

CAUSAL INFORMATION PRIORITIZATION FOR EFFICIENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement Learning (RL) methods often suffer from sample inefficiency, one of the underlying reasons is that blind exploration strategies may neglect causal relationships among states, actions, and rewards. Although recent causal approaches aim to address this problem, they lack grounded modeling of reward-guided causal understanding of states and actions for goal orientation, thus impairing learning efficiency. To tackle this issue, we propose a novel method named Causal Information Prioritization (**CIP**) that improves sample efficiency by leveraging factored MDPs to infer causal relationships between different dimensions of states and actions with respect to rewards, enabling the prioritization of causal information. Specifically, **CIP** identifies and leverages causal relationships between states and rewards to execute counterfactual data augmentation to prioritize high-impact state features under the causal understanding of the environments. Moreover, **CIP** integrates a causality-aware empowerment learning objective, which significantly enhances the agent’s execution of reward-guided actions for more efficient exploration in complex environments. To fully assess the effectiveness of **CIP**, we conduct extensive experiments across 39 tasks in 5 diverse continuous control environments, encompassing both locomotion and manipulation skills learning with pixel-based and sparse reward settings. Experimental results demonstrate that **CIP** consistently outperforms existing RL methods across a wide range of scenarios¹.

1 INTRODUCTION

Reinforcement Learning (RL) has emerged as a powerful paradigm for training intelligent decision-making agents to learn optimal behaviors by interacting with their environments, receiving reward feedback, and iteratively optimizing their decision-making policies (Harnoja et al., 2018; Ze et al., 2024; Sutton, 2018; Silver et al., 2017). Despite its notable successes, most RL approaches are faced with the sample-inefficiency problem, which means they typically necessitate an enormous number of interactions with the environment to learn policies, which can be impractical or costly in real-world scenarios (Savva et al., 2019; Kroemer et al., 2021). Inefficient policy learning often results from blind exploration strategies that neglect causal relationships, leading to spurious correlations and suboptimal solutions with high exploration costs (Zeng et al., 2023; Liu et al., 2024).

Causal reasoning captures essential information by analyzing causal relationships between different factors, filtering out irrelevant information, and avoiding interference from spurious correlations (Wang et al., 2022; Pitis et al., 2022; Zhang et al., 2024; Huang et al., 2022b). These approaches build internal causal structural models, enabling agents to strategically focus their exploration on the most pertinent aspects of the environment. They significantly reduce the number of samples required and demonstrate remarkable performance in single-task learning, generalization, and counterfactual reasoning (Richens & Everitt, 2024; Urpí et al., 2024; Deng et al., 2023; Huang et al., 2022a; Feng & Magliacane, 2023). However, most of these works overlook the reward-relevant causal relationships among different factors, or only partially consider the causal connections between states, actions, and rewards (Liu et al., 2024; Ji et al., 2024a), thus hindering efficient exploration.

In this work, we aim to identify and exploit task-specific causal relationships between states, actions, and rewards, enabling agents to discern relevant states and select actions that maximize rewards, ultimately facilitating precise and goal-oriented behaviors. Here we provide a motivating example in Figure 1, showing three trajectories for executing a manipulation soccer task, along with the

¹The anonymous project page is <https://sites.google.com/view/rl-cip/>.

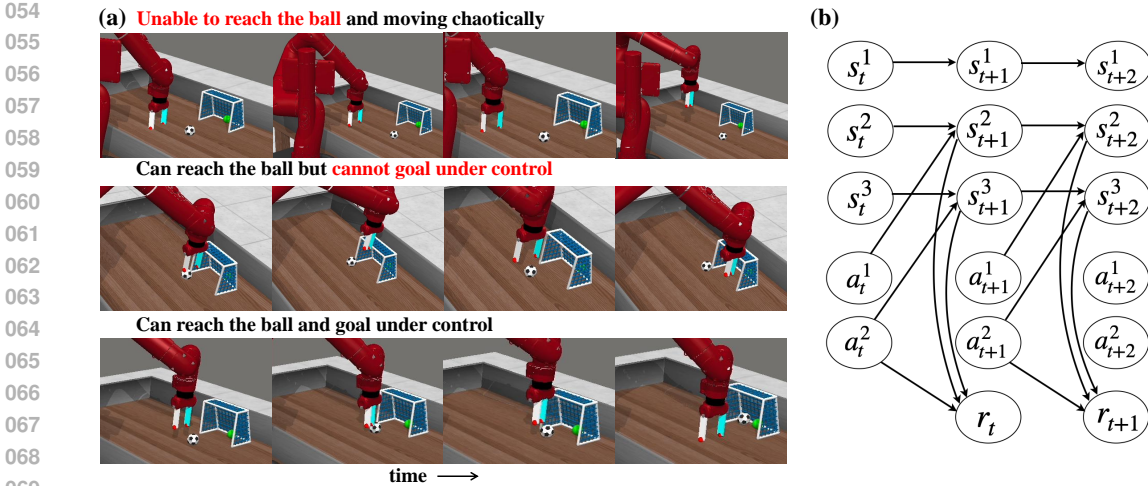


Figure 1: (a). An example of a robot manipulation soccer task with three trajectories, where the objective is to move the ball into the goal. (b). Underlying causal structure of this example in a factored MDP. Different nodes represent different dimensional states and actions.

underlying causal structure in a factored Markov Decision Process (MDP) (Kearns & Koller, 1999). In the first trajectory (row 1), when the agent fails to distinguish states with more intricate causal relationships of the task, the robotic arm exhibits chaotic moving and receives no rewards. The second trajectory (row 2) shows that even without chaotic movements, uncontrollable actions unrelated to the reward lead to an inability to guide the ball towards the goal. Only by filtering out irrelevant state features and executing more controllable actions can we guarantee that the ball is kicked into the goal like row 3. Quantifying the contribution of different factors to the reward can effectively help analyzing important causal relationships.

To address the limitation of sample-inefficiency and leverage the potential of causal reasoning, we propose a novel approach named Causal Information Prioritization (CIP) for efficient RL, improving learning efficiency from the perspective of rewards. Building upon the factored MDPs, CIP infers causal relationships between states, actions, and rewards across different dimensions, respectively. CIP employs counterfactual data augmentation based on the causality between states and rewards to generate transitions, prioritizing critical state transitions. Furthermore, CIP leverages the causality between actions and rewards to reweight actions, while utilizing empowerment to maximize mutual information between causally informed actions and future states, thereby enabling better control.

Specifically, CIP leverages collected data to construct a reward-guided structural model that explicitly reasons about state-reward causal influences, enabling the swapping of causally independent state features across observed trajectories to generate synthetic transitions without additional environment interactions. By swapping independent state features across different transitions (i.e., irrelevant state dimensions of chaotic movements in the soccer task), CIP accentuates causally dependent state information (i.e., relevant states to reach the ball), facilitating focused learning of critical state transitions. Subsequently, CIP constructs another structural model that incorporates actions and rewards to reweight actions of dimensions. To enhance the exploration efficiency, CIP integrates a causality-aware empowerment term, quantifying the agent’s capacity to exert controlled influence over its environment through the mutual information. This empowerment term, combined with causally weighted actions, is integrated into the learning objective, prioritizing actions with high causal influence. The synthesis of causal reasoning and action empowerment enables agents to focus on behaviors that are causally relevant to the task, leading to more efficient and effective policy learning. The main contributions of this work can be summarized as follows.

- To address limitations of blind exploration and sample-inefficiency, we introduce CIP, a novel efficient RL framework that prioritizes causal information through the lens of reward. CIP bridges the gap between causal reasoning and empowerment to facilitate efficient exploration.
- CIP constructs reward-guided structural models to uncover causal relationships between states, actions, and rewards across dimensions. By leveraging state-reward causality, it performs counterfactual data augmentation, eliminating the need for additional environment interactions, and

enabling learning on critical state transitions. Exploiting action-reward causality, it reweights actions to enhance exploration efficiency through empowerment. By prioritizing causal information, **CIP** enables agents to focus on behaviors that have causally significant effects on their tasks.

- To validate the effectiveness of **CIP**, we conduct extensive experiments in 39 tasks across 5 diverse continuous control environments, including manipulation and locomotion. These comprehensive evaluations demonstrate the effectiveness of **CIP** in pixel-based and sparse reward settings, underscoring its versatility and reliability.

2 RELATED WORK

2.1 CAUSAL RL

The application of causal reasoning in RL has shown significant potential to improve sample efficiency and generalization by effectively excluding irrelevant environmental factors through causal analysis (Huang et al., 2022a; Feng & Magliacane, 2023; Mutti et al., 2023; Sun et al., 2024; Sun & Wang). Wang (Wang et al., 2021) introduces a novel regularization-based method for causal dynamics learning, which explicitly identifies causal dependencies by regulating the number of variables used to predict each state variable. CDL (Wang et al., 2022) takes an innovative approach by using conditional mutual information to compute causal relationships between different dimensions of states and actions. IFactor (Liu et al., 2024) is a general framework to model four distinct categories of latent state variables, capturing various aspects of information. ACE (Ji et al., 2024a), an off-policy actor-critic method, integrates causality-aware entropy regularization. Table 2 provides a categorization of various causal RL methods, highlighting their focus on different reward-guided causal relationships. Existing approaches do not fully account for the causal relationships between both states and actions with rewards. Our goal is to explore these causal relationships from a reward-guided perspective to enhance sample efficiency across a broader range of tasks.

2.2 EMPOWERMENT IN RL

Empowerment, an information theory-based concept of intrinsic motivation, has emerged as a powerful paradigm for enhancing an agent’s environmental controllability (Mohamed & Jimenez Rezende, 2015; Klyubin et al., 2005; Cao et al., 2024). This framework conceptualizes actions and future states as information transmission channels, offering a novel perspective on agent-environment interactions. In RL, empowerment has been applied to uncover more controllable associations between states and actions, as well as to develop robust skill (Salge et al., 2014; Bharadhwaj et al., 2022; Choi et al., 2021; Eysenbach et al., 2018; Leibfried et al., 2019; Seitzer et al., 2021). Empowerment, expressed as maximizing mutual information $\max_{\pi} I$, serves as a learning objective in various RL frameworks, providing intrinsic motivation for exploration and potentially yielding more efficient and generalizable policies. Our approach extends empowerment in RL by examining the influence of state, actions, and rewards through a causal lens, integrating causal understanding with empowerment to enhance exploration strategy and learning efficiency.

2.3 OBJECT-CENTRIC RL AND OBJECT-ORIENTED RL

Recent advances in object-centric representation learning focus on acquiring and leveraging structured, object-wise representations from high-dimensional observations. Foundational works include Slot Attention (Locatello et al., 2020) and AIR (Eslami et al., 2016; Kosiorek et al., 2018), establishing basis for this field. Subsequent follow-ups have worked on these concepts by employing state-of-the-art architectures, including DINO-based approaches Zadaianchuk et al. (2023), transformer-based models (Wu et al., 2022), diffusion models (Jiang et al., 2023), and state-space models (Jiang et al., 2024). Notably, learning object-centric representations can enable compositional generalization across various domains, such as video and scene generation (Wu et al., 2023; 2024). Moreover, several theoretical studies have explored the mechanisms underlying compositional generalization and the causal identifiability (Kori et al., 2024; Brady et al., 2023; Lachapelle et al., 2024).

Object-centric representations have been effectively employed in world models to capture multi-object dynamics, as demonstrated by works (Jiang et al., 2019; Lin et al., 2020; Kossen et al., 2019). Building on these object-centric world models, various studies use them in RL by better modeling complex object-centric structures in partially observable MDPs (Kossen et al., 2019; Mambelli et al., 2022; Feng & Magliacane, 2023; Choi et al., 2024), identifying critical objects (Zadaianchuk et al., 2022; Park et al., 2021), and learning object-centric policies (Zadaianchuk et al., 2021; Yuan et al., 2022) and applications in robotic manipulation tasks (Li et al., 2020; Mitash et al., 2024; Haramati

et al., 2023; Li et al., 2024), as well as in learning intrinsic or curiosity-driven policies based on objects and their interactions (Watters et al., 2019; Wang et al., 2024c;b).

