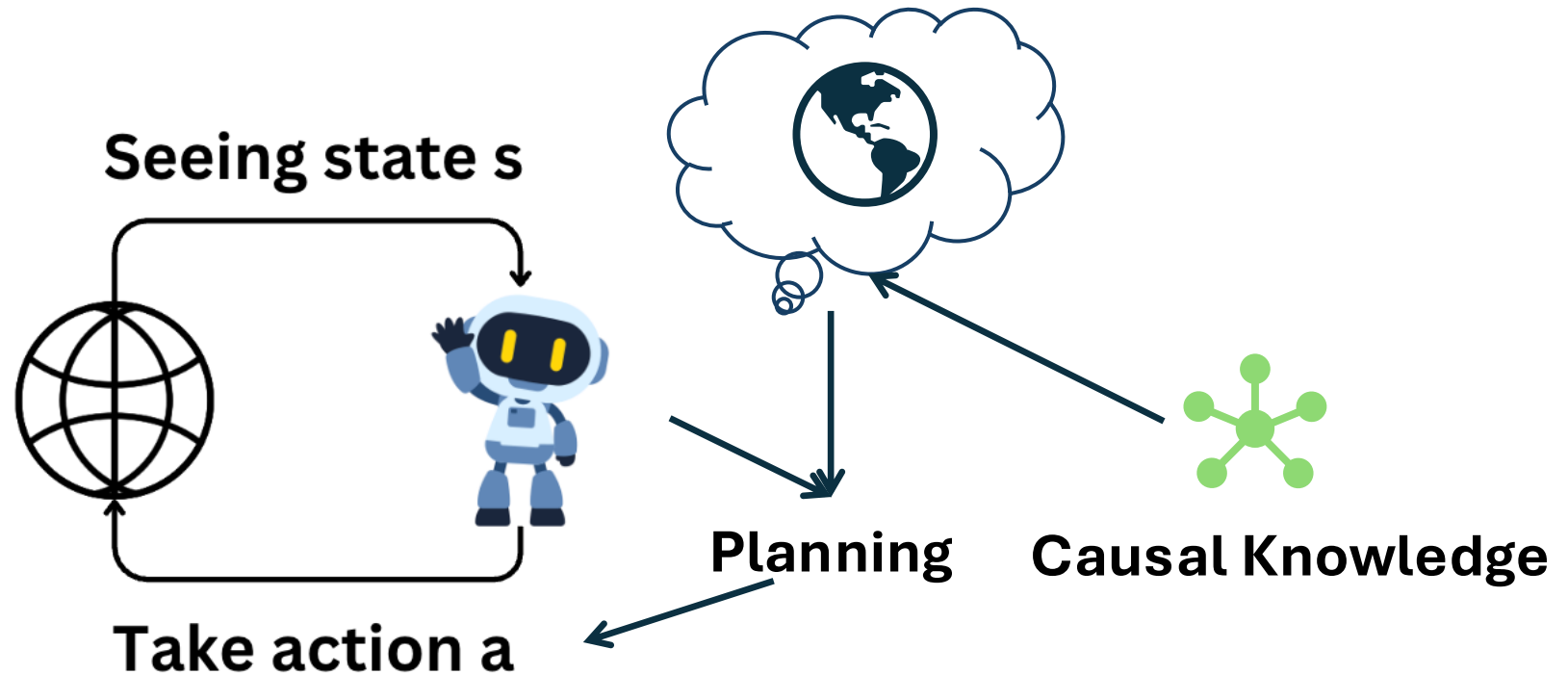


Towards Causal Foundation World Models

Mengyue Yang

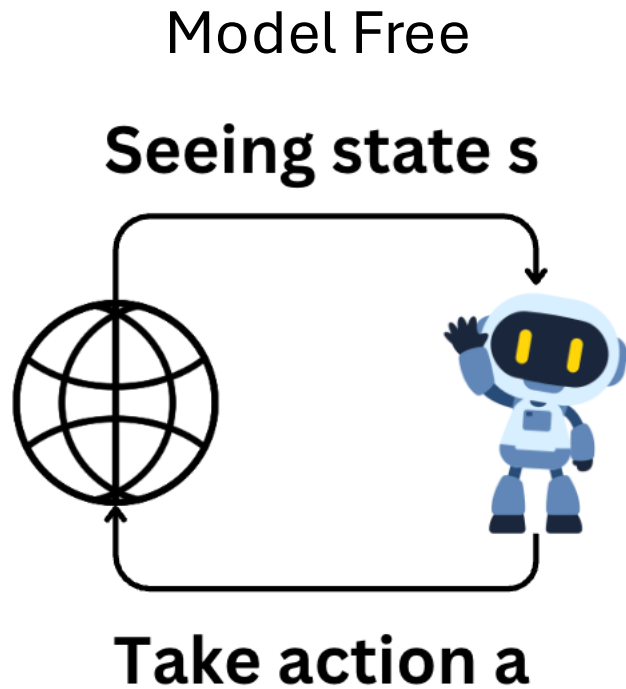
Lecturer in AI

University of Bristol

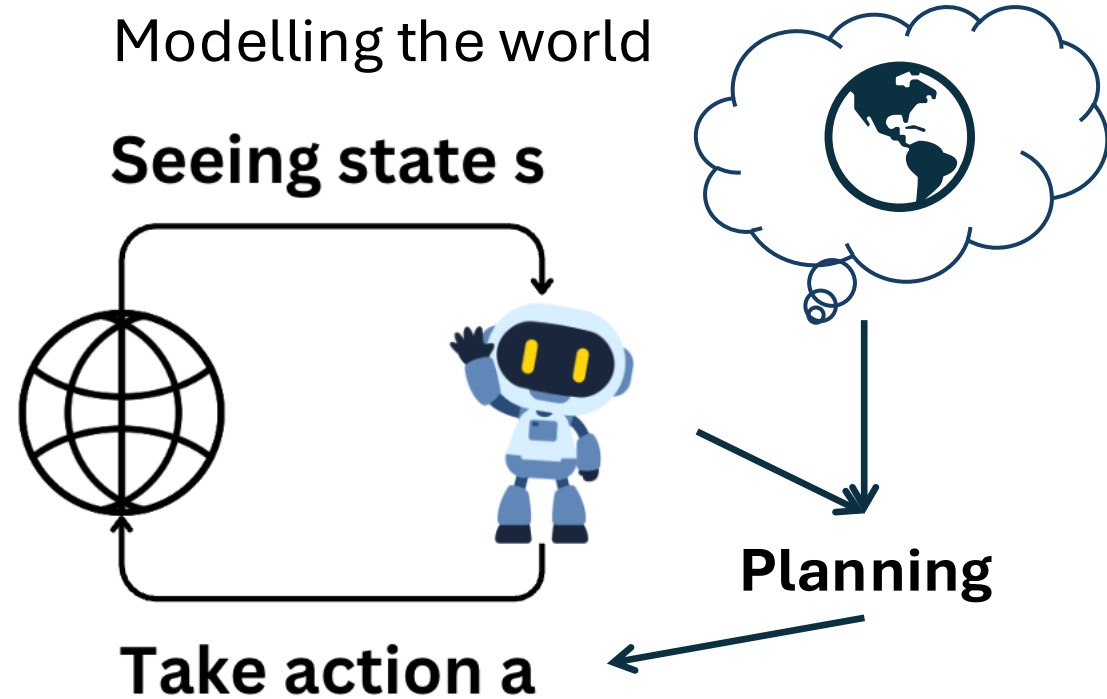


What is Agent Decision Making?

