



# Causal Reinforcement Learning

Empowering Agents with Causality



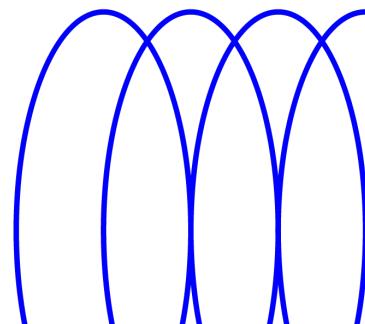
Jing Jiang (Associate Professor)

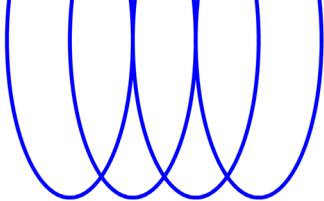
[Jing.Jiang@uts.edu.au](mailto:Jing.Jiang@uts.edu.au)



Zihong Deng (PhD Student)

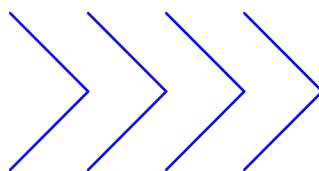
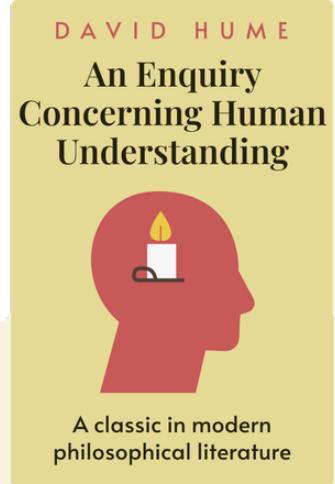
[Zhi-Hong.Deng@student.uts.edu.au](mailto:Zhi-Hong.Deng@student.uts.edu.au)





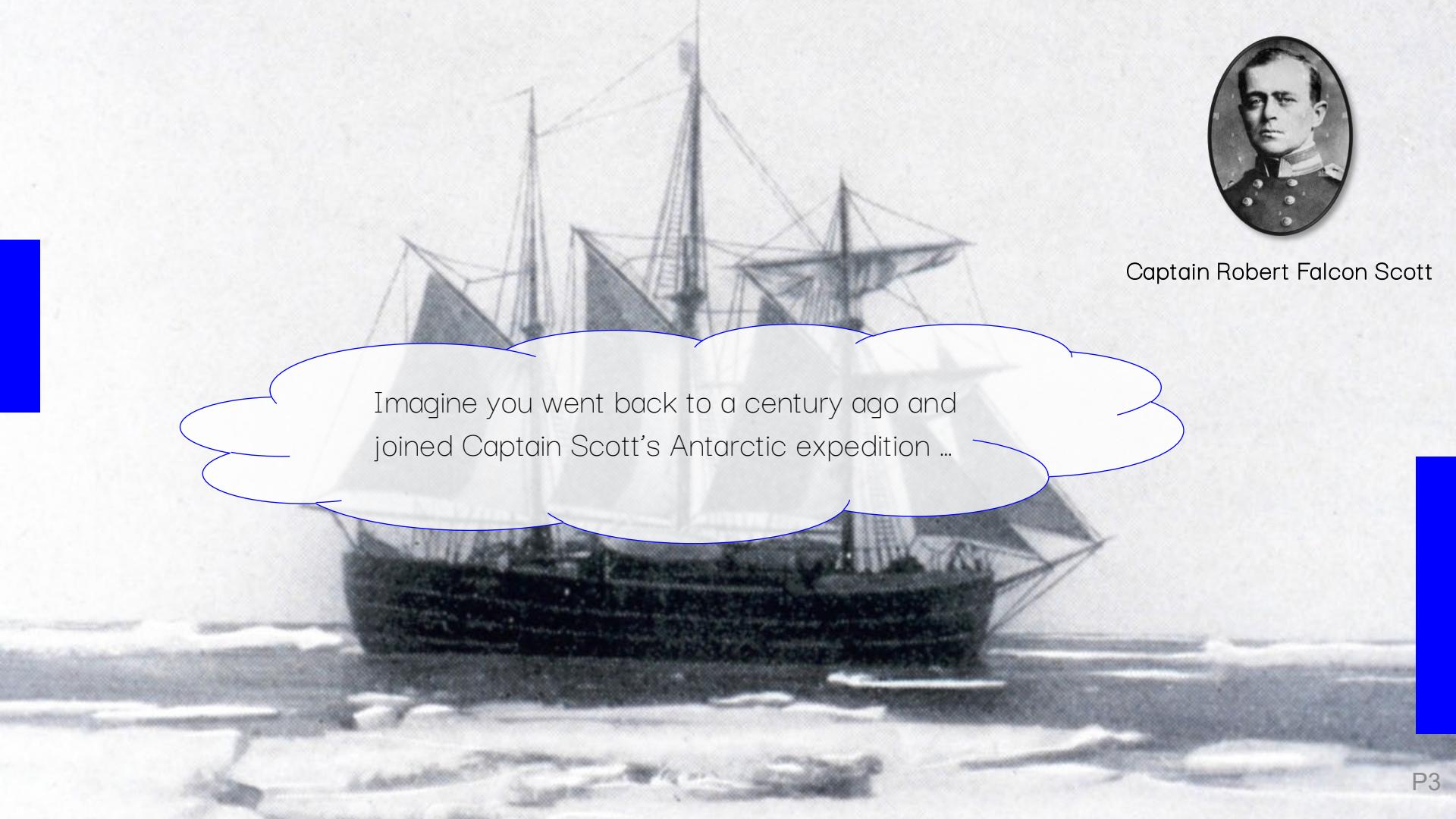
”All reasonings concerning matter of fact seem to be founded on the relation of cause and effect. By means of that relation alone we can go beyond the evidence of our memory and senses.”

-- David Hume, 1748





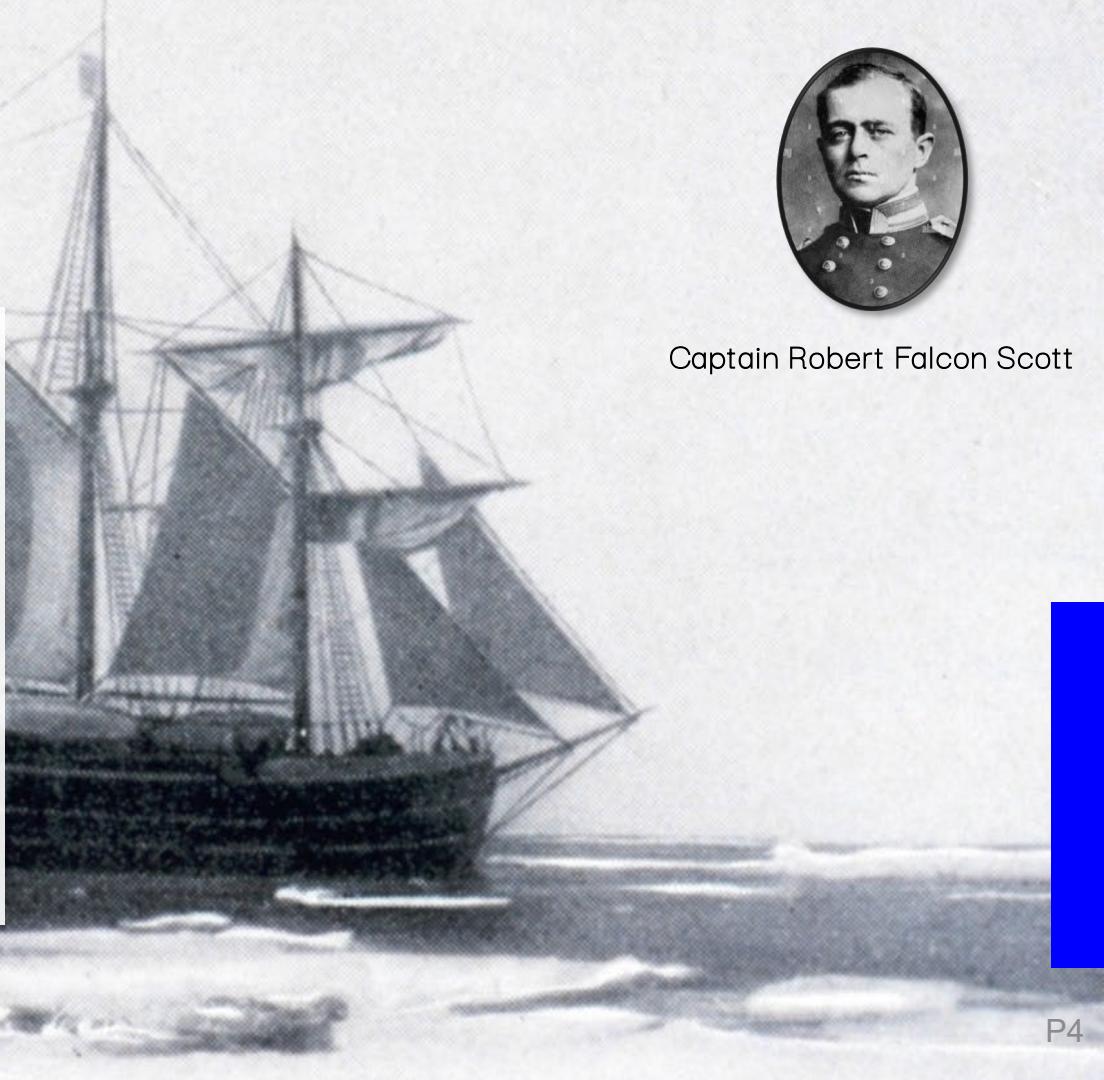
Captain Robert Falcon Scott

A black and white photograph of a three-masted sailing ship, likely the Terra Nova, navigating through a field of sea ice. The ship's hull is dark, and its masts are visible against a bright sky.

Imagine you went back to a century ago and joined Captain Scott's Antarctic expedition ...



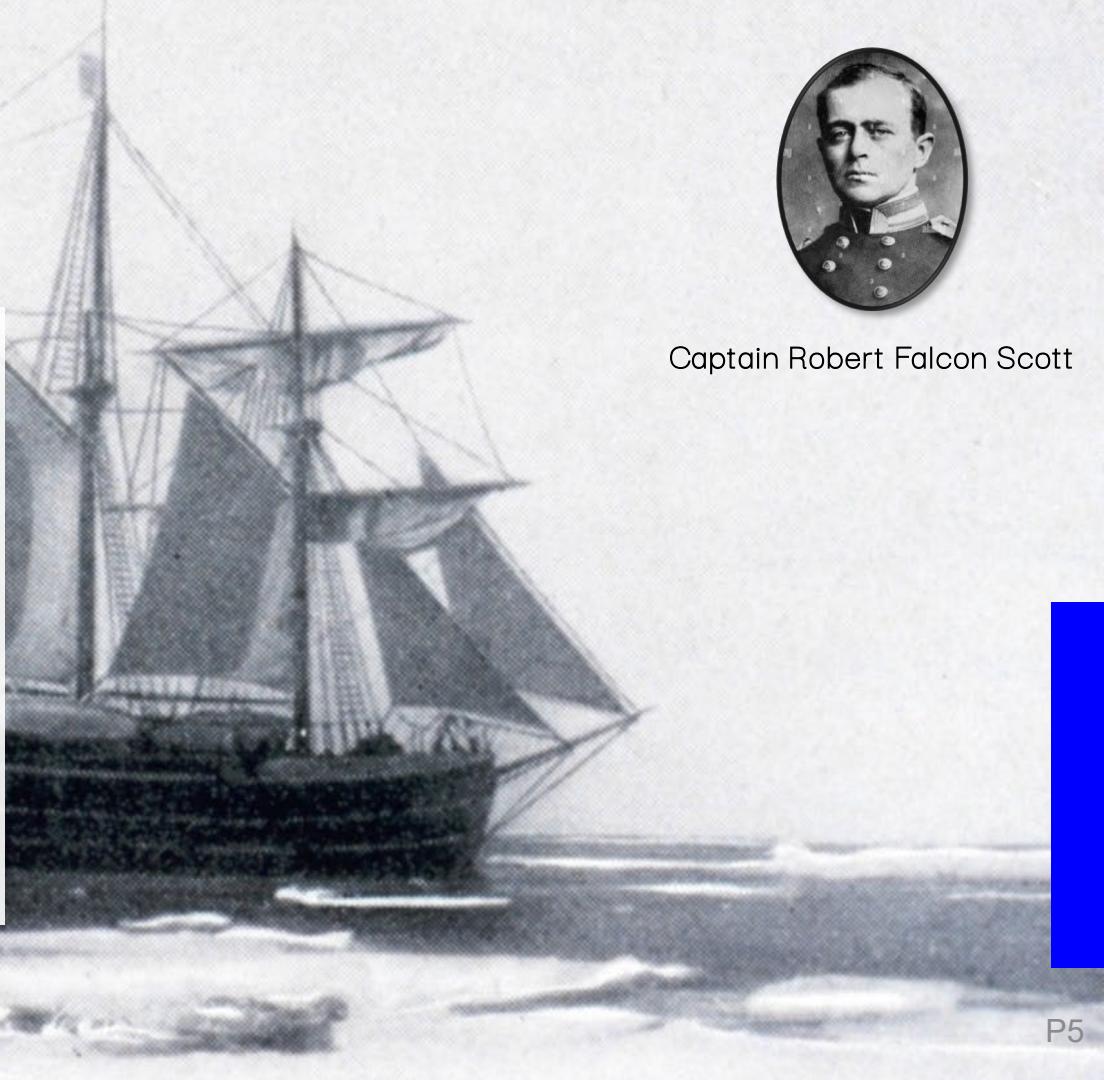
"Consuming rotten  
meat causes scurvy"



Captain Robert Falcon Scott



"Consuming rotten  
meat causes scurvy"



Captain Robert Falcon Scott



"Consuming rotten  
meat causes scurvy"



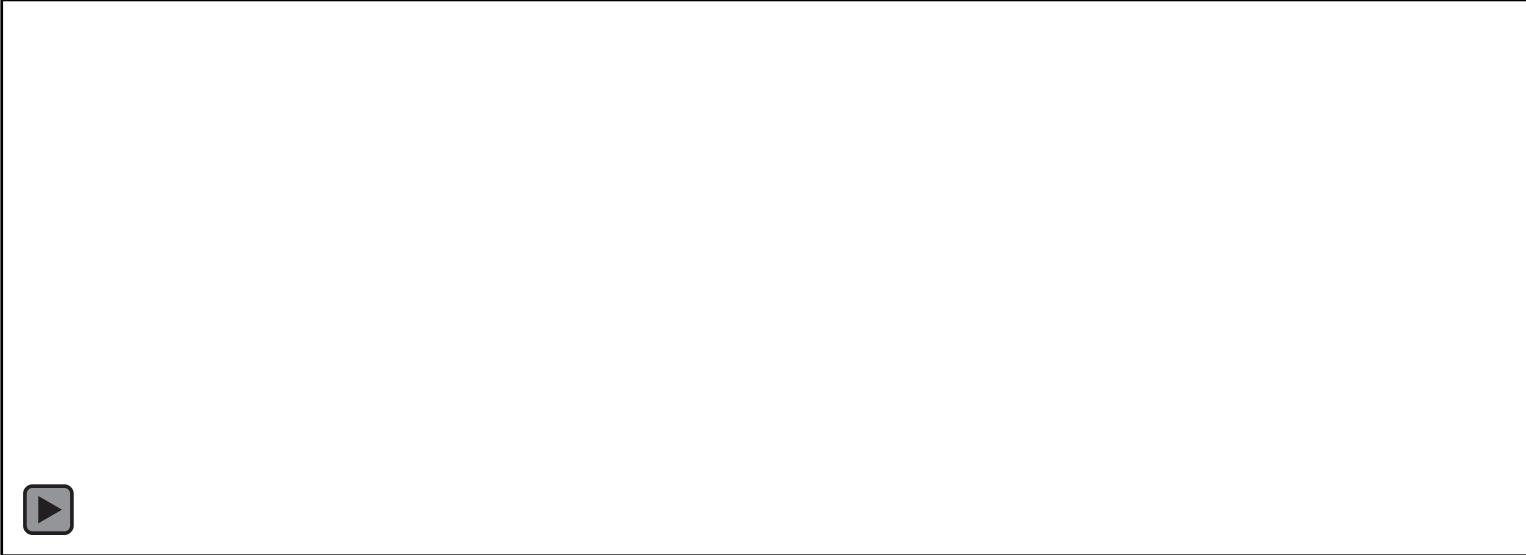
Lack of Vitamin C  
causes scurvy



Captain Robert Falcon Scott

# The AI Scaling Law

For more information, check  
“Pathways Language Model (PaLM):  
Scaling to 540 Billion Parameters for  
Breakthrough Performance”, 2022



<https://ai.googleblog.com/2022/04/pathways-language-model-palm-scaling-to.html>

The PaLM pretraining dataset consists of a high-quality corpus of **780 billion tokens**.

# What Can Go Wrong?



There is a strong correlation between rotten meat and a mysterious illness. During long voyages, the consumption of rotten meat and the onset of this illness usually coincide; however, sailors rarely experience it when they avoid eating spoiled meat. One day, Captain Scott led his sailors on an expedition to Antarctica, supplying them with an abundance of fresh meat instead of rotten meat. Would the sailors still be susceptible to this illness?



No, if the sailors are consuming only fresh meat and avoiding rotten meat, the correlation between rotten meat and the mysterious illness should not apply, and they would likely not be susceptible to the illness during the expedition.

A failure case of ChatGPT: ChatGPT, one of the best AI-powered chatbots trained using Reinforcement Learning from Human Feedback (RLHF) to align with human values.

# What Can Go Wrong?



There is a strong correlation between rotten meat and a mysterious illness. During long voyages, the consumption of rotten meat and the onset of the illness usually coincide; however, sailors rarely experience it when they avoid eating spoiled meat. One day, Captain Scott led his team's polar expedition to Antarctica, supplying them with an abundance of fresh meat instead of rotten meat. Would the sailors still be susceptible to this illness?

AI agents have not yet mastered causation in decision-making!



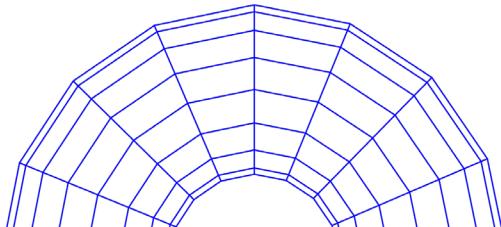
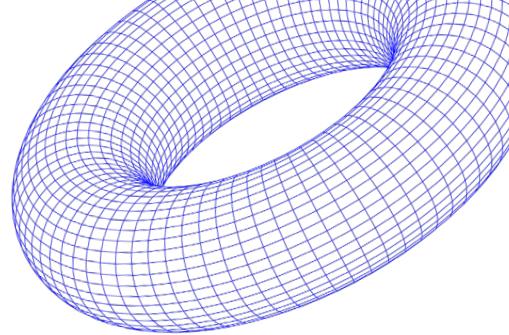
No, if the sailors are consuming only fresh meat and avoiding rotten meat, the correlation between rotten meat and the mysterious illness should not apply, and they would likely not be susceptible to the illness during the expedition.

A failure case for ChatGPT: ChatGPT, one of the best AI-powered chatbots trained using Reinforcement Learning from Human Feedback (RLHF) to align with human values.



"This idea that we're going to **just scale up** the current large language models and eventually human-level AI will emerge—**I don't believe this at all**, not for one second."

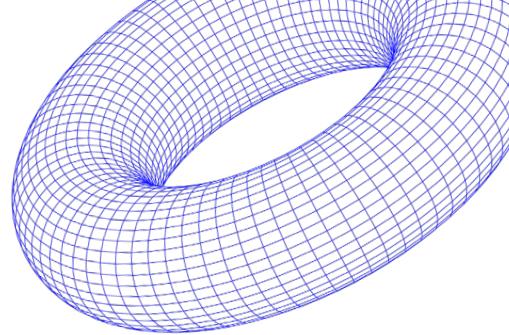
Yann Lecun





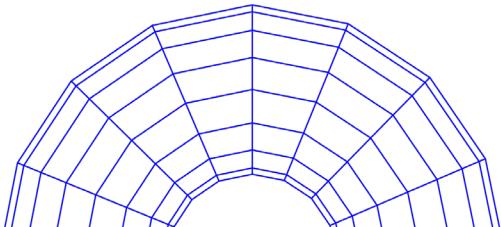
"This idea that we're going to **just scale up** the current large language models and eventually human-level AI will emerge—**I don't believe this at all**, not for one second."

Yann Lecun



Judea Pearl

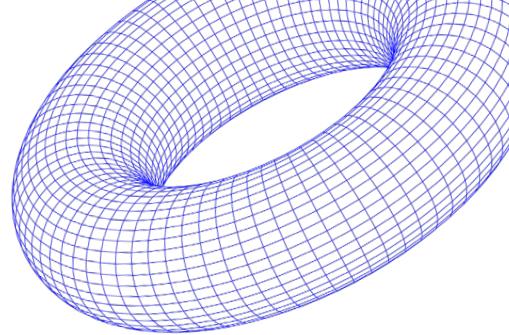
"Much of this data-centric history still haunts us today. We live in an era that presumes Big Data to be the solution to all our problems. Courses in 'data science' are proliferating in our universities, and jobs for 'data scientists' are lucrative in the companies that participate in the 'data economy.' But I hope with this book to convince you that **data are profoundly dumb.** "





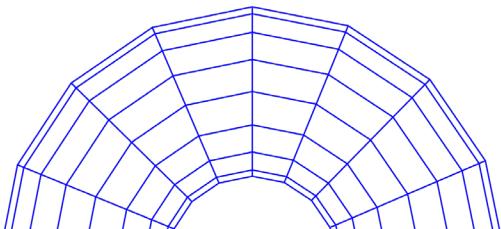
Yann Lecun

"This idea that we're going to **just scale up** the current large language models and eventually human-level AI will emerge—**I don't believe this at all**, not for one second."



