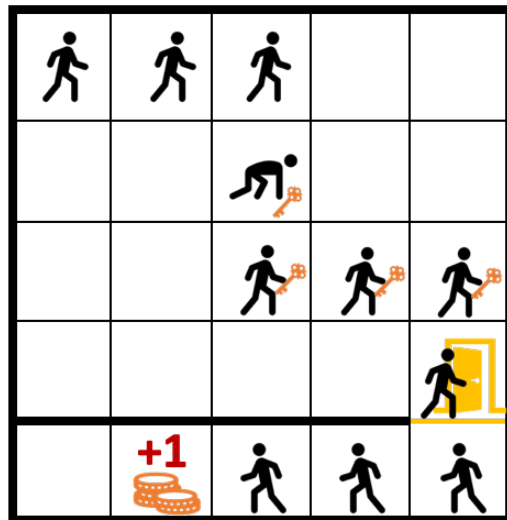# Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach

Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang,

Meng Fang, Mykola Pechenizkiy
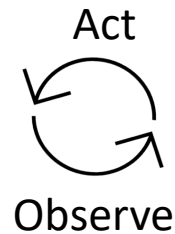
[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# A major challenge in RL: Delayed Reward

**Delayed Reward:**

Obtaining 🪙 : positive reward **+1**

Otherwise: reward **0**

Act

Observe

**+1**
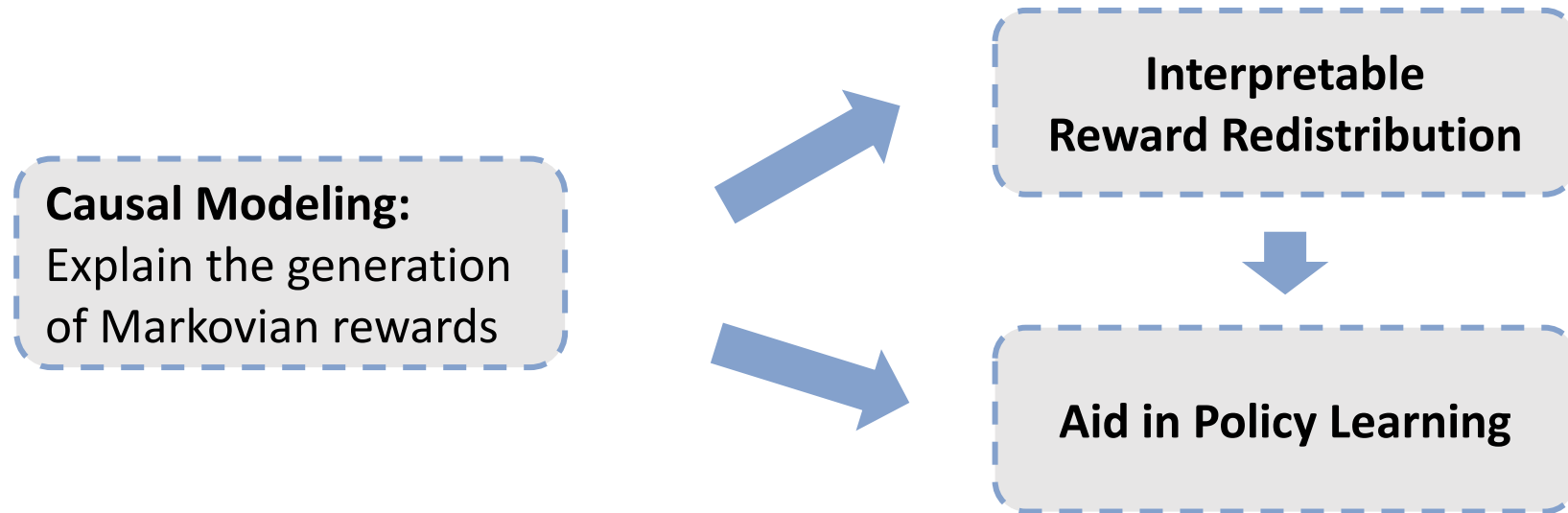
Lack of immediate feedback

Unstable policy optimization

**Reward Redistribution:**

Assign proxy rewards according to the contribution of each state-action pair

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Motivation

**Equally important for Interpretable Reward Redistribution:**
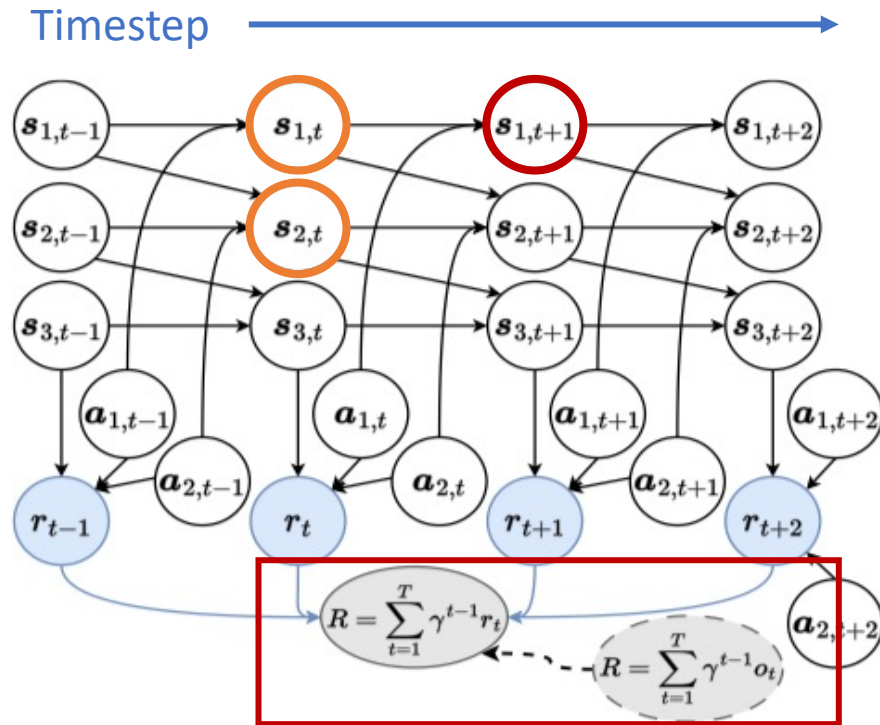
- Computing the contribution of each state-action pair towards delayed rewards?

- Explaining the reasons behind such contribution?

**Causal Modeling:** Explain the generation of Markovian rewards

**Interpretable Reward Redistribution**

**Aid in Policy Learning**

# Causal Reformulation of Reward Redistribution

Causality:

causes (which part of the state-action pair) → outcomes (Markovian reward)



Timestep

➢ **A graphical example** to model the generation of ,

Markovian rewards $r_t$,

long-term returns $R$,

by the causal structure over $s_t$, $a_t$, $r_t$ and $R$.

○ Observable