Agent policy $\pi(a|s)$: The way that agent to achieve the goal



Improve policy by exploration



Planning from the knowledge
and thinking

Why Modelling World is Necessary



Open World



New states • New mechanisms • Long tails



Agents



Must explore → learn → reuse



Data cannot cover all possibilities

World models enable “learning in imagination”, counterfactual prediction



Long-horizon control needs planning

Rollouts inside a world model support lookahead

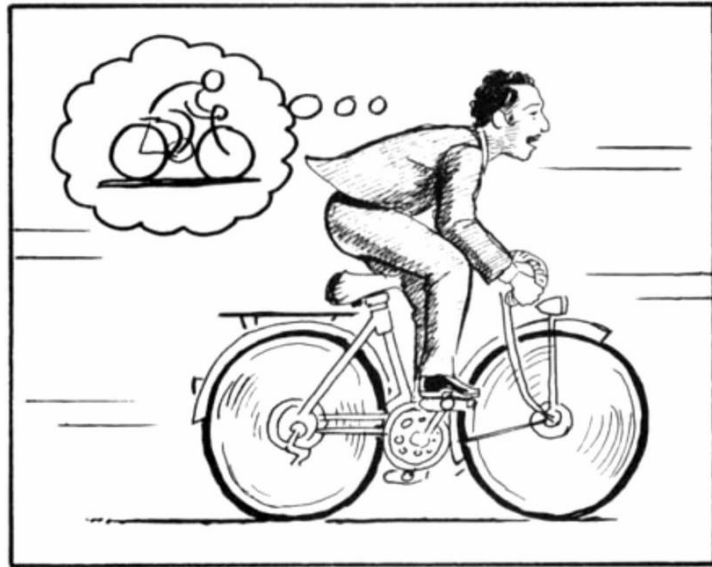


Exploration must be targeted, not random

“If you don’t explore it, you can’t learn it”, Model uncertainty and knowledge guides agent

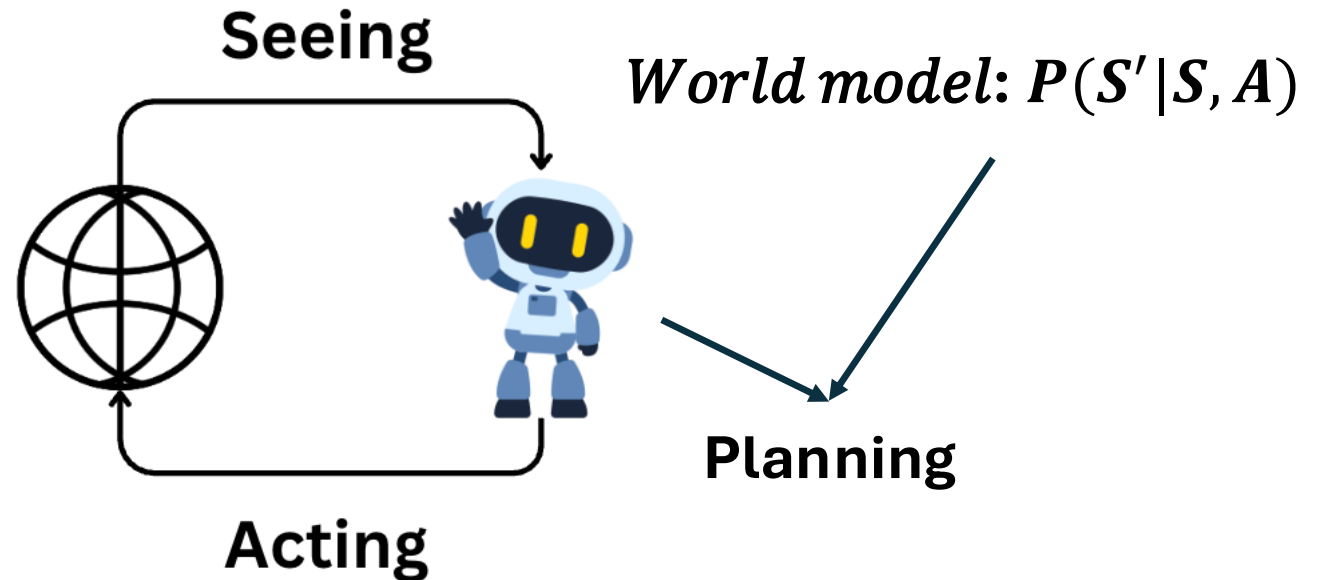
What is World Models

World Models modelling the **environment internal transition model**, makes agent **1. predict the future** and **2. reflect the past** possible.



What if I ride slower now? Can I achieve the goal?
Was it a good choice to start earlier?

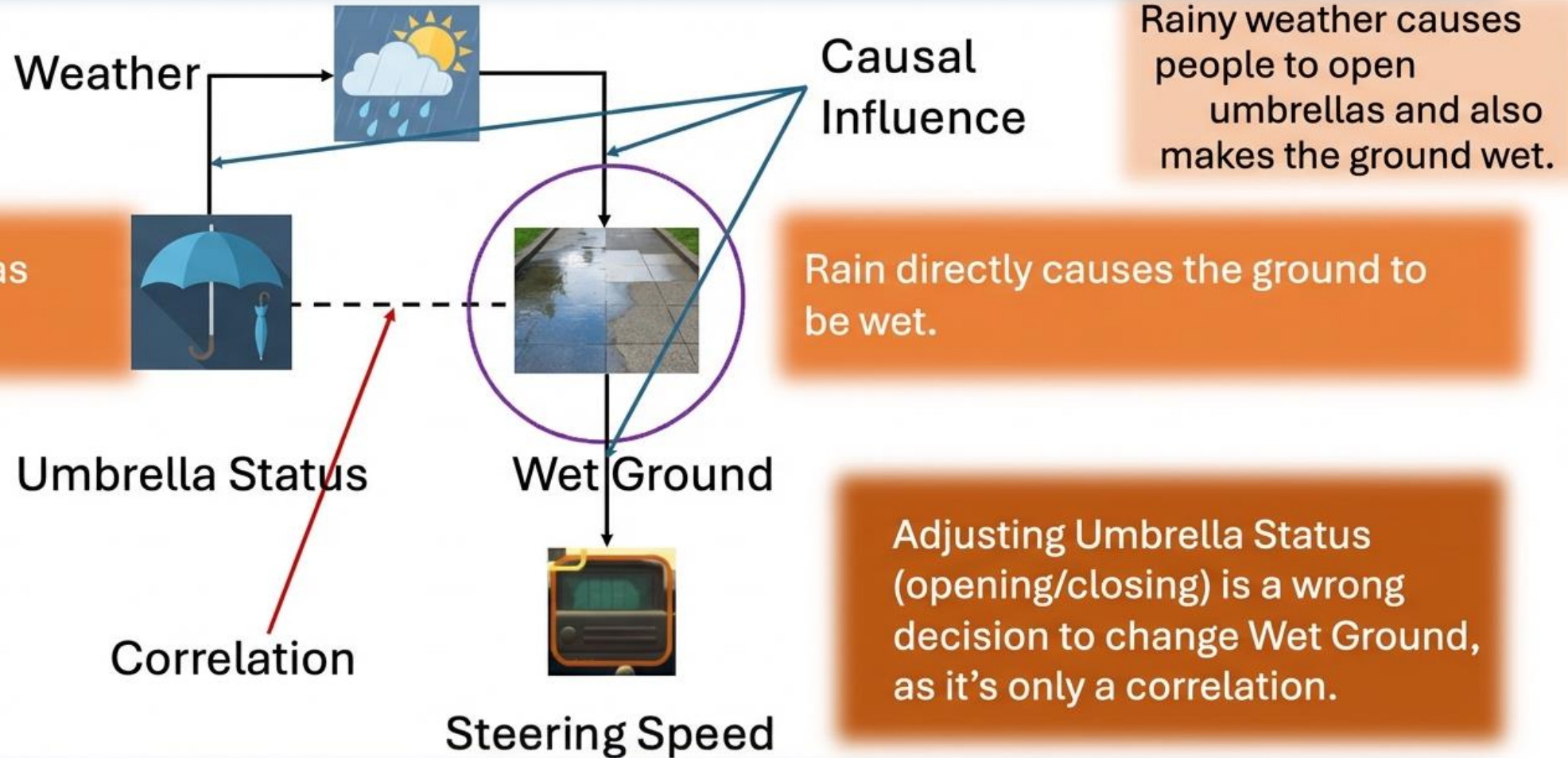
Picture borrow from Ha and Schmidhuber



Modelling the World \neq Understanding the World

Is current world model perfect?


