

515 A Broader Impacts

516 Our motivation is to strive to make decisions that are both understood and trusted by humans. By
517 increasing the credibility and transparency of the decision-making process, human users develop
518 a better understanding, validate and enhance the decisions made by algorithms. As a result, it is
519 possible to bridge the gap between humans and AI, fostering a symbiotic relationship that leverages
520 the strengths of both to enable more reliable and responsible decision-making in various filed. For
521 example, in finance, transparent algorithms can enhance risk assessment, investment strategies, and
522 fraud detection, while in transportation, trustworthy decision-making algorithms can contribute to
523 safe and efficient navigation and logistics. To advance this mission, our work on the GRD algorithm
524 represents a significant step forward. GRD explores facilitating policy learning in the presence of
525 delayed rewards. By decomposing the overall return into Markovian rewards, we provide a clearer
526 understanding of the contribution made by each state-action pair. Furthermore, we go beyond simply
527 explaining the rewards and delve into the causal view of reward generation. This approach allows
528 us to provide interpretable explanations of how Markovian rewards are generated, enabling a more
529 transparent decision-making process. This interpretability is vital for building trust with human
530 users who may need to understand and validate the decisions made by the algorithm. Moreover,
531 with an interpretable reward function, we can readily incorporate additional restrictions, such as
532 security constraints, into the decision-making process. This flexibility allows us to tailor the algorithm
533 to specific needs and requirements, further enhancing its trustworthiness. Additionally, since the
534 principles and techniques we are exploring can be applied across a wide range of domains and
535 industries, our collaboration with GRD not only contributes to the field of reinforcement learning but
536 also has broader implications beyond robotics.

537 B Proofs and Causal Background

538 B.1 Markov and faithfulness assumptions

539 A directed acyclic graph (DAG), $\mathcal{G} = (\mathbf{V}, \mathbf{E})$, can be deployed to represent a graphical criterion
540 carrying out a set of conditions on the paths, where \mathbf{V} and \mathbf{E} denote the set of nodes and the set of
541 directed edges, separately.

542 **Definition 1.** (*d-separation* [50]). *A set of nodes $\mathbf{Z} \subseteq \mathbf{V}$ blocks the path p if and only if (1) p contains*
543 *a chain $i \rightarrow m \rightarrow j$ or a fork $i \leftarrow m \rightarrow j$ such that the middle node m is in \mathbf{Z} , or (2) p contains a*
544 *collider $i \rightarrow m \leftarrow j$ such that the middle node m is not in \mathbf{Z} and such that no descendant of m is in*
545 *\mathbf{Z} . Let \mathbf{X}, \mathbf{Y} and \mathbf{Z} be disjoint sets of nodes. If and only if the set \mathbf{Z} blocks all paths from one node in*
546 *\mathbf{X} to one node in \mathbf{Y} , \mathbf{Z} is considered to d-separate \mathbf{X} from \mathbf{Y} , denoting as $(\mathbf{X} \perp_d \mathbf{Y} \mid \mathbf{Z})$.*

547 **Definition 2.** (*Global Markov Condition* [51, 50]). *If, for any partition $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, \mathbf{X} is d-separated*
548 *from \mathbf{Y} given \mathbf{Z} , i.e., $\mathbf{X} \perp_d \mathbf{Y} \mid \mathbf{Z}$. Then the distribution P over \mathbf{V} satisfies the global Markov*
549 *condition on graph G , and can be factorizes as, $P(\mathbf{X}, \mathbf{Y} \mid \mathbf{Z}) = P(\mathbf{X} \mid \mathbf{Z})P(\mathbf{Y} \mid \mathbf{Z})$. That is, \mathbf{X} is*
550 *conditionally independent of \mathbf{Y} given \mathbf{Z} , writing as $\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}$.*

551 **Definition 3.** (*Faithfulness Assumption* [51, 50]). *The variables, which are not entailed by the*
552 *Markov Condition, are not independent of each other.*

553 *Under the above assumptions, we can apply d-separation as a criterion to understand the conditional*
554 *independencies from a given DAG \mathcal{G} . That is, for any disjoint subset of nodes $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subseteq \mathbf{V}$,*
555 *$(\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z})$ and $\mathbf{X} \perp_d \mathbf{Y} \mid \mathbf{Z}$ are the necessary and sufficient condition of each other.*

556 B.2 Proofs

557 **Proposition 1** (Identifiability). *Suppose the state s_t , action \mathbf{a}_t , trajectory-wise long-term return R*
558 *are observable while Markovian rewards r_t are unobservable, and they form an MDP, as described*
559 *in Eq. 2. Then under the global Markov condition and faithfulness assumption, the reward function g*
560 *and the Markovian rewards r_t are identifiable, as well as the causal structure that is characterized by*
561 *binary masks \mathbf{c}^{\rightarrow} and \mathbf{C}^{\rightarrow} and the transition dynamics f .*

562 Below is the proof of Proposition 1. We begin by clarifying the assumptions we made and then
563 provide the mathematical proof.