Another research direction explores object-oriented MDPs, with the homomorphic object-oriented world model being a notable example that leverages MDP homomorphism to model object dynamics and enable efficient planning through symmetric equivalence in MDPs (Diuk et al., 2008; Scholz et al., 2014; Wandzel et al., 2019; Van der Pol et al., 2020; Rezaei-Shoshtari et al., 2022; Zhao et al., 2022), provides a powerful foundation for learning object-oriented MDPs and facilitates efficient planning (Wolfe & Barto, 2006). Our work, which focuses on uncovering general causal relationships among components in MDPs and empowerment optimization, is orthogonal to object-centric RL. However, object-centric RL could provide useful abstract object-based variables that could be useful for causal structure learning in complex environments ².

3 PRELIMINARIES

3.1 MARKOV DECISION PROCESS

In RL, the agent-environment interaction is formalized as an MDP. The standard MDP is defined by the tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mu_0, r, \gamma \rangle$, where \mathcal{S} denotes the state space, \mathcal{A} represents the action space, $\mathcal{P}(s'|s, a)$ is the transition dynamics, $r(s, a)$ is the reward function, and μ_0 is the distribution of the initial state s_0 . The discount factor $\gamma \in [0, 1)$ is also included. The objective of RL is to learn a policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ that maximizes the expected discounted cumulative reward $\eta_{\mathcal{M}}(\pi) := \mathbb{E}_{s_0 \sim \mu_0, s_t \sim \mathcal{P}, a_t \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$.

3.2 STRUCTURAL CAUSAL MODEL

A Structural Causal Model (SCM) (Pearl, 2009) is defined by a distribution over random variables, defined as $\mathcal{V} = \{s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^n, r_t, s_{t+1}^1, \dots, s_{t+1}^d\}$ and a Directed Acyclic Graph (DAG) $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with a conditional distribution $\mathcal{P}(v_i | \text{PA}(v_i))$ for node $v_i \in \mathcal{V}$. Then the distribution can be specified as:

$$p(v_1, \dots, v_{|\mathcal{V}|}) = \prod_{i=1}^{|\mathcal{V}|} p(v_i | \text{PA}(v_i)), \quad (1)$$

where $\text{PA}(v_i)$ is the set of parents of the node v_i in the graph \mathcal{G} .

Causal Structures in MDP We use a factored MDP (Kearns & Koller, 1999; Guestrin et al., 2003; 2001) to model the MDP and the underlying causal structures between states, actions, and rewards. In the factored MDP, nodes represent system variables (rewards and different dimensions of the states and actions), while the edges denote their relationships within the MDP. We employ causal discovery methods to learn the structures of \mathcal{G} .

We can identify the graph structure in \mathcal{G} , which can be represented as the adjacency matrix M . To integrate such relationships in MDP, we explicitly encode the causal mask over variables into the reward function. Hence, the reward function in MDP with the causal structure is defined as follows:

$$r_t = R(M^{s \rightarrow r} \odot s_t, M^{a \rightarrow r} \odot a_t, \epsilon_{r,t}) \quad (2)$$

where \odot denotes the element-wise product. $M^{s \rightarrow r} \in \mathbb{R}^{|\mathcal{S}| \times 1}$ and $M^{a \rightarrow r} \in \mathbb{R}^{|\mathcal{A}| \times 1}$ are the adjacency matrices indicating the influence of current states and actions on the reward, respectively, and $\epsilon_{r,t}$ represents i.i.d. Gaussian noise. Under the Markov condition and faithfulness assumption (Pearl, 2009; Spirtes et al., 2001), the structural vectors are identifiable. The detailed assumptions and propositions can be found in Appendix B. In this work, our objective is to discover and leverage these two causal matrices to prioritize causal information for efficient RL.

3.3 EMPOWERMENT IN RL

Empowerment quantifies an agent’s capacity to influence its environment and perceive the consequences of its actions (Klyubin et al., 2005; Bharadhwaj et al., 2022; Jung et al., 2011). In our framework, the empowerment is defined as the mutual information between the agent state s_{t+1} and action a_t , conditioned on the present state s_t and causal mask M , as shown follows:

$$\mathcal{E} := \max_{p(a_t)} \mathcal{I}(a_t; s_{t+1} | s_t, M), \quad (3)$$

²We provide a detailed discussion in Appendix C.1.

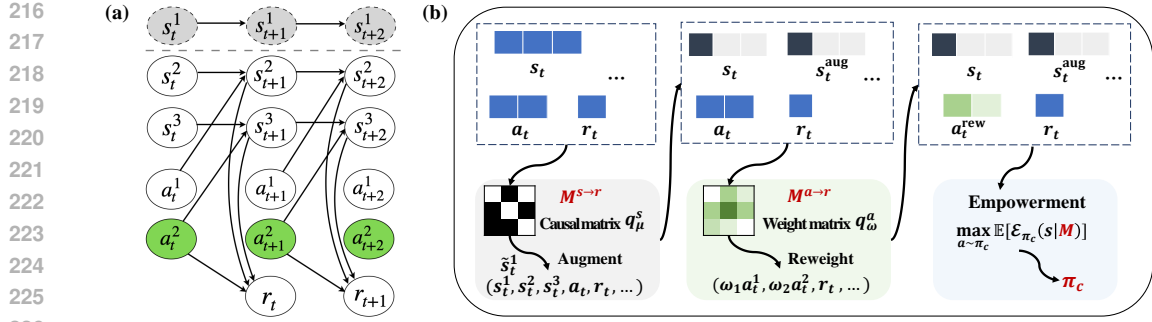


Figure 2: (a) Underlying causal structure of **CIP**. (b) The whole learning process of **CIP** includes counterfactual data augmentation, causal action reweight and causal action empowerment.

where \mathcal{E} denotes the channel capacity from action to state, and $p(a_t)$ is the distribution of actions. Unlike (Cao et al., 2024), which focuses on action-to-state empowerment effects, we leverage causal understanding and more accurate entropy calculation to analyze state-to-action influences, facilitating the development of more controllable behavioral policies.

4 CAUSAL INFORMATION PRIORITIZATION

In this section, we introduce the proposed framework Causal Information Prioritization (**CIP**), which implements causal information prioritization based on the causal relationships between states, actions, and rewards (as shown in Figure 2). First, we train a structural model based on the causal discovery method, DirectLiNGAM (Shimizu et al., 2011) using collected trajectories to obtain a causal matrix $M^{s \rightarrow r}$. Utilizing this matrix, **CIP** executes the swapping of causally independent state features, generating synthetic transitions (Section 4.1). This process of swapping independent state information accentuates causally dependent state information, enabling focused learning on critical state transitions. Subsequently, **CIP** constructs another structural model to get a weight matrix $M^{a \rightarrow r}$ that incorporates actions and rewards to reweight actions (Section 4.2). Furthermore, **CIP** integrates a causality-aware empowerment term $\mathcal{E}_{\pi_c}(s)$ combined with causally weighted actions into the learning objective to promote efficient exploration. This integration encourages the agent’s policy π_c to prioritize actions with high causal influence, thereby enhancing its goal-achievement capabilities.

4.1 COUNTERFACTUAL DATA AUGMENTATION

To discover the causal relationships between states and rewards, we initially collect trajectories to train a structural model by the DirectLiNGAM method, denoted as q_μ^s , to obtain the causal matrix $M^{s \rightarrow r}$. Subsequently, we infer the local factorization, which is utilized to generate counterfactual transitions. For each state s in the trajectories, we compute the uncontrollable set, defined as the set of variables in s for which the agent has no causal influence on rewards:

$$\mathcal{U}_s = \{s^i \mid M^{s \rightarrow r} \cdot (s_t^i, r_t) < \theta; i \in [1, N]\}, \quad (4)$$

where θ is a fixed threshold and N is the dimension of the state space. The set \mathcal{U}_s encompasses all dimensional state variables for which the causal relationship $s_t^i \rightarrow r_t$ does not exist in the causal matrix of states and rewards. Utilizing the learned causal matrix $M^{s \rightarrow r}$, we partition all state variables in the factored MDP into controllable and uncontrollable sets. These uncontrollable sets are then leveraged for counterfactual data augmentation, thereby prioritizing the causally-informed state information to improve learning efficiency.

To generate counterfactual samples, we perform a swap of variables that fall under the uncontrollable category (i.e., in set \mathcal{U}_s) sampled from the collect trajectories. Specifically, given two transitions (s_t, a_t, s_{t+1}, r_t) and $(\hat{s}_j, \hat{a}_j, \hat{s}_{j+1}, \hat{r}_j)$ sampled from trajectories, which share at least one uncontrollable sub-graph structure (i.e., $\mathcal{U}_s \cap \mathcal{U}_{\hat{s}} \neq \emptyset$), we construct a counterfactual transition $(\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}, \tilde{r}_t)$ by swapping the irrelevant state variables (s_t^i, s_{t+1}^i) with $(\hat{s}_j^i, \hat{s}_{j+1}^i)$ for each $i \in \mathcal{U}_s \cap \mathcal{U}_{\hat{s}}$. The augmented transitions will be added to the training data for causal reasoning during subsequent action empowerment, thus eliminating the need for additional environment interactions to prioritize causal information. Furthermore, we also consider directly using controllable state sets combined with

causal action empowerment to replace counterfactual data augmentation for policy learning. The comparative experimental results validating this approach are presented in Appendix D.3.1.

4.2 CAUSAL ACTION PRIORITIZATION THROUGH EMPOWERMENT

Causal action reweight Having analyzed the causal relationships between states and rewards to achieve efficient data augmentation, in this section, we further discover the causal relationships between actions and rewards to prioritize causally-informed decision-making behaviors. **CIP** constructs a reward-guided structural model, incorporating states (including augmented states), actions, and rewards. This model forms the foundation for action prioritization in policy learning, enabling action reweighting based on causality. Leveraging this structural model to delineate relationships between policy decisions and rewards, we evaluate the causal impact of different actions on reward outcomes. In this way, the agent focuses on pivotal actions with demonstrable causal links to desired reward outcomes, potentially accelerating learning and optimizing performance in complex environments.

Specifically, in **CIP**, we employ DirectLiNGAM method to train a causal structural model q_ω^a , which yields a weight matrix $M^{a \rightarrow r}$, delineating the relationships between actions and rewards, conditioned on the states. For a given set of actions $(a_t^1, a_t^2, a_t^3, \dots)$, we utilize the weight matrix $M^{a \rightarrow r}$ to reweight them as $(\omega_1 a_t^1, \omega_2 a_t^2, \omega_3 a_t^3, \dots)$, where ω represents the causal weights derived from the matrix $M^{a \rightarrow r}$. By leveraging this causal structure, we can prioritize the most pivotal actions, potentially leading to more efficient policy exploration and targeted policy improvements.

Causal action empowerment Based on the learned causal structure, we propose the causal action empowerment to incorporate the reweighted actions into the learning objective for efficient exploration in a controllable manner. To this end, we design a causality-aware empowerment term $\mathcal{E}_{\pi_c}(s)$ for policy optimization. We maximize the empowerment gain of the policy π_c , where π_c incorporates the learned causal structure. This approach allows us to quantify and maximize the empowerment that can be achieved by explicitly considering causal relationships, thereby bridging the gap between causal reasoning and empowerment.