Mimic the world data doesn't means model understand the world





Pearl's Causal Hierarchy

Predict something
haven't happened

Imagine based on something
already happened

Seeing
Association 

Doing
Intervention 

Imagining
Counterfactual 

Question:
What is?

Question:
What if?

Question:
Was it?

What does the
smoking tell us
about the lung
cancer.

What will happen
if someone keep
smoking

Would the lung
cancer got worse
if someone
smoking.

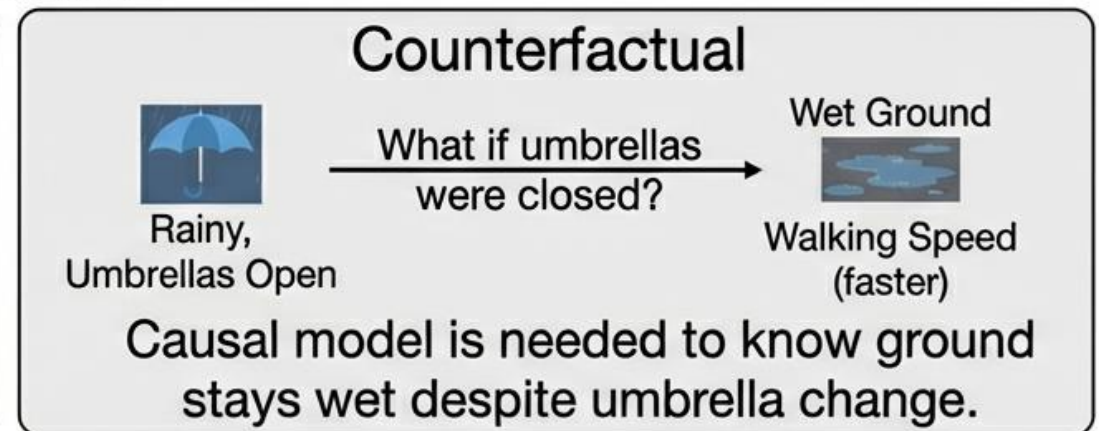
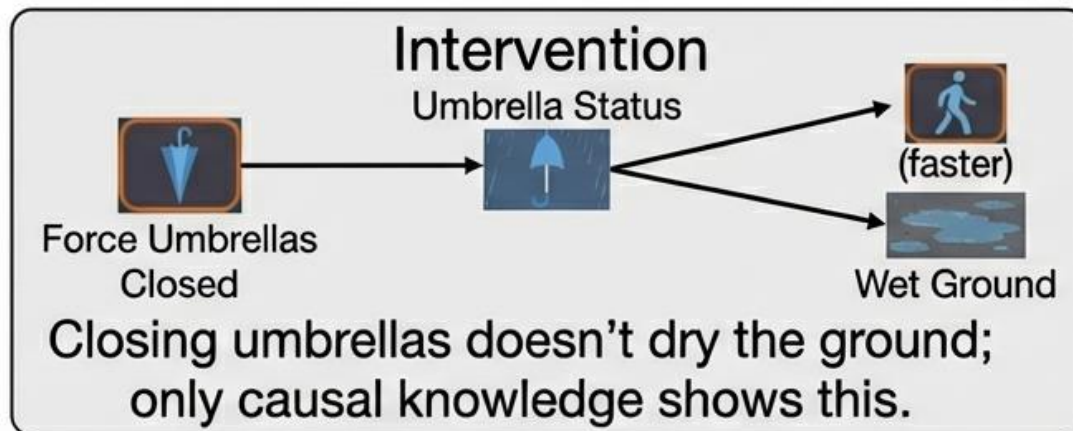
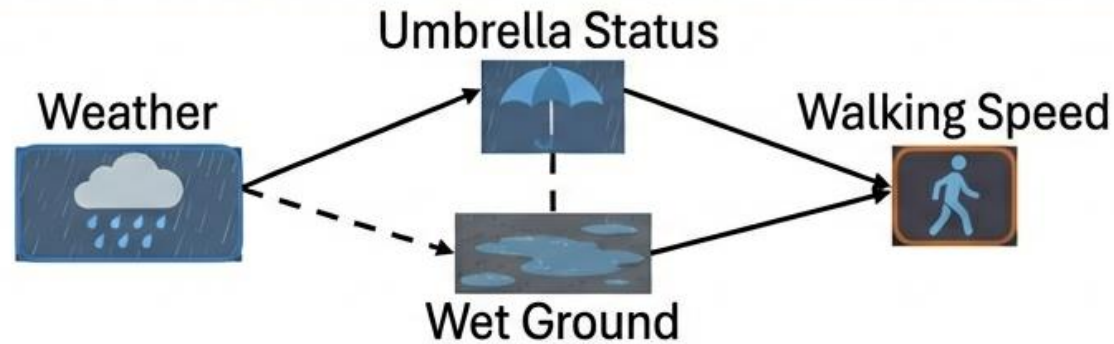
Predicting of future
Reflecting of past