564 **Assumption** We assume that, $\epsilon_{s,i,t}$ and $\epsilon_{r,t}$ in Eq. 2 are i.i.d additive noise. From the weight-space
565 view of Gaussian Process [52], equivalently, the causal models for $s_{i,t+1}$ and r_t can be represented as
566 follows, respectively,

$$s_{i,t+1} = f_i(s_t, \mathbf{a}_t) + \epsilon_{s,i,t} = W_{i,f}^T \phi_f(s_t, \mathbf{a}_t) + \epsilon_{s,i,t}, \quad (\text{A1})$$

567

$$r_t = g(s_t, \mathbf{a}_t) + \epsilon_{r,t} = W_g^T \phi_g(s_t, \mathbf{a}_t) + \epsilon_{r,t}, \quad (\text{A2})$$

568 where $\forall i \in [1, d^s]$, and ϕ_f and ϕ_g denote basis function sets.

569 Then we denote the variable set in the system by \mathbf{V} , with $\mathbf{V} = \{s_{1,t}, \dots, s_{d^s,t}, \mathbf{a}_{1,t}, \dots, \mathbf{a}_{d^a,t}, r_t\}_{t=1}^T \cup$
570 R , and the variables form a Bayesian network \mathcal{G} . Following AdaRL [23], there are possible edges
571 only from $s_{i,t-1} \in \mathbf{s}_{t-1}$ to $s_{i',t} \in \mathbf{s}_t$, from $\mathbf{a}_{j,t-1} \in \mathbf{a}_{t-1}$ to $s_{i',t} \in \mathbf{s}_t$, from $s_{i,t} \in \mathbf{s}_t$ to r_t , and from
572 $\mathbf{a}_{j,t} \in \mathbf{a}_t$ to r_t in \mathcal{G} . In particular, the r_t s are unobserved, while $R = \sum_{t=1}^T \gamma^{t-1} o_t$ is observed. Thus
573 there are deterministic edges from each r_t to R .

574 Below we omit the γ for simplicity.

575 *Proof.* **Given trajectory-wise long-term return R , the binary masks, $\mathbf{c}^{s \rightarrow r}$, $\mathbf{c}^{a \rightarrow r}$ and Markovian**
576 **reward function g and the rewards r_t are identifiable.** Following the above assumption, we first
577 rewrite the function to calculate trajectory-wise long-term return in Eq. 2 as,

$$\begin{aligned} R &= \sum_{t=1}^T r_t \\ &= \sum_{t=1}^T [W_g^T \phi_g(s_t, \mathbf{a}_t) + \epsilon_{r,t}] \\ &= W_g^T \sum_{t=1}^T \phi_g(s_t, \mathbf{a}_t) + \sum_{t=1}^T \epsilon_{r,t}. \end{aligned} \quad (\text{A3})$$

578 For simplicity, we replace the components in Eq. A3 by,

$$\begin{aligned} \zeta_g(X) &= \sum_{t=1}^T \phi_g(s_t, \mathbf{a}_t), \\ E_r &= \sum_{t=1}^T \epsilon_{r,t}, \end{aligned} \quad (\text{A4})$$

579 where $X := [s_t, \mathbf{a}_t]_{t=1}^T$ representing the concatenation of the covariates s_t and \mathbf{a}_t from $t = 1$ to T .
580 Consequently, we derive the following equation,

$$R = W_g^T \zeta_g(X) + E_r. \quad (\text{A5})$$

581 Then we can obtain a closed-form solution of W_g^T in Eq. A5 by modeling the dependencies between
582 the covariates X_τ and response variables R_τ , where both are continuous. One classical approach
583 to finding such a solution involves minimizing the quadratic cost and incorporating a weight-decay
584 regularizer to prevent overfitting. Specifically, we define the cost function as,

$$C(W_g) = \frac{1}{2} \sum_{X_\tau, R_\tau \sim \mathcal{D}} (R_\tau - W_g^T \zeta_g(X_\tau))^2 + \frac{1}{2} \lambda \|W_g\|^2. \quad (\text{A6})$$