Judea Pearl

"Much of this data-centric history still haunts us today. We live in an era that presumes Big Data to be the solution to all our problems. Courses in 'data science' are proliferating in our universities, and jobs for 'data scientists' are lucrative in the companies that participate in the 'data economy.' But I hope with this book to convince you that **data are profoundly dumb.** "



Yoshua Bengio

" One of the big debates these days is: **What are the elements of higher-level cognition? Causality is one element of it.**"

---

# Goals of This Tutorial

---

1. What is causal reinforcement learning and how is it different than traditional reinforcement learning?

---

# Goals of This Tutorial

---

1. What is causal reinforcement learning and how is it different than traditional reinforcement learning?
2. Different perspectives in the causal reinforcement learning literature.

---

# Goals of This Tutorial

---

1. What is causal reinforcement learning and how is it different than traditional reinforcement learning?
2. Different perspectives in the causal reinforcement learning literature.
3. Main results and techniques.

# Tutorial Outline

## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Correlation vs. Causation



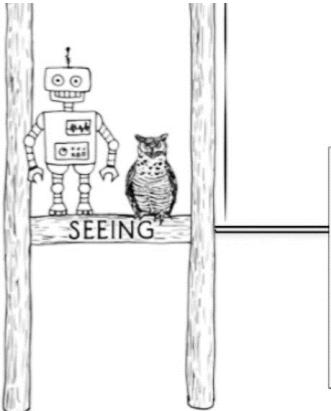
There is a strong [correlation](#) between rotten meat and scurvy in historical data, but eating rotten meat does not [cause](#) scurvy. The lack of vitamin C does.

# Ladder of Causation



Judea Pearl

“The Book of Why”, 2018



## 1. ASSOCIATION

ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see . . . ?*

(How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease?

What does a survey tell us about the election results?

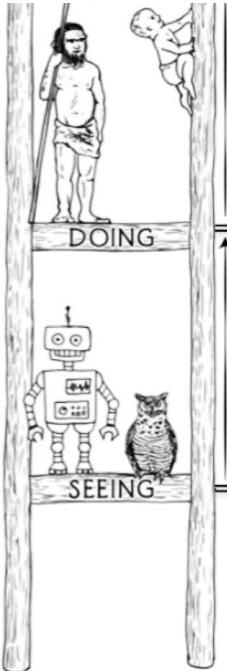


# Ladder of Causation



Judea Pearl

“The Book of Why”, 2018



## 2. INTERVENTION

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do . . . ? How?*  
(What would Y be if I do X?)

EXAMPLES: If I take aspirin, will my headache be cured?  
What if we ban cigarettes?



## 1. ASSOCIATION

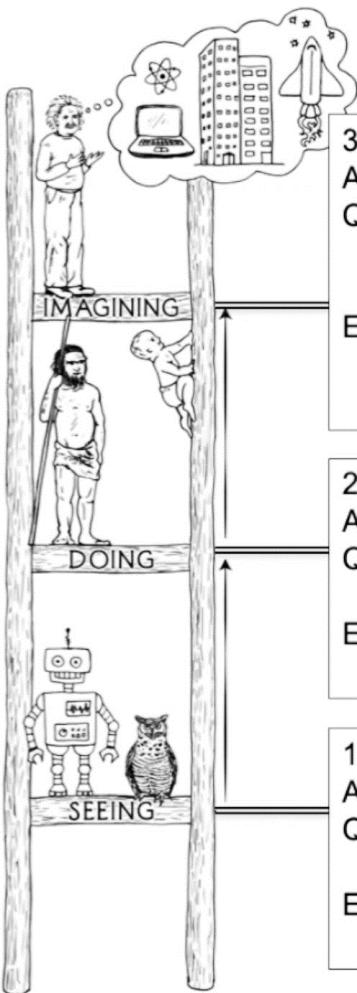
ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see . . . ?*  
(How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease?  
What does a survey tell us about the election results?



# Ladder of Causation



## 3. COUNTERFACTUALS

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done . . . ? Why?*

(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache?

Would Kennedy be alive if Oswald had not killed him? What if I had not smoked the last 2 years?



## 2. INTERVENTION

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do . . . ? How?*

(What would Y be if I do X?)

EXAMPLES: If I take aspirin, will my headache be cured?



What if we ban cigarettes?

## 1. ASSOCIATION

ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see . . . ?*

(How would seeing X change my belief in Y?)



EXAMPLES: What does a symptom tell me about a disease?

What does a survey tell us about the election results?



Judea Pearl

"The Book of Why", 2018

# Structural Causal Model (SCM)

**Definition.** An SCM is represented by a quadruple  $(V, U, F, P(U))$ , where

- $V = \{V_1, V_2, \dots, V_m\}$  is a set of **endogenous variables** that are of interest in a research problem,
- $U = \{U_1, U_2, \dots, U_n\}$  is a set of **exogenous variables** that represent the source of stochasticity in the model and are determined by external factors that are generally unobservable,
- $F = \{f_1, f_2, \dots, f_m\}$  is a set of **structural equations** that assign values to each of the variables in  $V$  such that  $f_i$  maps  $\text{PA}(V_i) \cup U_i$  to  $V_i$  where  $\text{PA}(V_i) \subseteq V \setminus V_i$  and  $U_i \subseteq U$ ,
- $P(U)$  is the joint probability distribution of the exogenous variables in  $U$ .

$$U_1 \perp\!\!\!\perp U_2 \perp\!\!\!\perp \cdots \perp\!\!\!\perp U_n$$

# Structural Causal Model (SCM)

Structural Equations

$$f_X: X = U_X$$

$$f_Z: Z = a \cdot X + U_Z$$

$$f_Y: Y = b \cdot X + c \cdot Z + U_Y$$

An example of SCM with structural equations

$$F = \{f_X, f_Z, f_Y\}$$



$X$ : food consumption,

$Z$ : intake of vitamin C,

$Y$ : occurrence of scurvy

# Structural Causal Model (SCM)

Endogenous Variables

$$f_X: \boxed{X} = U_X$$

$$f_Z: \boxed{Z} = a \cdot \boxed{X} + U_Z$$

$$f_Y: \boxed{Y} = b \cdot \boxed{X} + c \cdot \boxed{Z} + U_Y$$

An example of SCM with structural equations

$$F = \{f_X, f_Z, f_Y\}$$



**X**: food consumption,

**Z**: intake of vitamin C,

**Y**: occurrence of scurvy

# Structural Causal Model (SCM)

Exogenous Variables

$$f_X: X = U_X$$

$$f_Z: Z = a \cdot X + U_Z$$

$$f_Y: Y = b \cdot X + c \cdot Z + U_Y$$

An example of SCM with structural equations

$$F = \{f_X, f_Z, f_Y\}$$

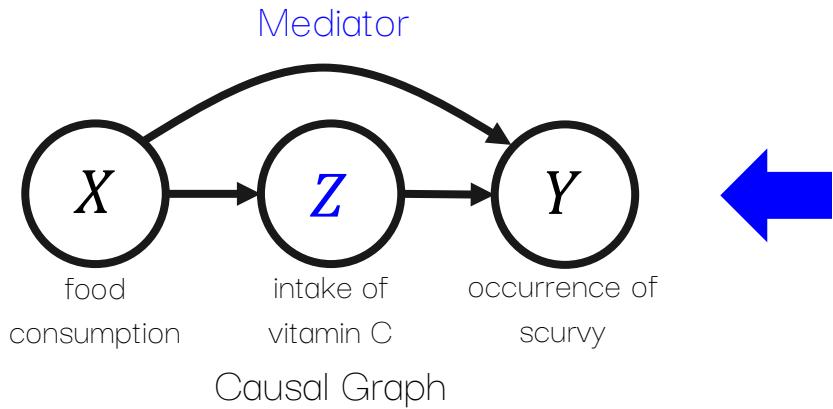


$X$ : food consumption,

$Z$ : intake of vitamin C,

$Y$ : occurrence of scurvy

# Causal Graph



$$f_X: X = U_X$$
$$f_Z: Z = a \cdot X + U_Z$$
$$f_Y: Y = b \cdot X + c \cdot Z + U_Y$$

An example of SCM with structural equations

$$\mathcal{F} = \{f_X, f_Z, f_Y\}$$



$X$ : food consumption,  
 $Z$ : intake of vitamin C,  
 $Y$ : occurrence of scurvy

# Observations vs. Interventions



Observations

Passively observe people with different food consumption.

# Observations vs. Interventions



## Observations

Q: What does consuming citrus fruits tell me about the possibility of getting scurvy?

$$P(Y \mid X=\text{orange})$$

# Observations vs. Interventions



Observations

Q: What does consuming citrus fruits tell me about the possibility of getting scurvy?

$$P(Y \mid X=\text{orange})$$



Interventions

Actively force all sailors to consume fresh citrus fruits.

# Observations vs. Interventions



Observations

Q: What does consuming citrus fruits tell me about the possibility of getting scurvy?

$$P(Y \mid X = \text{lemon})$$



Interventions

Q: What if all sailors consume fresh citrus fruits, will they get scurvy?

$$P(Y \mid \text{do}(X = \text{lemon}))$$

# Interventions

$f_X: X = \text{orange}$

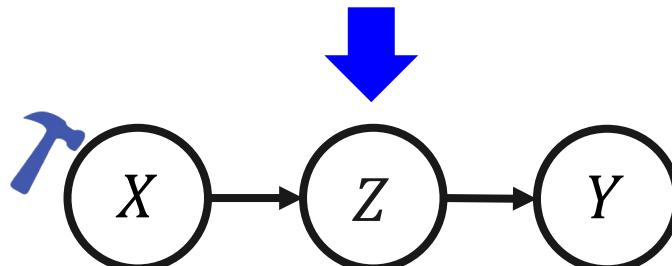
$f_Z: Z = a \cdot X + U_Z$

$f_Y: Y = c \cdot Z + U_Y$

New SCM



Interventions

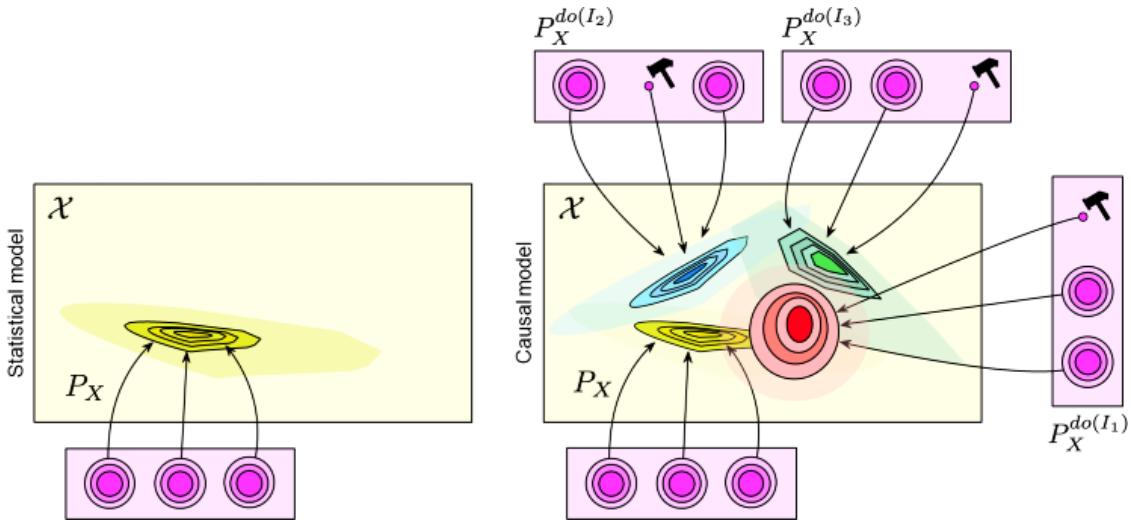


New Causal Graph

Q: What if all sailors consume fresh citrus fruits, will they get scurvy?

$P(Y | \text{do}(X=\text{orange}))$

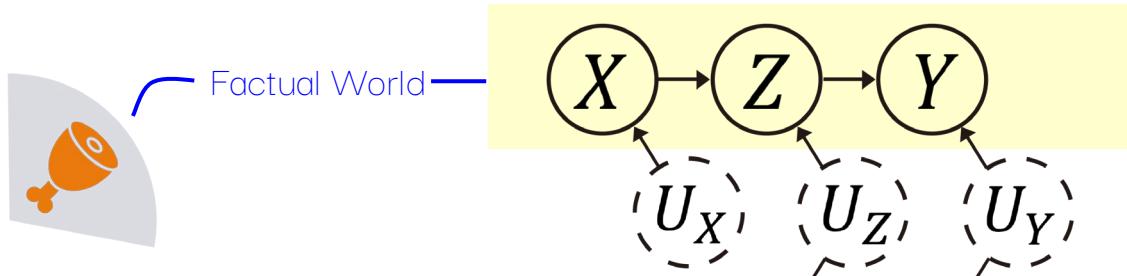
# Statistical vs. Causal Model



Difference between [statistical](#) (left) and [causal models](#) (right) on a given set of three variables. While a statistical model specifies a single probability distribution, a causal model represents a set of distributions, one for each possible intervention (indicated with a ↗ in the figure)

Causal models are inherently more powerful than statistical model!

# Counterfactuals



## Observations

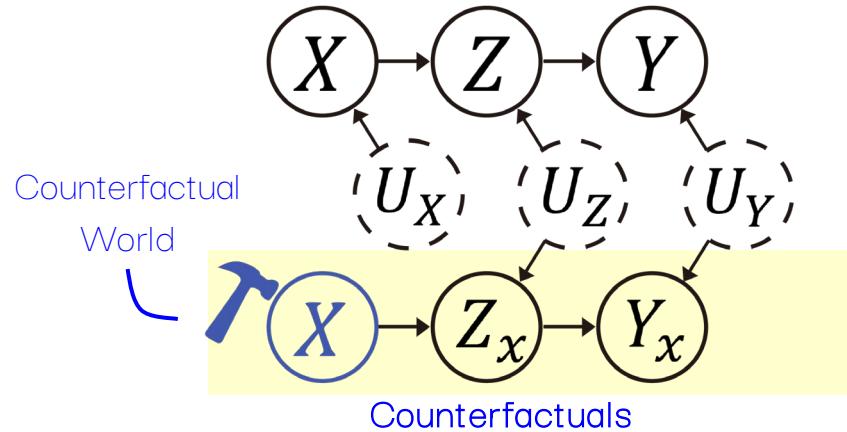
Passively observe people who consume rotten meat can also get scurvy.

# Counterfactuals



## Observations

Passively observe people who consume rotten meat can also get scurvy.



Imagine people who would have consumed rotten meat choosing to consume fresh citrus fruits instead.

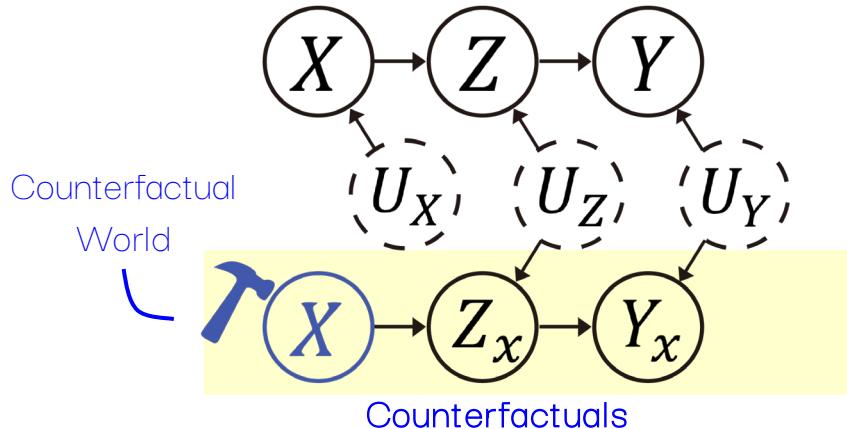
# Counterfactuals



## Observations

Passively observe people who consume rotten meat can also get scurvy.

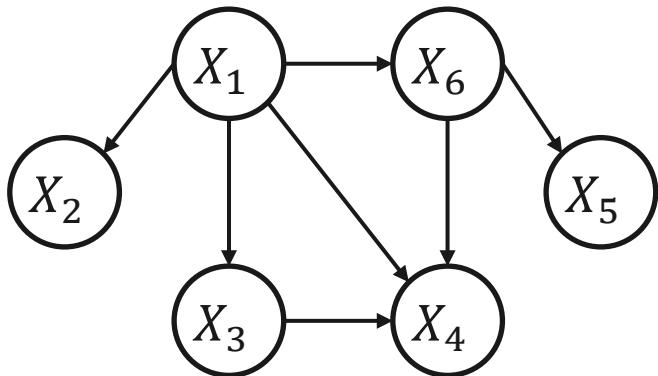
$$P(Y | X = \text{rotten meat})$$



Q: Considering that they consumed rotten meat in reality, would sailors have been protected from scurvy if they had consumed enough citrus fruit?

$$P(Y_{X=\text{orange}} | X = \text{rotten meat}, Y = \text{sailor})$$

# Non-Causal vs. Causal Factorization

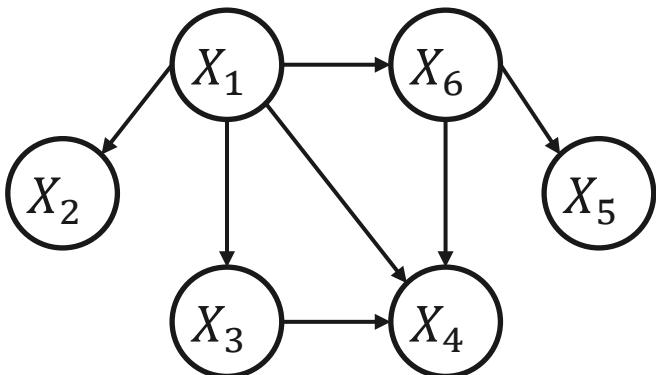


Non-Causal Factorization

e.g.,

$$\begin{aligned} P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | X_1, \dots, X_{i-1}) \\ &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdot P(X_4 | X_1, X_2, X_3) \\ &\quad \cdot P(X_5 | X_1, X_2, X_3, X_4) \cdot P(X_6 | X_1, X_2, X_3, X_4, X_5) \end{aligned}$$

# Non-Causal vs. Causal Factorization



Non-Causal Factorization

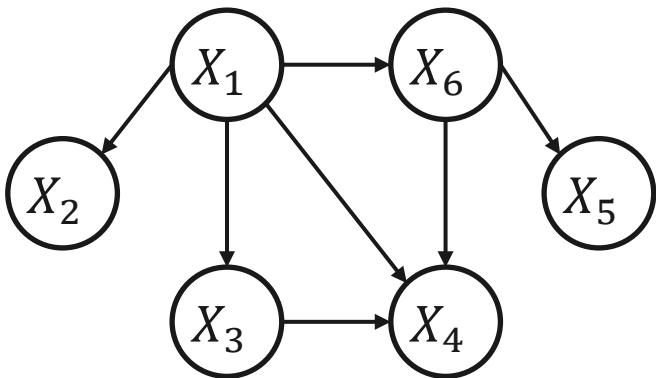
e.g.,

$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | X_1, \dots, X_{i-1}) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdot P(X_4 | X_1, X_2, X_3) \\
 &\quad \cdot P(X_5 | X_1, X_2, X_3, X_4) \cdot P(X_6 | X_1, X_2, X_3, X_4, X_5)
 \end{aligned}$$

Causal Factorization

$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | \text{PA}(X_i)) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1) \cdot P(X_4 | X_1, X_3, X_6) \\
 &\quad \cdot P(X_5 | X_6) \cdot P(X_6 | X_1)
 \end{aligned}$$

# Non-Causal vs. Causal Factorization



Non-Causal Factorization

e.g.,

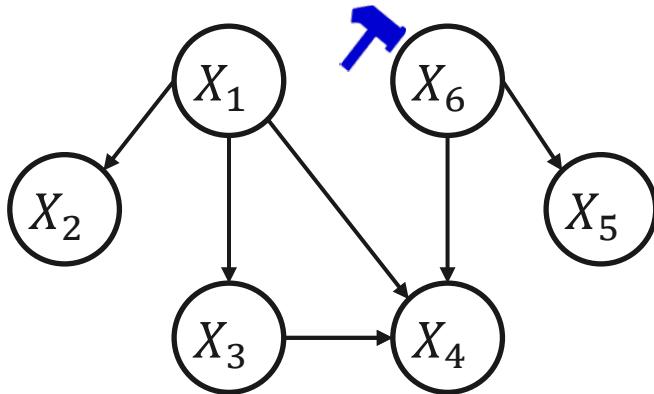
$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | X_1, \dots, X_{i-1}) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdot P(X_4 | X_1, X_2, X_3) \\
 &\quad \cdot P(X_5 | X_1, X_2, X_3, X_4) \cdot P(X_6 | X_1, X_2, X_3, X_4, X_5)
 \end{aligned}$$

Causal Factorization

$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | \text{PA}(X_i)) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1) \cdot P(X_4 | X_1, X_3, X_6) \\
 &\quad \cdot P(X_5 | X_6) \cdot P(X_6 | X_1)
 \end{aligned}$$

Causal factorization yields practical computational advantages during inference.