=

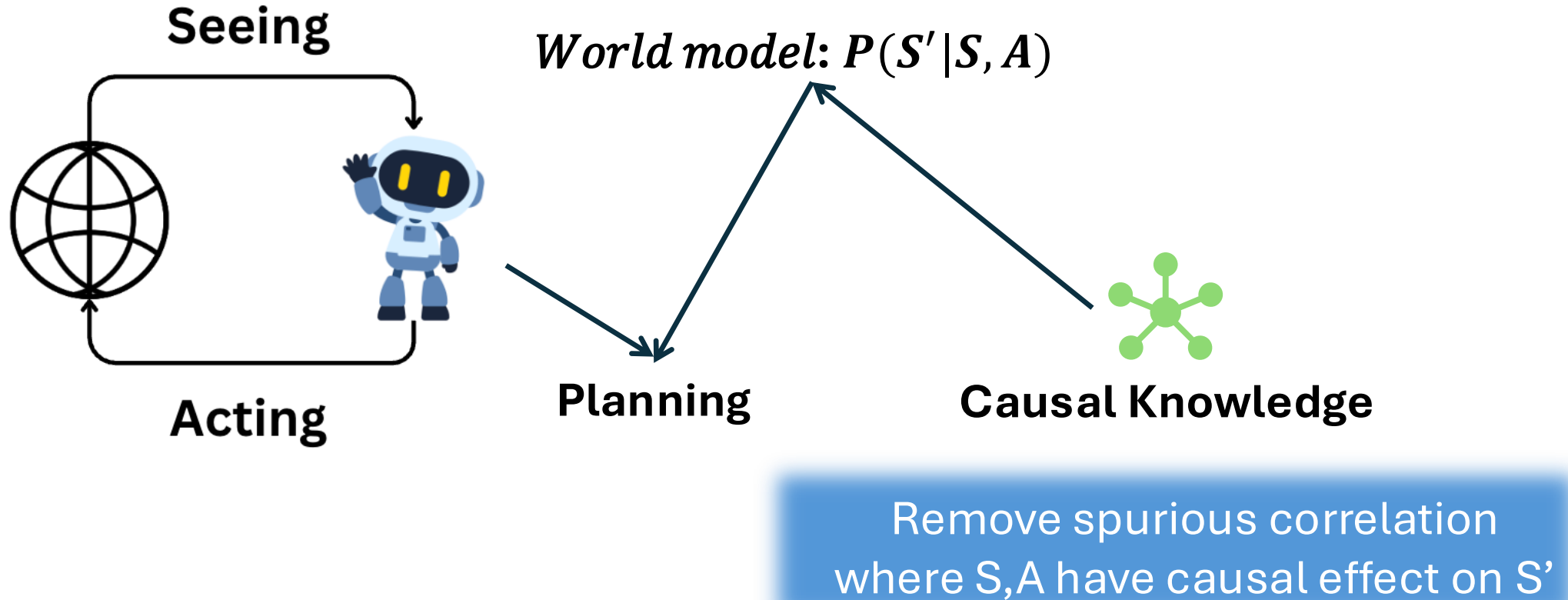
Intervention
Counterfactual

Causal Diagram

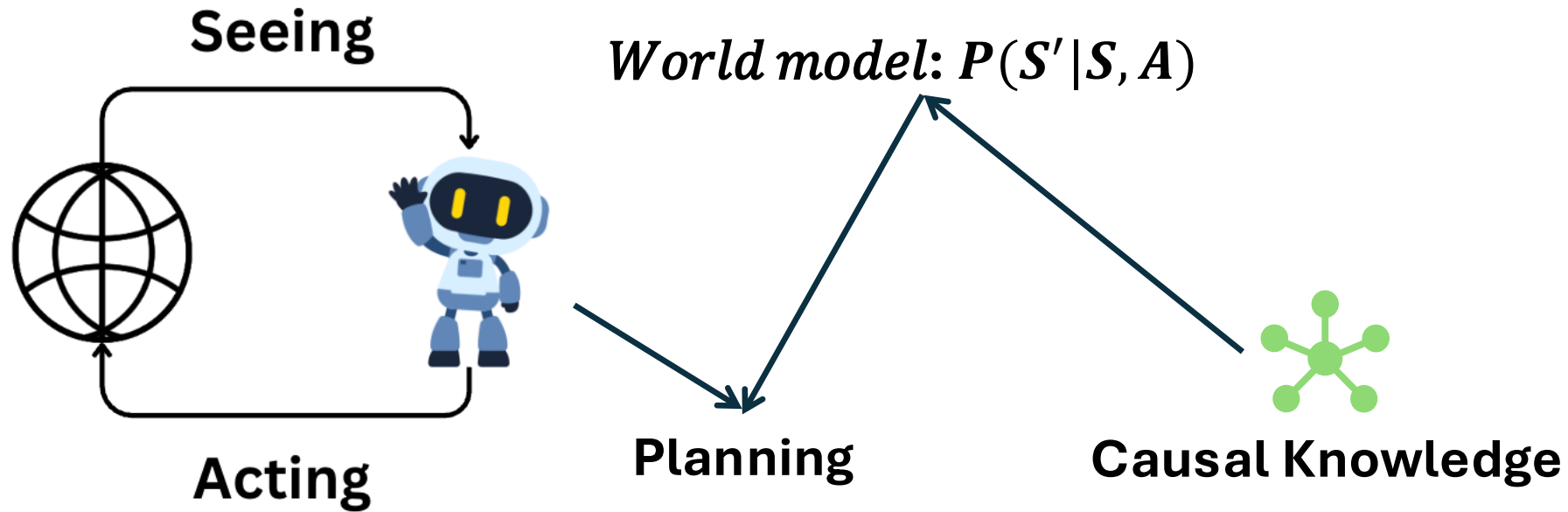
Understanding the underlying causal model is a prerequisite for inferring intervention and counterfactual