585 where τ represents trajectories consisting of state-action pairs X_τ and long-term returns R_τ , which
586 are sampled from the replay buffer \mathcal{D} . λ is the weight-decay regularization parameter. To find the
587 closed-form solution, we differentiate the cost function with respect to W_g and set the derivative to
588 zero:

$$\frac{\partial C(W_g)}{\partial W_g} = 0. \quad (\text{A7})$$

589 Solving this equation will yield the closed-form solution for W_g^T , *i.e.*,

$$W_g = (\lambda I_d + \zeta_g \zeta_g^T)^{-1} \zeta_g R = \zeta_g (\zeta_g^T \zeta_g + \lambda I_n)^{-1} R \quad (\text{A8})$$

590 Therefore, W_g , which indicates the causal structure and strength of the edge, can be identified from
 591 the observed data. In summary, given trajectory-wise long-term return R , the binary masks, $\mathbf{c}^{s \rightarrow r}$,
 592 $\mathbf{c}^{a \rightarrow r}$ and Markovian reward function g and the rewards r_t are identifiable.

593 **The binary masks, $\mathbf{C}^{s \rightarrow s}$, $\mathbf{C}^{a \rightarrow s}$ and the transition dynamics f are identifiable** In a similar
 594 manner, based on the assumption and Eq. 2, we can rewrite Eq. A1 to,

$$\mathbf{s}_{t+1} = W_{i,f}^T \phi_f(\mathbf{s}_t) + \epsilon_{s,i,t}. \quad (\text{A9})$$

595 To obtain a closed-form solution for $W_{i,f}$, f^T in Equation A9, we can model the dependencies between
 596 the covariates X_t and the response variables s_{t+1} , both of which are continuous. The closed-form
 597 solution can be represented as:

$$C(W_{i,f}) = \frac{1}{2} \sum_{\mathbf{s}_{i,t}, \mathbf{s}_{i,t+1} \sim \mathcal{D}} (\mathbf{s}_{i,t+1} - W_{i,f}^T \phi_{i,f}(\mathbf{s}_{i,t}))^2 + \frac{1}{2} \lambda \|W_{i,f}\|^2. \quad (\text{A10})$$

598 By taking derivatives of the cost function and setting them to zero, we can obtain the closed-form
 599 solution,

$$\begin{aligned} W_{i,f} &= (\lambda I_d + \phi_{i,f} \phi_{i,f}^T)^{-1} \phi_{i,f}^T \mathbf{s}_{i,t+1} \\ &= \phi_{i,f} (\phi_{i,f}^T \phi_{i,f} + \lambda I_n)^{-1} \mathbf{s}_{i,t+1}. \end{aligned} \quad (\text{A11})$$

600 Therefore, $W_{i,f}$ can be identified from the observed data. This conclusion applies to all dimensions
 601 of the state. As a result, the f , which indicates the parent nodes of the i -dimension of the state, as
 602 well as the strength of the causal edge, are identifiable. In summary, the binary masks, $\mathbf{C}^{s \rightarrow s}$, $\mathbf{C}^{a \rightarrow s}$
 603 and the transition dynamics f are identifiable.

604 Considering the Markov condition and faithfulness assumption, we can conclude that for any pair
 605 of variables $V_i, V_j \in \mathbf{V}$, V_i and V_j are not adjacent in the causal graph \mathcal{G} if and only if they are
 606 conditionally independent given some subset of $\{V_l \mid l \neq i, l \neq j\}$. Additionally, since there are no
 607 instantaneous causal relationships and the direction of causality can be determined if an edge exists,
 608 the binary structural masks $\mathbf{c}^{s \rightarrow r}$, $\mathbf{c}^{a \rightarrow r}$, $\mathbf{C}^{s \rightarrow s}$, and $\mathbf{C}^{a \rightarrow s}$ defined over the set \mathbf{V} are identifiable with
 609 conditional independence relationships [26]. Consequently, the functions f and g in Equation 2 are
 610 also identifiable.