We denote the empowerment of the causal policy as $\mathcal{E}_{\pi_c}(s) = \max_a \mathcal{I}(a_t; s_{t+1} | s_t; M)$. We then formulate the following objective empowerment function:

$$\begin{aligned} \mathcal{E}_{\pi_c}(s) &= \max_a \mathcal{I}(a_t; s_{t+1} | s_t; M) \\ &= \max_{a_t \sim \pi_c(\cdot | s)} \mathcal{H}(\pi_c(a_t | s_t)) - \mathcal{H}(\pi_c(a_t | s_t; s_{t+1})), \end{aligned} \quad (5)$$

where π_c is the policy under the causal weighted matrix $M^{a \rightarrow r}$. The first entropy term $\mathcal{H}(\pi_c(a_t | s_t))$ promotes action diversity within the constraints of the causal structure. It encourages the agent to explore a wide range of actions that are causally informed, while the second entropy term $-\mathcal{H}(\pi_c(a_t | s_t; s_{t+1}))$ enhances the action predictability in state transitions. It encourages the selection of actions that lead to predictable outcomes, given the current and subsequent states, thereby promoting controlled and goal-oriented behaviors. We train an inverse dynamics model to represent the policy $\pi_c(\cdot | s_t; s_{t+1})$. The detailed derivation proceeds as follows:

$$\mathcal{H}(\pi_c(\cdot | s_t)) = -\mathbb{E}_{a_t \in \mathcal{A}} \left[\sum_{i=1}^{d_A} M^{a^i \rightarrow r} \odot \pi_c(a_t^i | s_t) \log \pi(a_t^i | s_t) \right], \quad (6)$$

and

$$\mathcal{H}(\pi_c(\cdot | s_t; s_{t+1})) = -\mathbb{E}_{a_t \in \mathcal{A}} \left[\sum_{i=1}^{d_A} M^{a^i \rightarrow r} \odot \pi_c(a_t^i | s_t; s_{t+1}) \log \pi(a_t^i | s_t; s_{t+1}) \right], \quad (7)$$

where d_A is the dimension of the action space. Hence, the learning objective of the causal action empowerment can be defined as follows:

$$\begin{aligned} \mathcal{E}_{\pi_c}(s) &= \max_{a_t \sim \pi_c(\cdot | s)} \mathcal{H}(\pi_c(a_t | s_t)) - \mathcal{H}(\pi_c(a_t | s_t; s_{t+1})) \\ &= \max_{a_t \sim \pi_c(\cdot | s)} \mathbb{E}_{\pi_c(a_t | s_t) p_{\pi_c}(a_t | s_t, s_{t+1})} [\log \mathcal{P}_{\phi_c}(a_t | s_t, s_{t+1}; M) - \log \mathcal{P}_{\pi_c}(a_t | s_t; M)], \end{aligned} \quad (8)$$

where $\mathcal{P}_{\pi_c}(a_t | s_t; M)$ is the action distribution given current state of policy π_c with the causal structure, which can be denoted as $\pi_c(a_t | s_t)$. $\mathcal{P}_{\phi_c}(a_t | s_{t+1}, s_t; M)$ represents an inverse dynamics

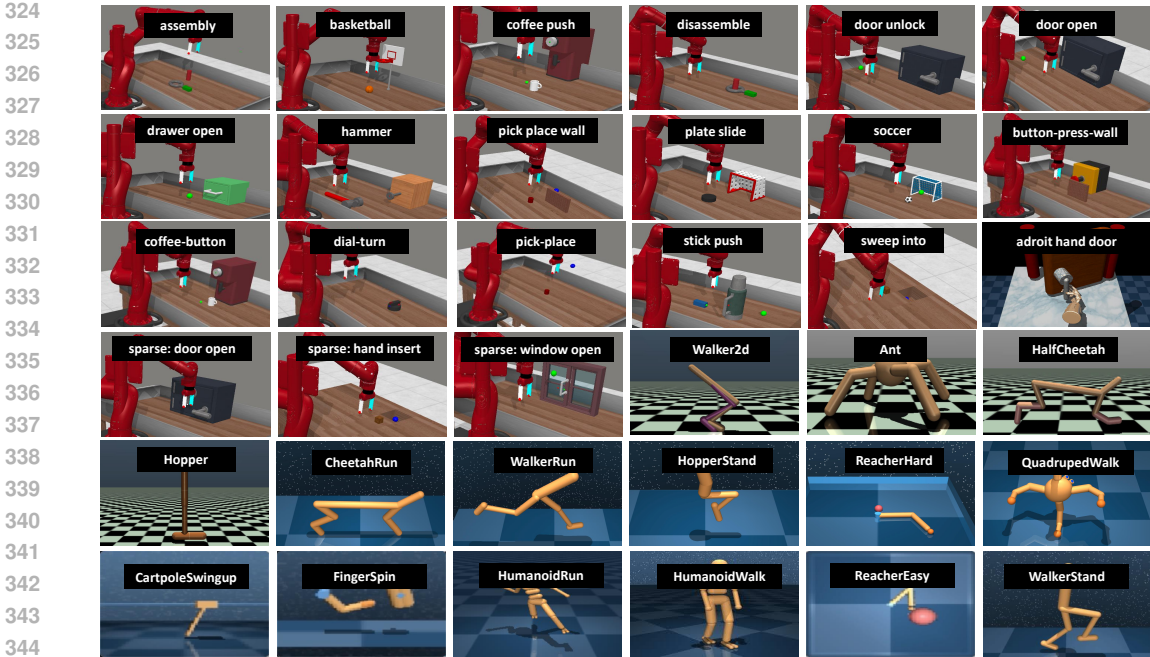


Figure 3: The 36 experimental tasks in 5 continuous control environments

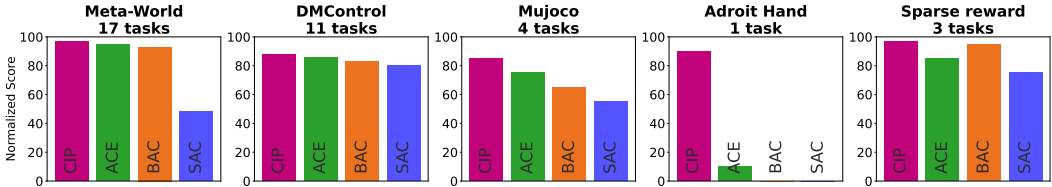


Figure 4: Experimental results with normalized score across all 36 tasks in 5 environments.

model trained on the collected transitions of state variables. Hence, we update the target policy π_c by maximizing the empowerment objective function derived in Eq. 8.

Adhering to the maximum entropy paradigm (Haarnoja et al., 2018), we calculate $\mathcal{E}_{\pi_c}(s)$ for maximization instead of standard entropy, thus prioritizing exploration of pivotal actions that are more likely to have significant causal effects on the reward. This targeted exploration strategy has the potential to accelerate learning by focusing on the most influential actions in current controllable states. Based on the causality-aware empowerment, the Q-value for policy π_c could be computed iteratively by applying a modified Bellman operator \mathcal{T}_c^π with $\mathcal{E}_{\pi_c}(s)$ term as stated below:

$$\begin{aligned} \mathcal{T}_c^\pi Q(s_t, a_t) &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_t \sim \pi_c} [Q(s_{t+1}, a_{t+1}) + \alpha \mathcal{E}_{\phi_c}(s)]] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{P}} [\mathbb{E}_{a_t \sim \pi_c} [Q(s_{t+1}, a_{t+1}) + \alpha (\mathcal{H}(\pi_c(a_t|s_t)) - \mathcal{H}(\pi_c(a_t|s_t; s_{t+1})))]]. \end{aligned} \quad (9)$$

Hence, we integrate the causality-aware empowerment term into the policy optimization objective function, $\hat{\eta}_{\mathcal{M}}(\pi_c) = \mathbb{E}_{s_0 \sim \mu_0, s_t \sim \mathcal{P}, a_t \sim \pi_c} [\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \alpha \mathcal{E}_{\pi_c}(s))]$.

In summary, **CIP** harnesses empowerment to integrate the causal understanding into decision-making. By maximizing the empowerment gain of the causally-informed policy, we guide the agent to prioritize actions that align with the environment’s underlying causal relationships. This approach enhances the agent’s exploration efficiency, focusing on actions with meaningful causal impacts and correlated with desired outcomes. Algorithm 1 illustrates the complete **CIP** pipeline.

5 EXPERIMENTS

Our experiments aim to address the following questions: (i) How does the performance of **CIP** compare to other RL approaches in diverse continuous control tasks, including manipulation and

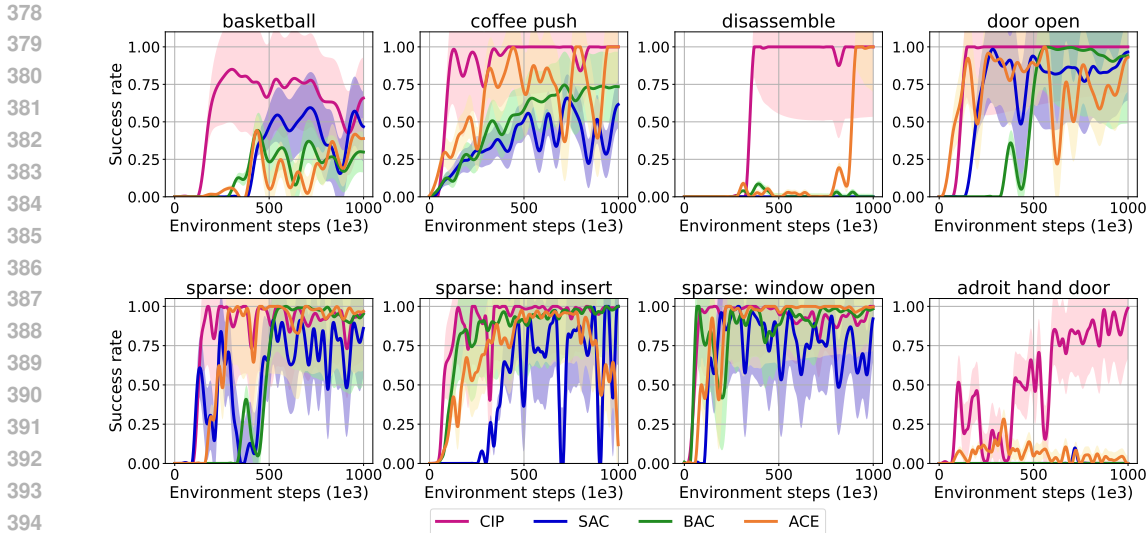


Figure 5: Experimental results of 8 manipulation skill learning tasks in Meta-World and adroit hand environments including sparse reward settings. For all tasks results, please refer to Appendix D.2.

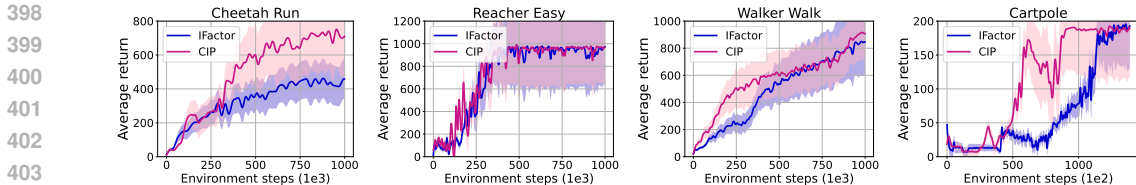


Figure 6: Experimental results of 4 pixel-based tasks in DMControl and Cartpole environments.

locomotion with sparse rewards, high-dimensional action spaces, and pixel-based challenges? (ii) Can **CIP**, through data augmentation and empowerment, improve sample efficiency and learn reliable policies? (iii) What are the effects of the components and hyperparameters in **CIP**?

5.1 EXPERIMENTAL SETUP

Environments. We evaluate **CIP** on 5 continuous control environments, including MuJoCo (Todorov et al., 2012), DMControl (Tassa et al., 2018), Meta-World (Yu et al., 2020), Adroit Hand (Rajeswaran et al., 2018), and sparse reward setting environments in Meta-World. This comprehensive evaluation encompasses 36 tasks, spanning both locomotion and manipulation skill learning, as illustrated in Figure 3. We also conduct experiments in 4 pixel-based tasks of the DMControl and Cartpole environment as shown in Figure 17. Our experimental tasks incorporate a wide range of challenges, including high-dimensional state and action spaces, sparse reward settings, pixel-based scenarios, and locomotion. For extensive experimental settings, please refer to Appendix D.1.

Baselines. We compare **CIP** with three popular RL baselines across all 36 tasks and against IFactor (Liu et al., 2024) in 3 pixel-based tasks: (1) SAC (Haarnoja et al., 2018), an off-policy actor-critic algorithm featuring maximum entropy regularization. (2) ACE (Ji et al., 2024a), a method employing causality-aware entropy regularization. (3) BAC (Ji et al., 2024b), a method that balances sufficient exploitation of past successes with exploration optimism. (4) IFactor (Liu et al., 2024), a causal framework modeling four distinct categories of latent state variables for pixel-based tasks. To ensure robustness and statistical significance, we conduct each experiment using 4 random seeds.