Causal World Models



Causal World Models



Focus planning by the true causal relations

Predicting of future
Reflecting of past

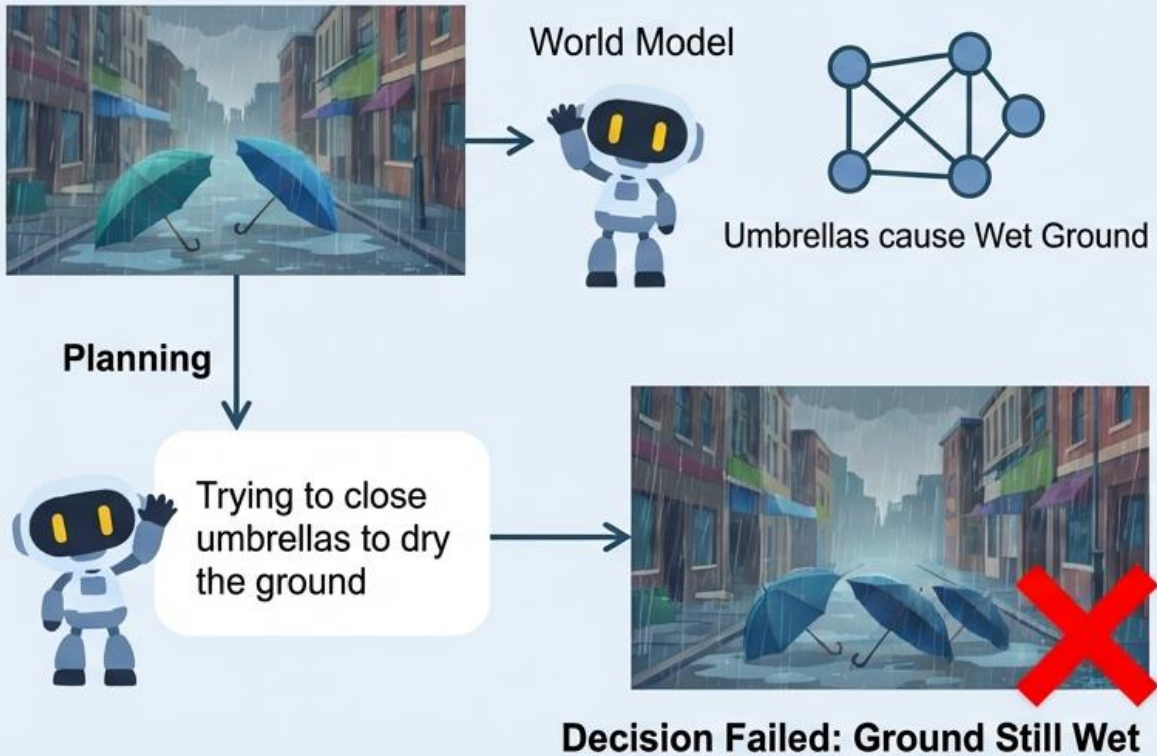
=

Intervention
Counterfactual

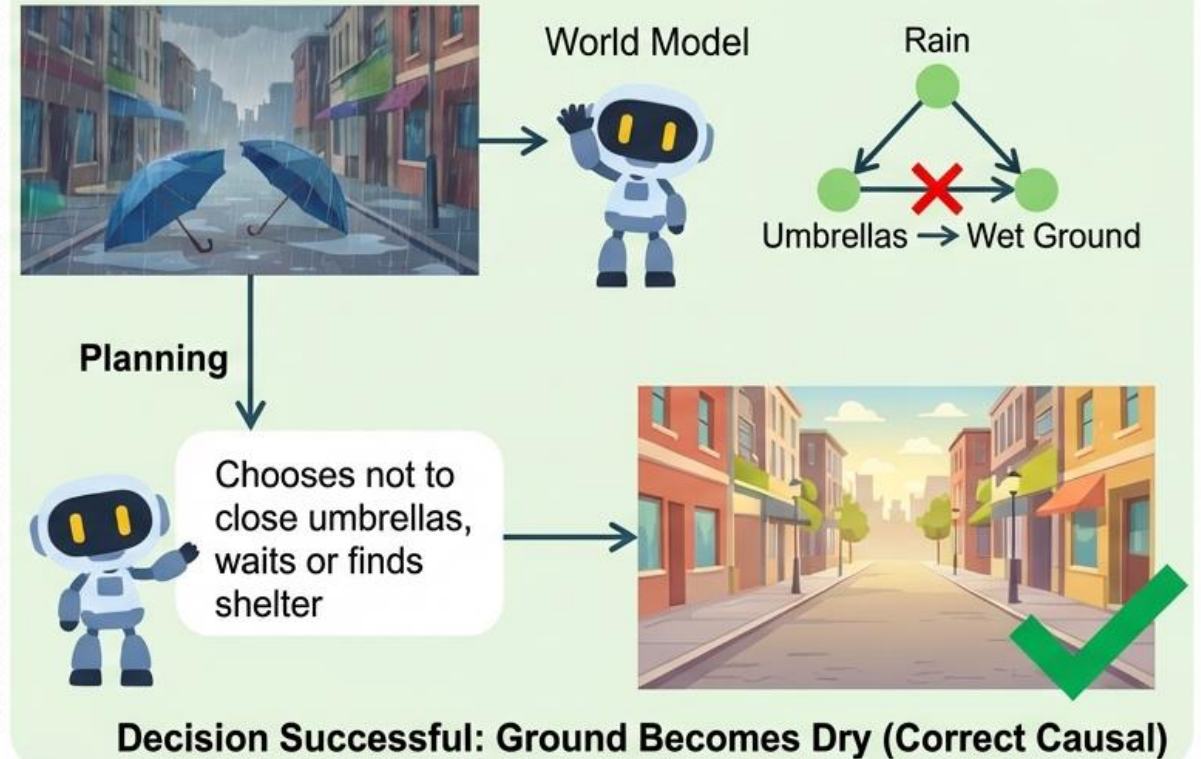
Causal World Models

World Model Comparison: How Causal Knowledge Removes Spurious Correlation and Improves Decision Making

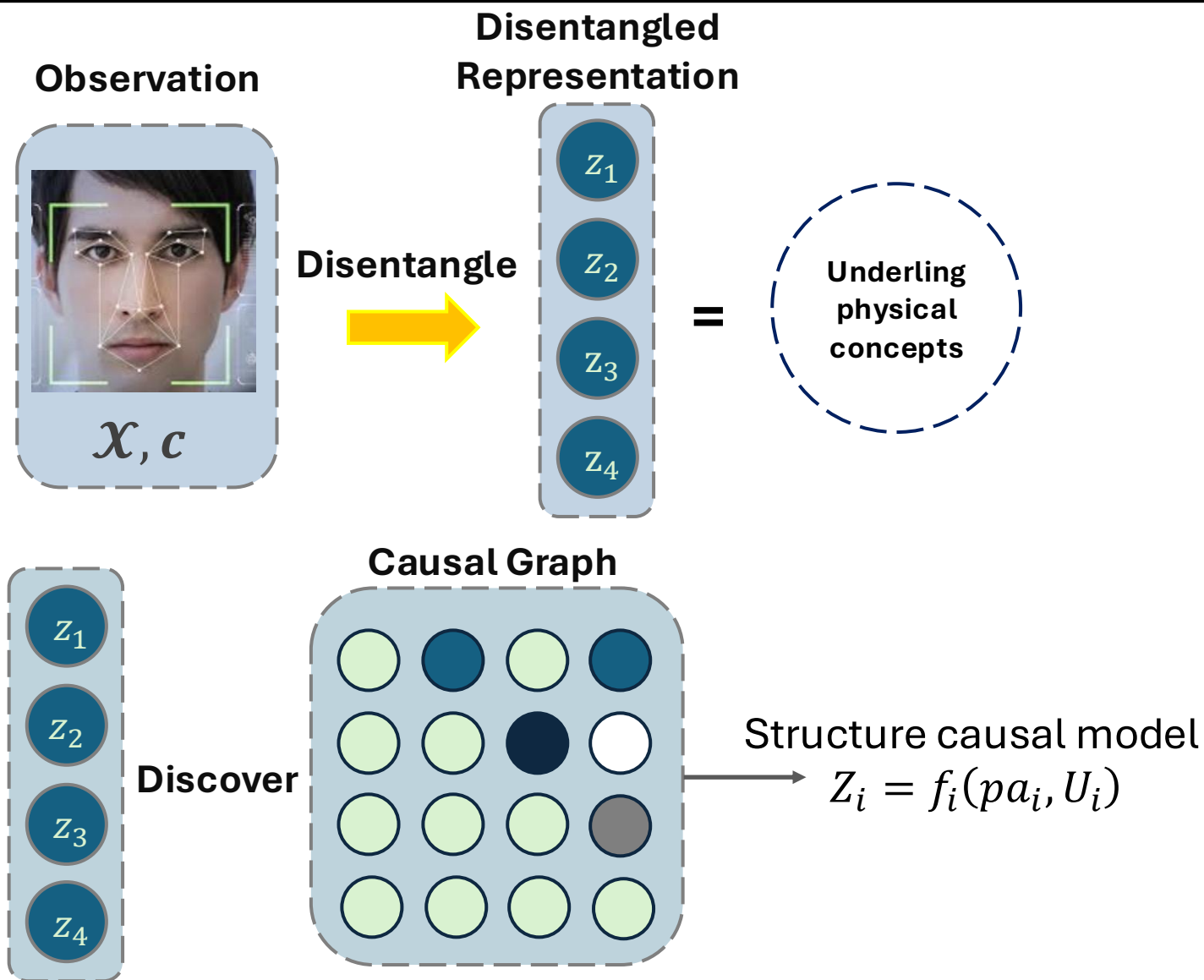
No Causal World Model (Spurious Correlation)



Causal World Model (Causal Knowledge)