# Independent Causal Mechanism



Non-Causal Factorization

e.g.,

$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | X_1, \dots, X_{i-1}) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdot P(X_4 | X_1, X_2, X_3) \\
 &\quad \cdot P(X_5 | X_1, X_2, X_3, X_4) \cdot P(X_6 | X_1, X_2, X_3, X_4, X_5)
 \end{aligned}$$

Causal Factorization

$$\begin{aligned}
 P(\mathbf{X}) &= \prod_{i=1}^N P(X_i | \text{PA}(X_i)) \\
 &= P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1) \cdot P(X_4 | X_1, X_3, X_6) \\
 &\quad \cdot P(X_5 | X_6) \cdot P(X_6 | X_1)
 \end{aligned}$$

Causal factorization is more robust to variations.

# Independent Causal Mechanism

## **Independent Causal Mechanisms (ICM) Principle.**

*The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other. In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other mechanisms.*

Applied to the causal factorization, this principle tells us that the factors should be independent in the sense that

- (a) changing (or performing an intervention upon) one mechanism  $P(X_i|\text{PA}(X_i))$  does not change any of the other mechanisms  $P(X_j|\text{PA}(X_j))$  ( $i \neq j$ ).
- (b) Knowing some other mechanisms  $P(X_j|\text{PA}(X_j))$  ( $i \neq j$ ) does not give us information about a mechanism  $P(X_i|\text{PA}(X_i))$ .

# Summary (so far)

- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- Structural Causal Model (SCM) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- Interventions refer to actively manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of imagination and retrospection.
- Causal Factorization decompose a joint distribution into independent causal mechanisms, yielding practical computational advantages and is robust to variations.

# Summary (so far)

- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- [Structural Causal Model \(SCM\)](#) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- Interventions refer to actively manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of imagination and retrospection.
- Causal Factorization decompose a joint distribution into independent causal mechanisms, yielding practical computational advantages and is robust to variations.

# Summary (so far)

- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- Structural Causal Model (SCM) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- **Interventions** refer to *actively* manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of imagination and retrospection.
- Causal Factorization decompose a joint distribution into independent causal mechanisms, yielding practical computational advantages and is robust to variations.

# Summary (so far)

- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- Structural Causal Model (SCM) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- Interventions refer to actively manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of **imagination** and **retrospection**.
- Causal Factorization decompose a joint distribution into independent causal mechanisms, yielding practical computational advantages and is robust to variations.

# Summary (so far)

- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- Structural Causal Model (SCM) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- Interventions refer to actively manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of imagination and retrospection.
- **Causal Factorization** decompose a joint distribution into **independent causal mechanisms**, yielding practical computational advantages and is **robust to variations**.

# Summary (so far)

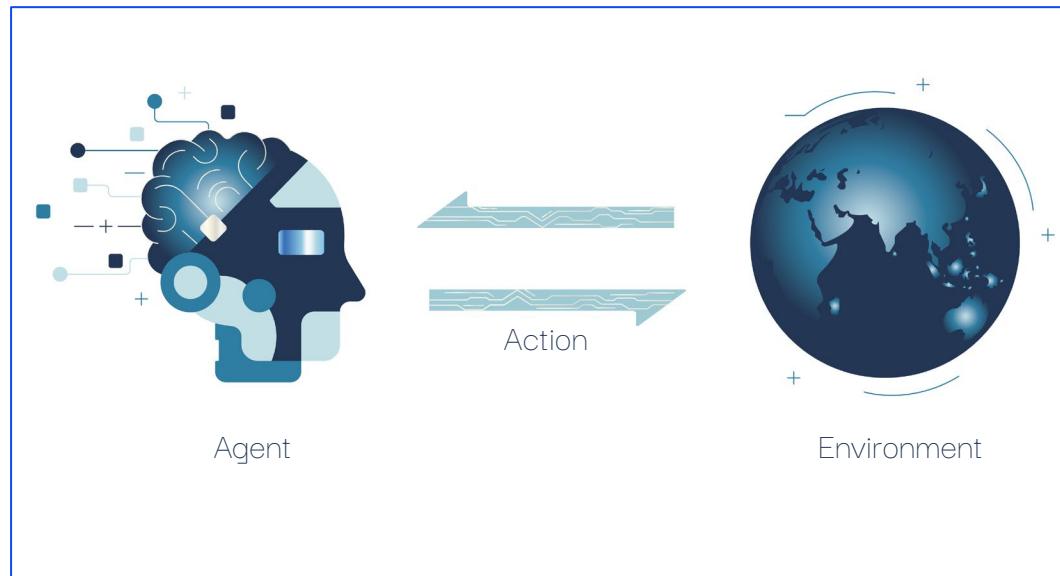
- Pearl's Ladder of Causation is a conceptual framework that categorizes levels of causal relationships, spanning from association, intervention, and counterfactuals.
- Structural Causal Model (SCM) provides a powerful framework for representing and analyzing causal relationships, offering a systematic approach to climb the Ladder of Causation.
- Interventions refer to actively manipulating a variable. Each intervention defines a new joint distribution but a statistical model can only captures one of them.
- Counterfactuals allow us to envision the outcomes of different decisions through the lens of imagination and retrospection.
- Causal Factorization decompose a joint distribution into independent causal mechanisms, yielding practical computational advantages and is robust to variations.

# Reinforcement Learning (RL)



Receive an initial observation

# Reinforcement Learning (RL)



Make a decision

# Reinforcement Learning (RL)



Receive and learning from feedback

# Reinforcement Learning (RL)



The agent-environment feedback loop

# Markov Decision Process (MDP)

**Definition** (Markov decision process). An MDP  $\mathcal{M}$  is specified by a tuple  $\{\mathcal{S}, \mathcal{A}, P, R, \mu_0, \gamma\}$ , where

- $\mathcal{S}$  denotes the state space and  $\mathcal{A}$  denotes the action space,
- $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability function that yields the probability of transitioning into the next states  $s_{t+1}$  after taking an action  $a_t$  at the current state  $s_t$ ,
- $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function that assigns the immediate reward for taking an action  $a_t$  at state  $s_t$ ,
- $\mu_0 : \mathcal{S} \rightarrow [0, 1]$  is the probability distribution that specifies the generation of the initial state, and
- $\gamma \in [0, 1]$  denotes the discount factor that accounts for how much future events lose their value as time passes.

RL aims to maximize the expected cumulative reward rather than the immediate one.

$$G_t = R_t + \gamma R_{t+1} + \dots + \gamma^T R_{t+T}$$

$R_t$

# Markov Decision Process (MDP)

**Definition** (Markov decision process). An MDP  $\mathcal{M}$  is specified by a tuple  $\{\mathcal{S}, \mathcal{A}, P, R, \mu_0, \gamma\}$

In essence, we want the agent to maximize

$$\mathbb{E}[G_t | S_t, \text{do}(A_t = a_t)],$$

- + the expected cumulative reward across a sequence of interventions.



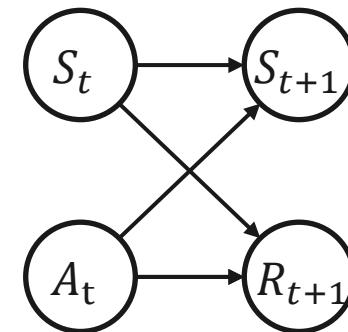
RL aims to maximize the expected cumulative reward rather than the immediate one.

# Markov Decision Process (MDP)

**Definition** (Markov decision process). An MDP  $\mathcal{M}$  is specified by a tuple  $\{\mathcal{S}, \mathcal{A}, P, R, \mu_0, \gamma\}$

We can always cast an MDP into an SCM without imposing any extra constraints.

- The state, action, and reward at each step correspond to endogenous variables.
- The state transition and reward functions are casted into structural equations  $\mathcal{F}$  in the SCM, represented by deterministic functions with independent exogenous variables.
- The initial state  $S_0$  can be considered an exogenous variable such that  $S_0 \in U$ .



Causal Graph

# Policy

**Definition** (Policy). A policy is defined as the probability distribution of actions at a give state:

$$\pi(A_t = a | S_t = s), \forall S_t \in \mathcal{S},$$

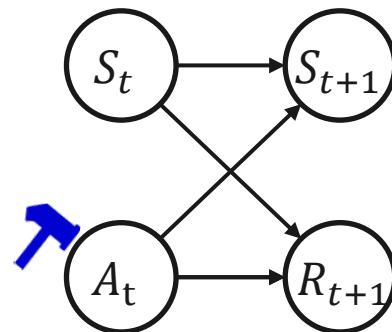
where  $A_t \in \mathcal{A}(s)$  is the state specific action space.

# Policy

**Definition (Policy).** A policy is defined as the probability distribution of actions at a give state:

$$\pi(A_t = a | S_t = s), \forall S_t \in \mathcal{S},$$

where  $A_t \in \mathcal{A}(s)$  is the state specific action space.

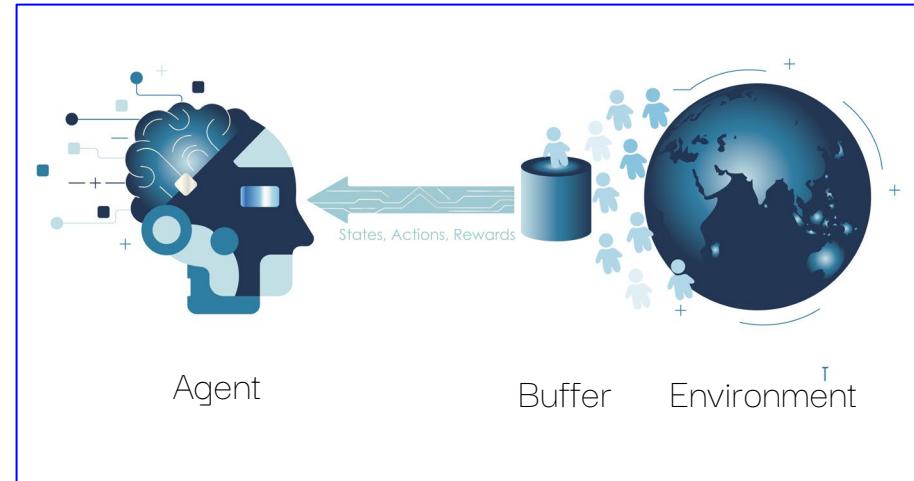


A policy  $\pi$  performs a soft intervention that preserves the dependency of the action on the state

# Categorizing RL Methods



Online Reinforcement Learning



Offline Reinforcement Learning

# Categorizing RL Methods

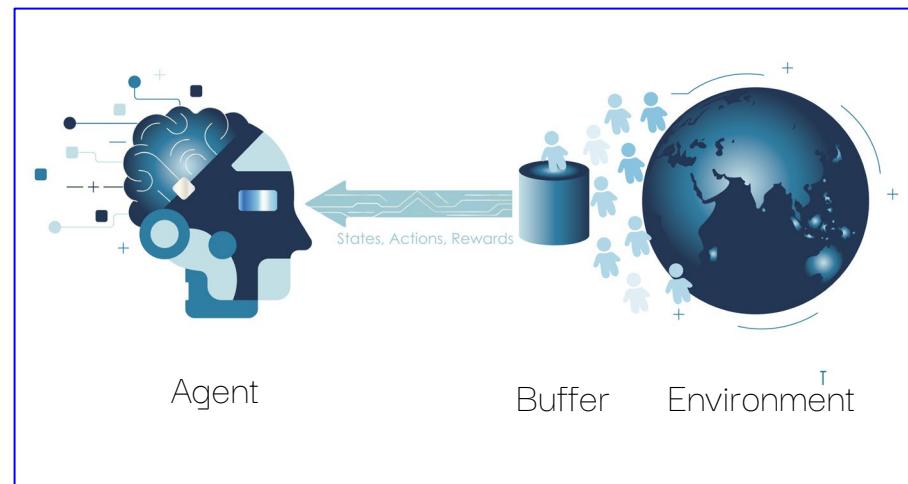


Agent

Environment

## Online Reinforcement Learning

The agent can **actively intervene** in the environment.



Agent

Buffer

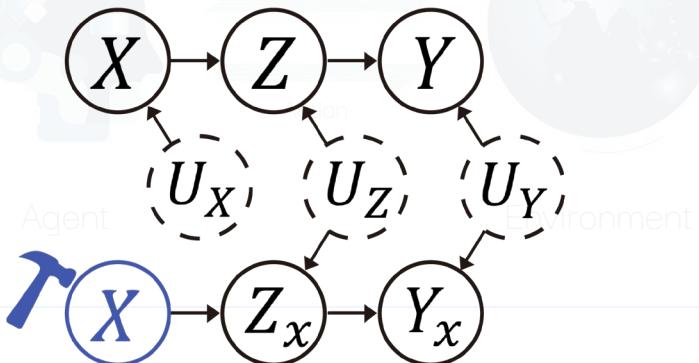
Environment

## Offline Reinforcement Learning

The agent can only **passively observe** the outcomes.

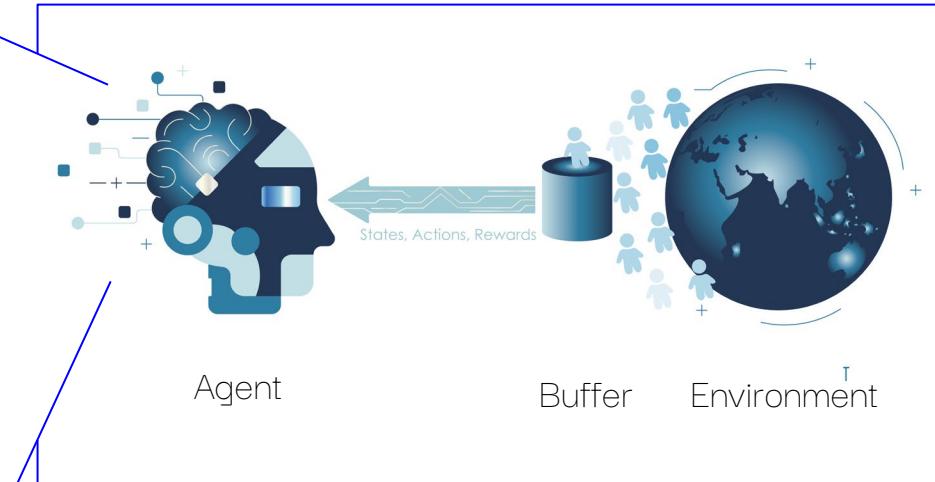
# Categorizing RL Methods

Q: Considering that they consumed rotten meat in reality, would sailors have been protected from scurvy if they had consumed enough citrus fruit?



$$P(Y_x=\text{orange} | X = \text{rotten meat}, Y = \text{sailor})$$

The agent cannot directly intervene in the environment.



Offline Reinforcement Learning

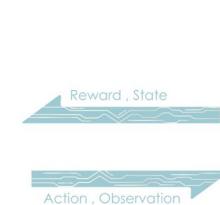
The agent can only **passively observe** the outcomes.

# Categorizing RL Methods

Model-free methods involve learning optimal policies or value functions directly from interaction with the environment without explicitly building a model of the environment's dynamics.



Model-based methods, on the other hand, revolve around creating and utilizing an explicit model of the environment to simulate and plan ahead for making informed decisions in reinforcement learning scenarios.

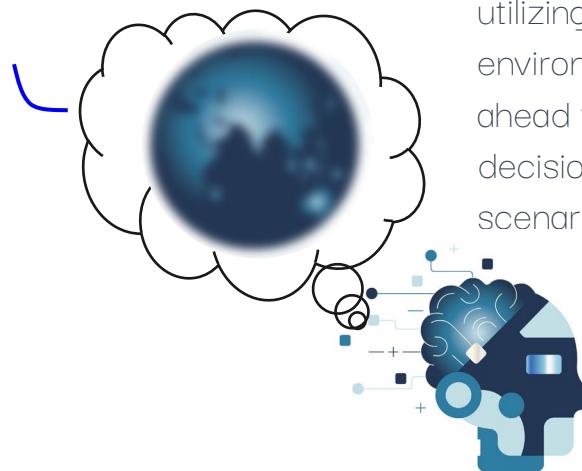


# Categorizing RL Methods



Yann Lecun

The ability to learn “world models” – internal models of **how the world works** – may be the key to build human-level AI.



**Model-based methods**, on the other hand, revolve around creating and utilizing an explicit model of the environment to simulate and plan ahead for making informed decisions in reinforcement learning scenarios.



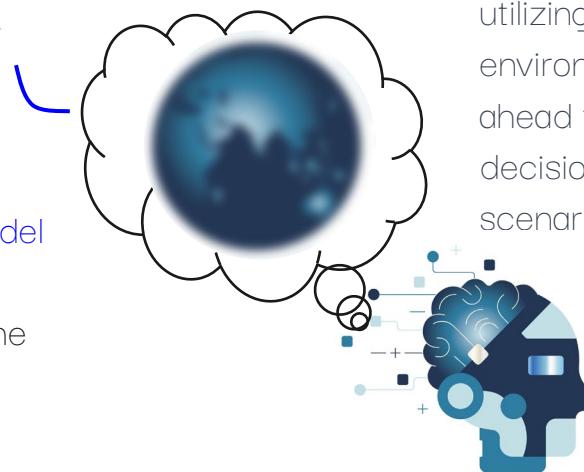
# Categorizing RL Methods



Yann LeCun

The ability to learn “world models” – internal models of **how the world works** – may be the key to build human-level AI.

Q: How to construct an **internal causal model** that describes the causal relationships between variables (concepts) governing the data generation process in our world?



**Model-based methods**, on the other hand, revolve around creating and utilizing an explicit model of the environment to simulate and plan ahead for making informed decisions in reinforcement learning scenarios.



---

# Causal Reinforcement Learning

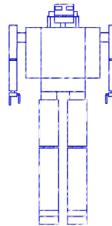
---

# Causal Reinforcement Learning

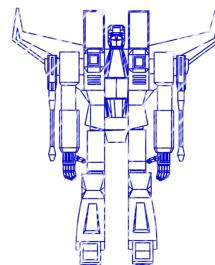
**Definition** (Causal Reinforcement Learning). Causal RL is an umbrella term for RL approaches that incorporate additional assumptions or prior knowledge to analyze and understand the causal mechanisms underlying actions and their consequences, enabling agents to make more informed and effective decisions.

# Causal Reinforcement Learning

**Definition** (Causal Reinforcement Learning). Causal RL is an umbrella term for RL approaches that incorporate **additional assumptions or prior knowledge** to analyze and understand the causal mechanisms underlying actions and their consequences, enabling agents to make more informed and effective decisions.



Traditional RL methods focus on learning the optimal policies through interactions with the environment, **without explicitly considering the causal relationships** between actions and outcomes.

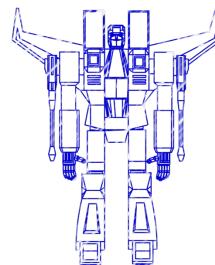


Causal RL, in contrast, go beyond the traditional framework by incorporating additional assumptions or prior knowledge about causality, empowering agents with a **deeper understanding of the underlying dynamics** of the world.

# Causal Reinforcement Learning

**Definition** (Causal Reinforcement Learning). Causal RL is an umbrella term for RL approaches that incorporate **additional assumptions or prior knowledge** to analyze and understand the causal mechanisms underlying actions and their consequences, enabling agents to make more informed and effective decisions.

Go beyond the evidence  
of memory and senses!



Causal RL, in contrast, go beyond the traditional framework by incorporating additional assumptions or prior knowledge about causality, empowering agents with a deeper understanding of the underlying dynamics of the world.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A Markov Decision Process (MDP) describes the dynamics of the environment during interaction, and it can also be represented as an SCM.
- A policy guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or knowledge regarding the underlying causal relationships within the data to inform decision-making.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A [Markov Decision Process \(MDP\)](#) describes the dynamics of the environment during interaction, and it can also be represented as an [SCM](#).
- A policy guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or knowledge regarding the underlying causal relationships within the data to inform decision-making.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A Markov Decision Process (MDP) describes the dynamics of the environment during interaction, and it can also be represented as an SCM.
- A [policy](#) guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or knowledge regarding the underlying causal relationships within the data to inform decision-making.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A Markov Decision Process (MDP) describes the dynamics of the environment during interaction, and it can also be represented as an SCM.
- A policy guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or knowledge regarding the underlying causal relationships within the data to inform decision-making.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A Markov Decision Process (MDP) describes the dynamics of the environment during interaction, and it can also be represented as an SCM.
- A policy guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or knowledge regarding the underlying causal relationships within the data to inform decision-making.

# Summary

- Reinforcement Learning (RL) focuses on sequential decision-making problems, where an agent intervenes in an environment with the goal of maximizing cumulative rewards.
- A Markov Decision Process (MDP) describes the dynamics of the environment during interaction, and it can also be represented as an SCM.
- A policy guides an agent's decision-making by mapping states to appropriate actions.
- RL methods can be categorized in various ways, such as online vs. offline and model-free vs. model-based methods.
- Causal RL aims to integrate assumptions or prior knowledge about the underlying causal relationships within the data to enhance decision-making.