5.2 MAIN RESULTS

Figure 4 presents the normalized scores of **CIP** compare to other methods across 36 tasks in 5 environments. In 17 Meta-World robot-arm tasks, **CIP** achieves a near-perfect score of 100, showcasing its exceptional performance in manipulation tasks. For locomotion tasks in DMControl and MuJoCo, **CIP** consistently attains scores exceeding 80, indicating robust performance across diverse locomotion challenges. Notably, **CIP** exhibits significant performance improvements in

Table 1: The experimental results of average return in 8 locomotion tasks. We bold the best scores, and underline second-best results, \pm is the standard deviation, w/o represents without. **• indicates CIP is statistically superior to compared method (pairwise t -test at 95% confidence interval).**

Method	Ant	HalfCheetah	Hopper	Walker2d	Cheetah Run	Hopper Stand	Quadruped Walk	Reacher Hard
CIP	<u>6418\pm81</u>	12594\pm210	2846\pm882	5624\pm91	893\pm12	936\pm17	<u>948\pm54</u>	991\pm11
w/o Aug	6231 \pm 81	<u>12225\pm102</u>	2308 \pm 785	5294 \pm 41	<u>885\pm13</u>	931 \pm 22	945 \pm 35	<u>989\pm13</u>
w/o Emp	6295 \pm 210	10986 \pm 572	2270 \pm 904	<u>5547\pm91</u>	876 \pm 21	785 \pm 114	924 \pm 23	971 \pm 13
SAC	6062 \pm 105•	10888 \pm 240•	2266 \pm 981	5251 \pm 106	767 \pm 16•	936\pm8	930 \pm 19•	980 \pm 8•
BAC	6511\pm30	10276 \pm 34•	2263 \pm 1063•	3316 \pm 702•	665 \pm 6•	932 \pm 4•	962\pm24	974 \pm 16
ACE	5922 \pm 106•	9390 \pm 25•	<u>2312\pm673</u> •	4922 \pm 96•	863 \pm 23•	912 \pm 16	933 \pm 57	973 \pm 17

challenging scenarios, such as adroit hand manipulation and 3 tasks with the sparse reward setting. These results underscore the effectiveness in tackling complex, high-dimensional control problems. In next sections, we present a comprehensive analysis of **CIP**'s performance across diverse tasks.

Robot-arm manipulation. Figure 5 presents the success rates across 7 Meta-World robot-arm manipulation tasks including sparse reward settings. **CIP** consistently outperforms all other methods across these tasks, demonstrating both faster policy learning and enhanced stability. In challenging tasks, such as disassemble, **CIP** achieves an impressive 100% success rate. The effectiveness of **CIP** can be attributed to focus on causally relevant information within the state and action spaces. In sparse reward settings, the efficient extraction of causal state information and the prioritization of controllable actions enable effective task completion. By systematically eliminating noise from non-causal factors, **CIP** allows the agent to construct a more controllable and efficient policy.

High-dimensional Adroit hand manipulation. To rigorously evaluate our method's efficacy in high-dimensional tasks, we conduct comparative experiments in the Adroit Hand environment of door open task. This challenging setup involves controlling a robotic hand with up to 28 actuated degrees of freedom ($\mathcal{A} \in \mathbb{R}^{28}$). Figure 5 illustrates the success rates achieved across all methods. Notably, while the three other comparative methods fail to demonstrate significant progress on this challenging task, **CIP** achieves a near 100% success rate after 700k environment steps.

Locomotion. We further evaluate **CIP** in another important category: locomotion. The part experimental results of average return in MuJoCo and DMControl environments are presented in Table 1. Learning curves are illustrated in Figure 7. **We observe that CIP achieves the best performance in six tasks and sub-optimal in other tasks, and shows statistically significant improvements in 5 out of 8 tasks.** Moreover, compared to the traditional method SAC, **CIP** demonstrates significant performance improvements in more challenging tasks such as CheetahRun and Hopper. Compared to the causality-based method ACE, **CIP** demonstrates improvements in all tasks. Overall, in locomotion tasks, **CIP** achieves superior performance and attains high sample efficiency. **Detailed performance and statistical analyses are provided in Appendix D.2 and D.3.5.**

Pixel-based task learning To further validate the performance in pixel-based tasks, we use 3 complex pixel-based DMControl tasks for evaluation, where video backgrounds serve as distractors. We apply the proposed counterfactual data augmentation and causal action empowerment to IFactor for comparison. As shown in Figure 6, **CIP** surpasses IFactor in terms of average return. These results underscore **CIP**'s efficacy in pixel-based tasks and its capacity to better overcome spurious correlations arising from video backgrounds, focusing on locomotion. **Moreover, the result of Cartpole task in Figure 6 demonstrate the effectiveness in discrete action space environment.** For visualization trajectories in pixel-based results, please refer to Appendix D.2.4.

5.3 ANALYSIS

Ablation study. We conduct ablation experiments involving **CIP**, **CIP** without (w/o) counterfactual data augmentation (Aug), and **CIP** w/o Empowerment (Emp). The results in 8 locomotion tasks are shown in Table 1. And the learning curves of all tasks are depicted in Appendix D.3. The experimental results reveal that the variant without the empowerment learning objective performs poorly, underscoring the critical role of empowerment maximization in enhancing control capabilities.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

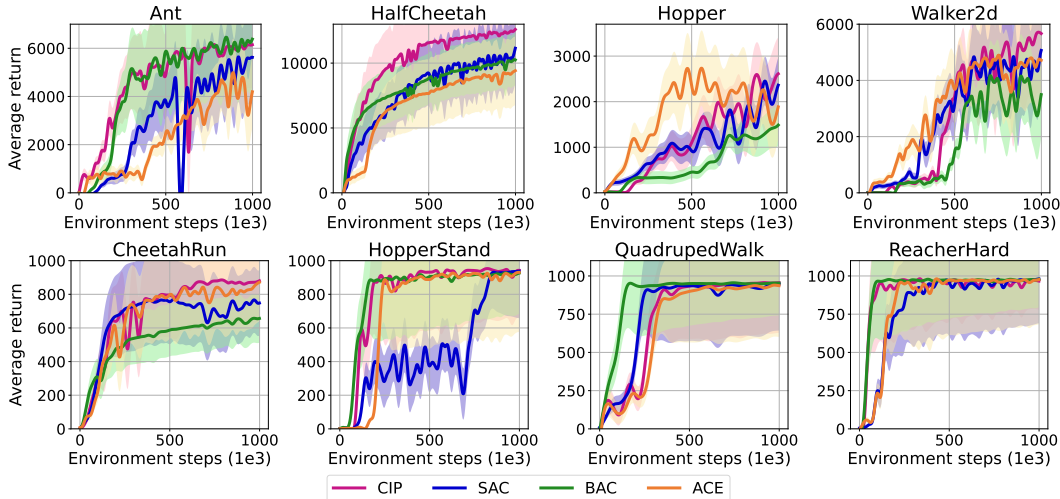


Figure 7: Experimental results with average return across 8 tasks in locomotion tasks.

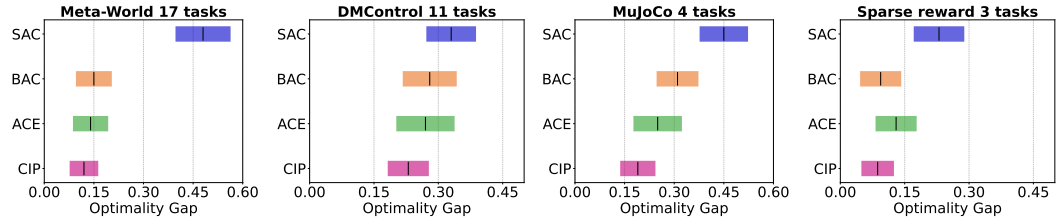


Figure 8: Experimental results of reliability evaluation by the metric Optimality Gap (lower values are better) on 4 diverse environments across 35 tasks.

Additionally, **CIP** without counterfactual data augmentation is less sample efficient than **CIP**, highlighting the importance of augmentation.

Reliability evaluation. We evaluate **CIP**’s reliability across 35 tasks in 4 environments, excluding the Adroit Hand door task due to **CIP**’s exceptional performance there. Figure 8 illustrates the experimental results using the Optimality Gap metric (Agarwal et al., 2021). **CIP** consistently achieves the lowest values across all tasks in four environments, with lower values indicating superior performance. This consistent excellence across diverse scenarios underscores the robustness and reliability of our proposed method.

6 CONCLUSION

This study introduces an efficient RL framework, designed to enhance sample efficiency. This approach begins by counterfactual data augmentation using the causality between states and rewards, effectively mitigating interference from irrelevant states without additional environmental interactions. We then develop a reward-guided structural model that leverages causal awareness to prioritize causal actions through empowerment. We conduct extensive experiments across 39 tasks spanning 5 diverse continuous control environments which demonstrate the exceptional performance of our proposed method, showcasing its robustness and adaptability across challenging scenarios.

Limitation and Future Work The current limitations of our work are twofold. First, **CIP** has not yet been extended to complex scenarios, such as real-world 3D robotics tasks. Potential approaches to address this limitation include leveraging object-centric models (Wu et al., 2023), 3D perception models (Wang et al., 2024a), and robotic foundation models (Team et al., 2024; Firoozi et al., 2023) to construct essential variables for causal world modeling. Second, **CIP** does not adequately consider non-stationarity and heterogeneity, which are critical challenges in causal discovery. Future work could integrate method designed to handle such complexities, such as CD-NOD (Huang et al., 2020).

540 REPRODUCIBILITY STATEMENT

541 We provide the core code of **CIP** in the supplementary material. The implementation details are
542 shown in Appendix D.1.

543
544
545 REFERENCES

- 546 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare.
547 Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information*
548 *processing systems*, 34:29304–29320, 2021.
- 549
550 Homanga Bharadhwaj, Mohammad Babaeizadeh, Dumitru Erhan, and Sergey Levine. Information
551 prioritization through empowerment in visual model-based rl. In *International Conference on*
552 *Learning Representations*, 2022.
- 553 Jack Brady, Roland S Zimmermann, Yash Sharma, Bernhard Schölkopf, Julius Von Kügelgen, and
554 Wieland Brendel. Provably learning object-centric representations. In *International Conference on*
555 *Machine Learning*, pp. 3038–3062. PMLR, 2023.
- 556
557 Hongye Cao, Fan Feng, Meng Fang, Shaokang Dong, Jing Huo, and Yang Gao. Towards em-
558 powerment gain through causal structure learning in model-based rl. In *ICML 2024 Workshop:*
559 *Foundations of Reinforcement Learning and Control–Connections and Perspectives*, 2024.
- 560 David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine*
561 *learning research*, 3(Nov):507–554, 2002.
- 562
563 Jongwook Choi, Archit Sharma, Honglak Lee, Sergey Levine, and Shixiang Shane Gu. Variational
564 empowerment as representation learning for goal-based reinforcement learning. *arXiv preprint*
565 *arXiv:2106.01404*, 2021.
- 566
567 Jongwook Choi, Sungtae Lee, Xinyu Wang, Sungryull Sohn, and Honglak Lee. Unsupervised object
568 interaction learning with counterfactual dynamics models. In *Proceedings of the AAAI Conference*
569 *on Artificial Intelligence*, volume 38, pp. 11570–11578, 2024.
- 570
571 ZH Deng, J Jiang, G Long, and C Zhang. Causal reinforcement learning: A survey. *Transactions on*
Machine Learning Research, 2023.
- 572
573 Carlos Diuk, Andre Cohen, and Michael L Littman. An object-oriented representation for efficient
574 reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*,
575 pp. 240–247, 2008.
- 576
577 SM Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Geoffrey E Hinton,
578 et al. Attend, infer, repeat: Fast scene understanding with generative models. *Advances in neural*
information processing systems, 29, 2016.
- 579
580 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you
581 need: Learning skills without a reward function. In *International Conference on Learning*
Representations, 2018.
- 582
583 Fan Feng and Sara Magliacane. Learning dynamic attribute-factored world models for efficient
584 multi-object reinforcement learning. *Advances in Neural Information Processing Systems*, 36,
585 2023.
- 586
587 Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu,
588 Yuke Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. Foundation models in robotics:
589 Applications, challenges, and the future. *The International Journal of Robotics Research*, pp.
02783649241281508, 2023.
- 590
591 Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. *Advances*
in neural information processing systems, 14, 2001.
- 592
593 Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms
for factored mdps. *Journal of Artificial Intelligence Research*, 19:399–468, 2003.