Causal World Models – Representation + Causal Structure



Intervene causes **Smile**,
Mouth Open will change



Smile = -0.75

But Intervene effect
concept **Mouth Open**,
Smile will **not** be influenced

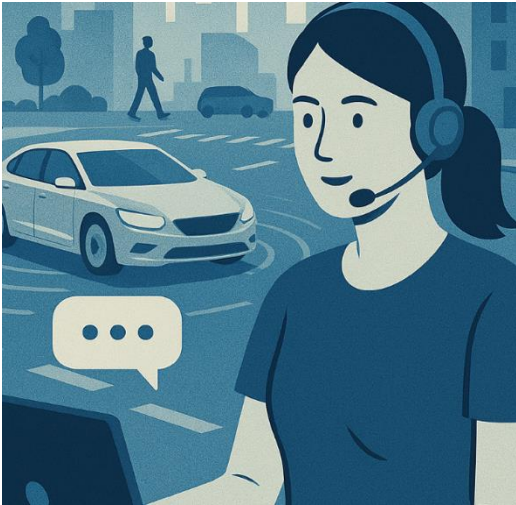


Mouth Open = -0.75

The Challenge of Scaling in Open-ended Environments

Open-ended world

Infinite state, action, multi-agent dynamic - exploding of strategy



Chat Agent



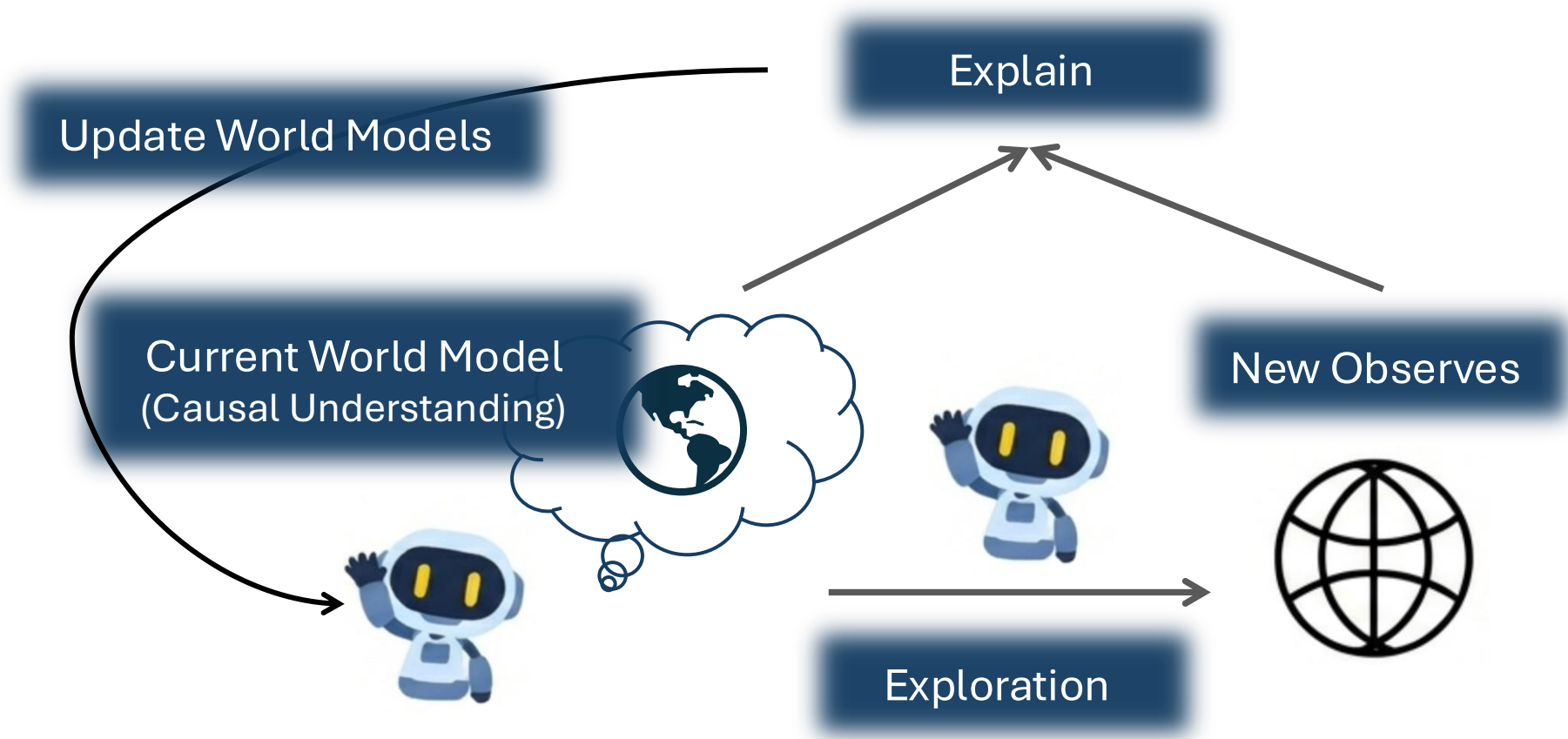
Automatic driving in real-world



Football AI

Partial observation, causality changing through the observation window....

Continual Causal Learning in Open-Ended Worlds



Causal Discovery in Open-Ended World

[NeurIPS 2025]

Curious Causality-Seeking Agents in Open-Ended World



Zhiyu Zhao,



Haoxuan Li, Haifeng Zhang,



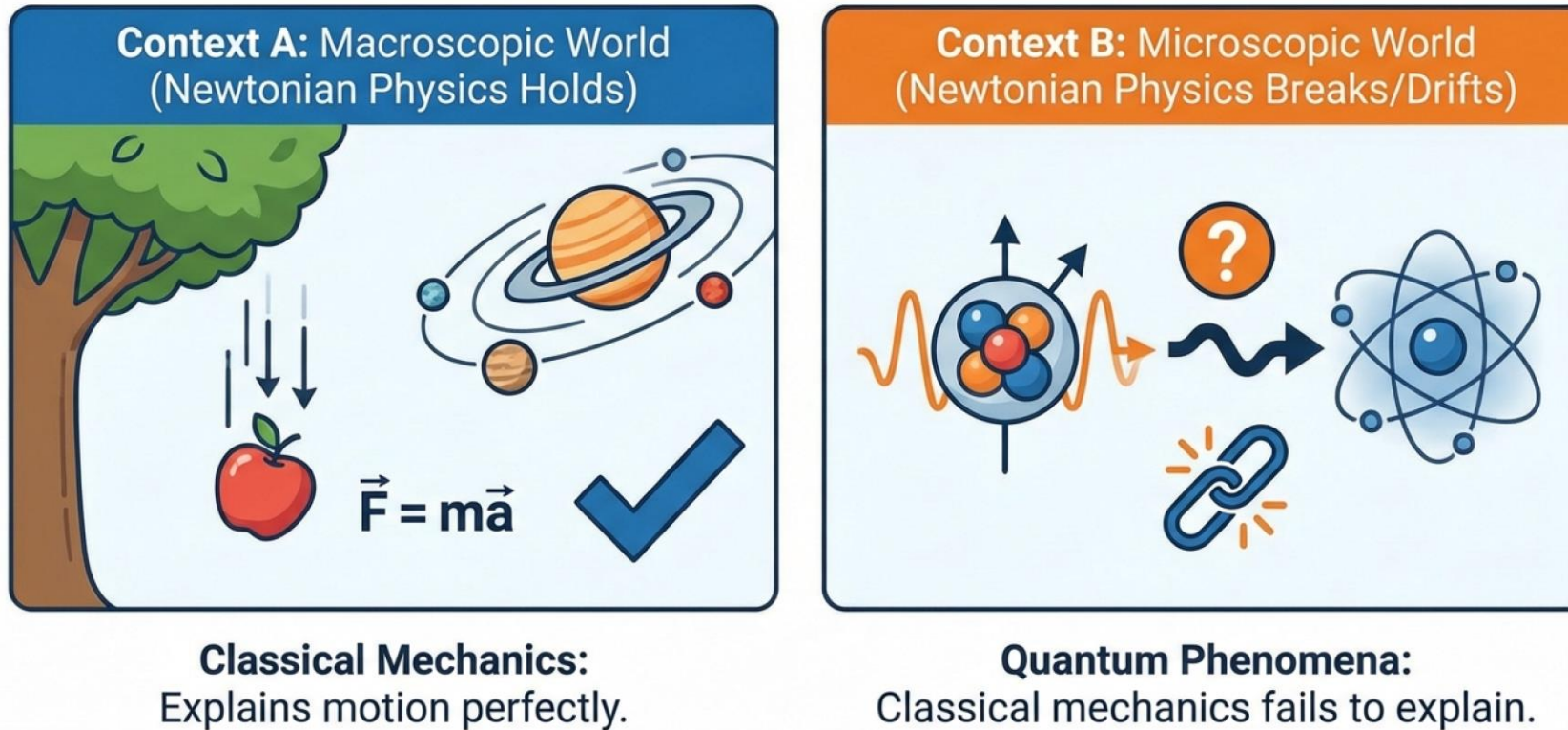
Jun Wang, Francesco Faccio, Jürgen Schmidhuber, Mengyue Yang