# Tutorial Outline

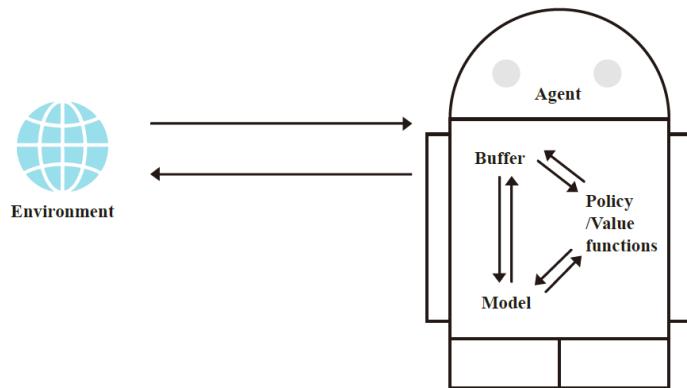
## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Causal RL



# Causal RL

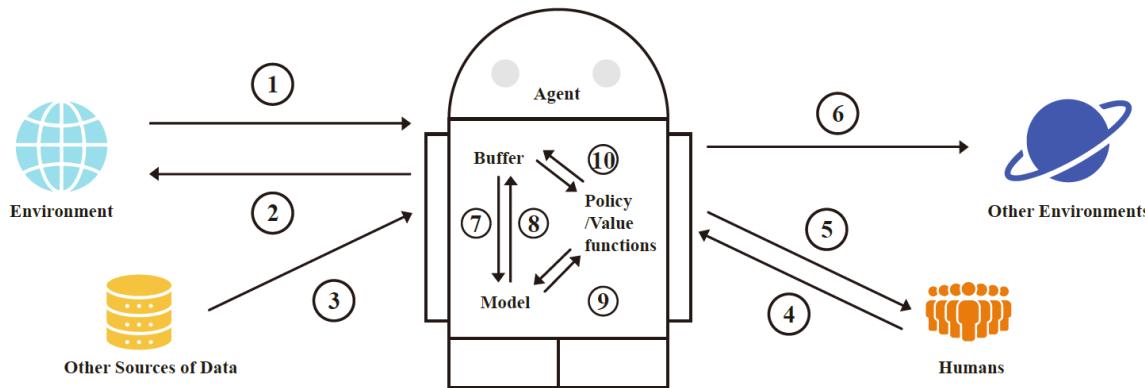


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans. 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Sample Efficiency



AlphaGo

$3 \times 10^7$  games of self-play



AlphaStar

$2 \times 10^2$  years of self-play



OpenAI Rubik's Cube

$10^4$  years of simulation

# Sample Efficiency



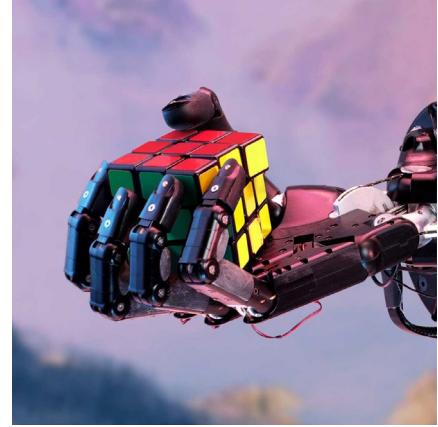
AlphaGo

$3 \times 10^7$  games of self-play



AlphaStar

$2 \times 10^2$  years of self-play

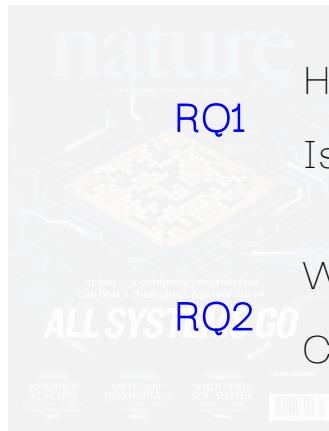


OpenAI Rubik's Cube

$10^4$  years of simulation

Does human learning also require such a large sample size?

# Sample Efficiency



How can exploration be made more efficient?

Is all unexplored area in the state space equally important?

What are the causal variables that govern the environmental dynamics?

Can we accelerate the learning process utilizing these factors?

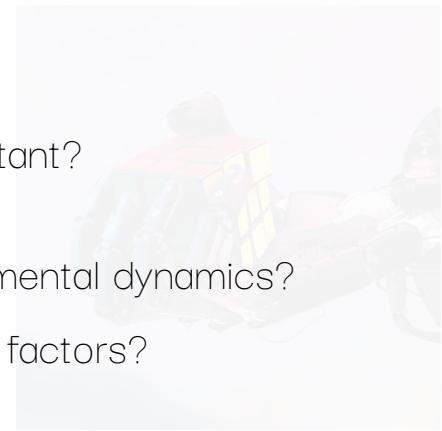
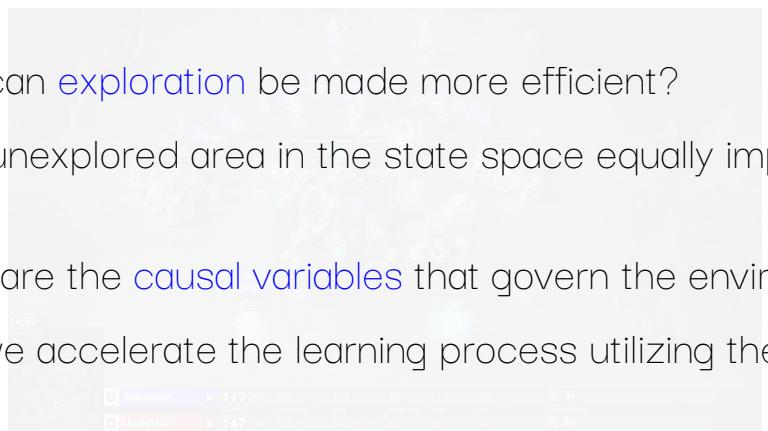
How can agents be equipped with introspective capabilities?

Can agents effectively learn from imaginative experiences?

$3 \times 10^7$  games of self-play

$2 \times 10^2$  years of self-play

$10^4$  years of simulation



RQ3  
AlphaGo

AlphaStar

OpenAI Rubik's Cube

# Directed Exploration

RQ1

How can exploration be made more efficient?

Is all unexplored area in the state space equally important?

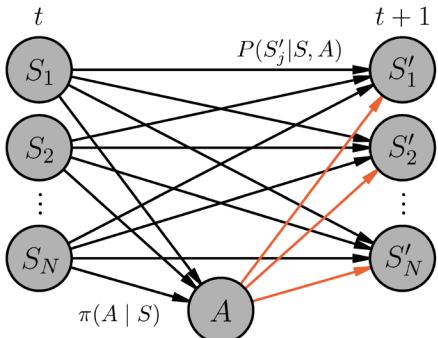
# Directed Exploration

Causal Graph

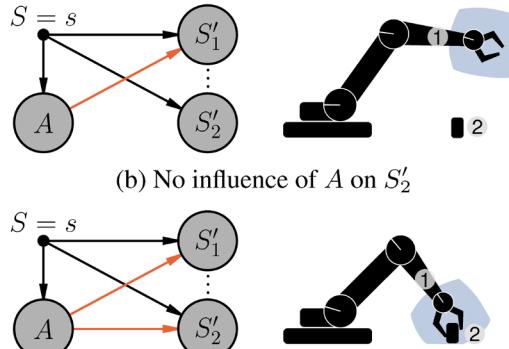
For more information, check  
 "Causal Influence Detection  
 for Improving Efficiency in  
 Reinforcement Learning".  
 NeurIPS 2021

RQ1 How can exploration be made more efficient?

Is all unexplored area in the state space equally important?



(a) Causal Graph  $\mathcal{G}$



(b) No influence of  $A$  on  $S'_2$   
 (c) Influence of  $A$  on  $S'_1$  and  $S'_2$

Figure 1: Causal graphical model capturing the environment transition from state  $S$  to  $S'$  by action  $A$ , factorized into state components. (a): Viewed globally over all time steps, all components of the state and the action can influence all state components at the next time step. (b, c): Given a situation  $S = s$ , some influences may or may not hold in the *local causal graph*  $\mathcal{G}_{S=s}$ . In this paper, our aim is to detect which influence the action has on  $S'$ , i.e. the presence of the red arrows.

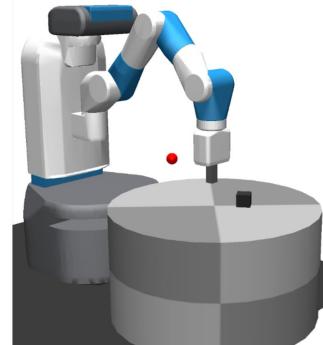


Figure 8: FETCHROT-TABLE. The table rotates periodically.

# Directed Exploration

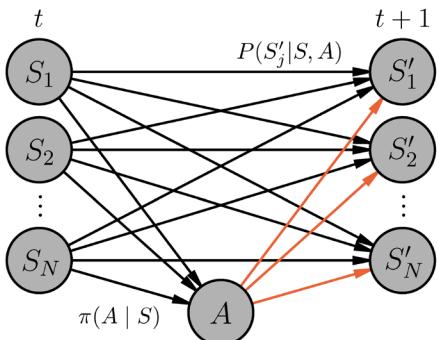
Causal Graph

For more information, check  
“Causal Influence Detection  
for Improving Efficiency in  
Reinforcement Learning”,  
NeurIPS 2021

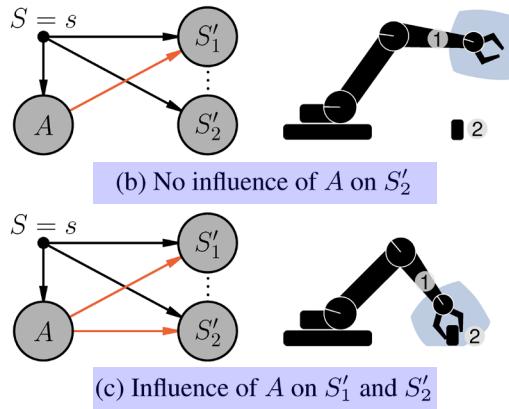
RQ1

How can exploration be made more efficient?

Is all unexplored area in the state space equally important?



(a) Causal Graph  $\mathcal{G}$



(b) No influence of  $A$  on  $S'_2$

(c) Influence of  $A$  on  $S'_1$  and  $S'_2$

How to infer the influence the action has in a specific state configuration  $S = s$ ?

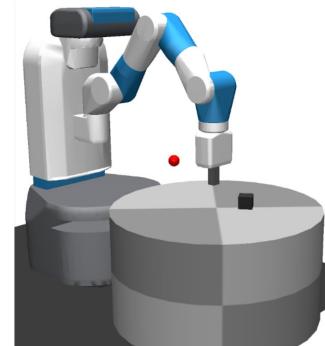


Figure 8: FETCHROT-TABLE. The table rotates periodically.

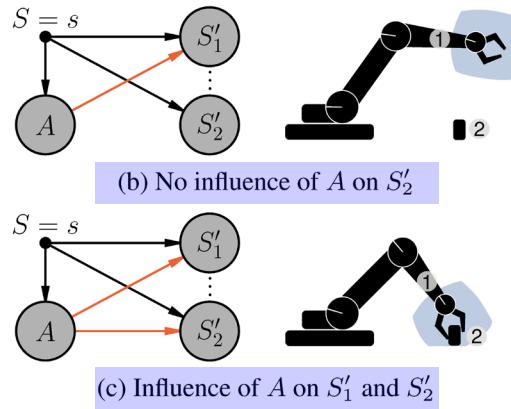
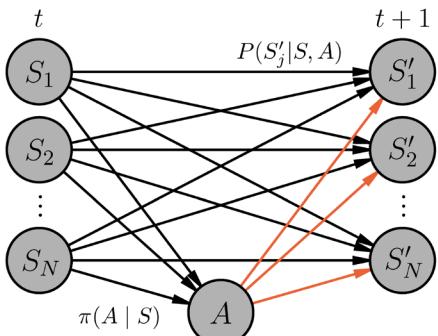
# Directed Exploration

Causal Graph

For more information, check  
 "Causal Influence Detection  
 for Improving Efficiency in  
 Reinforcement Learning".  
 NeurIPS 2021

RQ1 How can exploration be made more efficient?

Is all unexplored area in the state space equally important?



How to infer the influence the action has in a specific state configuration  $S = s$ ?

Conditional Action Influence (CAI)

$$C^j(s) := I(S'_j; A \mid S = s) = \mathbb{E}_{a \sim \pi} [\text{D}_{\text{KL}}(P_{S'_j|s,a} \parallel P_{S'_j|s})]$$

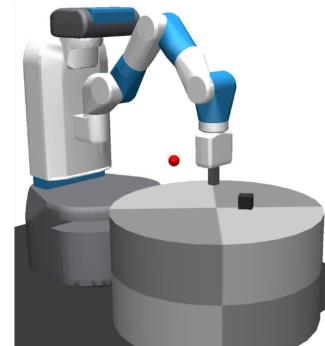
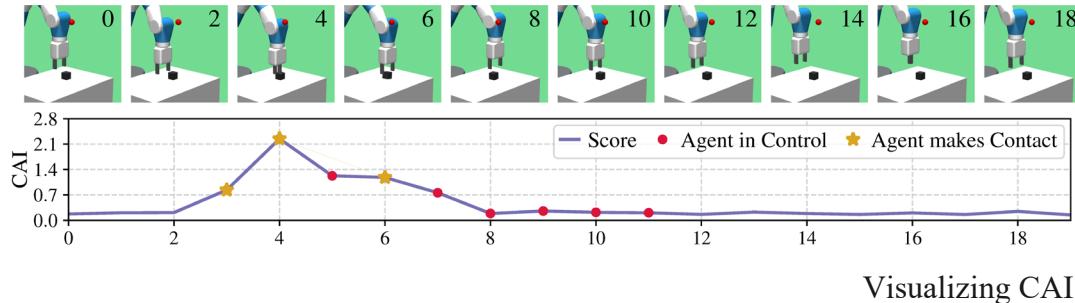


Figure 8: FETCHROT-TABLE. The table rotates periodically.

# Directed Exploration

Causal Graph

- RQ1 How can *exploration* be made more efficient?  
Is all unexplored area in the state space equally important?



For more information, check  
“Causal Influence Detection  
for Improving Efficiency in  
Reinforcement Learning”,  
NeurIPS 2021

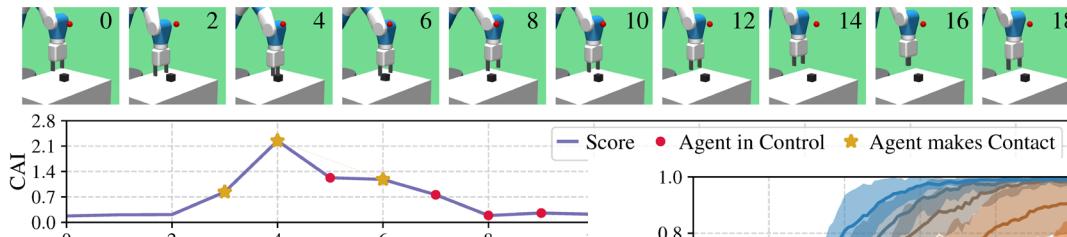
# Directed Exploration

Causal Graph

For more information, check  
“Causal Influence Detection  
for Improving Efficiency in  
Reinforcement Learning”,  
NeurIPS 2021

RQ1 How can exploration be made more efficient?

Is all unexplored area in the state space equally important?



- Causal influence as intrinsic motivation
- Greedy w.r.t action influence
- Causal influence as replay priority

Figure 5: Performance of *active exploration* in `FETCHPICKANDPLACE` depending on the fraction of exploratory actions chosen actively (Eq. 6) from a total of 30% exploratory actions.

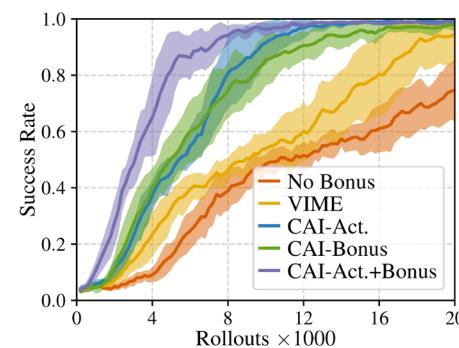


Figure 6: Experiment comparing exploration strategies on `FETCHPICKANDPLACE`. The combination of active exploration and reward bonus yields the largest sample efficiency.

# Causal RL

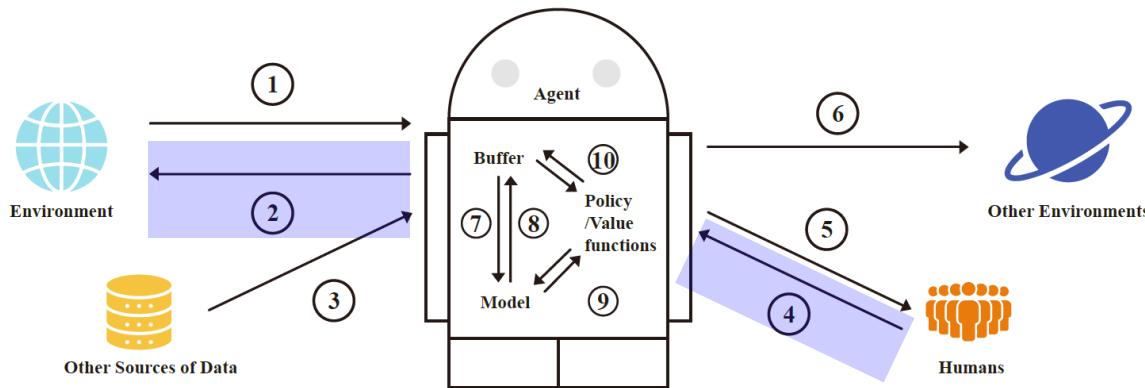


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans; 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Causal Representation

RQ2

What are the causal variables that govern the environmental dynamics?

Can we accelerate the learning process utilizing these factors?

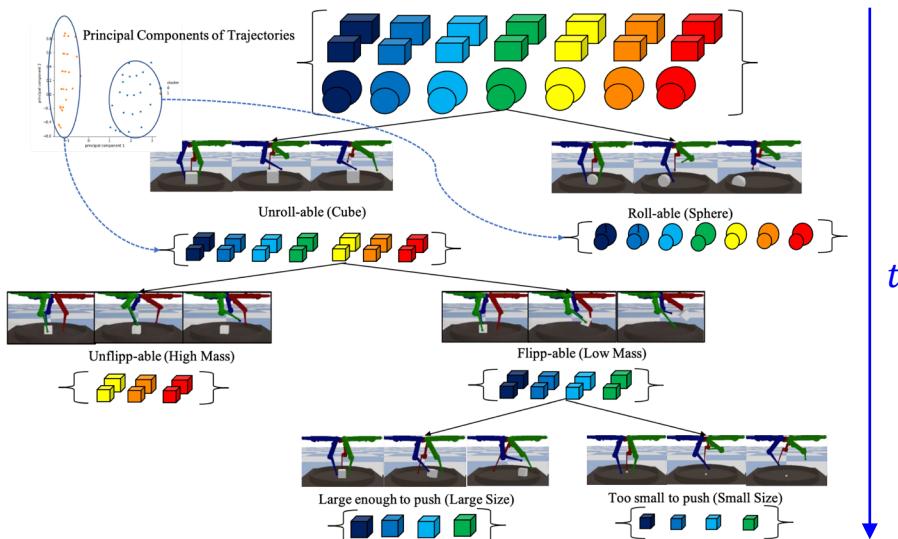
For more information, check  
 “Causal Curiosity: RL Agents  
 Discovering Self-supervised  
 Experiments for  
 Causal Representation  
 Learning”, ICML 2021

# Causal Representation

What are the *causal variables* that govern the environmental dynamics?

RQ2

Can we accelerate the learning process utilizing these factors?



*Figure 3.* Discovered hierarchical latent space. The agent learns experiments that differentiate the full set of blocks in ShapeSizeMass into hierarchical binary clusters. At each level, the environments are divided into 2 clusters on the basis of the value of a single causal factor. We also show the principal components of the trajectories in the top left. For brevity, the full extent of the tree is not depicted here. For each level of hierarchy  $k$ , there are  $2^k$  number of clusters.

# Causal Representation

What are the *causal variables* that govern the environmental dynamics?

RQ2

Can we accelerate the learning process utilizing these factors?

For more information, check  
“Causal Curiosity: RL Agents  
Discovering Self-supervised  
Experiments for  
Causal Representation  
Learning”, ICML 2021

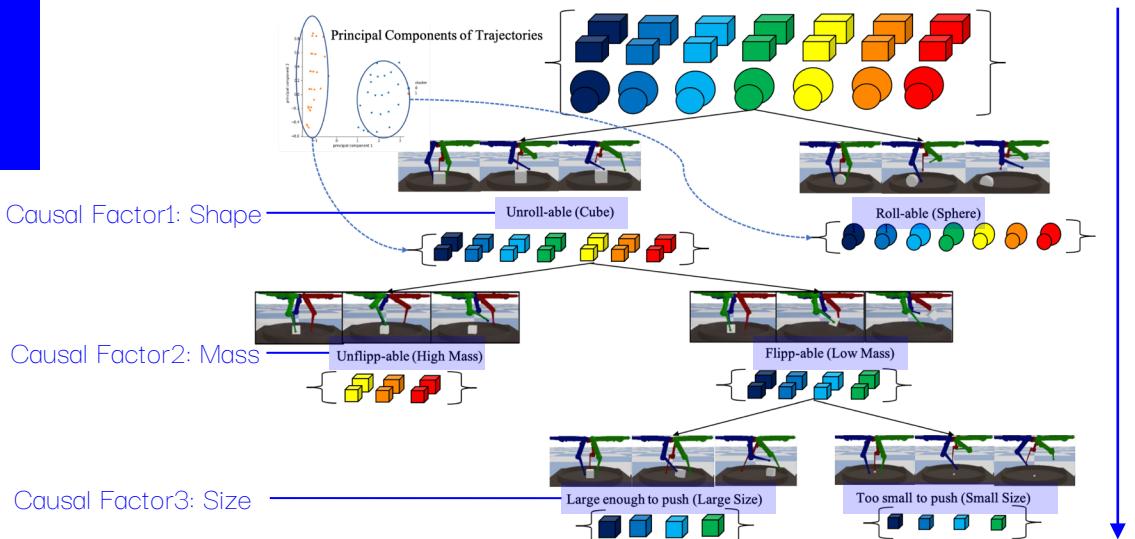


Figure 3. Discovered hierarchical latent space. The agent learns experiments that differentiate the full set of blocks in ShapeSizeMass into hierarchical binary clusters. At each level, the environments are divided into 2 clusters on the basis of the value of a single causal factor. We also show the principal components of the trajectories in the top left. For brevity, the full extent of the tree is not depicted here. For each level of hierarchy  $k$ , there are  $2^k$  number of clusters.