● Unobservable &
Goal of Reward Redistribution

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Causal Reformulation of Reward Redistribution

A generative process in MDP:

$$\begin{cases} s_{i,t+1} = f(\boldsymbol{C}^{\boldsymbol{s}\rightarrow\boldsymbol{s}}_{\cdot,i} \odot \boldsymbol{s}_t, \boldsymbol{C}^{\boldsymbol{a}\rightarrow\boldsymbol{s}}_{\cdot,i} \odot \boldsymbol{a}_t, \epsilon_{s,i,t}) \\ r_t = g(\boldsymbol{C}^{\boldsymbol{s}\rightarrow r} \odot \boldsymbol{s}_t, \boldsymbol{C}^{\boldsymbol{a}\rightarrow r} \odot \boldsymbol{a}_t, \epsilon_{r,t}) \\ R = \sum_{t=1}^{T} \gamma^{t-1} r_t \end{cases}$$

Dynamics function

Markovian reward function

Return Equivalence

Causal structure $\boldsymbol{C}^{\cdot\rightarrow\cdot}$:

$$\boldsymbol{C}^{s\rightarrow r} \in \{0,1\}^{|s|}, \boldsymbol{C}^{a\rightarrow r} \in \{0,1\}^{|a|},$$
$$\boldsymbol{C}^{s\rightarrow s} \in \{0,1\}^{|s|\times|s|}, \boldsymbol{C}^{a\rightarrow s} \in \{0,1\}^{|a|\times|s|};$$

$\epsilon_{s,i,t}$ and $\epsilon_{r,t}$: i.i.d random noises , $\odot$: element-wise product

**Identifiability Result:** Given observed state $\boldsymbol{s}_t$, action $\boldsymbol{a}_t$, long-term return $R$ , under the global Markov condition and faithfulness assumption, the causal structure $\boldsymbol{C}^{\cdot\rightarrow\cdot}$ , unknown functions, $f, g$ and the rewards $r_t$ are identifiable.

