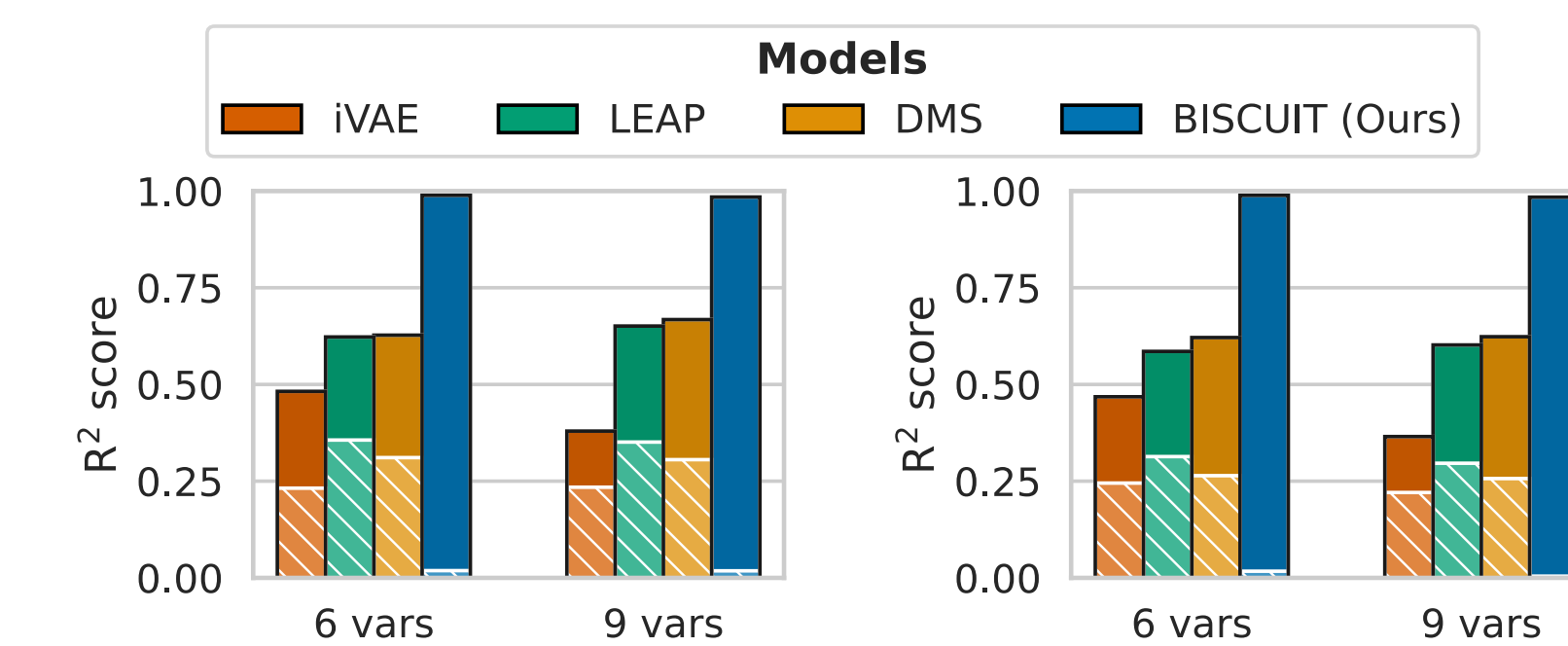
- BISCUIT learns **causal representations** from videos of **interactive systems**
 - Example: identify the causal variables (e.g. microwave state, plate position, etc.) of the kitchen environment
- Key assumption: interactions between agent and a causal variable can be described by a **binary interaction variable**
 - Interventional (e.g. open microwave) vs observational



Experiments

- Evaluating accuracy of identifying causal variables from high-dimensional videos
- Actions being clicks or robotic input

Synthetic Environment (additive Gaussian noise)