Motivation: Causal “Drift” in Open Worlds

Causal Discovery in Open-Ended World

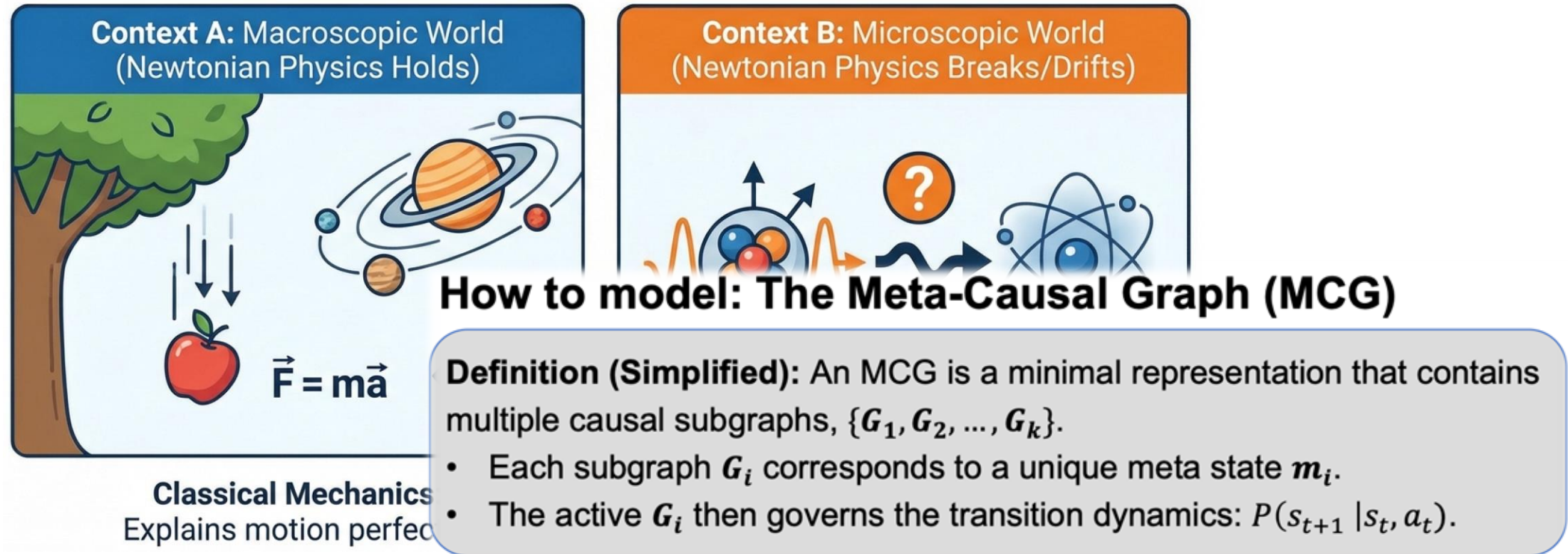
Causal “Drift” in Open Worlds



Curious Causality-Seeking Agents in Open-Ended World. NeurIPS 2025

Causal Discovery in Open-Ended World

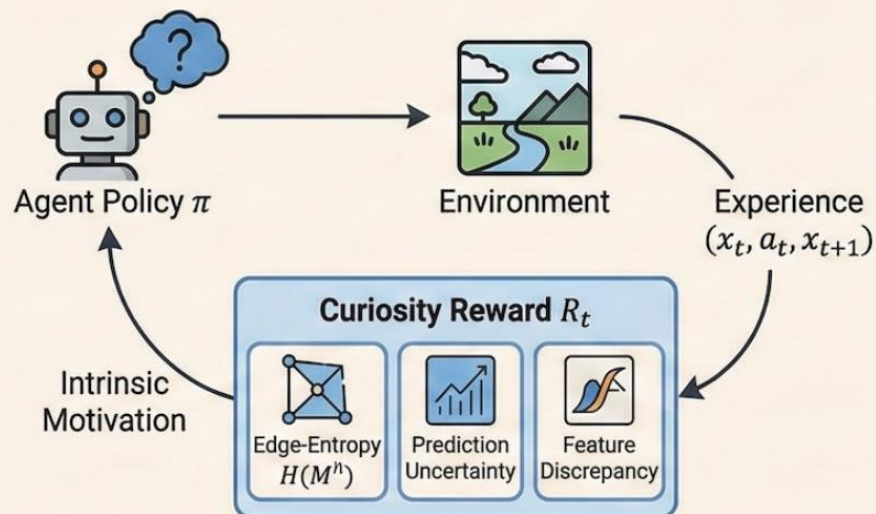
Causal “Drift” under Different Condition



Learning in Open-Ended World

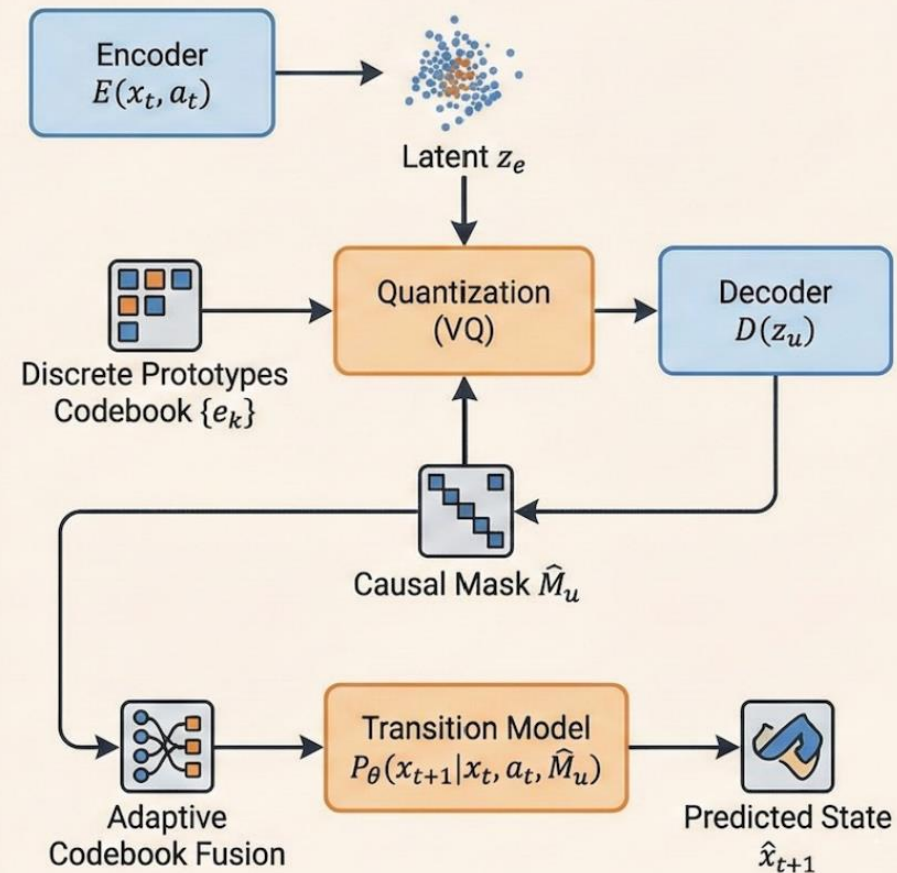
Causal Understanding =
Statistical Learning + Active Exploration