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Generative Return Decomposition

> **How do we estimate the generative model?**

Overall objective to optimize parameterized generative model:

$$L_{\mathrm{m}} = L_{\mathrm{rew}} + L_{\mathrm{cau}} + L_{\mathrm{reg}}$$

Minimize MSE for reward function:

$$L_{\mathrm{rew}}(\phi_{\mathrm{rew}}, \phi_{\mathrm{cau}}^{s\rightarrow r}, \phi_{\mathrm{cau}}^{a\rightarrow r}) = \mathbb{E}_{\tau\sim\mathcal{D}} \left\| R - \sum_{t=1}^{T} \gamma^{t-1}\hat{r}_t \right\|^2 = \mathbb{E}_{\tau\sim\mathcal{D}} \left\| \sum_{t=1}^{T} \gamma^{t-1}o_t - \sum_{t=1}^{T} \gamma^{t-1}\hat{r}_t \right\|^2$$
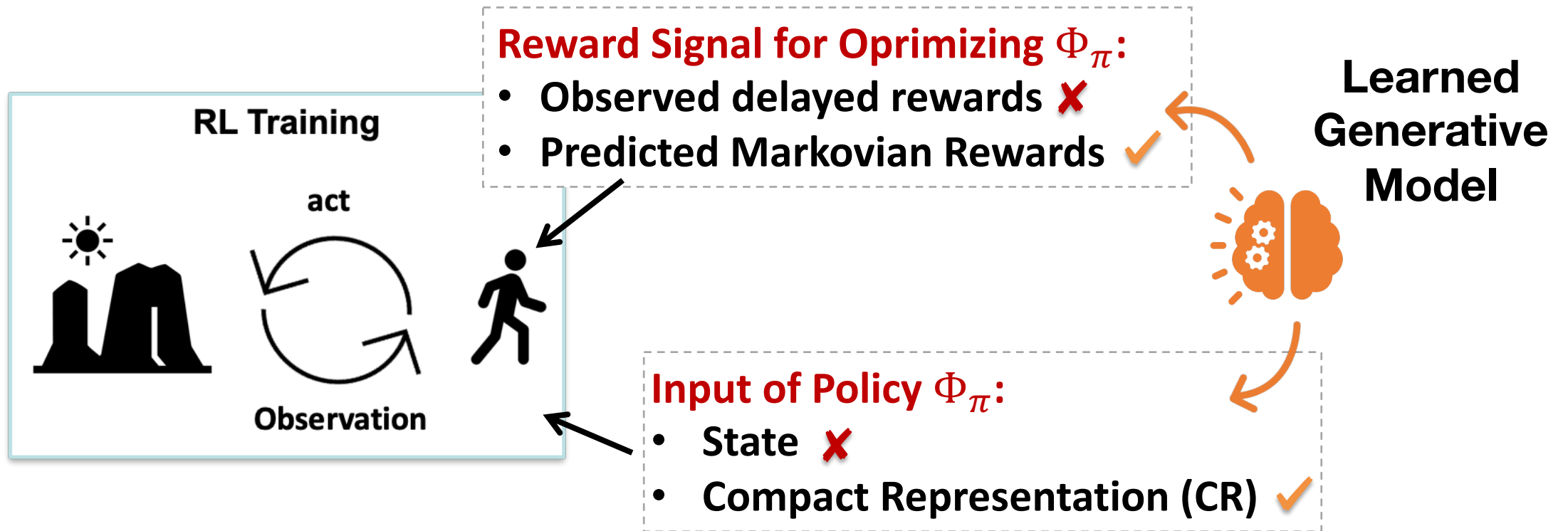
Maximize the likelihood for dynamic function:

$$L_{\mathrm{dyn}}\left(\phi_{\mathrm{dyn}}, \phi_{\mathrm{cau}}^{s\rightarrow s}, \phi_{\mathrm{cau}}^{a\rightarrow s}\right) = \mathbb{E}_{s_t,a_t,s_{t+1}\sim\mathcal{D}} \left[ -\sum_{i=1}^{|s|} \log P\left(s_{i,t+1}\middle| s_t, a_t, \boldsymbol{C}^{s\rightarrow s}, \boldsymbol{C}^{a\rightarrow s}\right) \right]$$