- 594 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
595 maximum entropy deep reinforcement learning with a stochastic actor. In *International conference*
596 *on machine learning*, pp. 1861–1870. PMLR, 2018.
- 597 Dan Haramati, Tal Daniel, and Aviv Tamar. Entity-centric reinforcement learning for object manip-
598 ulation from pixels. In *NeurIPS 2023 Workshop on Goal-Conditioned Reinforcement Learning*,
599 2023.
- 600 Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour,
601 and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of*
602 *Machine Learning Research*, 21(89):1–53, 2020.
- 603 Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, and Kun Zhang. Adarl: What, where,
604 and how to adapt in transfer reinforcement learning. In *International Conference on Learning*
605 *Representations*, 2022a.
- 606 Biwei Huang, Chaochao Lu, Liu Leqi, José Miguel Hernández-Lobato, Clark Glymour, Bernhard
607 Schölkopf, and Kun Zhang. Action-sufficient state representation learning for control with
608 structural constraints. In *International Conference on Machine Learning*, pp. 9260–9279. PMLR,
609 2022b.
- 610 Tianying Ji, Yongyuan Liang, Yan Zeng, Yu Luo, Guowei Xu, Jiawei Guo, Ruijie Zheng, Furong
611 Huang, Fuchun Sun, and Huazhe Xu. Ace: Off-policy actor-critic with causality-aware entropy
612 regularization. In *Forty-first International Conference on Machine Learning*, 2024a.
- 613 Tianying Ji, Yu Luo, Fuchun Sun, Xianyuan Zhan, Jianwei Zhang, and Huazhe Xu. Seizing serendip-
614 ity: Exploiting the value of past success in off-policy actor-critic. In *Forty-first International*
615 *Conference on Machine Learning*, 2024b.
- 616 Jindong Jiang, Sepehr Janghorbani, Gerard De Melo, and Sungjin Ahn. Scalor: Generative world mod-
617 els with scalable object representations. In *International Conference on Learning Representations*,
618 2019.
- 619 Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. *arXiv*
620 *preprint arXiv:2303.10834*, 2023.
- 621 Jindong Jiang, Fei Deng, Gautam Singh, Minseung Lee, and Sungjin Ahn. Slot state space models.
622 *arXiv preprint arXiv:2406.12272*, 2024.
- 623 Tobias Jung, Daniel Polani, and Peter Stone. Empowerment for continuous agent—environment
624 systems. *Adaptive Behavior*, 19(1):16–39, 2011.
- 625 Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*,
626 volume 16, pp. 740–747, 1999.
- 627 Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. Empowerment: A universal
628 agent-centric measure of control. In *2005 IEEE congress on evolutionary computation*, volume 1,
629 pp. 128–135. IEEE, 2005.
- 630 Avinash Kori, Francesco Locatello, Ainkaran Santhirasekaram, Francesca Toni, Ben Glocker, and
631 Fabio De Sousa Ribeiro. Identifiable object-centric representation learning via probabilistic slot
632 attention. *arXiv preprint arXiv:2406.07141*, 2024.
- 633 Adam Kosiorek, Hyunjik Kim, Yee Whye Teh, and Ingmar Posner. Sequential attend, infer, repeat:
634 Generative modelling of moving objects. *Advances in Neural Information Processing Systems*, 31,
635 2018.
- 636 Jannik Kossen, Karl Stelzner, Marcel Hussing, Claas Voelcker, and Kristian Kersting. Structured
637 object-aware physics prediction for video modeling and planning. In *International Conference on*
638 *Learning Representations*, 2019.
- 639 Oliver Kroemer, Scott Niekum, and George Konidaris. A review of robot learning for manipulation:
640 Challenges, representations, and algorithms. *Journal of machine learning research*, 22(30):1–82,
641 2021.

- 648 Sébastien Lachapelle, Divyat Mahajan, Ioannis Mitliagkas, and Simon Lacoste-Julien. Additive
649 decoders for latent variables identification and cartesian-product extrapolation. *Advances in Neural*
650 *Information Processing Systems*, 36, 2024.
- 651 Felix Leibfried, Sergio Pascual-Diaz, and Jordi Grau-Moya. A unified bellman optimality principle
652 combining reward maximization and empowerment. *Advances in Neural Information Processing*
653 *Systems*, 32, 2019.
- 654 Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. Towards practical multi-object manipu-
655 lation using relational reinforcement learning. In *2020 IEEE International Conference on Robotics*
656 *and Automation (ICRA)*, pp. 4051–4058. IEEE, 2020.
- 657 Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming
658 Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric
659 robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
660 *Pattern Recognition*, pp. 18061–18070, 2024.
- 661 Yulong Li and Deepak Pathak. Object-aware gaussian splatting for robotic manipulation. In *ICRA*
662 *2024 Workshop on 3D Visual Representations for Robot Manipulation*, 2024. URL <https://openreview.net/forum?id=gdRI43hDgo>.
- 663 Zhixuan Lin, Yi-Fu Wu, Skand Peri, Bofeng Fu, Jindong Jiang, and Sungjin Ahn. Improving
664 generative imagination in object-centric world models. In *International conference on machine*
665 *learning*, pp. 6140–6149. PMLR, 2020.
- 666 Yuren Liu, Biwei Huang, Zhengmao Zhu, Honglong Tian, Mingming Gong, Yang Yu, and Kun
667 Zhang. Learning world models with identifiable factorization. *Advances in Neural Information*
668 *Processing Systems*, 36, 2024.
- 669 Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold,
670 Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention.
671 *Advances in neural information processing systems*, 33:11525–11538, 2020.
- 672 Davide Mambelli, Frederik Träuble, Stefan Bauer, Bernhard Schölkopf, and Francesco Locatello.
673 Compositional multi-object reinforcement learning with linear relation networks. In *ICLR2022*
674 *Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022.
- 675 Chaitanya Mitash, Mostafa Hussein, Jeroen Vanbaar, Vikedo Terhuja, and Kapil Katyal. Scaling
676 object-centric robotic manipulation with multimodal object identification. 2024.
- 677 Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically
678 motivated reinforcement learning. *Advances in neural information processing systems*, 28, 2015.
- 679 Mirco Mutti, Riccardo De Santi, Emanuele Rossi, Juan Felipe Calderon, Michael Bronstein, and
680 Marcello Restelli. Provably efficient causal model-based reinforcement learning for systematic
681 generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp.
682 9251–9259, 2023.
- 683 Jongjin Park, Younggyo Seo, Chang Liu, Li Zhao, Tao Qin, Jinwoo Shin, and Tie-Yan Liu. Object-
684 aware regularization for addressing causal confusion in imitation learning. *Advances in Neural*
685 *Information Processing Systems*, 34:3029–3042, 2021.
- 686 Judea Pearl. *Causality*. Cambridge university press, 2009.
- 687 Silviu Pitis, Elliot Creager, and Animesh Garg. Counterfactual data augmentation using locally
688 factored dynamics. *Advances in Neural Information Processing Systems*, 33:3976–3990, 2020.
- 689 Silviu Pitis, Elliot Creager, Ajay Mandlekar, and Animesh Garg. Mocoda: Model-based counterfac-
690 tual data augmentation. *Advances in Neural Information Processing Systems*, 35:18143–18156,
691 2022.
- 692 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
693 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
694 learning and demonstrations. *Robotics: Science and Systems XIV*, 2018.

- 702 Sahand Rezaei-Shoshtari, Rosie Zhao, Prakash Panangaden, David Meger, and Doina Precup. Con-
703 tinuous mdp homomorphisms and homomorphic policy gradient. *Advances in Neural Information*
704 *Processing Systems*, 35:20189–20204, 2022.
- 705 Jonathan Richens and Tom Everitt. Robust agents learn causal world models. *arXiv preprint*
706 *arXiv:2402.10877*, 2024.
- 707
708 Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment—an introduction. *Guided*
709 *Self-Organization: Inception*, pp. 67–114, 2014.
- 710
711 Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain,
712 Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied
713 ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.
714 9339–9347, 2019.
- 715 Jonathan Scholz, Martin Levihn, Charles Isbell, and David Wingate. A physics-based model prior for
716 object-oriented mdps. In *International Conference on Machine Learning*, pp. 1089–1097. PMLR,
717 2014.
- 718
719 Maximilian Seitzer, Bernhard Schölkopf, and Georg Martius. Causal influence detection for improv-
720 ing efficiency in reinforcement learning. *Advances in Neural Information Processing Systems*, 34:
721 22905–22918, 2021.
- 722 Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvarinen, Yoshinobu Kawahara,
723 Takashi Washio, Patrik O Hoyer, Kenneth Bollen, and Patrik Hoyer. Directlingam: A direct
724 method for learning a linear non-gaussian structural equation model. *Journal of Machine Learning*
725 *Research-JMLR*, 12(Apr):1225–1248, 2011.
- 726
727 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
728 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without
729 human knowledge. *nature*, 550(7676):354–359, 2017.
- 730 Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press,
731 2001.
- 732
733 Hao Sun and Taiyi Wang. Toward causal-aware rl: State-wise action-refined temporal difference. In
734 *Deep Reinforcement Learning Workshop NeurIPS 2022*.
- 735 Yüewen Sun, Erli Wang, Biwei Huang, Chaochao Lu, Lu Feng, Changyin Sun, and Kun Zhang.
736 Acamda: Improving data efficiency in reinforcement learning through guided counterfactual data
737 augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp.
738 15193–15201, 2024.
- 739 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 740
741 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,
742 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint*
743 *arXiv:1801.00690*, 2018.
- 744
745 Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep
746 Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot
747 policy. *arXiv preprint arXiv:2405.12213*, 2024.
- 748 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
749 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033.
750 IEEE, 2012.
- 751 Núría Armengol Urpí, Marco Bagatella, Marin Vlastelica, and Georg Martius. Causal action
752 influence aware counterfactual data augmentation. In *Forty-first International Conference on*
753 *Machine Learning*, 2024.
- 754
755 Elise Van der Pol, Thomas Kipf, Frans A Oliehoek, and Max Welling. Plannable approximations to
mdp homomorphisms: Equivariance under actions. *arXiv preprint arXiv:2002.11963*, 2020.