611

□

612 C Implementation Details

613 C.1 Baselines

614 We compare our method against the following baselines,

- 615 • RRD (biased). This baseline utilizes a surrogate objective called randomized return decom-
 616 position loss for reducing the consumption of estimating the Markovian reward function. It
 617 applies Monte-Carlo sampling to get a biased estimation of the Mean Square Error (MSE)
 618 between the observed episodic reward and the sum of Markovian reward predictions in a
 619 sequence. We keep the same setting and hyper-parameters with its official implementation
 620 to reproduce the results, in which the policy module is optimized by soft actor-critic (SAC)
 621 algorithm [44].
- 622 • RRD (unbiased). This variant of RRD (biased) provides an unbiased estimation of MSE by
 623 sampling short sub-sequences. It offers a computationally efficient approach to optimize
 624 MSE. According to (author?) [14], RRD (biased) and RRD (unbiased) achieve state-of-the-
 625 art performance in episodic MuJoCo tasks.
- 626 • This baseline performs non-parametric uniform reward redistribution. At each time step,
 627 the proxy reward is set to the normalized value of the trajectory return. IRCR is a simple
 628 and efficient approach, and except for RRD, it achieves state-of-the-art performance in the
 629 literature. The implementation is from RRD [14].

Algorithm 1 Learning the generative process and policy jointly.

```
1: Initialize: Environment  $\mathcal{E}$ , trajectory  $\tau \leftarrow \emptyset$ , buffer  $\mathcal{D} \leftarrow \emptyset$ 
2: Initialize: Generative Model  $\Phi_m := [\phi_{\text{cau}}, \phi_{\text{dyn}}, \phi_{\text{rew}}]$ ; Policy Model  $\Phi_\pi$ 
3: for  $i = 1, 2, \dots, 3 \times 10^4$  do
4:    $\tau \leftarrow \emptyset$ , reset  $\mathcal{E}$ 
5:   for  $n_{\text{step}} = 1, 2, \dots, 100$  do
6:     sample data  $\langle s_t, \mathbf{a}_t, o_t \rangle$  from  $\mathcal{E}$ , and store them to trajectory  $\tau$ 
7:     if  $\mathcal{E}$  done then
8:       store trajectory  $\tau = \{s_{1:T}, \mathbf{a}_{1:T}, R\}$  to buffer  $\mathcal{D}$ , where  $R = \sum_{i=1}^T \gamma^{t-1} o_t$ 
9:        $\tau \leftarrow \emptyset$ , reset  $\mathcal{E}$ 
10:    end if
11:    for  $n_{\text{batch}} = 1, 2, \dots, \text{train\_batches}$  do
12:      Sample  $D_1$  consisting of  $M$  trajectories from  $\mathcal{D}$ :  $D_1 = \{\langle s_t^m, \mathbf{a}_t^m \rangle \mid_{t=1}^T, R^m\} \mid_{m=1}^M$ 
13:      Sample  $D_2$  consisting of  $N$  samples from  $\mathcal{D}$ :  $D_2 = \{s_{t_n}, \mathbf{a}_{t_n}, s_{t_n+1}\} \mid_{n=1}^N$ 
14:      Sample binary masks by Gumbel-Softmax from  $\phi_{\text{cau}}$ :  $\hat{\mathbf{c}}^{s \rightarrow r}$ ,  $\hat{\mathbf{c}}^{a \rightarrow r}$ ,  $\hat{\mathbf{C}}^{s \rightarrow s}$  and  $\hat{\mathbf{C}}^{a \rightarrow s}$ 
15:      Sample binary masks deterministically from  $\phi_{\text{cau}}$ :  $\tilde{\mathbf{c}}^{s \rightarrow r}$ ,  $\tilde{\mathbf{c}}^{a \rightarrow r}$ , and  $\tilde{\mathbf{C}}^{s \rightarrow s}$ 
16:      Calculate  $\tilde{\mathbf{e}}^{s \rightarrow \pi}$  based on  $\tilde{\mathbf{c}}^{s \rightarrow r}$  and  $\tilde{\mathbf{C}}^{s \rightarrow s}$ 
17:      Update  $D_2$ :  $D_2 \leftarrow \{\tilde{\mathbf{e}}^{s \rightarrow \pi} \odot s_{t_n}, \mathbf{a}_{t_n}, s_{t_n+1}, \phi_{\text{rew}}(s_{t_n}, \mathbf{a}_{t_n}, \tilde{\mathbf{c}}^{s \rightarrow r}, \tilde{\mathbf{c}}^{a \rightarrow r})\} \mid_{n=1}^N$ 
18:      Optimize  $\phi_{\text{rew}}$  with  $D_1$  (Using  $\hat{\mathbf{c}}^{s \rightarrow r}$  and  $\hat{\mathbf{c}}^{a \rightarrow r}$ ):  $\phi_{\text{rew}} \leftarrow \phi_{\text{rew}} - \alpha \nabla_{\phi_{\text{rew}}} L_{\text{rew}}$  (Eq. 4)
19:      Optimize  $\phi_{\text{dyn}}$  with  $D_2$  (Using  $\hat{\mathbf{C}}^{s \rightarrow s}$  and  $\hat{\mathbf{C}}^{a \rightarrow s}$ ):  $\phi_{\text{dyn}} \leftarrow \phi_{\text{dyn}} - \alpha \nabla_{\phi_{\text{dyn}}} L_{\text{dyn}}$  (Eq. 5)
20:      Optimize  $\phi_{\text{cau}}$ :  $\phi_{\text{cau}} \leftarrow \phi_{\text{cau}} - \alpha \nabla_{\phi_{\text{cau}}} (L_{\text{sp}} + L_{\text{rew}} + L_{\text{dyn}})$  (Eq. 6)
21:      Optimize  $\Phi_\pi$ :  $\Phi_\pi \leftarrow \Phi_\pi - \alpha \nabla_{\Phi_\pi} J_\pi$  (Eq. 8)
22:    end for
23:  end for
24: end for
```