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

# Causal Representation

What are the *causal variables* that govern the environmental dynamics?

RQ2

Can we accelerate the learning process utilizing these factors?

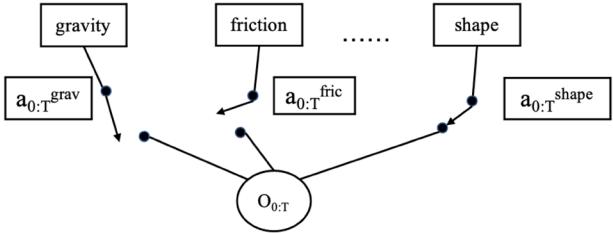


Figure 2. Gated Causal Graph. A subset of the unobserved parent causal variables influence the observed variable  $O$ . The action sequence  $a_{0:T}$  serves a gating mechanism, allowing or blocking particular edges of the causal graph using the implicit Causal Selector Function (Equation 4).

For more information, check  
“Causal Curiosity: RL Agents Discovering Self-supervised Experiments for Causal Representation Learning”, ICML 2021

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

# Causal Representation

ICM

- What are the *causal variables* that govern the environmental dynamics?
- RQ2  
Can we accelerate the learning process utilizing these factors?

Independent Causal Mechanism



The information in an observed trajectory is *the sum of information* “injected” into it from the multiple causes

The information content will be greater for a larger *number of causal parents* in the graph.

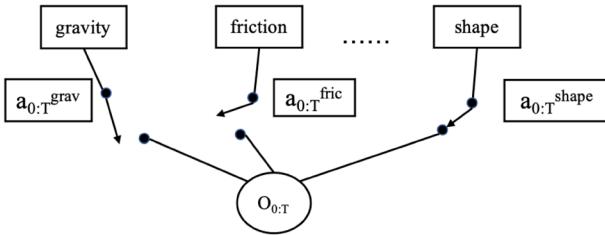


Figure 2. Gated Causal Graph. A subset of the unobserved parent causal variables influence the observed variable  $O$ . The action sequence  $a_{0:T}$  serves a gating mechanism, allowing or blocking particular edges of the causal graph using the implicit Causal Selector Function (Equation 4).

For more information, check  
“Causal Curiosity: RL Agents Discovering Self-supervised Experiments for Causal Representation Learning”, ICM 2021

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

# Causal Representation

ICM

- What are the *causal variables* that govern the environmental dynamics?
- RQ2  
Can we accelerate the learning process utilizing these factors?

One-Factor-at-A-Time



We want to search for one causal factor at a time.



Find an action sequence for which the number of causal parents of  $\mathbf{O}$  is low, i.e., the complexity of  $\mathbf{O}$ .



Independent Causal Mechanism

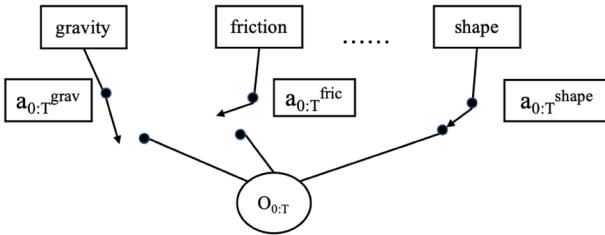


Figure 2. Gated Causal Graph. A subset of the unobserved parent causal variables influence the observed variable  $\mathbf{O}$ . The action sequence  $\mathbf{a}_{0:T}$  serves a gating mechanism, allowing or blocking particular edges of the causal graph using the implicit Causal Selector Function (Equation 4).

For more information, check  
“Causal Curiosity: RL Agents Discovering Self-supervised Experiments for Causal Representation Learning”, ICM 2021

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

# Causal Representation

ICML

- What are the *causal variables* that govern the environmental dynamics?  
**RQ2** Can we accelerate the learning process utilizing these factors?

For more information, check  
“Causal Curiosity: RL Agents  
Discovering Self-supervised  
Experiments for  
Causal Representation  
Learning”, ICML 2021

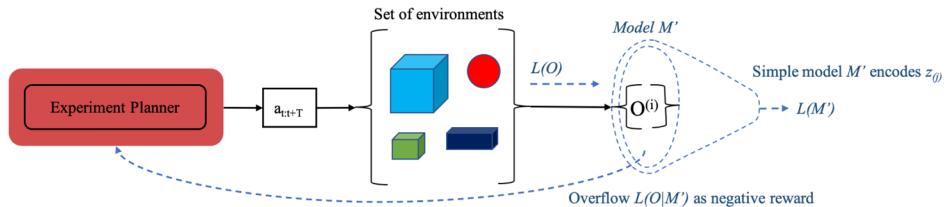


Figure 9. Overview of training. The experiment planner generates a trajectory of actions which is applied to each of the environments with varying causal factors namely mass, shape and size of blocks. For each environment, an observation trajectory or state  $\mathbf{O}^i \in \mathcal{O}$  is obtained. A simple model with fixed low expressive power is used to approximate the generative model for  $\mathbf{O}$ . The "information overflow"  $L(\mathbf{O}|M)$  is returned as negative reward forcing  $O$  to be caused by few causal factors.

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

Minimizing the complexity



Maximizing the causal curiosity

$$a_{0:T}^* = \underset{a_{0:T}}{\operatorname{argmin}} L(O|M)$$

$$a_{0:T}^* = \underset{a_{0:T}}{\operatorname{argmax}} -L(O|M)$$

# Causal Representation

ICM

- What are the *causal variables* that govern the environmental dynamics?  
**RQ2**  
 Can we accelerate the learning process utilizing these factors?

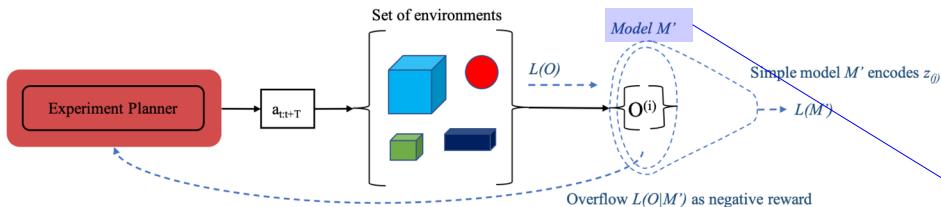


Figure 9. Overview of training. The experiment planner generates a trajectory of actions which is applied to each of the environments with varying causal factors namely mass, shape and size of blocks. For each environment, an observation trajectory or state  $\mathbf{O}^i \in \mathcal{O}$  is obtained. A simple model with fixed low expressive power is used to approximate the generative model for  $\mathbf{O}$ . The "information overflow"  $L(\mathbf{O}|M)$  is returned as negative reward forcing  $O$  to be caused by few causal factors.

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?

Minimizing the complexity



Maximizing the causal curiosity

$$a_{0:T}^* = \underset{a_{0:T}}{\operatorname{argmin}} L(O|M)$$

$$a_{0:T}^* = \underset{a_{0:T}}{\operatorname{argmax}} -L(O|M)$$

Assume  $M$  is a bimodal clustering model.

# Causal Representation

ICM

What are the *causal variables* that govern the environmental dynamics?  
**RQ2**  
 Can we accelerate the learning process utilizing these factors?

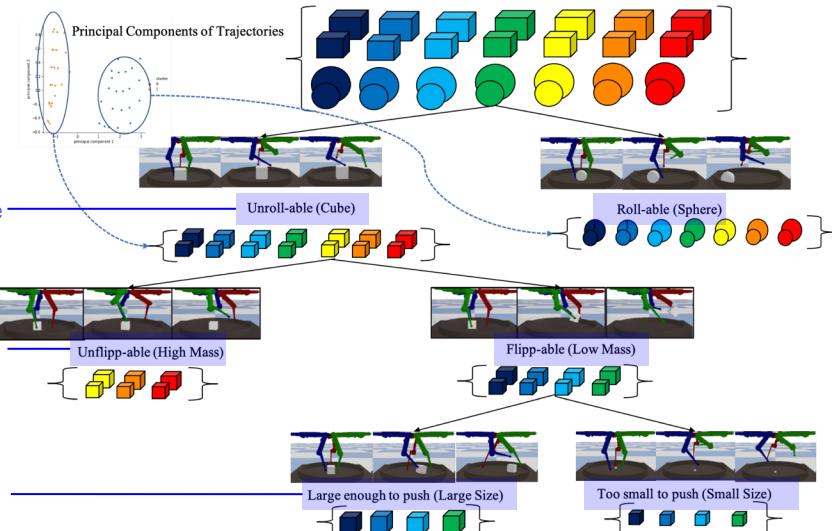
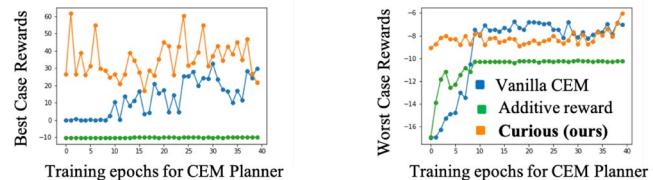


Figure 3. Discovered hierarchical latent space. The agent learns experiments that differentiate the full set of blocks in  $\text{Shape} \times \text{Size} \times \text{Mass}$  into hierarchical binary clusters. At each level, the environments are divided into 2 clusters on the basis of the value of a single causal factor. We also show the principal components of the trajectories in the top left. For brevity, the full extent of the tree is not depicted here. For each level of hierarchy  $k$ , there are  $2^k$  number of clusters.

For more information, check  
 “Causal Curiosity: RL Agents  
 Discovering Self-supervised  
 Experiments for  
 Causal Representation  
 Learning”, ICML 2021

How to empower agents to discover semantically meaningful experimental behaviors rather than maximizing reward for a particular task?



The behaviors discovered by the agents while optimizing causal curiosity show high zero-shot Generalization Ability and converge to the same performance as conventional planners for downstream tasks.

# Causal RL

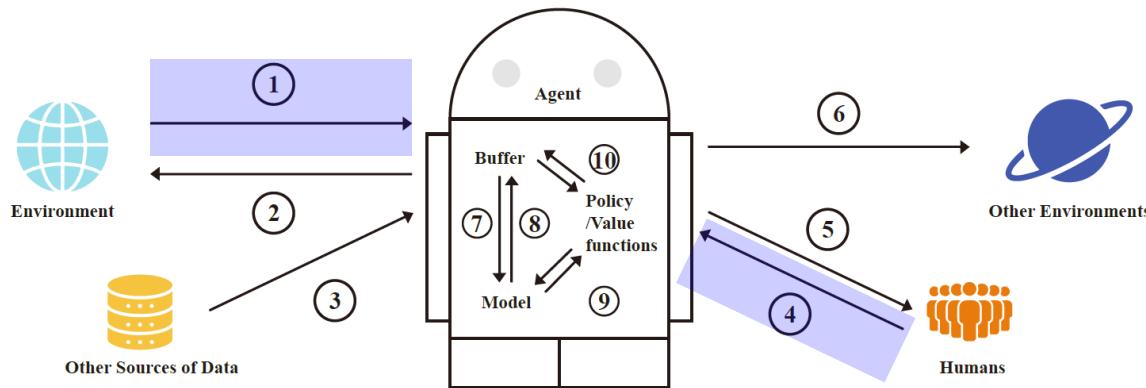


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans. 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Counterfactual Reasoning

RQ3

How can agents be equipped with introspective capabilities?

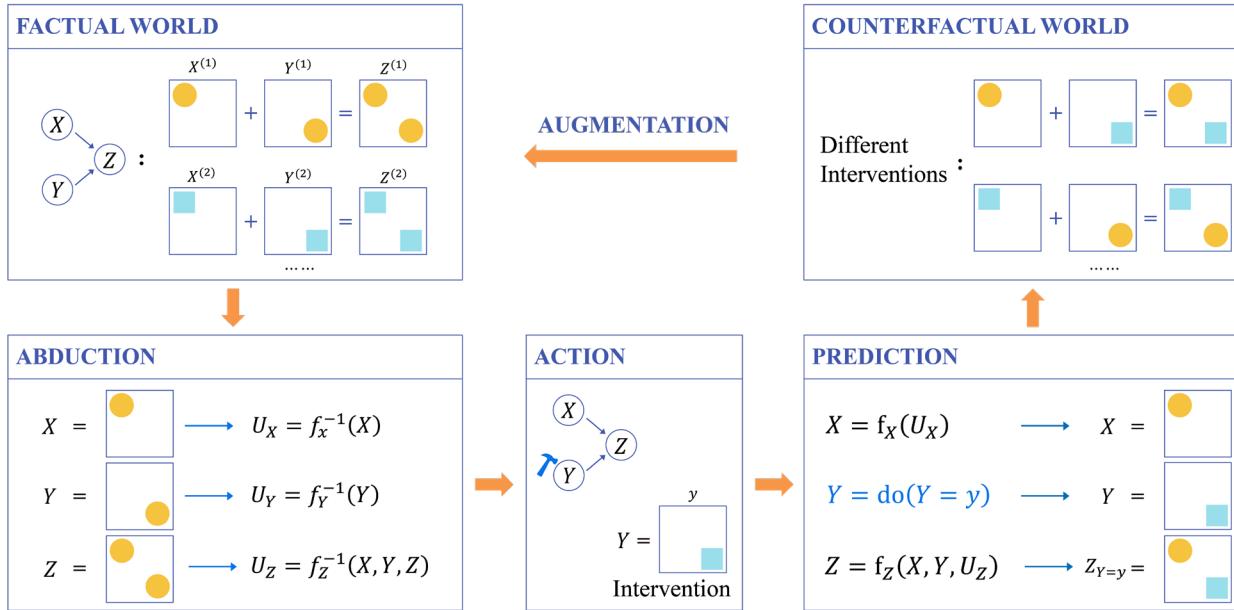
Can agents effectively learn from imaginative experiences?

# Counterfactual Reasoning

How can agents be equipped with introspective capabilities?

RQ3

Can agents effectively learn from imaginative experiences?



# Counterfactual Reasoning

SCM

- How can agents be equipped with introspective capabilities?
- RQ3**  
 Can agents effectively learn from imaginative experiences?

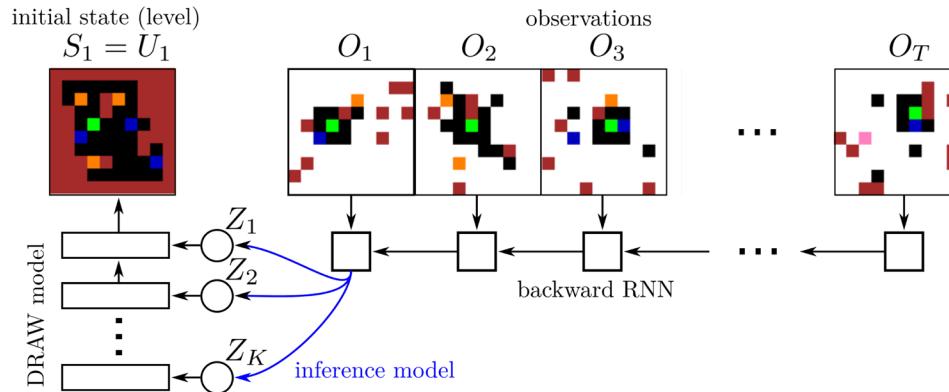
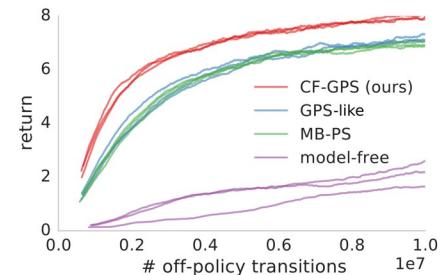


Figure 3: **Top: PO-SOKOBAN.** Shown on the left is a procedurally generated initial state. The agent is shown in green, boxes in yellow, targets in blue and walls in red. The agent does not observe this state but a sequence of observations, which are masked by iid noise with 0.9 probability, except a 3x3 window around the agent. **Bottom: Inference model.** For counterfactual inference in PO-SOKOBAN, we need the (approximate) inference distribution  $p(U_{s1}|\hat{h}_T)$  over the initial state  $U_{s1} = S_1$ , conditioned on the history of observations  $\hat{h}_T$ . We model this distribution using a DRAW generative model with latent variables  $Z$ , which are conditioned on the output of a backward RNN summarizing the observation history.



Counterfactually-Guided Policy Search (CF-GPS) outperforms a naive model-based RL (MB-PS) algorithm as well as model-free methods

# Causal RL

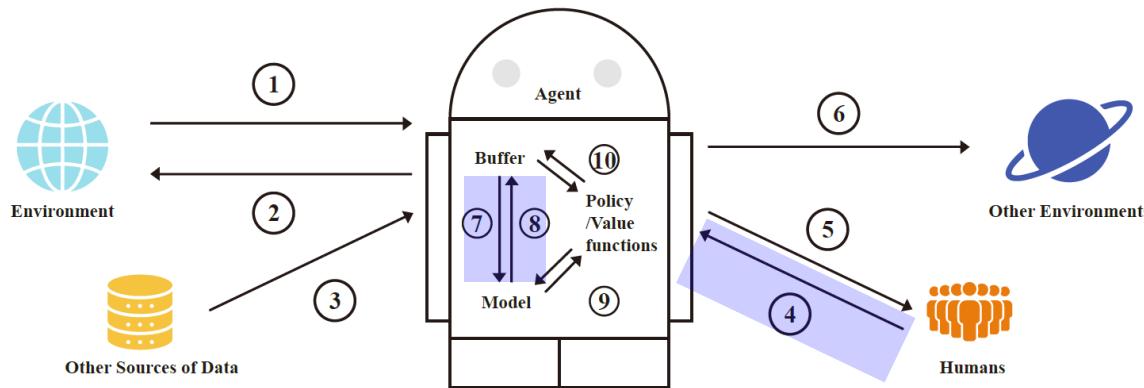


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans; 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Tutorial Outline

## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Generalization Ability



Jacob Andreas  
@jacobandreas

Deep RL is popular because it's the only area in ML where it's socially acceptable to train on the test set.

6:27 AM · Oct 29, 2017

---

111 Reposts   10 Quotes   627 Likes   4 Bookmarks

# Generalization Ability

Jacob Andreas [Submitted on 19 Sep 2017 (v1), last revised 30 Jan 2019 (this version, v3)]

## Deep Reinforcement Learning that Matters

Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, David Meger

In recent years, significant progress has been made in solving challenging problems across various domains using deep reinforcement learning (RL). Reproducing existing work and accurately judging the improvements offered by novel methods is vital to sustaining this progress. Unfortunately, reproducing results for state-of-the-art deep RL methods is seldom straightforward. In particular, non-determinism in standard benchmark environments, combined with variance intrinsic to the methods, can make reported results tough to interpret. Without significance metrics and tighter standardization of experimental reporting, it is difficult to determine whether improvements over the prior state-of-the-art are meaningful. In this paper, we investigate challenges posed by reproducibility, proper experimental techniques, and reporting procedures. We illustrate the variability in reported metrics and results when comparing against common baselines and suggest guidelines to make future results in deep RL more reproducible. We aim to spur discussion about how to ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted.

# Generalization Ability

Jacob Andreas [Submitted on 19 Sep 2017 (v1), last revised 30 Jan 2019 (this version, v3)] ...

## Deep Reinforcement Learning that Matters

Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, David Meger

In recent reinforcement learning, it has been vital to succeed in complex tasks. In particular, it has been reported that deep learning models can be used to determine the best actions to take in complex environments. However, this work has been largely based on reproducing existing results without providing any guarantees about the quality of the learned policies.

**Reinforcement Learning  
never worked, and “deep”  
only helped a bit.**

results when comparing against common baselines and suggest guidelines to make future results in deep RL more reproducible. We aim to spur discussion about how to ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted.

# Generalization Ability

Jacob Andreas  
[Submitted on 19 Sep 2017 (v1), last revised 30 Jan 2019 (this version, v3)]

## Deep Reinforcement Learning that Matters

Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, David Meger

In recent  
reinforcer  
vital to su  
In particula  
reported t  
to determ  
by reproduc

Reinforcement Learning  
never works, and “deep”  
Reproducibility Crisis  
only helped a bit.

results when comparing against common baselines and suggest guidelines to make future results in deep RL more reproducible. We aim to spur discussion about how to ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted.

# Generalization Ability

Jacob Andreas

[Submitted on 19 Sep 2017 (v1), last revised 30 Jan 2019 (this version, v3)]

## Deep Reinforcement Learning that Matters

What does generalization mean for agents?

RQ1

How can agents achieve reliable generalization despite unknown variations?