- 756 Arthur Wandzel, Yoonseon Oh, Michael Fishman, Nishanth Kumar, Lawson LS Wong, and Stefanie
757 Tellex. Multi-object search using object-oriented pomdps. In *2019 International Conference on*
758 *Robotics and Automation (ICRA)*, pp. 7194–7200. IEEE, 2019.
- 759
- 760 Chenxi Wang, Hongjie Fang, Hao-Shu Fang, and Cewu Lu. Rise: 3d perception makes real-world
761 robot imitation simple and effective. *arXiv preprint arXiv:2404.12281*, 2024a.
- 762
- 763 Zizhao Wang, Xuesu Xiao, Yuke Zhu, and Peter Stone. Task-independent causal state abstraction.
764 In *Proceedings of the 35th International Conference on Neural Information Processing Systems,*
765 *Robot Learning workshop*, 2021.
- 766
- 767 Zizhao Wang, Xuesu Xiao, Zifan Xu, Yuke Zhu, and Peter Stone. Causal dynamics learning for task-
768 independent state abstraction. In *International Conference on Machine Learning*, pp. 23151–23180.
769 PMLR, 2022.
- 770
- 771 Zizhao Wang, Jiaheng Hu, Caleb Chuck, Stephen Chen, Roberto Martín-Martín, Amy Zhang, Scott
772 Niekum, and Peter Stone. Skild: Unsupervised skill discovery guided by factor interactions. In
773 *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b.
- 774
- 775 Zizhao Wang, Jiaheng Hu, Peter Stone, and Roberto Martín-Martín. Elden: exploration via local
776 dependencies. *Advances in Neural Information Processing Systems*, 36, 2024c.
- 777
- 778 Zizhao Wang, Caroline Wang, Xuesu Xiao, Yuke Zhu, and Peter Stone. Building minimal and
779 reusable causal state abstractions for reinforcement learning. *arXiv preprint arXiv:2401.12497*,
780 2024d.
- 781
- 782 Nicholas Watters, Loic Matthey, Matko Bosnjak, Christopher P Burgess, and Alexander Lerchner.
783 Cobra: Data-efficient model-based rl through unsupervised object discovery and curiosity-driven
784 exploration. *arXiv preprint arXiv:1905.09275*, 2019.
- 785
- 786 Alicia Peregrin Wolfe and Andrew G Barto. Defining object types and options using mdp homomor-
787 phisms. In *Proceedings of the ICML-06 Workshop on Structural Knowledge Transfer for Machine*
788 *Learning*, 2006.
- 789
- 790 Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. Slotformer: Unsupervised
791 visual dynamics simulation with object-centric models. In *The Eleventh International Conference*
792 *on Learning Representations*, 2022.
- 793
- 794 Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. Slotdiffusion: Object-centric
795 generative modeling with diffusion models. *Advances in Neural Information Processing Systems*,
796 36:50932–50958, 2023.
- 797
- 798 Ziyi Wu, Yulia Rubanova, Rishabh Kabra, Drew A Hudson, Igor Gilitschenski, Yusuf Aytar, Sjoerd
799 van Steenkiste, Kelsey R Allen, and Thomas Kipf. Neural assets: 3d-aware multi-object scene
800 synthesis with image diffusion models. *arXiv preprint arXiv:2406.09292*, 2024.
- 801
- 802 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
803 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
804 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- 805
- 806 Wentao Yuan, Chris Paxton, Karthik Desingh, and Dieter Fox. Sornet: Spatial object-centric
807 representations for sequential manipulation. In *Conference on Robot Learning*, pp. 148–157.
808 PMLR, 2022.
- 809
- 810 Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Self-supervised visual reinforcement
811 learning with object-centric representations. In *International Conference on Learning Representa-*
812 *tions*, 2021.
- 813
- 814 Andrii Zadaianchuk, Georg Martius, and Fanny Yang. Self-supervised reinforcement learning with
815 independently controllable subgoals. In *Conference on Robot Learning*, pp. 384–394. PMLR,
816 2022.

810 Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Object-centric learning for real-
811 world videos by predicting temporal feature similarities. In *Thirty-seventh Conference on Neural*
812 *Information Processing Systems (NeurIPS 2023)*, 2023.

813 Yanjie Ze, Yuyao Liu, Ruizhe Shi, Jiabin Qin, Zhecheng Yuan, Jiashun Wang, and Huazhe Xu.
814 H-index: Visual reinforcement learning with hand-informed representations for dexterous manipu-
815 lation. *Advances in Neural Information Processing Systems*, 36, 2024.

816 Yan Zeng, Ruichu Cai, Fuchun Sun, Libo Huang, and Zhifeng Hao. A survey on causal reinforcement
817 learning. *arXiv preprint arXiv:2302.05209*, 2023.

818 Yudi Zhang, Yali Du, Biwei Huang, Ziyang Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy.
819 Interpretable reward redistribution in reinforcement learning: a causal approach. *Advances in*
820 *Neural Information Processing Systems*, 36, 2024.

821 Linfeng Zhao, Lingzhi Kong, Robin Walters, and Lawson LS Wong. Toward compositional general-
822 ization in object-oriented world modeling. In *International Conference on Machine Learning*, pp.
823 26841–26864. PMLR, 2022.

824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864	CONTENTS	
865		
866	1 Introduction	1
867		
868	2 Related Work	3
869		
870	2.1 Causal RL	3
871		
872	2.2 Empowerment in RL	3
873		
874	2.3 Object-centric RL and Object-Oriented RL	3
875		
876	3 Preliminaries	4
877		
878	3.1 Markov Decision Process	4
879		
880	3.2 Structural Causal Model	4
881		
882	3.3 Empowerment in RL	4
883		
884	4 Causal Information Prioritization	5
885		
886	4.1 Counterfactual Data Augmentation	5
887		
888	4.2 Causal Action Prioritization Through Empowerment	6
889		
890	5 Experiments	7
891		
892	5.1 Experimental setup	8
893		
894	5.2 Main Results	8
895		
896	5.3 Analysis	9
897		
898	6 Conclusion	10
899		
900	A Broader Impact	19
901		
902	B Assumptions and Propositions	19
903		
904	C Extended Related Work	20
905		
906	C.1 Extended Discussion on object-centric RL and 3D world models	20
907		
908	D Details on Experimental Design and Results	21
909		
910	D.1 Experimental setup	21
911		
912	D.2 Full Results	21
913		
914	D.2.1 Effectiveness in robot arm manipulation	21
915		
916	D.2.2 Effectiveness in spare reward settings	22
917		
	D.2.3 Effectiveness in locomotion	22
	D.2.4 Effectiveness in pixel-based tasks	27
	D.3 Property Analysis	27
	D.3.1 Analysis for replacing counterfactual data augmentation	27
	D.3.2 Extensive ablation study	28

918	D.3.3	Hyperparameter analysis	29
919	D.3.4	Computation cost analysis	30
920	D.3.5	Statistical performance analysis	31
921	D.3.6	Generalization analysis	31
922	D.3.7	Causal discovery analysis	32
923			
924			
925			
926	E	Details on the Proposed Framework	36
927			
928	F	Experimental Platforms and Licenses	36
929			
930	F.1	Experimental platforms	36
931	F.2	Licenses	36
932			
933			
934			
935			
936			
937			
938			
939			
940			
941			
942			
943			
944			
945			
946			
947			
948			
949			
950			
951			
952			
953			
954			
955			
956			
957			
958			
959			
960			
961			
962			
963			
964			
965			
966			
967			
968			
969			
970			
971			

972 A BROADER IMPACT

973
974 To avoid blind exploration and improve sample efficiency, we propose **CIP** for efficient reinforcement
975 learning. **CIP** leverages the causal relationships among states, actions, and rewards to prioritize causal
976 information for efficient policy learning. **CIP** first learns a causal matrix between states and rewards
977 to execute counterfactual data augmentation, prioritizing important state features without additional
978 environmental interactions. Subsequently, it learns a causal reweight matrix between actions and
979 rewards to prioritize causally-informed behaviors. We then introduce a causal action empowerment
980 term into the learning objective to enhance the controllability. By prioritizing the causal information,
981 **CIP** enables agents to focus on behaviors that have causally significant effects on their tasks. **CIP**
982 offers substantial broader impact by prioritizing causal information through individual assessment
983 of how different factors contribute to rewards. Our novel empowerment learning objective achieves
984 efficient policy optimization by leveraging entropy via the policy and learned inverse dynamics model.
985 This approach shows promise for extension into research frameworks centered on maximum entropy
986 algorithms.

987 Despite its strengths, **CIP** has limitations beyond its reliance on the method DirectLiNGAM. There’s
988 potential to explore alternative causal discovery techniques for more robust relationship mapping.
989 Moreover, analyzing inter-entity causal connections could lead to better disentanglement of diverse
990 behaviors. Our future work will investigate a range of causal discovery methods to refine our approach.
991 We aim to extend **CIP** to model-based RL frameworks, focusing on building causal world models to
992 enhance generalization.

993 B ASSUMPTIONS AND PROPOSITIONS

994
995 **Assumption 1** (*d-separation (Pearl, 2009)*) *d-separation is a graphical criterion used to determine,*
996 *from a given causal graph, if a set of variables X is conditionally independent of another set Y , given*
997 *a third set of variables Z . In a directed acyclic graph (DAG) \mathcal{G} , a path between nodes n_1 and n_m is*
998 *said to be blocked by a set S if there exists a node n_k , for $k = 2, \dots, m - 1$, that satisfies one of the*
999 *following two conditions:*

1000
1001 (i) $n_k \in S$, and the path between n_{k-1} and n_{k+1} forms $(n_{k-1} \rightarrow n_k \rightarrow n_{k+1})$, $(n_{k-1} \leftarrow n_k \leftarrow$
1002 $n_{k+1})$, or $(n_{k-1} \leftarrow n_k \rightarrow n_{k+1})$.

1003 (ii) Neither n_k nor any of its descendants is in S , and the path between n_{k-1} and n_{k+1} forms
1004 $(n_{k-1} \rightarrow n_k \leftarrow n_{k+1})$.

1005 In a DAG, we say that two nodes n_a and n_b are *d-separated* by a third node n_c if every path between
1006 nodes n_a and n_b is blocked by n_c , denoted as $n_a \perp\!\!\!\perp n_b | n_c$.

1007
1008 **Assumption 2** (*Global Markov Condition (Spirtes et al., 2001; Pearl, 2009)*) *The state is fully*
1009 *observable and the dynamics is Markovian. The distribution p over a set of variables $\mathcal{V} =$*
1010 *$(s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^d, r_t)^T$ satisfies the global Markov condition on the graph if for any parti-*
1011 *tion $(\mathcal{S}, \mathcal{A}, \mathcal{R})$ in \mathcal{V} such that if \mathcal{A} *d-separates* \mathcal{S} from \mathcal{R} , then $p(\mathcal{S}, \mathcal{R} | \mathcal{A}) = p(\mathcal{S} | \mathcal{A}) \cdot p(\mathcal{R} | \mathcal{A})$*

1012
1013 **Assumption 3** (*Faithfulness Assumption (Spirtes et al., 2001; Pearl, 2009)*) *For a set of variables*
1014 *$\mathcal{V} = (s_t^1, \dots, s_t^d, a_t^1, \dots, a_t^d, r_t)^T$, there are no independencies between variables that are not*
1015 *implied by the Markovian Condition.*

1016
1017 **Assumption 4** *Under the assumptions that the causal graph is Markov and faithful to the observa-*
1018 *tions, the edge $s_t^i \rightarrow s_{t+1}^i$ exists for all state variables s^i .*

1019 **Assumption 5** *No simultaneous or backward edges in time.*

1020
1021 **Proposition 1** *Under the assumptions that the causal graph is Markov and faithful to the observa-*
1022 *tions, there exists an edge from $a_t^i \rightarrow r_t$ if and only if $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$.*

1023 *Proof:* We proceed by proving both directions of the if and only if statement.

1024
1025 (\Rightarrow) Suppose there exists an edge from a_t^i to r_t . We prove that $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$ by contradiction.
Assume $a_t^i \perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$. By the faithfulness assumption, this independence must be reflected in

the graph structure. However, this implies the absence of a directed path from a_t^i to r_t , contradicting the existence of the edge. Thus, $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$.

(\Leftarrow) Now, suppose $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$. We prove the existence of an edge from a_t^i to r_t by contradiction. Assume no such edge exists. By the Markov assumption, the absence of this edge implies $a_t^i \perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$, contradicting our initial supposition. Therefore, an edge from a_t^i to r_t must exist. Thus, we have shown that an edge from a_t^i to r_t exists if and only if $a_t^i \not\perp\!\!\!\perp r_t | a_t \setminus a_t^i, s_t$, completing the proof.

Proposition 2 *Under the assumptions that the causal graph is Markov and faithful to the observations, there exists an edge from $s_t^i \rightarrow r_t$ if and only if $s_t^i \not\perp\!\!\!\perp r_t | \{a_t, s_t \setminus r_t\}$.*

The proof of Proposition 2 follows a similar line of reasoning as that of Proposition 1.

Theorem 1 *Based on above 5 assumptions and 2 propositions, suppose s_t, a_t, s_t follow the factored MDP reward function Eq. 2, the causal matrices $M^{s \rightarrow r}$ and $M^{a \rightarrow r}$ are identifiable.*

C EXTENDED RELATED WORK

We categorize existing causal RL approaches based on problem domains and task types, providing a systematic analysis of how different methods explore causal relationships between states, actions, and rewards, as illustrated in Table 2.