Curiosity-Driven Intervention

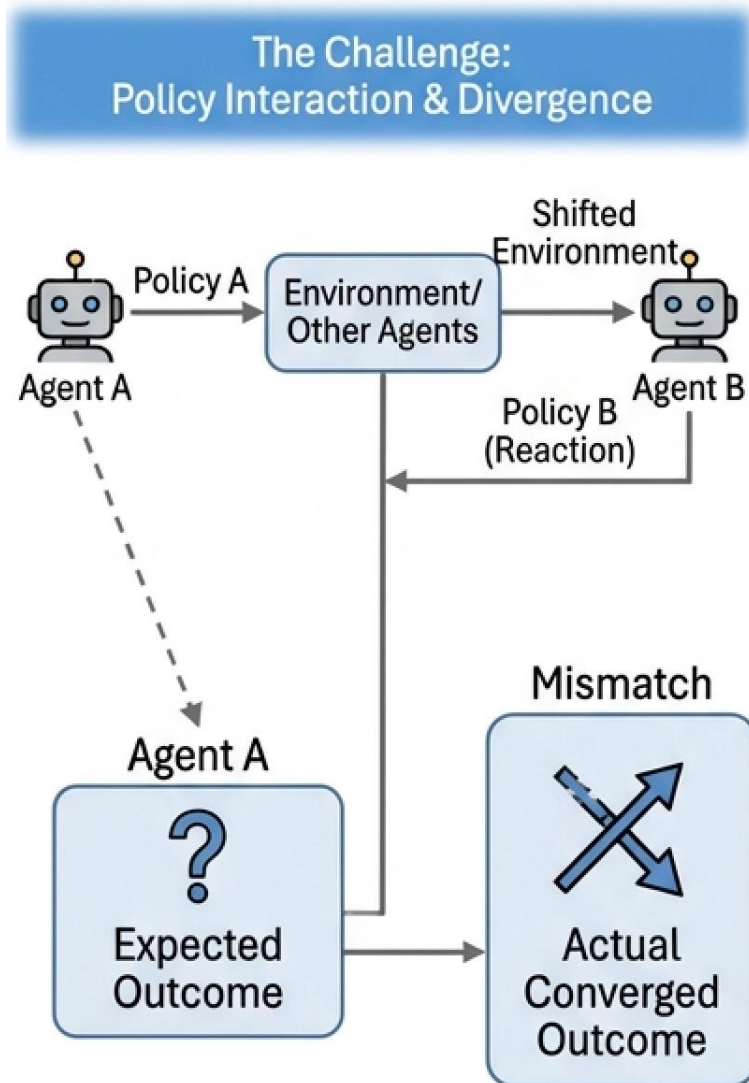


Updating causal graph in the learning loop

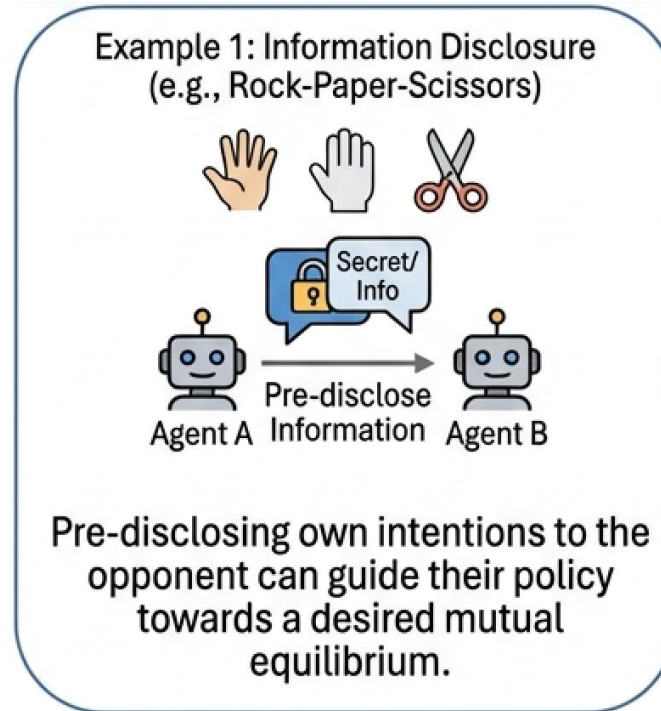
VQ-VAE Based Discovery Architecture



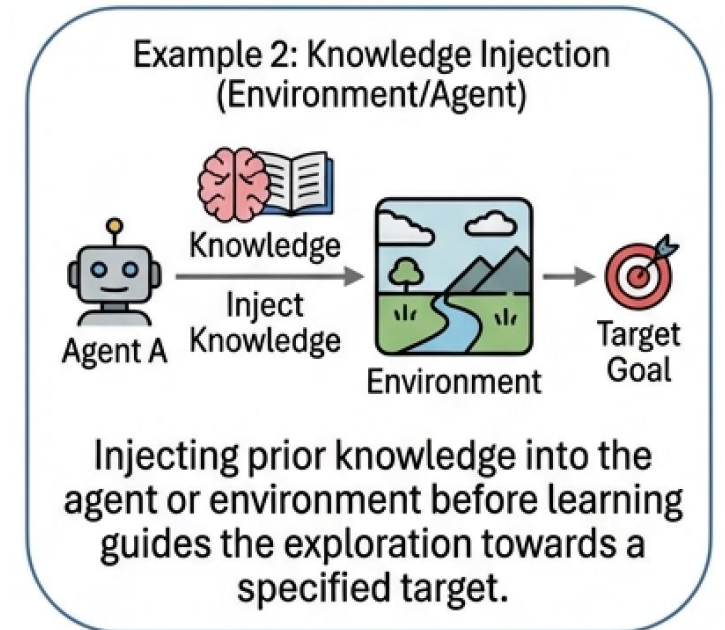
What If the World can be Changed by Policy



Our Explorations: Pre-policy Intervention for Convergence



Leverage knowledge from foundation models



Future Outlook: Towards Generalizable Causal Representations

Current: Task-Specific Causal Learning (RL Exploration)



Football AI



Chat Agent/
Driving

Limited transferability.
Re-learns basic physics
from scratch for each task.



Goal: Learn
Universal
Causal
Primitives

Future: Generalizable Causal Foundation



Generic Warehouse
Robot



Household
Kitchen



Sci-Fi Planet
Rover



Learn invariant
mechanisms
across domains



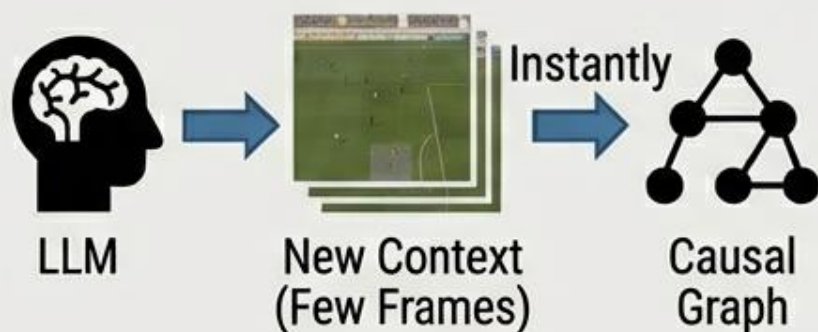
Enable zero-shot
generalization to
new
environments.

Moving from learning 'how to play football' to learning 'how objects interact physically'.

Future Outlook: Causal Foundation World Models & Scaling



In-Context Causal Learning



Leverage LLM reasoning for fast, in-context causal discovery

The Scaling Hypothesis



Causal understanding emerges and robustifies with large-scale pre-training.

Final Goal: A unified, scalable world model with robust causal reasoning capabilities.