Regularizer:

$$L_{\mathrm{reg}}(\phi_{\mathrm{cau}}) = \lambda_1 \sum_i D_i(C^{s\rightarrow r}) + \lambda_2 \sum_i D_i(C^{a\rightarrow r}) + \lambda_3 \sum_{j\neq i} D_{i,j}(C^{s\rightarrow s})$$
$$+\lambda_4 \sum_{i=j} D_{i,j}(C^{s\rightarrow s}) + \lambda_5 \sum_{i=i,j} D_{i,j}(C^{a\rightarrow s}), \text{ where } D_i(C) = \log P(C_i = 1)$$

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Generative Return Decomposition

> **How do we exploit the estimated generative model?**



**Reward Signal for Oprimizing** $\Phi_\pi$**:**
- **Observed delayed rewards** ✗
- **Predicted Markovian Rewards** ✓

**Learned Generative Model**

**Input of Policy** $\Phi_\pi$**:**
- **State** ✗
- **Compact Representation (CR)** ✓

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Compact Representation (CR)

**> a minimally sufficient subset of all state components for policy learning**

**CR**: All the state components influence *rewards*.

$s_{i,t}$ is selected into **CR** if,

  1) $C_i^{s \to r} = 1$:

       $s_{i,t}$ directly impacts $r_t$ ($s_{2,t}$)

  2) $C_{i,j}^{s \to s} = C_j^{s \to r} = 1$:
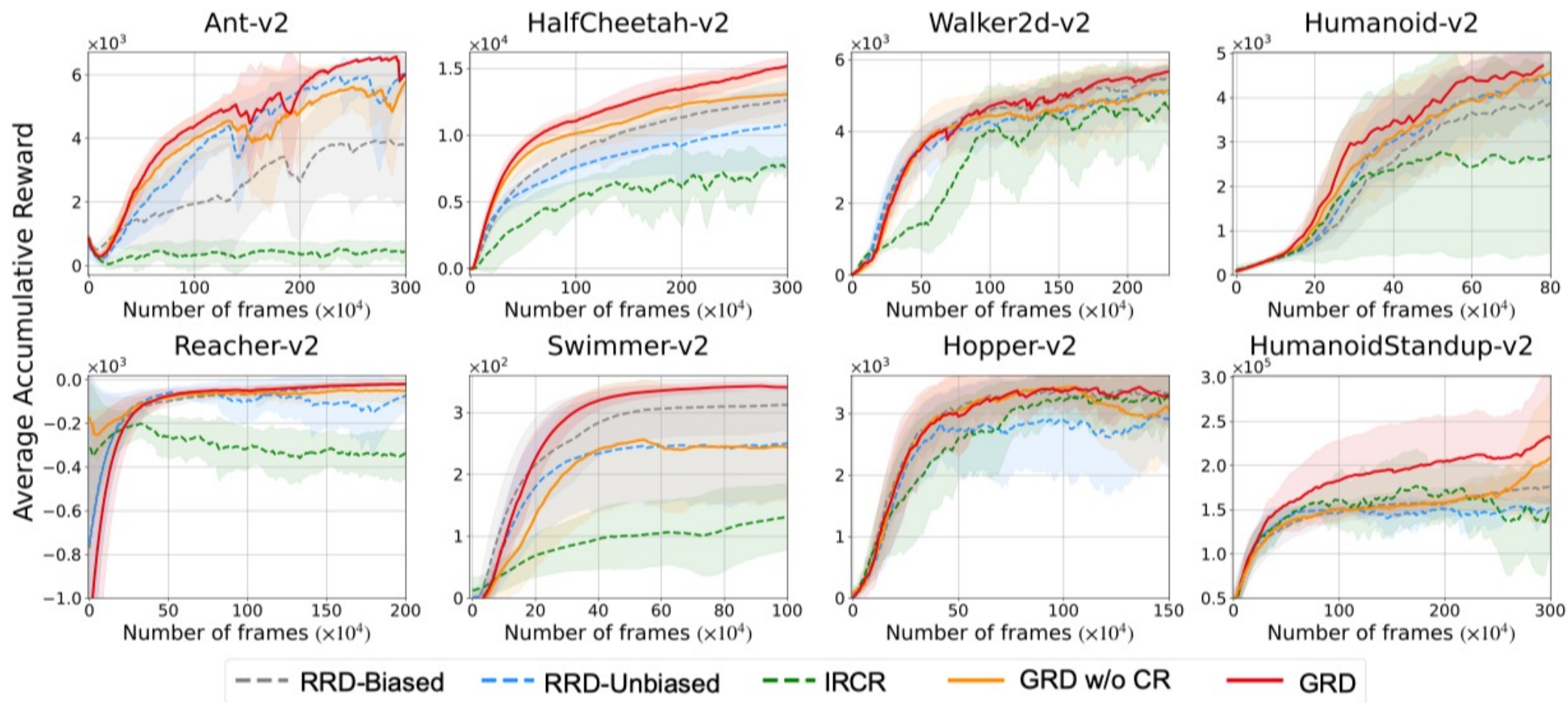
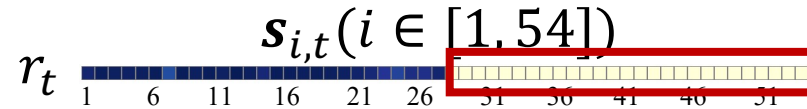       $s_{i,t}$ indirectly impacts $r_{t+1}$, ($s_{3,t}$)