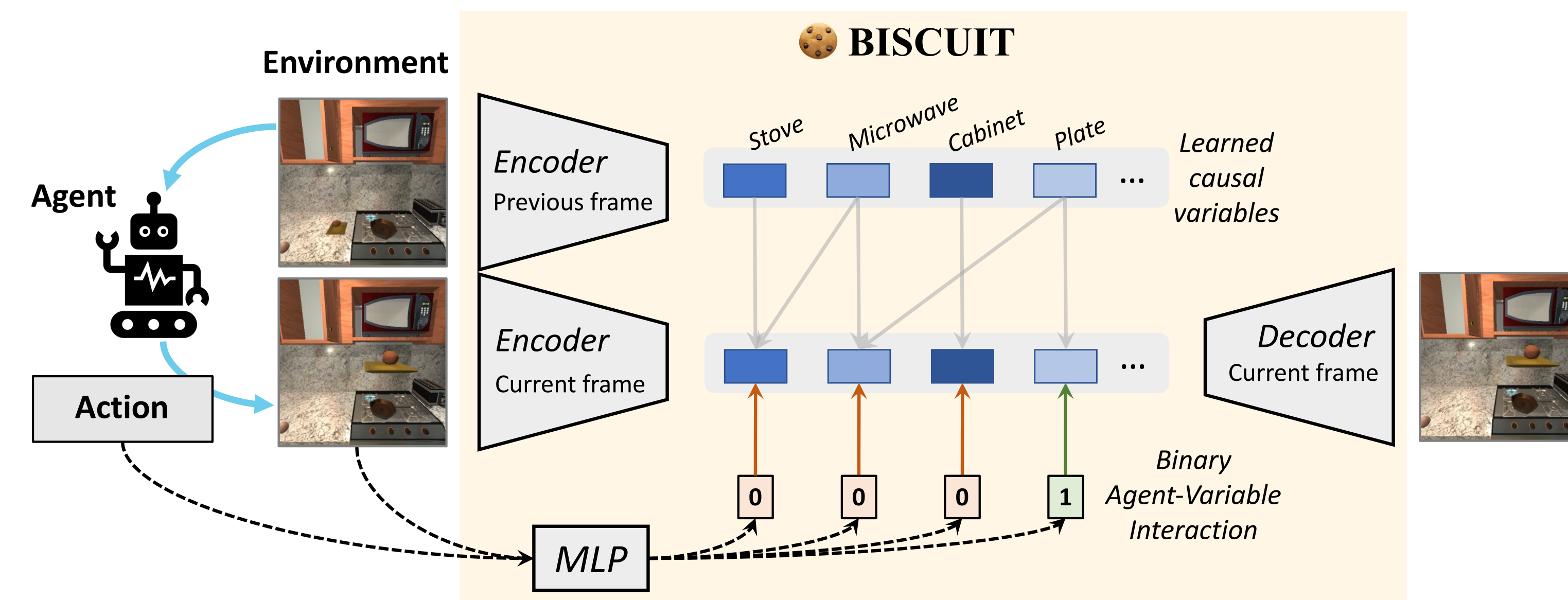
Robotic Environments

Table 1: R^2 scores (diag \uparrow / sep \downarrow) for the identification of the causal variables on CausalWorld and iTHOR.

Models	CausalWorld	iTHOR
iVAE (Khemakhem et al., 2020a)	0.28 / 0.00	0.48 / 0.35
LEAP (Yao et al., 2022b)	0.30 / 0.00	0.63 / 0.45
DMS (Lachapelle et al., 2022b)	0.32 / 0.00	0.61 / 0.40
BISCUIT-NF (Ours)	0.97 / 0.01	0.96 / 0.15