In the single-task learning domain, methods such as ACE (Ji et al., 2024a) and IFactor (Liu et al., 2024) have shown success in learning policies for manipulation and locomotion tasks. However, both approaches are limited by focusing on a single reward-guided causal relationship. Regarding generalization, AdaRL (Huang et al., 2022a) effectively leverages both state-reward and action-reward causal relationships. However, AdaRL focuses primarily on applying causal inference to address generalization challenges in locomotion tasks. Its application is limited to locomotion tasks, leaving more complex manipulation tasks unaddressed. Since our work focuses on the single-task problem domain, we do not provide a direct comparison with AdaRL. Conversely, CBM (Wang et al., 2024d) considers the causal relationship between states and rewards but overlooks the causal link between actions and rewards. In the problem domain of counterfactual data augmentation, current causal RL methods (Urpí et al., 2024; Pitis et al., 2020; 2022) have not yet explored the inference and utilization of both causal relationships.

In summary, current research on reward-guided causal discovery remains incomplete and lacks validation across a broader spectrum of tasks. This gap underscores the need for more comprehensive investigation and application in the field of causal reinforcement learning.

C.1 EXTENDED DISCUSSION ON OBJECT-CENTRIC RL AND 3D WORLD MODELS

The main similarity lies between our framework and object-centric RL is both are learning and using factored MDPs (Kearns & Koller, 1999), but they differ in granularity: our framework operates at the component level (e.g., raw state variables), whereas object-centric RL factors states based on objects.

Although our work is orthogonal to object-centric RL, we believe certain elements of object-centric RL could complement our framework in specific applications, particularly in real-world robotic manipulation tasks. Potential future work include:

- **Using object-centric representation as input:** Object-centric models can help identify object-factored variables, such as object attributes, geometry, and physical states, which are useful for planning (Jiang et al., 2019; Lin et al., 2020; Kossen et al., 2019; Mambelli et al., 2022; Feng & Magliacane, 2023; Choi et al., 2024; Zadaianchuk et al., 2022; Park et al., 2021; Zadaianchuk et al., 2021; Yuan et al., 2022; Li et al., 2020; Mitash et al., 2024; Haramati et al., 2023; Li et al., 2024). In this case, states are factored as objects, and we can learn causal graphs over these variables. This is useful in robotic environments involving numerous objects. We will leave this as a future work for adapting our current framework to the applications of the object-centric robotic task.

Table 2: Categorization of different causal RL methods with two different causal relationship of state-to-reward (state-reward) and action-to-reward (action-reward).

Problem domain	Task type	Method	Causal relationship	
			state-reward	action-reward
Single-task	manipulation; locomotion	ACE (Ji et al., 2024a)	✗	✓
	manipulation; locomotion	IFactor (Liu et al., 2024)	✓	✗
	manipulation	CAI (Seitzer et al., 2021)	✗	✗
Generalization	manipulation	CDL (Wang et al., 2022)	✗	✗
	locomotion	AdaRL (Huang et al., 2022a)	✓	✓
	manipulation; locomotion	CBM (Wang et al., 2024d)	✓	✗
Augmentation	manipulation	CAIAC (Urpí et al., 2024)	✗	✗
	manipulation	CoDA (Pitis et al., 2020)	✗	✗
	manipulation	MoCoDA (Pitis et al., 2022)	✗	✗

- **Learning more compact factored object representations with our framework:** Our structure learning approach could benefit object-centric RL by disentangling the internal representations of individual objects to the reward-relevant and reward-irrelevant groups by learning the causal structures. This can enhance the compactness and interpretability of object-centric representations.
- **Using object-aware 3D world models for applications:** In 3D environments, object-aware 3D world models (Li & Pathak, 2024) can provide essential representations of objects. Our framework could then build causal structures on top of these factored 3D-object representations.

While these directions are promising and could advance the applicability of our framework in certain domains, they are outside the primary focus of this work. We plan to explore these ideas as part of future work.

D DETAILS ON EXPERIMENTAL DESIGN AND RESULTS

D.1 EXPERIMENTAL SETUP

We present the detailed hyperparameter settings of the proposed method **CIP** across all 5 environments in Table 3. Additionally, the Q-value and V-value networks are used MLP with 512 hidden size. And the policy network is the Gaussian MLP with 512 hidden size. Moreover, we set the target update interval of 2. For fair comparison, the hyperparameters of the baseline methods (SAC (Haarnoja et al., 2018), BAC (Ji et al., 2024b), ACE (Ji et al., 2024a)) follow the same settings in the experiments.

For pixel-based DMControl environments, we employ IFactor (Liu et al., 2024) to encode latent states and integrate the **CIP** framework for policy learning. We utilize the $s_t^{\bar{}}$ state features in IFactor as uncontrollable states unrelated to rewards to execute counterfactual data augmentation. Furthermore, for simplicity, we maximize the mutual information between future states and actions to facilitate empowerment. All parameter settings in these three tasks adhere to those specified in IFactor. Additionally, We use the same background video for the comparison.

D.2 FULL RESULTS

D.2.1 EFFECTIVENESS IN ROBOT ARM MANIPULATION

Figure 9 presents the learning curves for all 17 manipulation skill tasks within the Meta-World environment. The **CIP** framework demonstrates superior learning outcomes and efficiency compared to the three baseline methods, despite exhibiting minor instabilities in the basketball and dial-turn tasks. Notably, **CIP** achieves a 100% success rate in more complex tasks, such as pick-place-wall and assembly. The visualization results presented in Figures 11 and 12 further demonstrate **CIP**'s

Table 3: Hyperparameter settings of **CIP** in 5 environments

Hyperparameter	Environment				
	Meta-World	Sparse	MuJoCo	DMControl	Adroit Hand
batch size	512	512	256	512	256
hidden size	1024	1024	256	1024	256
Q-value network hidden size			512		
V-value network hidden size			512		
policy network hidden size			512		
learning step			1000000		
replay size			1000000		
causal sample size			10000		
gamma			0.99		
learning rate			0.0003		
update interval			2		

ability to effectively and efficiently complete tasks, even in high-dimensional action spaces such as the Adroit Hand environment.

In the hammer task, **CIP** allows the robot arm to execute reach and pick actions with precision, enabling it to accurately identify the nail’s position and successfully perform the hammering action. In the Adroit Hand door task, **CIP** effectively controls the complex joints to grasp the doorknob and applies the appropriate force to twist it, thereby opening the door.

These findings affirm the effectiveness of **CIP** in robot arm manipulation skill learning, highlighting its capacity to enhance sample efficiency while mitigating the risks associated with blind exploration.

Visualization. We employ trajectory visualization to comparatively validate the efficacy of our method. As depicted in Figure 10, the light-shaded regions delineate the policy exploration space, while the point clustering area indicates the area of frequent interaction. Our analysis reveals that **CIP**, leveraging counterfactual data augmentation, achieves substantially broader exploration compared to ACE and SAC. Concurrently, the causal information prioritization framework facilitates more focused execution in critical state regions. These visual findings provide robust empirical support for the effectiveness of our proposed augmentation framework.

D.2.2 EFFECTIVENESS IN SPARE REWARD SETTINGS

Figure 13 presents the learning curves for all three sparse reward setting tasks within the Meta-World environment, while Figure 14 showcases their corresponding visualization trajectories. These findings reveal that **CIP** not only achieves superior learning efficiency but also adeptly executes critical actions necessary for task completion, such as opening the door and window and maneuvering the node to the target place.

These results substantiate the effectiveness of **CIP** in sparse reward scenarios. The counterfactual data augmentation process prioritizes salient state information, effectively filtering out irrelevant factors that could hinder learning. Meanwhile, causal action empowerment enhances policy controllability by focusing on actions that are causally linked to desired outcomes. This dual approach not only accelerates the learning process but also fosters a more robust policy capable of navigating the complexities inherent in sparse reward settings. Overall, these findings underscore **CIP**’s potential to significantly improve performance in challenging environments characterized by limited feedback.

D.2.3 EFFECTIVENESS IN LOCOMOTION

We further evaluate **CIP** in 15 locomotion tasks in DMControl and MuJoCo environments. Figure 15 presents the learning curves, while Figure 16 showcases the corresponding visualization trajectories in 4 specific tasks. A comprehensive analysis indicates that **CIP** achieves faster learning efficiency and greater stability compared to ACE and SAC, while demonstrating comparable policy learning

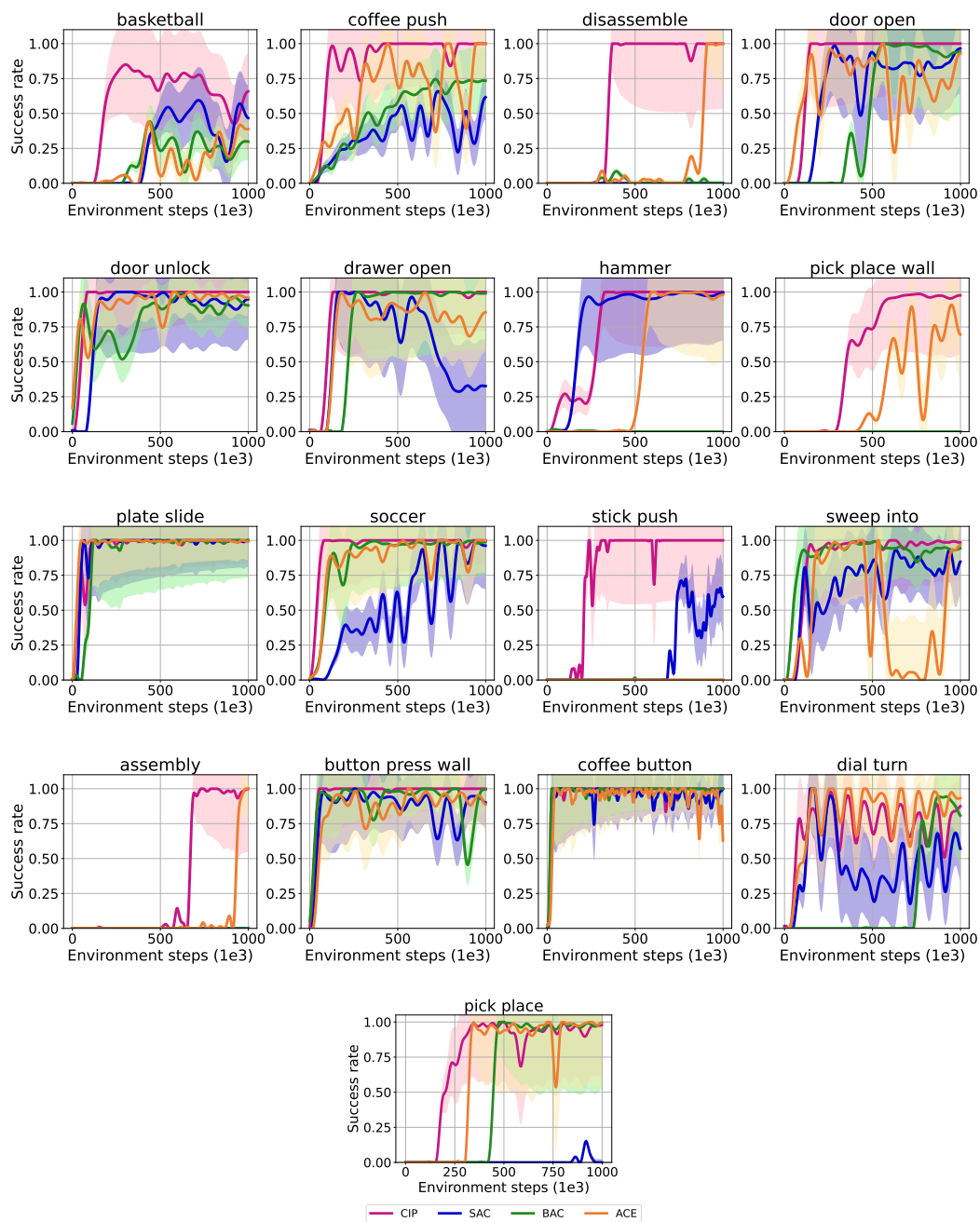
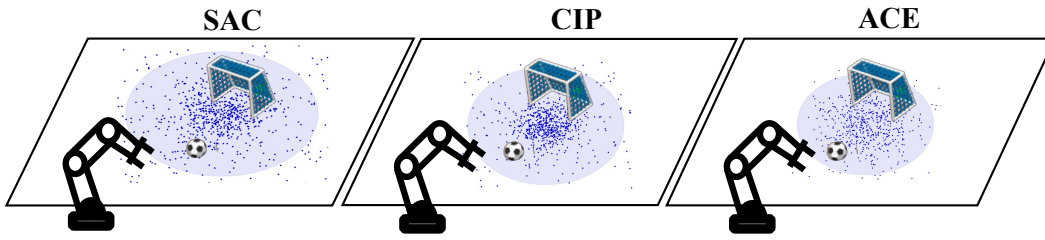