[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Experimental Results – Episodic MuJoCo

At time step $t$, the agent is marked as $r_t$.

The observed sparse and delayed rewards are,

$$o_t = \begin{cases} 0, & \text{if } t \neq T \\ \gamma^{t-1} r_t, & \text{if } t = T \end{cases}$$

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.
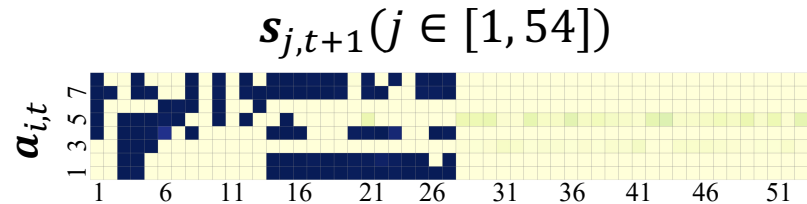
# Experimental Results – Episodic MuJoCo



[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

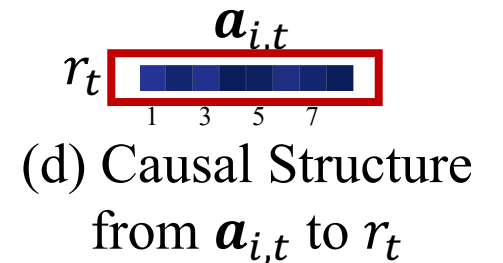# Visualization of Learned Causal Structure



(a) Causal Structure from $s_{i,t}$ to $s_{j,t+1}$

(b) Causal Structure from $a_{i,t}$ to $s_{j,t+1}$
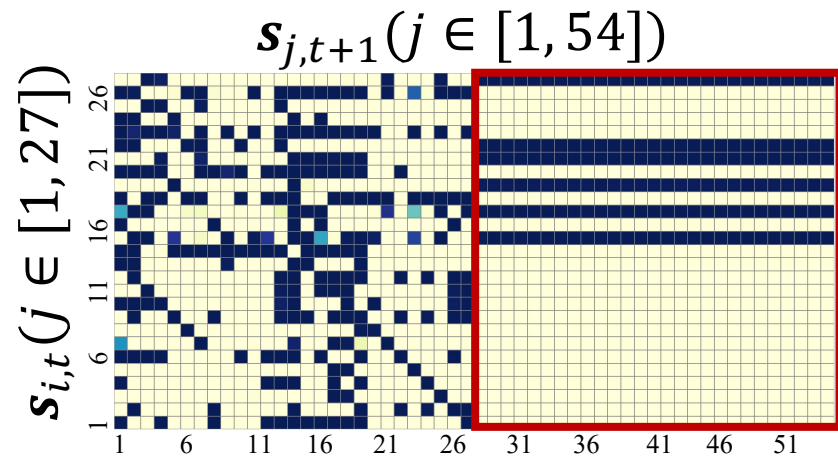
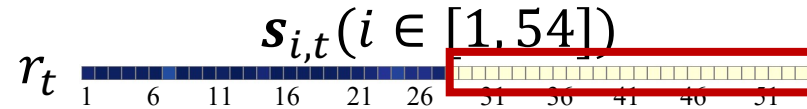(c) Causal Structure from $s_{i,t}$ to $r_t$