630 C.2 Detailed Generative Model

631 The parametric generative model Φ_m used in the MDP environment consists of three components:
632 ϕ_{cau} , ϕ_{rew} , and ϕ_{dyn} . We provide a detailed description of their model structures below.

633 **ϕ_{cau} for predicting the causal structure** ϕ_{cau} comprises a set of free parameters without input. We
634 divide ϕ_{cau} into four parts, each corresponding to the binary masks in Equation 2. Specifically, we
635 have

- 636 • $\phi_{\text{cau}}^{s \rightarrow s} \in \mathbb{R}^{d^s \times d^s \times 2}$ for $\mathbf{C}^{s \rightarrow s} \in \{0, 1\}^{d^s \times d^s}$,
- 637 • $\phi_{\text{cau}}^{a \rightarrow s} \in \mathbb{R}^{d^a \times d^s \times 2}$ for $\mathbf{C}^{a \rightarrow s} \in \mathbb{R}^{d^a \times d^s}$,
- 638 • $\phi_{\text{cau}}^{s \rightarrow r} \in \mathbb{R}^{d^s \times 2}$ for $\mathbf{c}^{s \rightarrow r} \in \mathbb{R}^{d^s}$,
- 639 • $\phi_{\text{cau}}^{a \rightarrow r} \in \mathbb{R}^{d^a \times 2}$ for $\mathbf{c}^{a \rightarrow r} \in \mathbb{R}^{d^a}$.

640 Below we explain the shared workflows in ϕ_{cau} using the example of predicting the causal edge
641 from the i -th dimension of state $s_{i,t}$ to the j -th dimension of the next state $s_{j,t+1}$, by part of the free
642 parameters, $\phi_{\text{cau},i,j}^{s \rightarrow s}$.

643 For simplicity, we denote $\phi_{\text{cau},i,j}^{s \rightarrow s}$ as ψ . The shape of ψ is now easy to be determined. That is
644 $\psi \in \mathbb{R}^2$ and we write it as $\psi = [\psi_0, \psi_1]$. With this 2-element vector, we can characterize a Bernoulli
645 distribution, where each element corresponds to the unnormalized probability of classifying the
646 edge as existing (ψ_0) or not existing (ψ_1), respectively. Therefore, the probability of the causal edge
647 existing from the i -th dimension of state $s_{i,t}$ to the j -th dimension of the next state $s_{j,t+1}$ can be
648 calculated as:

$$P(\mathbf{C}_{i,j}^{s \rightarrow s}) = \frac{\exp(\psi_0)}{\exp(\psi_0) + \exp(\psi_1)} \quad (\text{A12})$$

649 **Obtain $\hat{\mathbf{C}}_{i,j}^{s \rightarrow s}$ through Gumbel-Softmax sampling in the training phases.** During training, it is
650 crucial to maintain the gradient flow for backpropagation. To achieve this, we sample the binary

Layer#	1	2	3
ϕ_{rew}	FC256	FC256	FC1
ϕ_{dyn}	FC256	FC256	FC9
ϕ_{π}	FC256	FC256	FC2d ^a
ϕ_{ν}	FC256	FC256	FC1

Table A1: The network structures of ϕ_{rew} , ϕ_{dyn} , ϕ_{π} and ϕ_{ν} . FC256 denotes a fully-connected layer with an output size of 256. Each hidden layer is followed by an activation function, ReLU. d^a is the number of dimensions of the action in a specific task.