RQ2

What knowledge can be transferred?

How can algorithms be designed to facilitate efficient adaptation?

by reproducibility, proper experimental techniques, and reporting procedures. We illustrate the variability in reported metrics and results when comparing against common baselines and suggest guidelines to make future results in deep RL more reproducible. We aim to spur discussion about how to ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted.

# Generalization Ability

RQ1

What does generalization mean for agents?

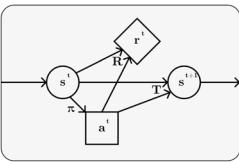
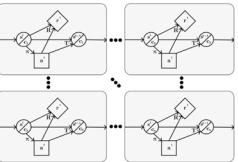
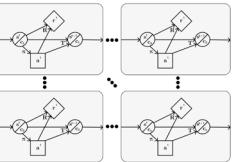
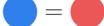
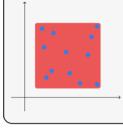
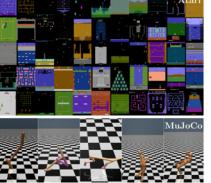
How can agents achieve **reliable generalization** despite unknown variations?

# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

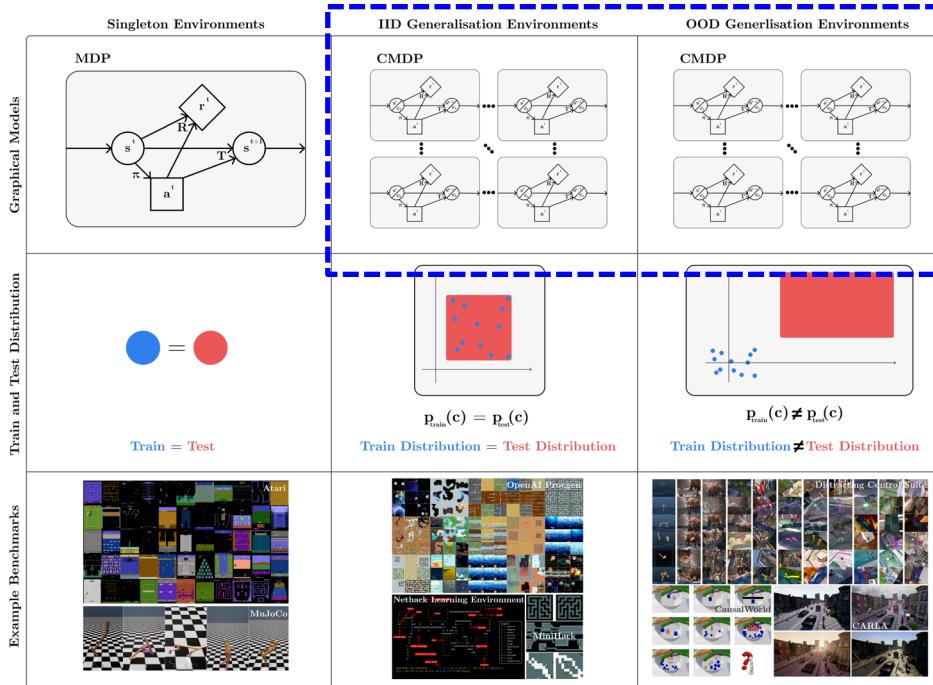
	Singleton Environments	IID Generalisation Environments	OOD Generalisation Environments
Graphical Models			
Train and Test Distribution	 $p_{\text{train}}(c) = p_{\text{test}}(c)$ <b>Train = Test</b>	 $p_{\text{train}}(c) = p_{\text{test}}(c)$ <b>Train Distribution = Test Distribution</b>	 $p_{\text{train}}(c) \neq p_{\text{test}}(c)$ <b>Train Distribution <math>\neq</math> Test Distribution</b>
Example Benchmarks			

# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

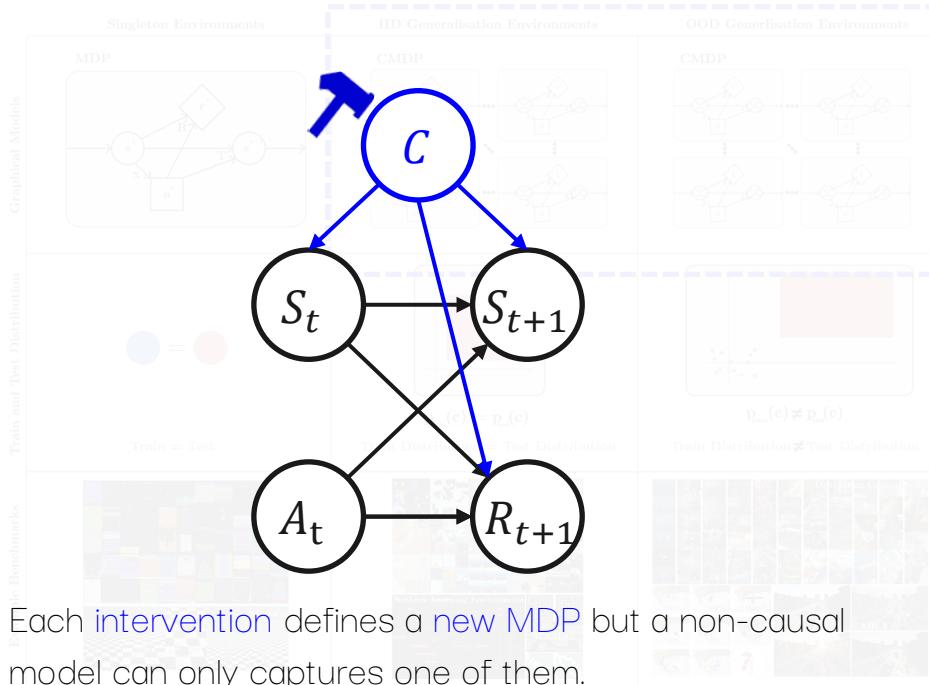


# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

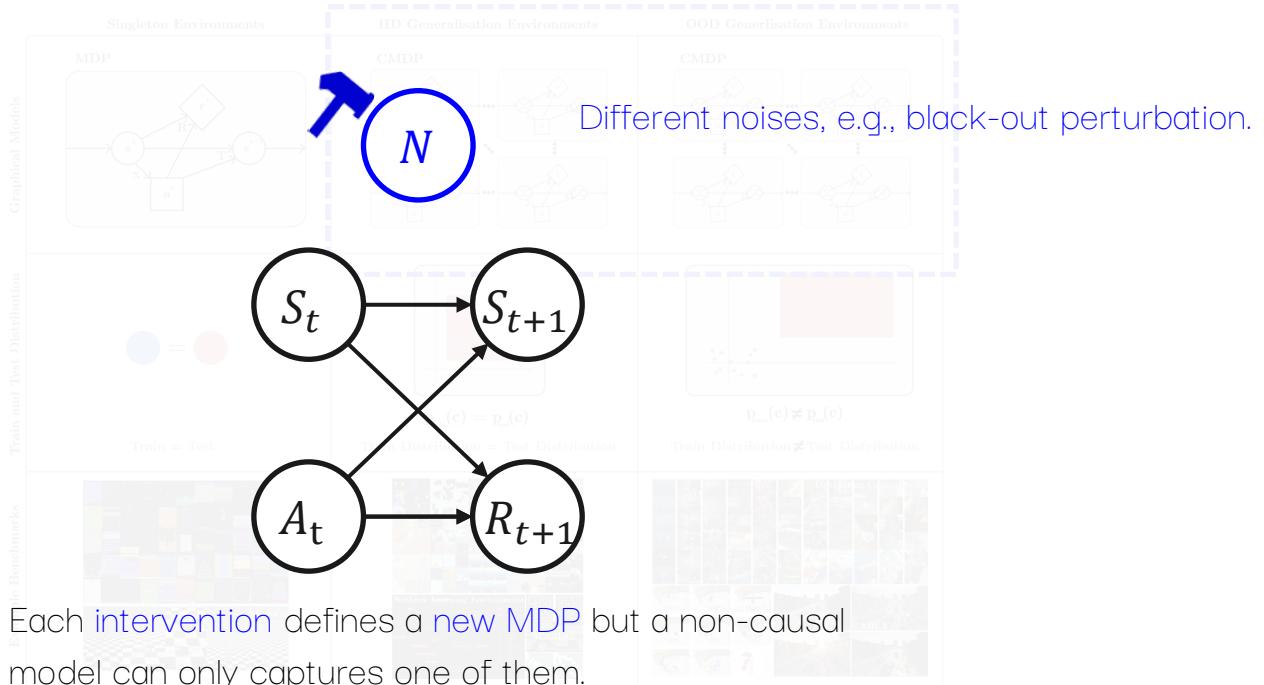


# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

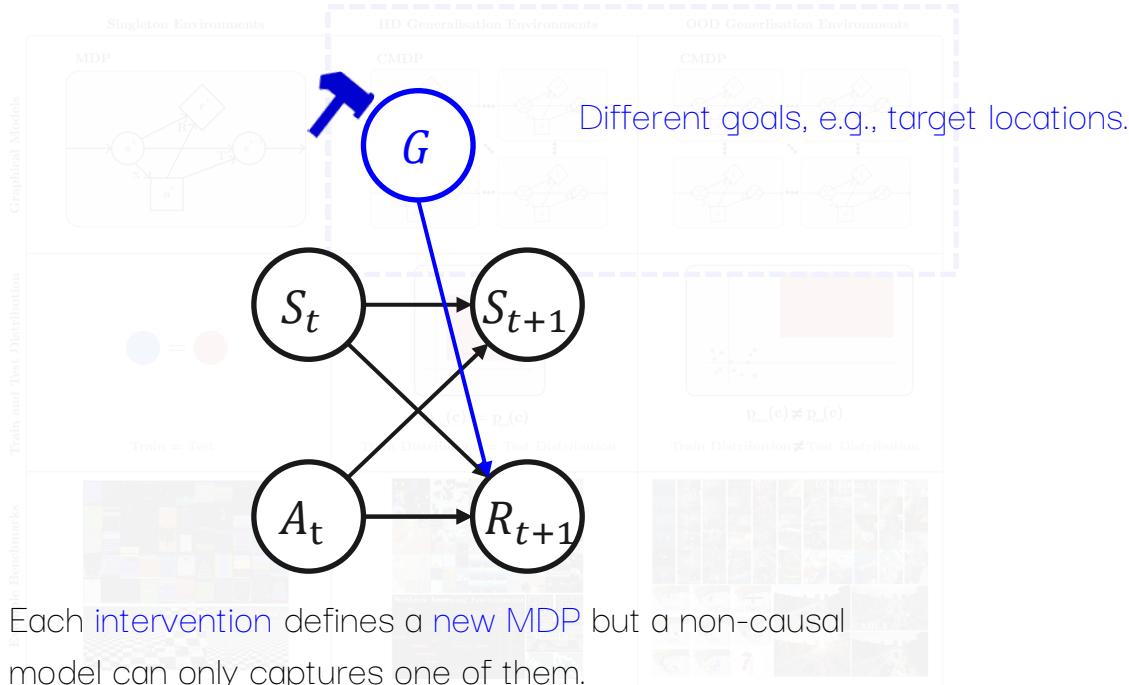


# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

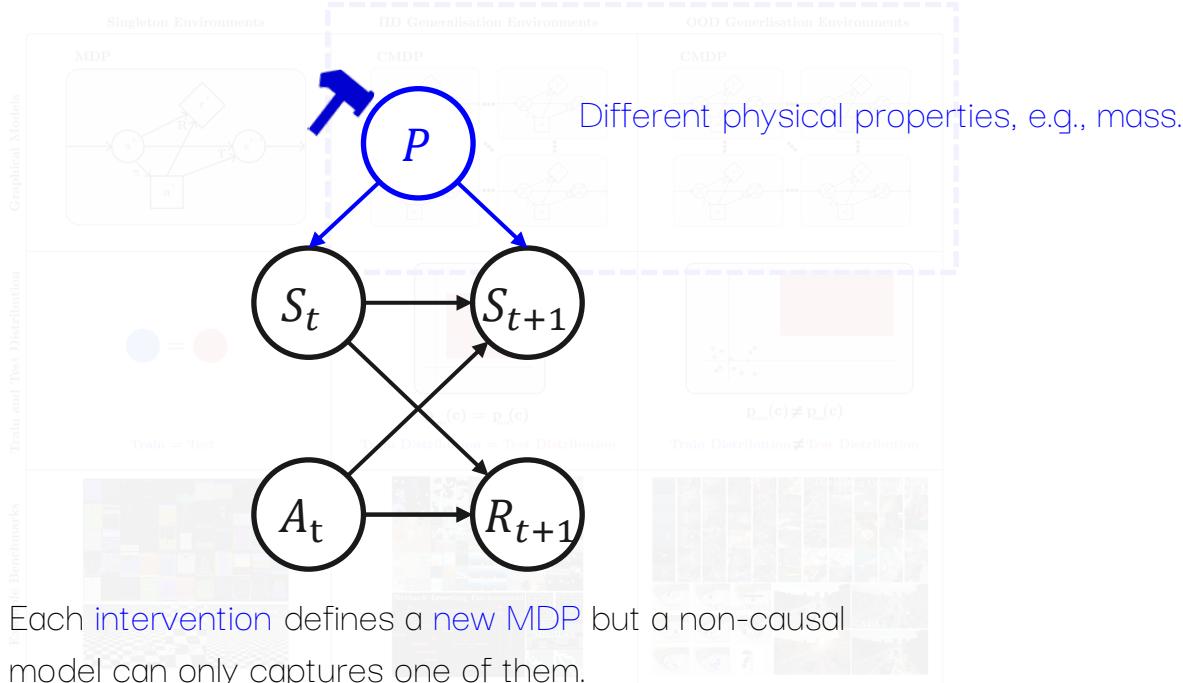


# Generalization Ability

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?



# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

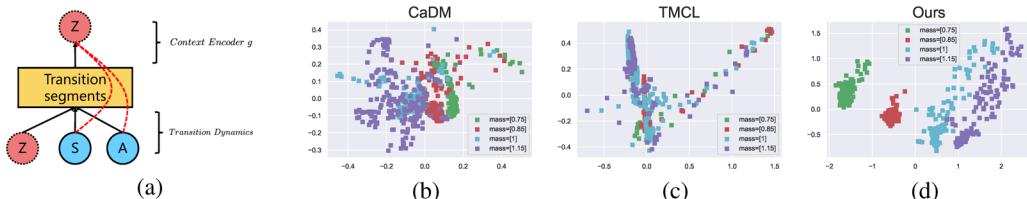


Figure 1: (a) The illustration of why historical states and actions are encoded in environment-specified factor  $Z$ , (b)(c)(d) The PCA visualization of estimated context (environmental-specific) vectors in **Pendulum** task, where the dots with different colors denote that the context vector (after PCA) estimated from different environments. More visualization results are given at Appendix A.13.

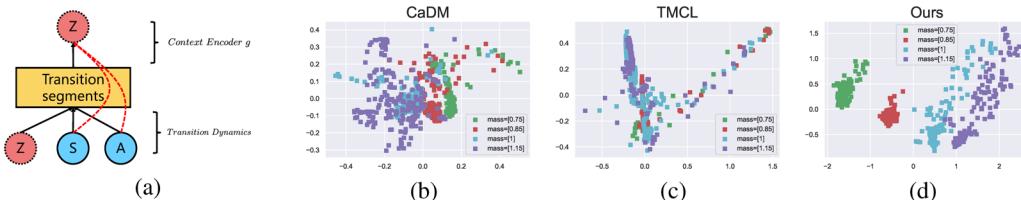
For more information, check  
“A Relational Intervention Approach  
for Unsupervised Dynamics  
Generalization in Model-Based  
Reinforcement Learning”, ICLR 2022

# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?



For more information, check  
“A Relational Intervention Approach  
for Unsupervised Dynamics  
Generalization in Model-Based  
Reinforcement Learning”, ICLR 2022

How can the extraction of **contextual variable  $Z$**  from past transition segments be improved to ensure that  $Z$  maintains its crucial property of being **similar in the same environment** and **dissimilar in different ones** for effective model generalization?

# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?

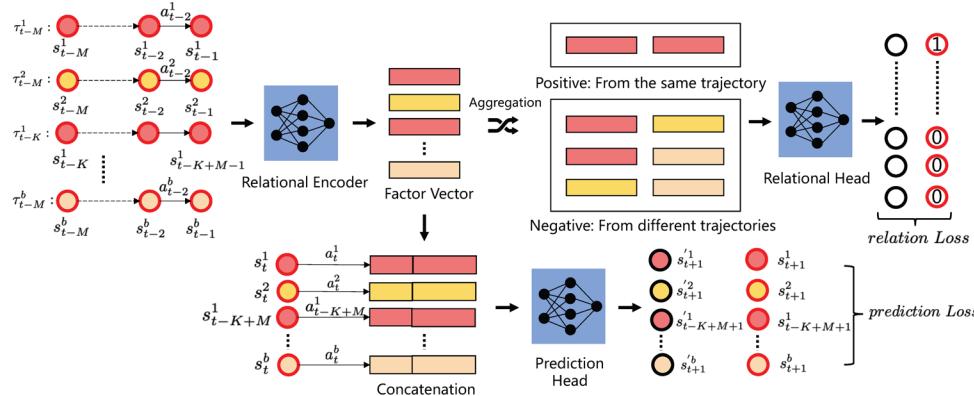


Figure 2: An overview of our Relational Intervention approach, where Relational Encoder, Prediction Head and Relational Head are three learnable functions, and the circles denote states (Ground-Truths are with red boundary, and estimated states are with black boundary), and the rectangles denote the estimated vectors. Specifically, *prediction Loss* enables the estimated environmental-specified factor can help the Prediction head to predict the next states, and the *relation Loss* aims to enforce the similarity between factors estimated from the same trajectory or similar trajectories.

For more information, check  
 "A Relational Intervention Approach  
 for Unsupervised Dynamics  
 Generalization in Model-Based  
 Reinforcement Learning", ICLR 2022

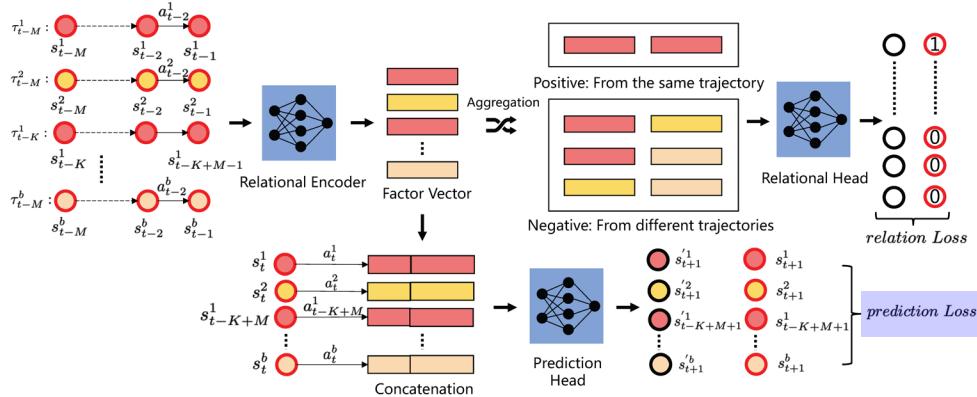
How can the extraction of **contextual variable Z** from past transition segments be improved to ensure that **Z** maintains its crucial property of being **similar in the same environment** and **dissimilar in different ones** for effective model generalization?

# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?



$$\mathcal{L}_{\theta, \phi}^{pred} = -\frac{1}{N} \sum_{i=1}^N \log \hat{f}(s_{t+1}^i | s_t^i, a_t^i, g(\tau_{t-k:t-1}^i; \phi); \theta)$$

For more information, check  
 "A Relational Intervention Approach  
 for Unsupervised Dynamics  
 Generalization in Model-Based  
 Reinforcement Learning", ICLR 2022

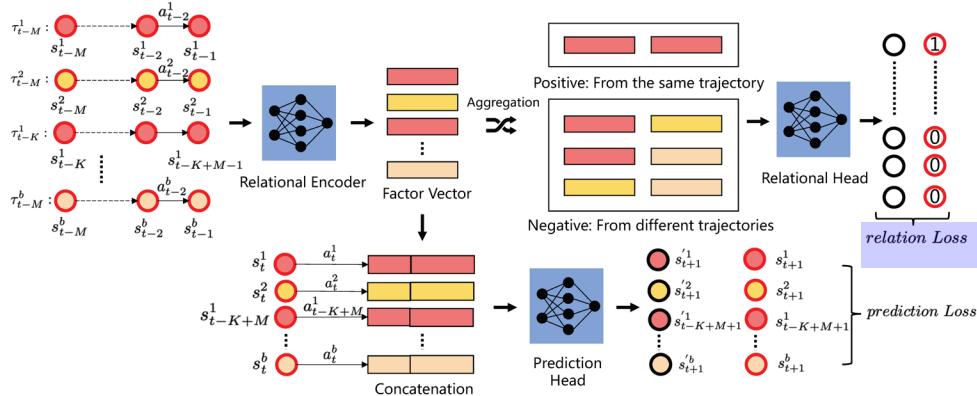
How can the extraction of **contextual variable Z** from past transition segments be improved to ensure that **Z** maintains its crucial property of being **similar in the same environment** and **dissimilar in different ones** for effective model generalization?

# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?



$$\mathcal{L}_{\varphi, \phi}^{relation} = -\frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j=1}^N \left[ y^{i,j} \cdot \log h([\hat{z}^i, \hat{z}^j]; \varphi) + (1-y^{i,j}) \cdot \log (1-h([\hat{z}^i, \hat{z}^j]; \varphi)) \right]$$

For more information, check  
 "A Relational Intervention Approach  
 for Unsupervised Dynamics  
 Generalization in Model-Based  
 Reinforcement Learning", ICLR 2022

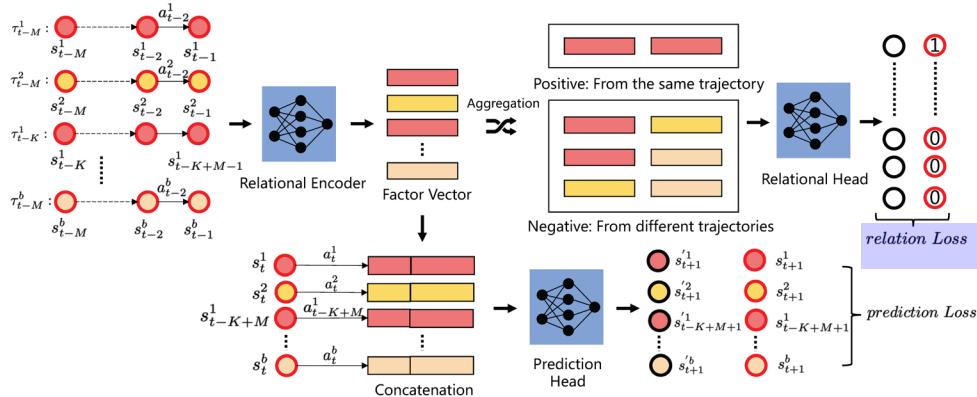
How can the extraction of **contextual variable Z** from past transition segments be improved to ensure that **Z** maintains its crucial property of being **similar in the same environment** and **dissimilar in different ones** for effective model generalization?

# Generalization

What does generalization mean for agents?

RQ1

How can agents achieve **reliable generalization** despite unknown variations?



$$\mathcal{L}_{\varphi, \phi}^{relation} = -\frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j=1}^N \left[ y^{i,j} \cdot \log h([\hat{z}^i, \hat{z}^j]; \varphi) + (1-y^{i,j}) \cdot \log (1-h([\hat{z}^i, \hat{z}^j]; \varphi)) \right]$$

For more information, check  
 "A Relational Intervention Approach  
 for Unsupervised Dynamics  
 Generalization in Model-Based  
 Reinforcement Learning", ICLR 2022

How can the extraction of **contextual variable Z** from past transition segments be improved to ensure that **Z** maintains its crucial property of being **similar in the same environment** and **dissimilar in different ones** for effective model generalization?

Estimating **trajectory invariant** information is insufficient because **the estimated  $\hat{Z}$ s in the same environment will also be pushed away**, which may undermine the cluster compactness for estimated  $\hat{Z}$ s.

# Generalization

The causal effects induced by the same context are similar.

For more information, check  
 "A Relational Intervention Approach for Unsupervised Dynamics Generalization in Model-Based Reinforcement Learning", ICLR 2022

What does generalization mean for agents?

RQ1

How can agents achieve reliable generalization despite unknown variations?

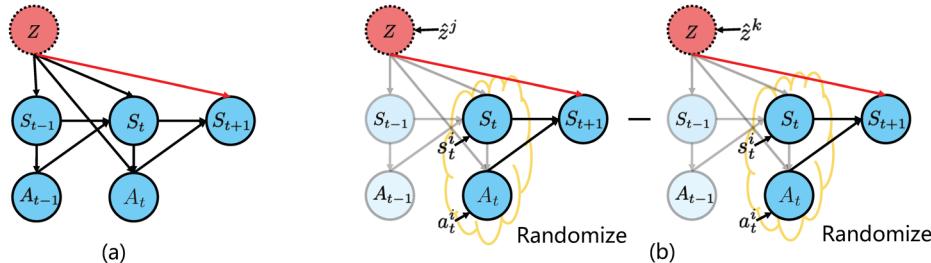


Figure 3: (a) The illustration of causal graph, and the red line denotes the direct causal effect from  $Z$  to  $S_{t+1}$ . (b) The illustration of estimating the controlled causal effect.

$$\begin{aligned} \mathcal{L}_{\varphi, \phi}^{i\text{-relation}} = & -\frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j=1}^N \left[ [y^{i,j} + (1-y^{i,j}) \cdot w^{i,j}] \cdot \log h([\hat{z}^i, \hat{z}^j]; \varphi) \right. \\ & \left. + (1-y^{i,j}) \cdot (1-w^{i,j}) \cdot \log (1-h([\hat{z}^i, \hat{z}^j]; \varphi)) \right], \end{aligned}$$

similarity

$$ACDE_{\hat{z}^j, \hat{z}^k} = \frac{1}{N} \sum_{i=1}^N |CDE_{\hat{z}^j, \hat{z}^k}(s_t^i, a_t^i)|$$

The estimated  $\hat{Z}$ s in the same environment should have similar causal effect on  $S_{t+1}$ .

# Generalization

The causal effects induced by the same context are similar.

What does generalization mean for agents?

RQ1

How can agents achieve reliable generalization despite unknown variations?

Table 3: The prediction errors of methods on test environments

	CaDM (Lee et al., 2020)	TMCL (Seo et al., 2020)	Ours
Hopper	$0.0551 \pm 0.0236$	$0.0316 \pm 0.0138$	<b><math>0.0271 \pm 0.0011</math></b>
Ant	$0.3850 \pm 0.0256$	$0.1560 \pm 0.0106$	<b><math>0.1381 \pm 0.0047</math></b>
C_Halfcheetah	$0.0815 \pm 0.0029$	$0.0751 \pm 0.0123$	<b><math>0.0525 \pm 0.0061</math></b>
HalfCheetah	$0.6151 \pm 0.0251$	$1.0136 \pm 0.6241$	<b><math>0.4513 \pm 0.2147</math></b>
Pendulum	$0.0160 \pm 0.0036$	$0.0130 \pm 0.0835$	<b><math>0.0030 \pm 0.0012</math></b>
Slim_Humanoid	$0.8842 \pm 0.2388$	$0.3243 \pm 0.0027$	<b><math>0.3032 \pm 0.0046</math></b>

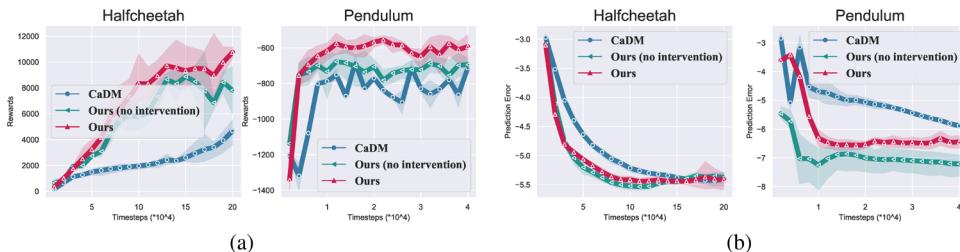


Figure 6: (a) The average rewards of trained model-based RL agents on unseen environments. The results show the mean and standard deviation of returns averaged over three runs. (b) The average prediction errors over the training procedure.

For more information, check  
“A Relational Intervention Approach  
for Unsupervised Dynamics  
Generalization in Model-Based  
Reinforcement Learning”, ICLR 2022

How can the extraction of contextual variable  $Z$  from past transition segments be improved to ensure that  $Z$  maintains its crucial property of being similar in the same environment and dissimilar in different ones for effective model generalization?

# Causal RL

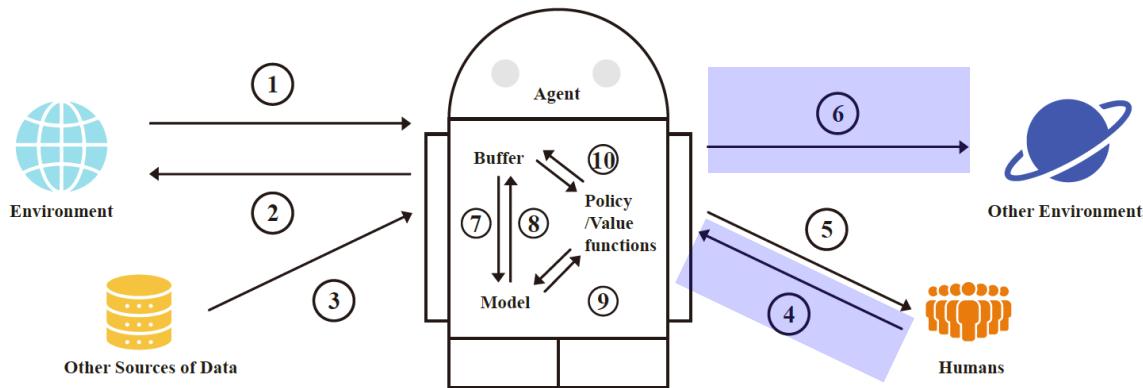


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans; 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Knowledge Transfer

RQ2

What knowledge can be transferred?

How can algorithms be designed to facilitate efficient adaptation?

For more information, check  
“AdaRL: What, Where, And  
How to Adapt in Transfer  
Reinforcement Learning”,  
ICLR 2022

# Knowledge Transfer

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?

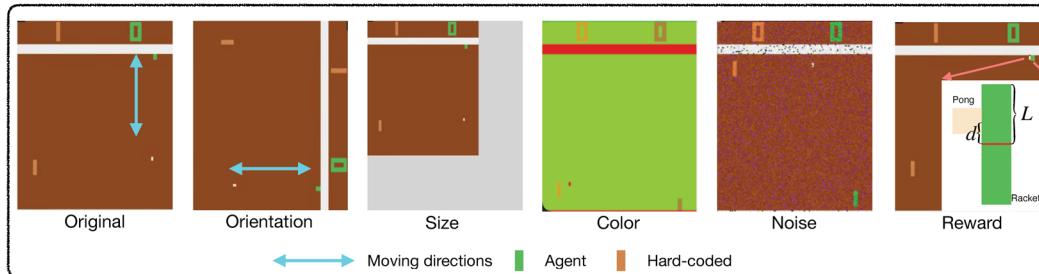


Figure A8: Visual example of the original Pong game and the various change factors. The light blue arrows are added to show the direction in which the agent can move.

For more information, check  
“AdaRL: What, Where, And  
How to Adapt in Transfer  
Reinforcement Learning”,  
ICLR 2022

# Knowledge Transfer

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?

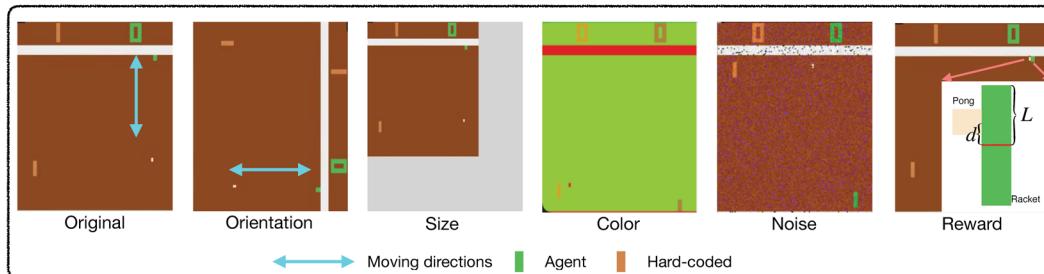


Figure A8: Visual example of the original Pong game and the various change factors. The light blue arrows are added to show the direction in which the agent can move.

How to adapt reliably and efficiently to changes across domains **with a few samples from the target domain**, even in partially observable environments?

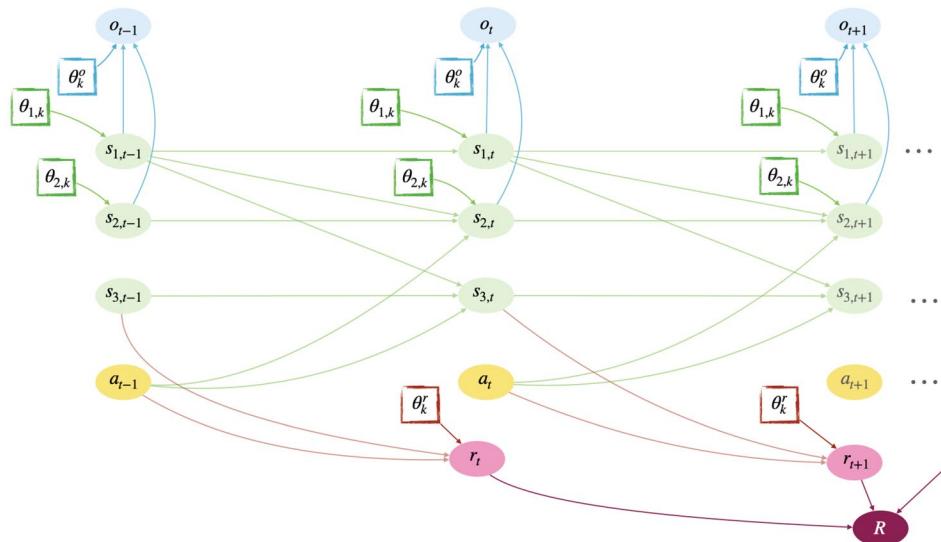
# Knowledge Transfer

Sparse Mechanisms Shift

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?



How to adapt reliably and efficiently to changes across domains with a few samples from the target domain, even in partially observable environments?

# Knowledge Transfer

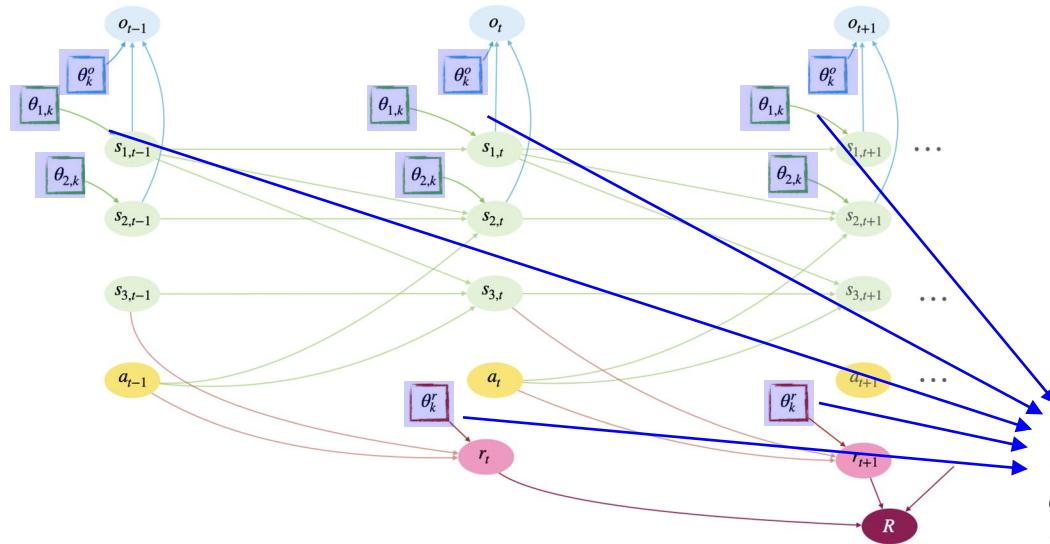
Sparse Mechanisms Shift

For more information, check  
“AdaRL: What, Where, And  
How to Adapt in Transfer  
Reinforcement Learning”,  
ICLR 2022

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?



How to adapt reliably and efficiently to changes across domains with a few samples from the target domain, even in partially observable environments?

Introduce a low-dimensional vector  $\theta_k$  to characterize the domain-specific information in a compact way.

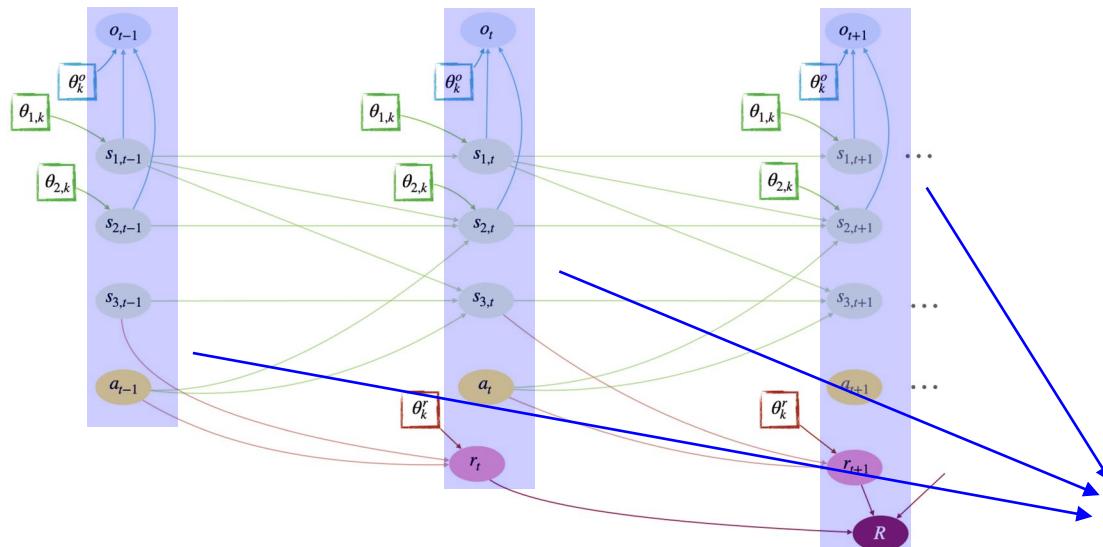
# Knowledge Transfer

Sparse Mechanisms Shift

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?



How to adapt reliably and efficiently to changes across domains **with a few samples from the target domain**, even in partially observable environments?

Shared across different domains.

# Knowledge Transfer

Sparse Mechanisms Shift

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?

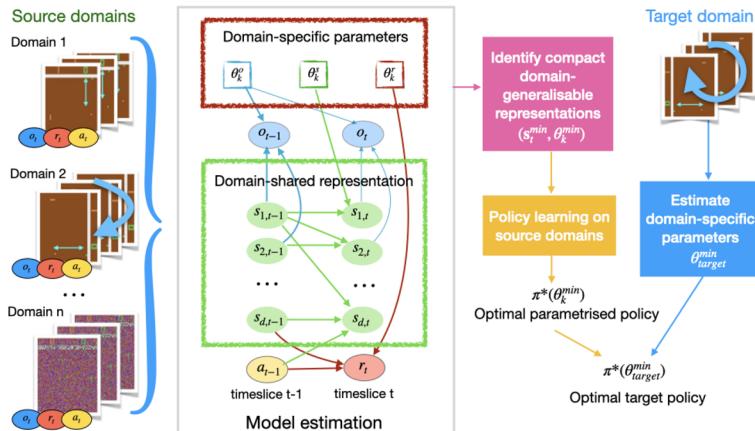


Figure 1: The overall AdaRL framework. We learn a Dynamic Bayesian Network (DBN) over the observations, latent states, reward, actions and domain-specific change factors that is shared across the domains. We then characterize a minimal set of representations that suffice for policy transfer, so that we can quickly adapt the optimal source policy with only a few samples from the target domain.

How to adapt reliably and efficiently to changes across domains with a few samples from the target domain, even in partially observable environments?

All we need to update in the target domain is the low-dimensional  $\theta_k$ .

# Knowledge Transfer

Sparse Mechanisms Shift

What knowledge can be transferred?

RQ2

How can algorithms be designed to facilitate efficient adaptation?

	Oracle Upper bound	Non-t lower bound	PNN (Rusu et al., 2016)	PSM (Agarwal et al., 2021)	MTQ (Fakoor et al., 2020)	AdaRL* Ours w/o masks	AdaRL Ours
O_in	18.65 ( $\pm 2.43$ )	6.18 • ( $\pm 2.43$ )	9.70 • ( $\pm 2.09$ )	11.61 • ( $\pm 3.85$ )	15.79 • ( $\pm 3.26$ )	14.27 • ( $\pm 1.93$ )	<b>18.97</b> ( $\pm 2.00$ )
O_out	19.86 ( $\pm 1.09$ )	6.40 • ( $\pm 3.17$ )	9.54 • ( $\pm 2.78$ )	10.82 • ( $\pm 3.29$ )	10.82 • ( $\pm 4.13$ )	12.67 • ( $\pm 2.49$ )	<b>15.75</b> ( $\pm 3.80$ )
C_in	19.35 ( $\pm 0.45$ )	8.53 • ( $\pm 2.08$ )	14.44 • ( $\pm 2.37$ )	19.02 ( $\pm 1.17$ )	16.97 • ( $\pm 2.02$ )	18.52 • ( $\pm 1.41$ )	<b>19.14</b> ( $\pm 1.05$ )
C_out	19.78 ( $\pm 0.25$ )	8.26 • ( $\pm 3.45$ )	14.84 • ( $\pm 1.98$ )	17.66 • ( $\pm 2.46$ )	15.45 • ( $\pm 3.30$ )	17.92 ( $\pm 1.83$ )	<b>19.03</b> ( $\pm 0.97$ )
S_in	18.32 ( $\pm 1.18$ )	6.91 • ( $\pm 2.02$ )	11.80 • ( $\pm 3.25$ )	12.65 • ( $\pm 3.72$ )	13.68 • ( $\pm 3.49$ )	14.23 • ( $\pm 3.19$ )	<b>16.65</b> ( $\pm 1.72$ )
S_out	19.01 ( $\pm 1.04$ )	6.60 • ( $\pm 3.11$ )	9.07 • ( $\pm 4.58$ )	8.45 • ( $\pm 4.51$ )	11.45 • ( $\pm 2.46$ )	12.80 • ( $\pm 2.62$ )	<b>17.82</b> ( $\pm 2.35$ )
N_in	18.48 ( $\pm 1.25$ )	5.51 • ( $\pm 3.88$ )	12.73 • ( $\pm 3.67$ )	11.30 • ( $\pm 2.58$ )	12.67 • ( $\pm 3.84$ )	13.78 • ( $\pm 2.15$ )	<b>16.84</b> ( $\pm 3.13$ )
N_out	18.26 ( $\pm 1.11$ )	6.02 • ( $\pm 3.19$ )	13.24 • ( $\pm 2.55$ )	11.26 • ( $\pm 3.15$ )	15.77 • ( $\pm 2.12$ )	14.65 • ( $\pm 3.01$ )	<b>18.30</b> ( $\pm 2.24$ )

Table 3: Average final scores on modified Pong (POMDP) with  $N_{target} = 50$ . The best non-oracle are marked in red. O, C, S, and N denote the orientation, color, size, and noise factors, respectively.