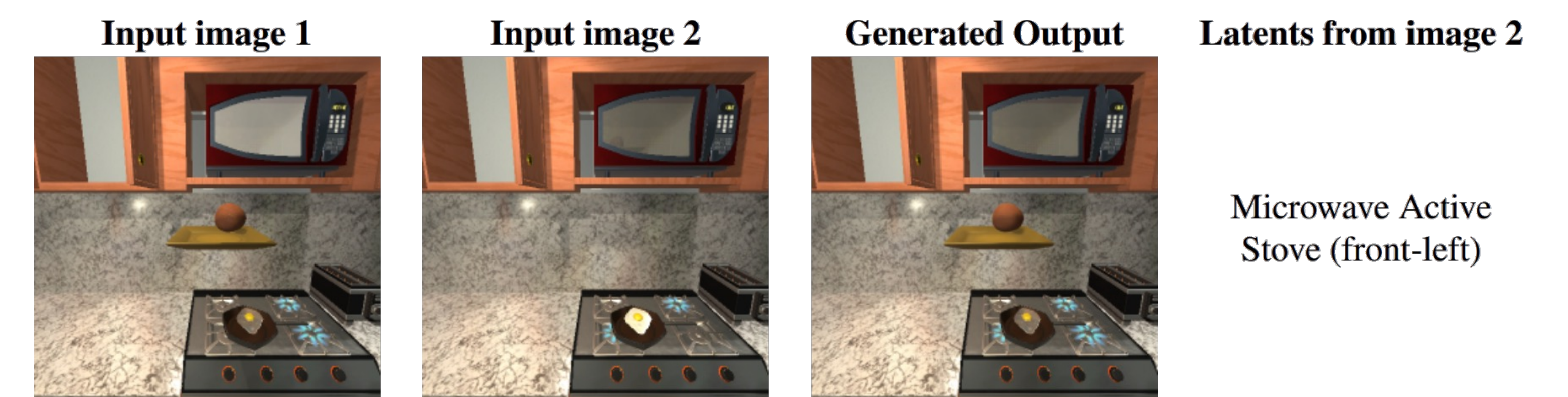
How does BISCUIT work?

- Temporal VAE with causal vars in latent space and MLPs learning interaction vars
- Alternative setup: normalizing flow applied on autoencoder representation