651 values of $\hat{\mathbf{C}}_{i,j}^{s \rightarrow s}$ by applying Gumbel-Softmax [41],

$$\hat{\mathbf{C}}_{i,j}^{s \rightarrow s} = \text{GS}(\psi) \quad (\text{A13})$$

652 where GS denotes the Gumbel-Softmax sampling, which allows us to obtain binary discrete samples
653 from the Bernoulli distribution. By applying Gumbel Softmax sampling allows us to randomly
654 sample from the Bernoulli distribution in a stochastic manner, rather than simply selecting the class
655 with the highest probability. This introduces some randomness, enabling the model to explore the
656 balance and uncertainty between different classifications more flexibly.

657 **Obtain $\tilde{\mathbf{C}}_{i,j}^{s \rightarrow s}$ by deterministic sampling in the inference phases.** During inference, including data
658 sampling and policy learning, we get the prediction of $\mathbf{C}_{i,j}^{s \rightarrow s}$ through a deterministic sampling,

$$\tilde{\mathbf{C}}_{i,j}^{s \rightarrow s} = \begin{cases} 1, \psi_0 \geq \psi_1 \\ 0, \psi_0 < \psi_1. \end{cases} \quad (\text{A14})$$

659 This is a greedy sampling to avoid introducing randomness during the Gumble-Softmax sampling.

660 The above explanation of the workflow in ϕ_{cau} for predicting a single causal edge provides insight
661 into the overall implementation of the entire module ϕ_{cau} and can be applicable for all the causal
662 edges. Therefore, we can obtain $\hat{\mathbf{C}}^{a \rightarrow s}$, $\hat{\mathbf{c}}^{s \rightarrow r}$, $\hat{\mathbf{c}}^{a \rightarrow r}$, $\tilde{\mathbf{c}}^{s \rightarrow r}$ and $\tilde{\mathbf{c}}^{a \rightarrow r}$, using similar procedures.

663 **ϕ_{rew} for predicting the Markovian rewards** ϕ_{rew} is a stacked fully-connected network, and the
664 details of the network structure are provided in Table A1.

665 During training, the prediction of Markovian reward can be written as,

$$\hat{r} = \phi_{\text{rew}}(s_t, \mathbf{a}_t, \hat{\mathbf{c}}^{s \rightarrow r}, \hat{\mathbf{c}}^{a \rightarrow r}) = \text{FCs}([\hat{\mathbf{c}}^{s \rightarrow r} \odot s_t, \hat{\mathbf{c}}^{a \rightarrow r} \odot \mathbf{a}_t]), \quad (\text{A15})$$

666 where $[\cdot, \cdot]$, \odot denotes concatenation and element-wise multiply operations, respectively. FCs denotes
667 the stacked fully-connected network. $\hat{\mathbf{c}}^{s \rightarrow r}$ and $\hat{\mathbf{c}}^{a \rightarrow r}$ are derived from ϕ_{cau} by Gumbel-Softmax.

668 During inference, including policy learning and data sampling, the predicted Markovian reward is

$$\tilde{r} = \phi_{\text{rew}}(s_t, \mathbf{a}_t, \tilde{\mathbf{c}}^{s \rightarrow r}, \tilde{\mathbf{c}}^{a \rightarrow r}) = \text{FCs}([\tilde{\mathbf{c}}^{s \rightarrow r} \odot s_t, \tilde{\mathbf{c}}^{a \rightarrow r} \odot \mathbf{a}_t]), \quad (\text{A16})$$

669 where $\tilde{\mathbf{c}}^{s \rightarrow r}$ and $\tilde{\mathbf{c}}^{a \rightarrow r}$ are derived from ϕ_{cau} greedily by deterministic sampling.

670 **ϕ_{dyn} for modeling the environment dynamics** In our experiment, we do not directly utilize ϕ_{dyn}
671 in policy learning. Instead, this module serves as a bridge to optimize $\phi_{\text{cau}}^{s \rightarrow s}$ and $\phi_{\text{cau}}^{a \rightarrow s}$. Subsequently,
672 $\phi_{\text{cau}}^{s \rightarrow s}$ can be utilized in the calculation of $\tilde{\mathbf{c}}^{s \rightarrow \pi}$.