(d) Causal Structure from $a_{i,t}$ to $r_t$

**Two characters of Ant:**
- 28~111 dimensions in the state are not used.
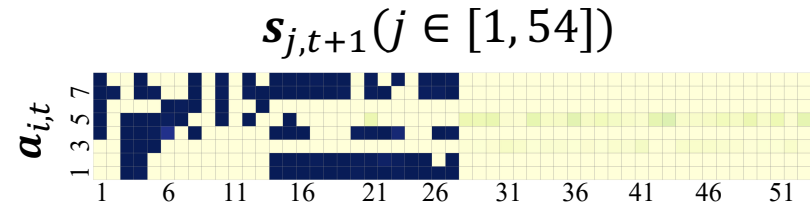- Low-cost control: all dimensions of action cause rewards.

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Visualization of Learned Causal Structure



$s_{j,t+1}(j \in [1,54])$

$s_{i,t}(j \in [1,27])$

(a) Causal Structure from $s_{i,t}$ to $s_{j,t+1}$

$s_{i,t}(i \in [1,54])$

$r_t$

(c) Causal Structure from $s_{i,t}$ to $r_t$

$s_{j,t+1}(j \in [1,54])$

$a_{i,t}$

(b) Causal Structure from $a_{i,t}$ to $s_{j,t+1}$

$a_{i,t}$

$r_t$

(d) Causal Structure from $a_{i,t}$ to $r_t$

- The edges related to no-used states do not exist (in yellow);

- Although learned redundant edges, they are not in CR;

- The edges from all dimensions of action to reward exist (in blue).

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.
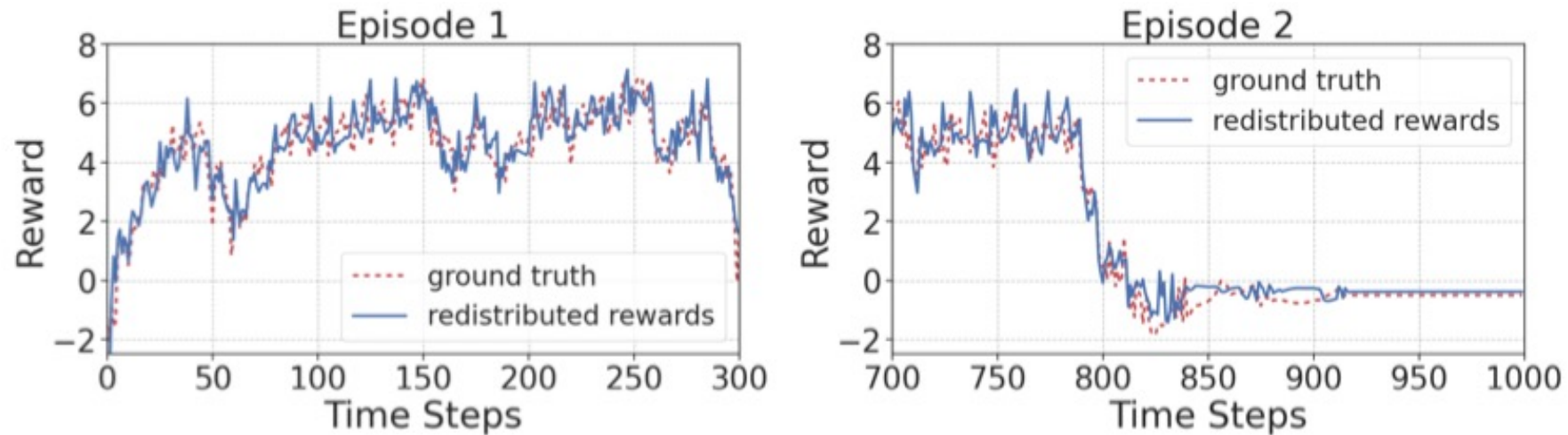
# Visualization of Redistributed Rewards



Figure 5: The visualization of redistributted rewards (blue solid lines) and the grounded rewards (red dotted lines).

[1] Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. "Interpretable Reward Redistribution in Reinforcement Learning: A Causal Approach." NeurIPS 2023.

# Thanks!

Project

Find Me

Paper : https://arxiv.org/abs/2305.18427
Project Page: https://reedzyd.github.io/GenerativeReturnDecomposition/