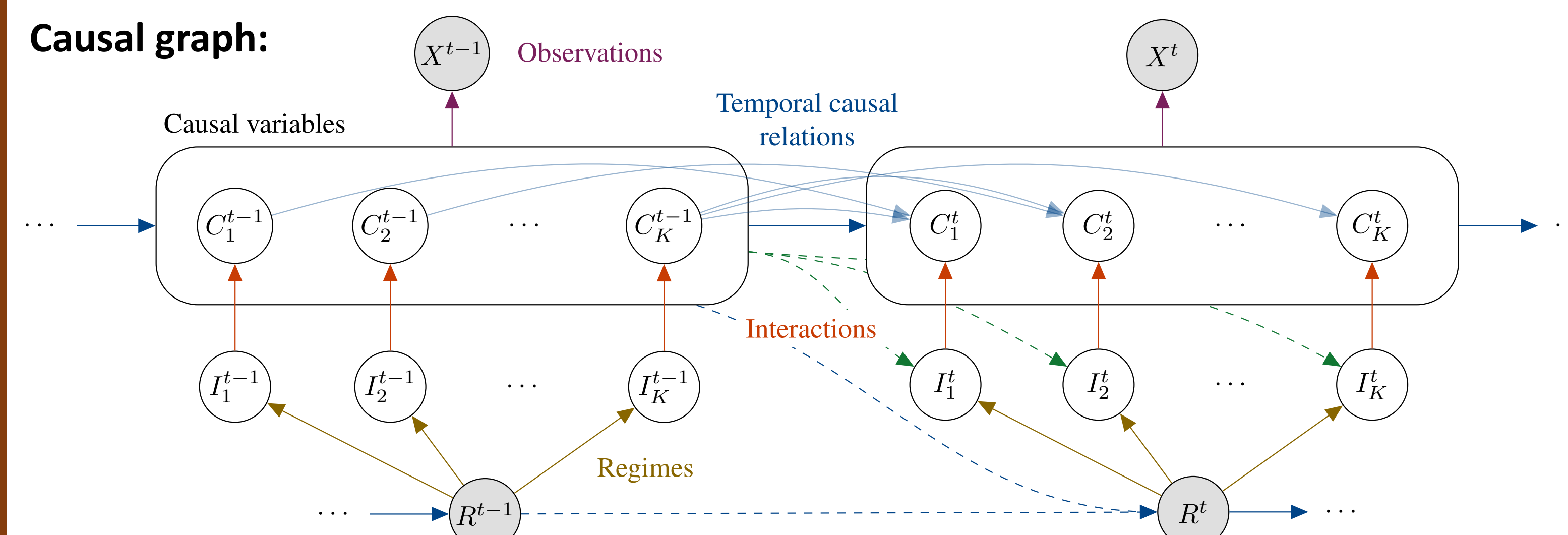


Simulating Interventions in Latent Space

- Latent interventions by (1) encoding two images, (2) replacing latents of first image by latents of second image for respective causal variables, (3) decoding new latents
- Achieves novel combinations of causal vars, e.g. uncooked egg on burning stove



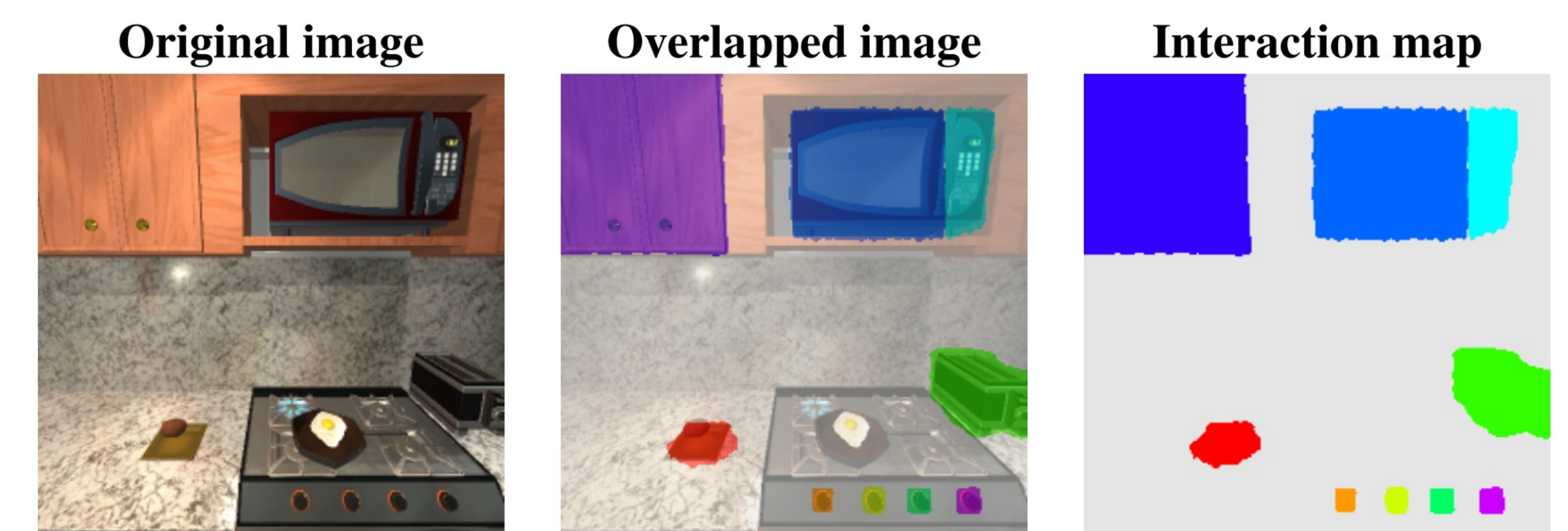
BISCUIT provably identifies causal variables



- Assumption 1: each causal variable has a distinct interaction pattern
- Assumption 2: mechanisms sufficiently vary on intervention or over time
- Allows for additive Gaussian noise models if mean changes over time

Interaction Maps

- In iTHOR, an action is a random x-y position of object interacted with
- Visualizing learned interaction variables for each causal variable segments objects



Paper and Demo