673 During training, we initially sample $\hat{\mathbf{C}}^{s \rightarrow s}$ and $\hat{\mathbf{C}}^{a \rightarrow s}$ using Gumbel-Softmax. The prediction for the
674 i -th dimension of the next state can be represented as follows,

$$\hat{s}_{i,t} = \text{MDN}([\hat{\mathbf{C}}_{\cdot,i}^{s \rightarrow s} \odot s_t, \hat{\mathbf{C}}_{\cdot,i}^{a \rightarrow s} \odot \mathbf{a}_t]), \quad (\text{A17})$$

675 where $[\cdot, \cdot]$ denotes concatenation and MDN denotes the Mixture Density Network which outputs the
676 means, variances, and probabilities for N_{Gau} Gaussian cores. The parameters of MDN are shared
677 across the predictions of different dimensions of the next state. We set $N_{\text{cau}} = 3$ in our experiments.
678 More details about ϕ_{dyn} can be found in Table A1.

679 C.3 Detailed Policy Model

680 Considering the specific requirements of the employed RL algorithm, Soft Actor-Critic (SAC), our
681 Policy Model Φ_{π} comprises two components, the actor ϕ_{π} and the critic ϕ_{ν} . Detailed network
682 structures for both components can be found in Table A1.

Table A2: The table of the hyper-parameters used in the experiments for GRD.

Envs	λ_1	λ_2	λ_3	λ_4	λ_5
<i>Ant</i>	10^{-5}	0	10^{-7}	10^{-8}	10^{-8}
<i>HalfCheetah</i>	10^{-5}	10^{-5}	10^{-5}	10^{-6}	10^{-5}
<i>Walker2d</i>	10^{-5}	10^{-5}	10^{-6}	10^{-6}	10^{-7}
<i>Humanoid</i>	10^{-5}	10^{-8}	10^{-5}	10^{-7}	10^{-8}
<i>Reacher</i>	5×10^{-7}	10^{-8}	10^{-8}	10^{-8}	10^{-8}
<i>Swimmer</i>	10^{-7}	10^{-9}	10^{-9}	0	10^{-9}
<i>Hopper</i>	10^{-6}	10^{-6}	10^{-6}	10^{-7}	10^{-6}
<i>HumanStandup</i>	10^{-5}	10^{-4}	10^{-6}	10^{-7}	10^{-7}

Table A3: The hyper-parameters.

hyperparameters	value	hyperparameters	value
epochs	3	optimizer	Adam
cycles	100	learning rate	3×10^{-4}
iteration	100	N	256
train_batches	100	M	4
replay buffer size	10^6	γ	1.00
evaluation episodes	10	Polyak-averaging coefficient	0.0005

683 C.4 Training Process.

684 We follow the line of joint learning in (author?) [14], which avoids learning a return decomposition
 685 model in advance using data sampled by optimal or sub-optimal policies [13]. During each mini-batch
 686 training iteration, we sample two sets of data separately from the replay buffer \mathcal{D} :

- 687 • $D_1 = \{\langle s_t^m, \mathbf{a}_t^m \rangle |_{t=1}^T, R^m\} |_{m=1}^M$ consists of M trajectories. Provided with the trajectory-wise
 688 long-term returns $R^m |_{m=1}^M$, D_1 is utilized to optimize $\phi_{\text{cau}}^{s \rightarrow r}$, $\phi_{\text{cau}}^{a \rightarrow r}$ and ϕ_{rew} , with L_{rew} .
- 689 • $D_2 = \{s_{t_n}, \mathbf{a}_{t_n}, s_{t_n+1}\} |_{n=1}^N$ consists of N state-action pairs. D_2 are used for policy optimiza-
 690 tion and optimize the parts of causal structure, $\phi_{\text{cau}}^{s \rightarrow s}$ and $\phi_{\text{cau}}^{a \rightarrow s}$, ϕ_{dyn} . With such a D_2 , GRD
 691 breaks the temporal cues in the training data to learn the policy and dynamics function.

692 Please refer to Algorithm 1 for a detailed training process.

693 C.5 Hyper-Parameters.

694 The network is trained from scratch using the Adam optimizer, without any pre-training. The initial
 695 learning rate for both model estimation and policy learning is set to 3×10^{-4} . The hyperparameters
 696 for policy learning are shared across all tasks, with a discount factor of 1.00 and a Polyak-averaging
 697 coefficient of 5×10^{-4} . The target entropy is set to the negative value of the dimension of the robot
 698 action. To facilitate training, we utilize a replay buffer with a size of 1×10^6 time steps. The warmup
 699 size of the buffer for training is set to 1×10^4 . The model is trained for 3 epochs, with each epoch
 700 consisting of 100 training cycles. In each cycle, we repeat the process of data collection and model
 701 training for 100 iterations. During each iteration, we collect data from 100 time steps of interaction
 702 with the MuJoCo simulation, which is then stored in the replay buffer. For training the ϕ_{rew} , we
 703 sample 4 episodes, each containing 5×10^3 steps. As for policy learning and the optimization of ϕ_{dyn} ,
 704 we use data from 256 time steps. ϕ_{cau} is trained together with ϕ_{rew} and ϕ_{dyn} . Validation is performed
 705 after every cycle, and the average metric is computed based on 10 test rollouts. The hyperparameters
 706 for learning the GRD model can be found in Table A2. All experiments were conducted on an
 707 HPC system equipped with 128 Intel Xeon processors operating at a clock speed of 2.2 GHz and 5
 708 terabytes of memory.