How to adapt reliably and efficiently to changes across domains with a few samples from the target domain, even in partially observable environments?

# Causal RL

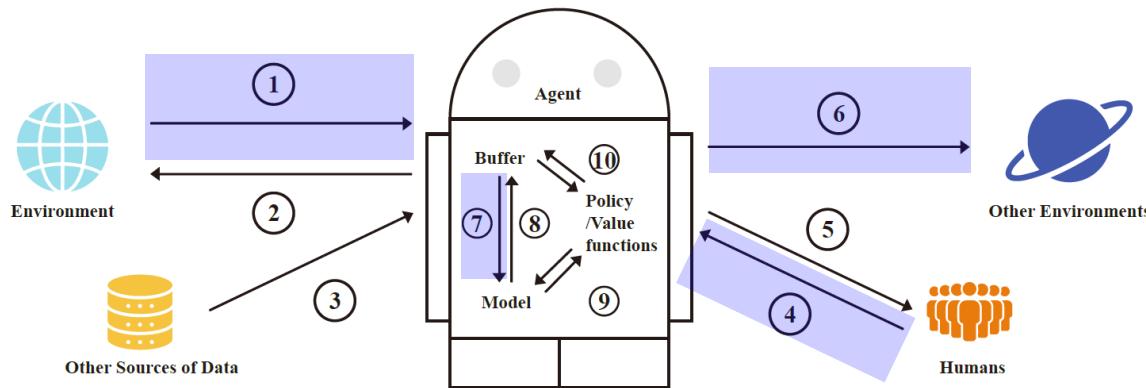


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans. 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Tutorial Outline

## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Spurious Correlation

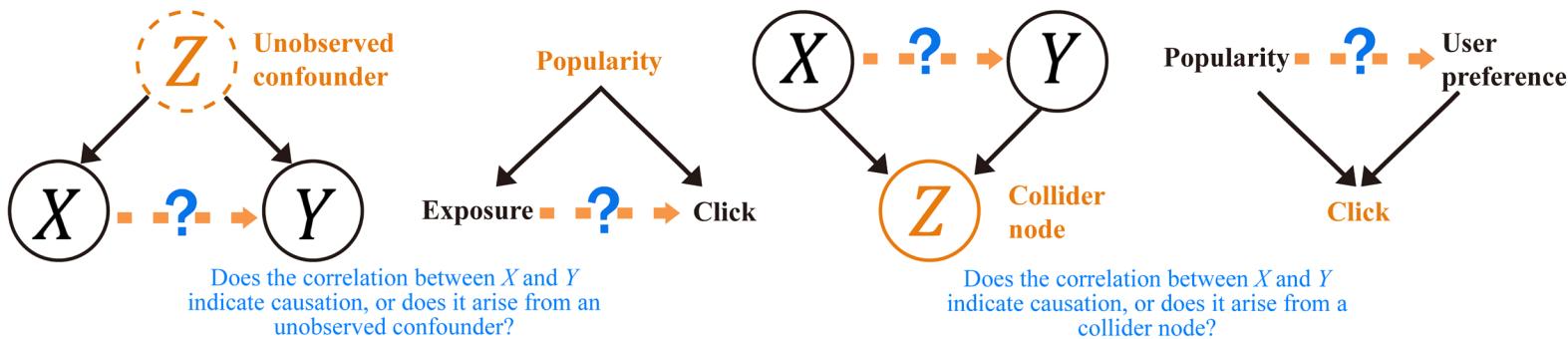


Figure 8: Causal graph illustrating the two types of spurious correlations, with examples from real-world applications.

Correlation does not imply causation.

# Spurious Correlation



People are biased.

Does the correlation between  $X$  and  $Y$  indicate causation, or does it arise from an unobserved confounder?

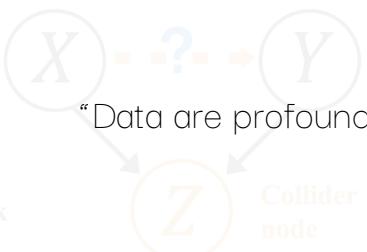
Data is biased, in part because people are biased.

Causal graph illustrating the two types of spurious correlations, with examples from real-world applications.

Yann Lecun



"Data are profoundly dumb."



Judea Pearl

Does the correlation between  $X$  and  $Y$  indicate causation, or does it arise from a collider node?

# Spurious Correlation



Yann Lecun

RQ

What types of spurious correlation exist in reinforcement learning?  
How to eliminate or mitigate spurious correlations?

# Spurious Correlation

RQ

What types of spurious correlation exist in reinforcement learning?

# Spurious Correlation

Causal Graph

RQ What types of spurious correlation exist in reinforcement learning?

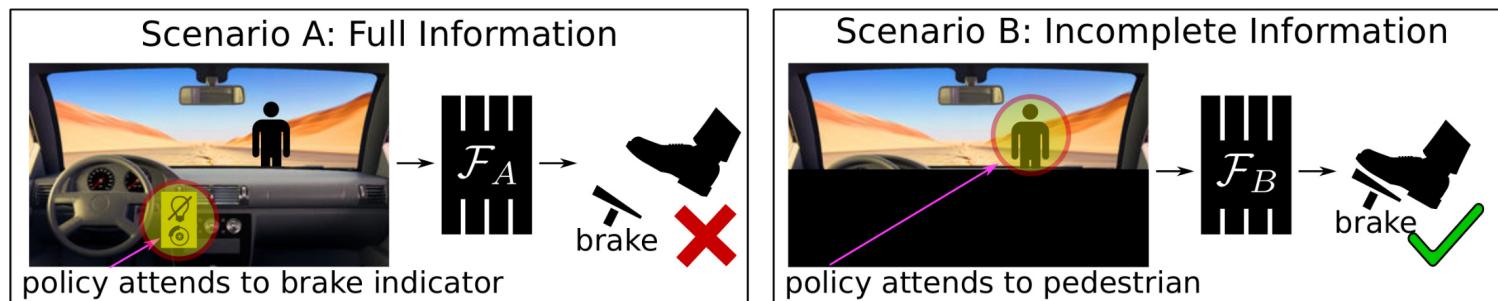


Figure 1: Causal misidentification: *more* information yields worse imitation learning performance. Model A relies on the braking indicator to decide whether to brake. Model B instead correctly attends to the pedestrian.

# Spurious Correlation

Causal Graph

For more information, check  
“Spurious Correlation  
Reduction for Offline  
Reinforcement Learning”,  
2021

RQ What types of spurious correlation exist in reinforcement learning?

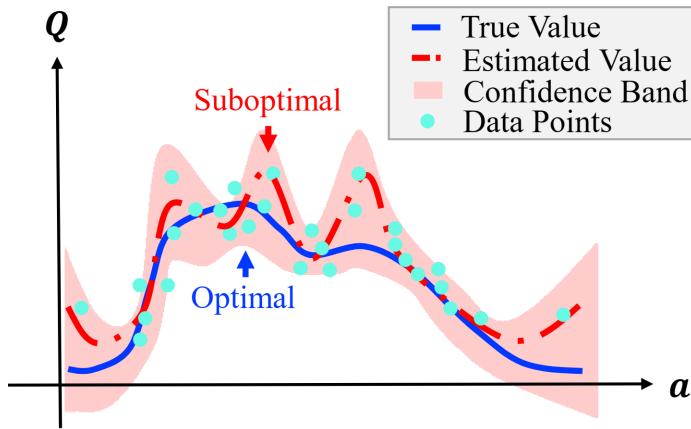
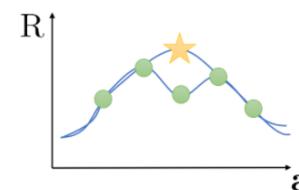
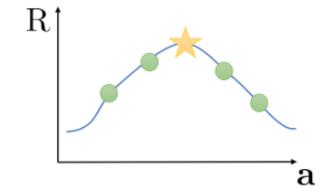


Figure 1. An example of false correlation: the epistemic uncertainty is correlated with the value, making a suboptimal action with high uncertainty appear to be better than the optimal one.

online RL setting



offline RL setting



Credit: NeurIPS 2020 Tutorial - Offline Reinforcement Learning:  
From Algorithms to Practical Challenges

# Causal RL

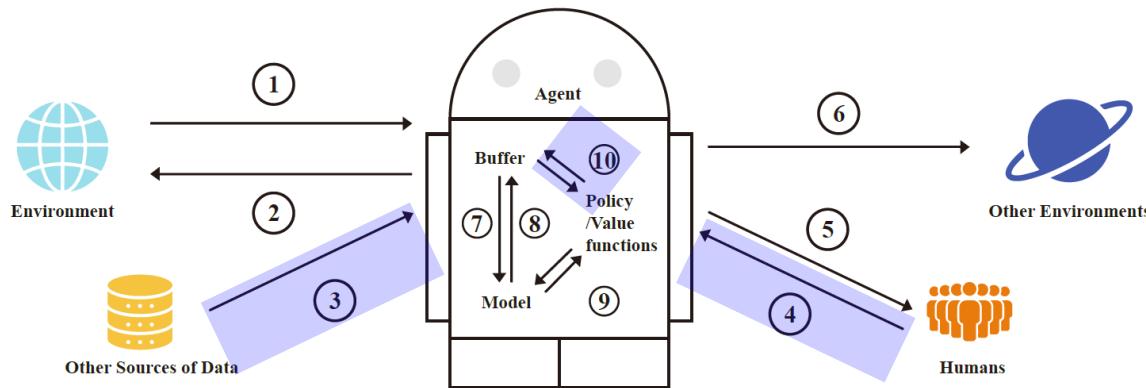


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans; 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Tutorial Outline

## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Beyond Return

Article • Machine Learning

## Report calls for transparency in AI automation

Centre for Data Ethics and Innovation calls for more transparency in algorithms for machine learning in the public sector...

# Beyond Return

Article • Machine Learning

Article • AI Strategy

## UK launch national standards for algorithmic transparency

The UK Government has announced one of the world's first national standards for algorithm transparency

# Beyond Return

Article • AI Strategy

## Need for responsible AI in some of the world's largest banks

Research shows one-third of North America and Europe's largest banks lack transparency and are not publicly reporting on their AI development

for algorithm transparency

# Beyond Return

The image shows a screenshot of a New York Times article. At the top left, there's a vertical sidebar with the text "Article • A" at the top, followed by "Need the", "Research", "transpa", and "for algo". The main title "Is an Algorithm Less Racist Than a Loan Officer?" is prominently displayed in large, bold, black font. Below the title is a subtitle: "Digital mortgage platforms have the potential to reduce discrimination. But automated systems provide rich opportunities to perpetuate bias, too." The New York Times logo is visible at the top center of the article area.

Article • A

Need the

Research

transpa

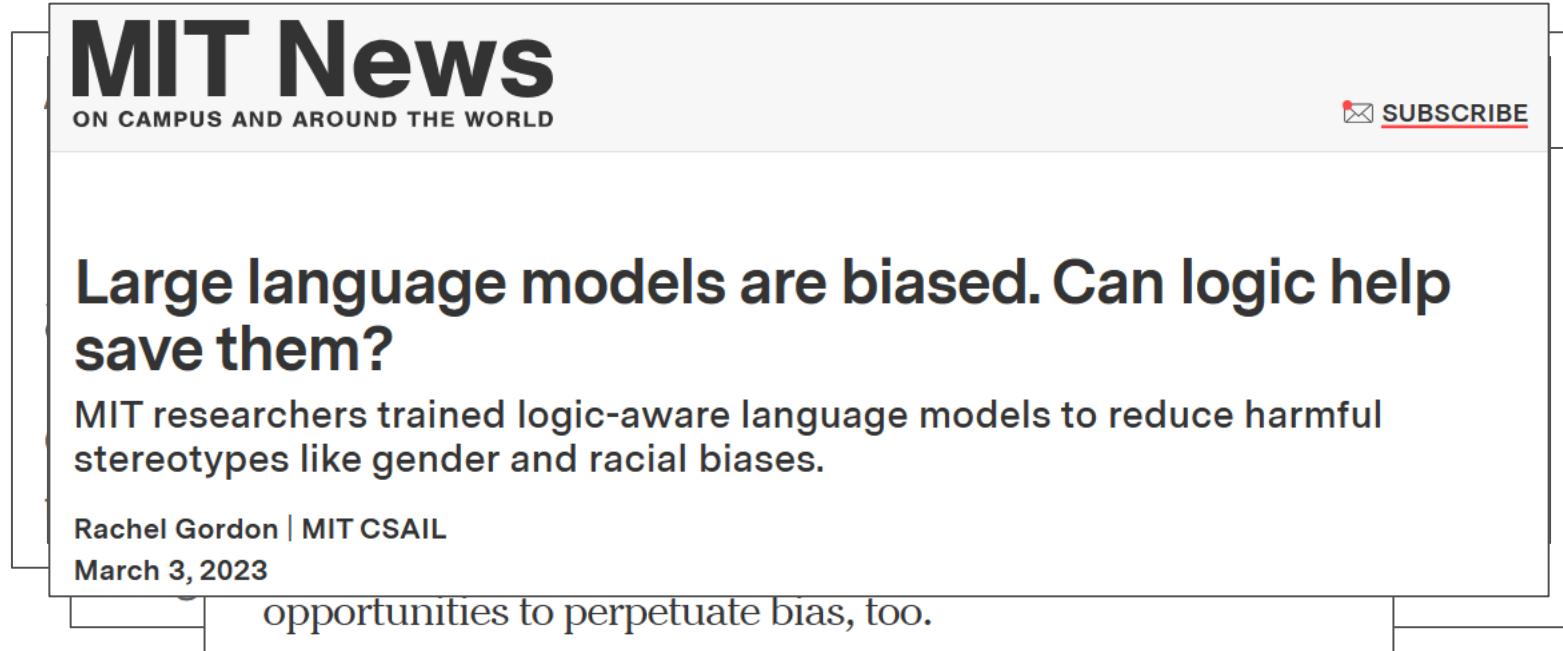
for algo

**The New York Times**

# Is an Algorithm Less Racist Than a Loan Officer?

Digital mortgage platforms have the potential to reduce discrimination. But automated systems provide rich opportunities to perpetuate bias, too.

# Beyond Return



**MIT News**  
ON CAMPUS AND AROUND THE WORLD

[!\[\]\(34d41167ceb06c65352cd3f761cde0f6\_img.jpg\) SUBSCRIBE](#)

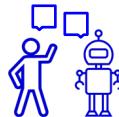
## Large language models are biased. Can logic help save them?

MIT researchers trained logic-aware language models to reduce harmful stereotypes like gender and racial biases.

Rachel Gordon | MIT CSAIL  
March 3, 2023

opportunities to perpetuate bias, too.

# Beyond Return



**Explainability** - the ability to understand and interpret the decisions of an agent.



**Fairness** - agents should strive to genuinely benefit humans and promote social good, avoiding any form of discrimination or harm towards specific individuals or groups.



**Safety** - agents should not prioritize higher returns over safety.

# Tutorial Outline

## Part 1

- Introduction
- Causality
- Reinforcement Learning
- Causal Reinforcement Learning

## Part 2

- Sample Efficiency
- Generalization
- Spurious Correlation
- Beyond Return

# Takeaway

## 1. Enhancing Sample Efficiency through Causal Reinforcement Learning

- 👉 Explore the areas that agents can causally influence the environment.
- 👉 Extract the causal factors to simplify the learning problem.
- 👉 Generate counterfactual rollouts for data augmentation.

## 2. Advancing Generalization Ability and Knowledge Transfer through Causal Reinforcement Learning

- 👉 Different interventions induce different MDPs and causal model can capture such variations.
- 👉 Transfer the domain-invariant information and only adapt the changed causal mechanisms.

## 3. Addressing Spurious Correlations through Causal Reinforcement Learning

- 👉 Identify the factors that exhibit spurious correlations and address them with causal reasoning.

## 4. Consideration Beyond Return

- 👉 Explanability, Fairness, Safety, ...

# Open Problems

## 1. Causal Learning in Reinforcement Learning



Causal representation learning, causal discovery, causal dynamics learning, ...

## 2. Causality-aware Multitask, Meta, and Lifelong Reinforcement Learning



Organize knowledge using causal structures.

## 3. Human-in-the-loop Learning and Reinforcement Learning from Human Feedback

## 4. Theoretical Advances in Causal Reinforcement Learning



Identifiability, convergence, suboptimality, ...

## 5. Benchmarking Causal Reinforcement Learning



How to design a reliable benchmark for causal RL methods? What metrics should be considered?

## 6. Real-world Causal Reinforcement Learning



Applications, e.g., robotics, self-driving, healthcare, finance, ...

# Causal RL Survey

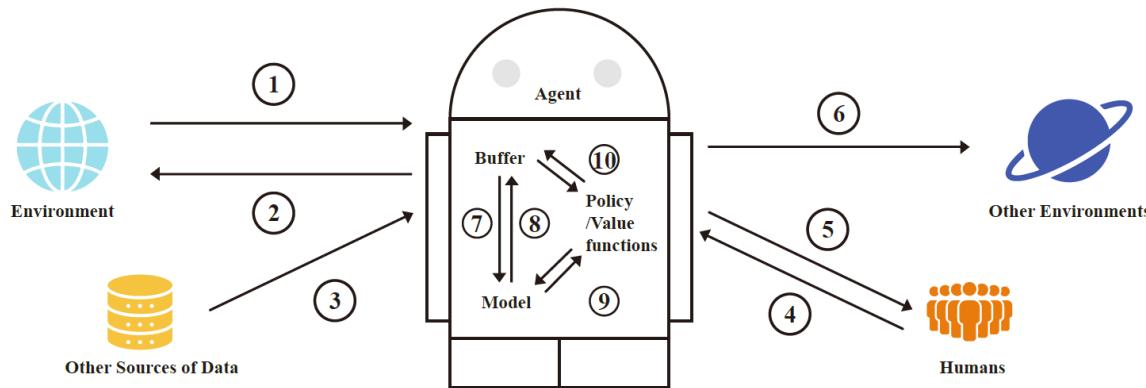


Figure 10: A schematic diagram illustrating the integration of causality into the reinforcement learning process. The numbered edges represent some key components: 1) Abstraction and extraction of causal representations from raw observations; 2) Directed exploration guided by causal knowledge; 3) Fusing (possibly confounded) data; 4) Incorporating causal assumptions or knowledge from humans. 5) Providing causality-based explanations; 6) Generalization and knowledge transfer; 7) Learning causal world models; 8) Counterfactual data generation; 9) Planning with world models; 10) Enhanced training of policies and value functions with causal reasoning.

# Causal Reinforcement Learning: A Survey

Zhihong Deng

Australian Artificial Intelligence Institute, University of Technology Sydney

[zhi-hong.deng@student.uts.edu.au](mailto:zhi-hong.deng@student.uts.edu.au)

Jing Jiang

Australian Artificial Intelligence Institute, University of Technology Sydney

[jing.jiang@uts.edu.au](mailto:jing.jiang@uts.edu.au)

Guodong Long

Australian Artificial Intelligence Institute, University of Technology Sydney

[guodong.long@uts.edu.au](mailto:guodong.long@uts.edu.au)

Chengqi Zhang

Australian Artificial Intelligence Institute, University of Technology Sydney

[chengqi.zhang@uts.edu.au](mailto:chengqi.zhang@uts.edu.au)

## Abstract

Reinforcement learning is an essential paradigm for solving sequential decision problems under uncertainty. Despite many remarkable achievements in recent decades, applying reinforcement learning methods in the real world remains challenging. One of the main obstacles is that reinforcement learning agents lack a fundamental understanding of the world and must therefore learn from scratch through numerous trial-and-error interactions. They may also face challenges in providing explanations for their decisions and generalizing the acquired knowledge. Causality, however, offers a notable advantage as it can formalize knowledge in a systematic manner and leverage invariance for effective knowledge transfer. This has led to the emergence of causal reinforcement learning, a subfield of reinforcement learning that seeks to enhance existing algorithms by incorporating causal relationships into the learning process. In this survey, we comprehensively review the literature on causal reinforcement learning. We first introduce the basic concepts of causality and reinforcement learning, and then explain how causality can address core challenges in non-causal reinforcement learning. We categorize and systematically review existing causal reinforcement learning approaches based on their target problems and methodologies. Finally, we outline open issues and future directions in this emerging field.

## 1 Introduction

*“All reasonings concerning matter of fact seem to be founded on the relation of cause and effect. By means of that relation alone we can go beyond the evidence of our memory and senses.”*

—David Hume, *An Enquiry Concerning Human Understanding*.

# Q & A

Thank you!



Zhihong Deng

Zhi-Hong.Deng@student.uts.edu.au