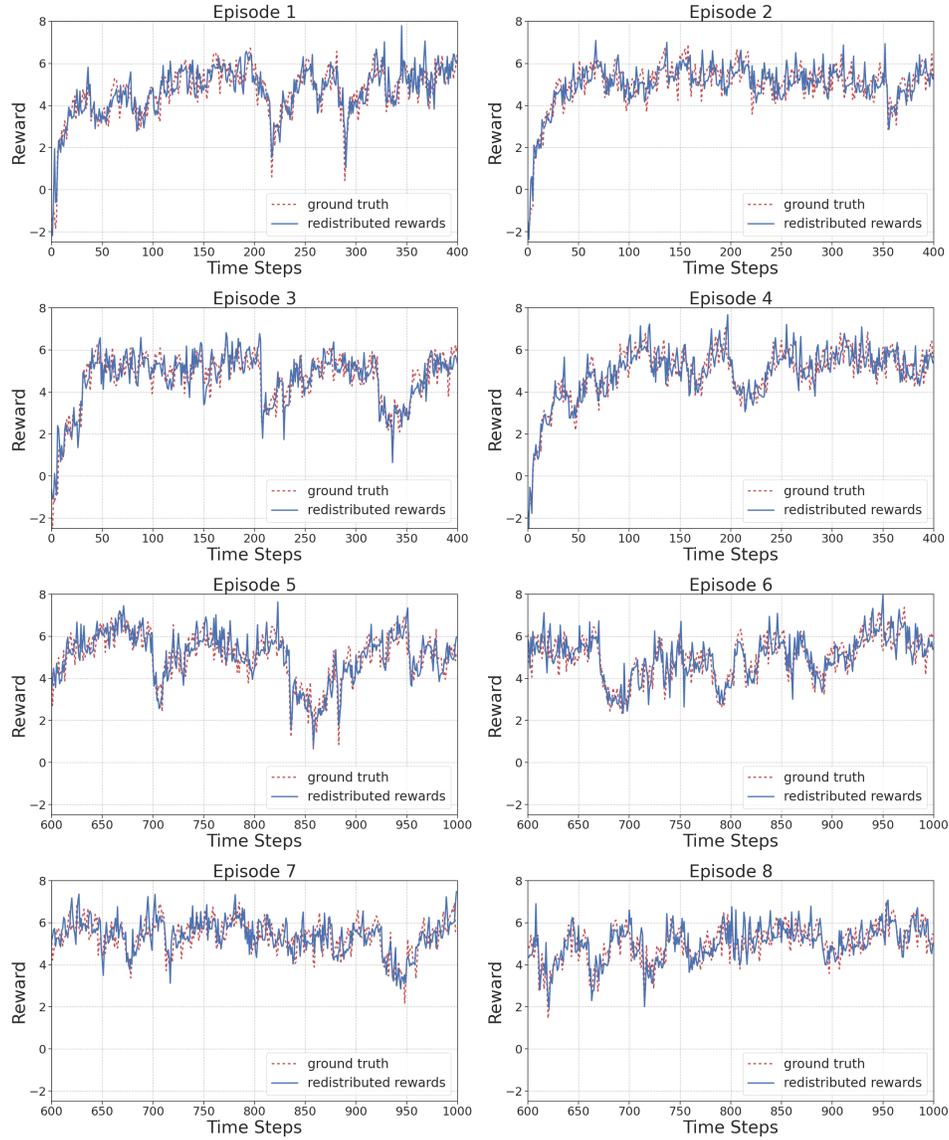


Figure A1: The visualization of redistributed rewards and grounded rewards in *Ant*. The results are produced by the GRD model trained after 1×10^6 steps. The redistributed rewards are shown in red, and the grounded rewards are shown in blue.

709 **D Visualization**

710 As shown in Figure A1, we visualize the redistributed rewards in *Ant* by GRD, as well as the grounded
 711 rewards provided by the environment.