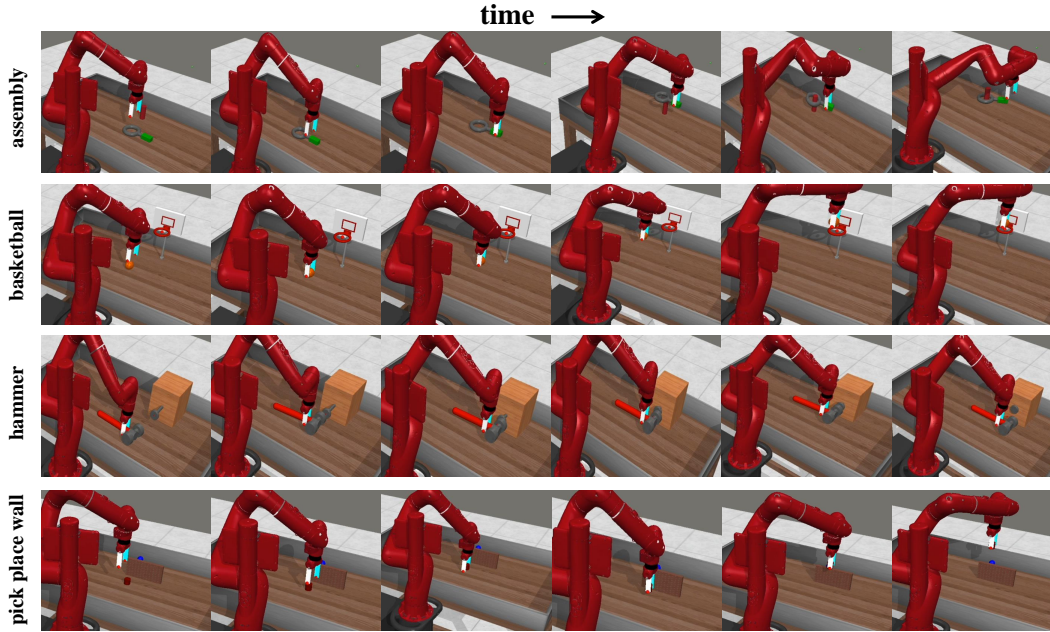
Figure 9: Experimental results across 17 manipulation skill learning tasks in Meta-World.

performance to BAC, which is known for its proficiency in control tasks. The visualization results reveal that **CIP** effectively executes running and walking actions in complex humanoid scenarios.

These findings collectively underscore the efficacy of **CIP** in locomotion tasks, highlighting its potential to advance the state-of-the-art in reinforcement learning for intricate motor control problems. The method’s success across varied environments suggests a robust framework that could generalize effectively to other challenging domains within robotics and control systems, paving the way for future research and applications in these areas.



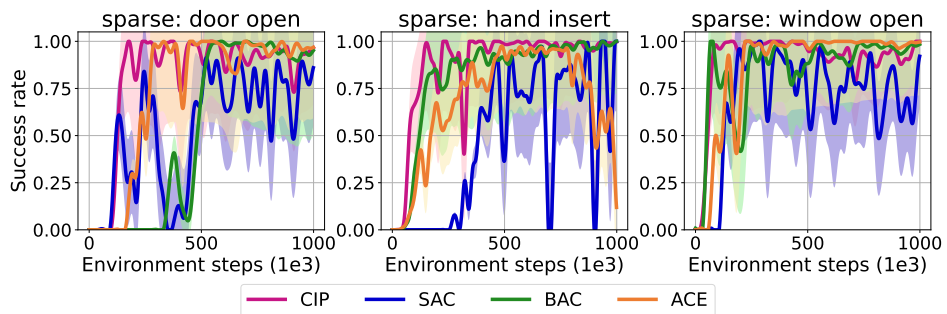
1250 Figure 10: Visualization of the trajectories in soccer task.



1273 Figure 11: Visualization trajectories of 4 manipulation skill learning tasks in Meta-World environment.



1281 Figure 12: Visualization trajectory of Adroit Hand door open task.



1294 Figure 13: Experimental results across 3 manipulation skill learning tasks in sparse reward settings
1295 of Meta-World environment.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

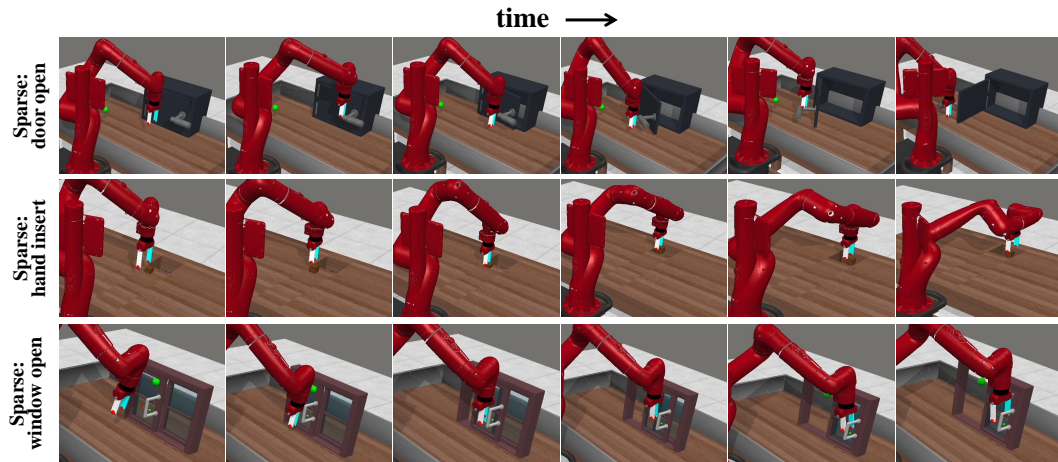


Figure 14: Visualization trajectories of 3 manipulation skill learning tasks in sparse reward settings of Meta-World environment.

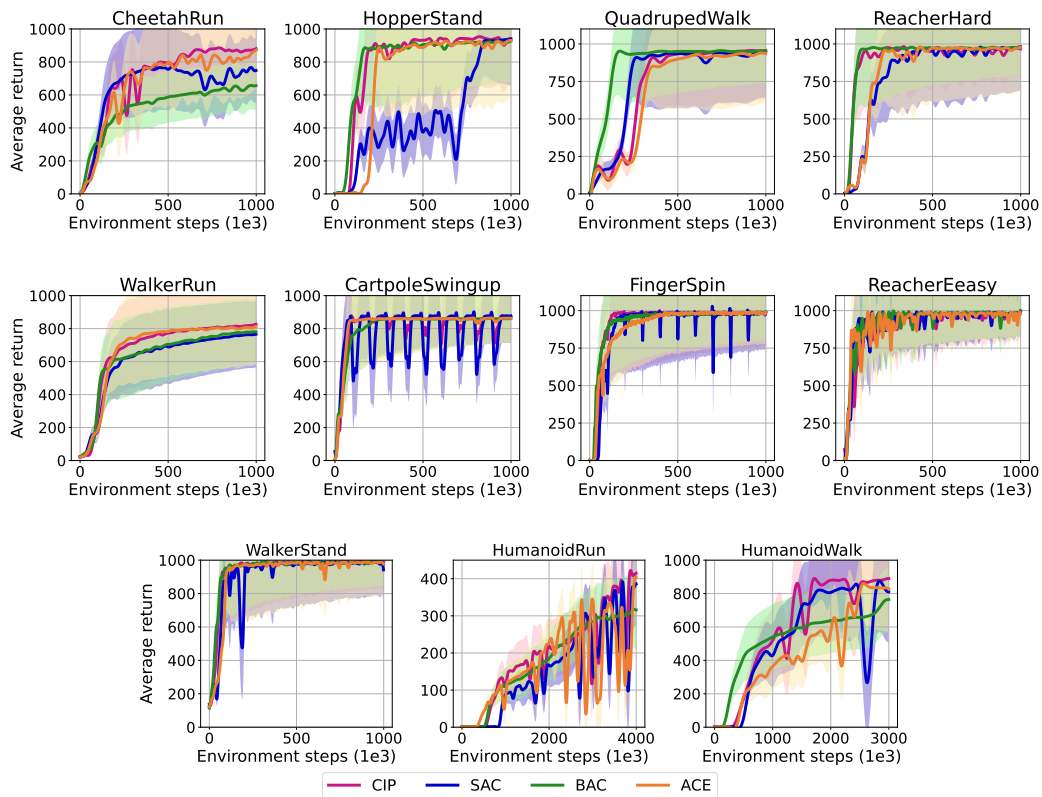


Figure 15: Experimental results across 11 locomotion tasks in DMControl environment.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

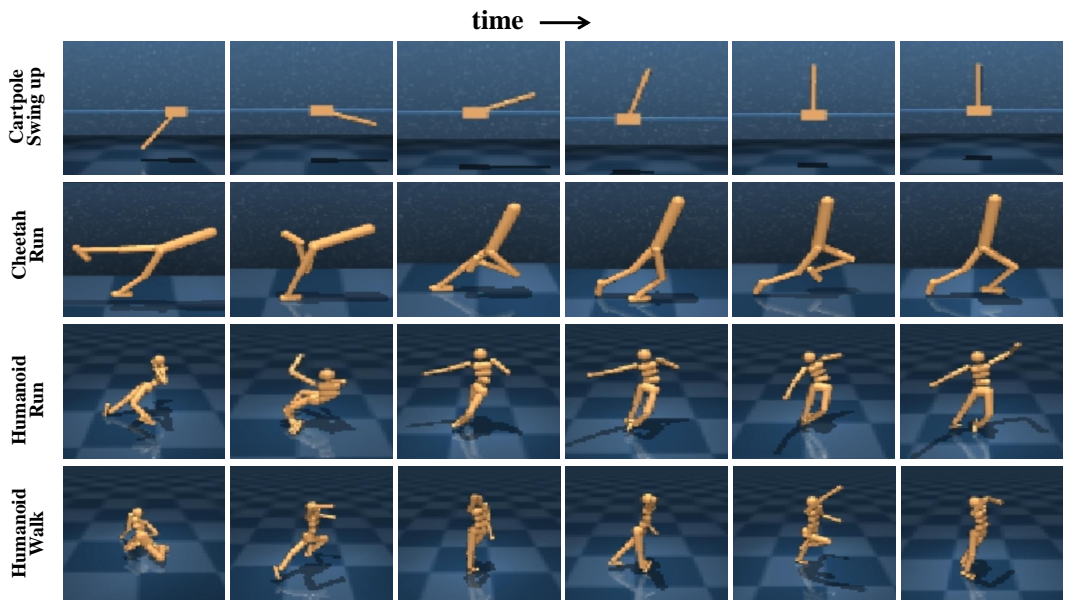


Figure 16: Visualization trajectories of 4 locomotion tasks in DMControl environment.

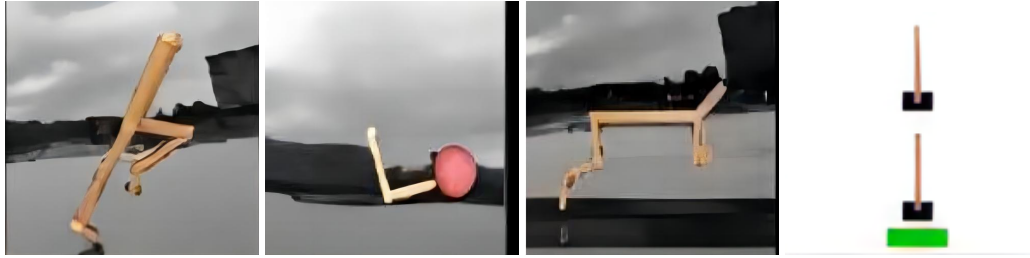


Figure 17: The DMControl environment of 3 pixel-based tasks (Walker Walk, Cheetah Run, Reacher Easy) and 1 task in Cartpole environment (Liu et al., 2024).

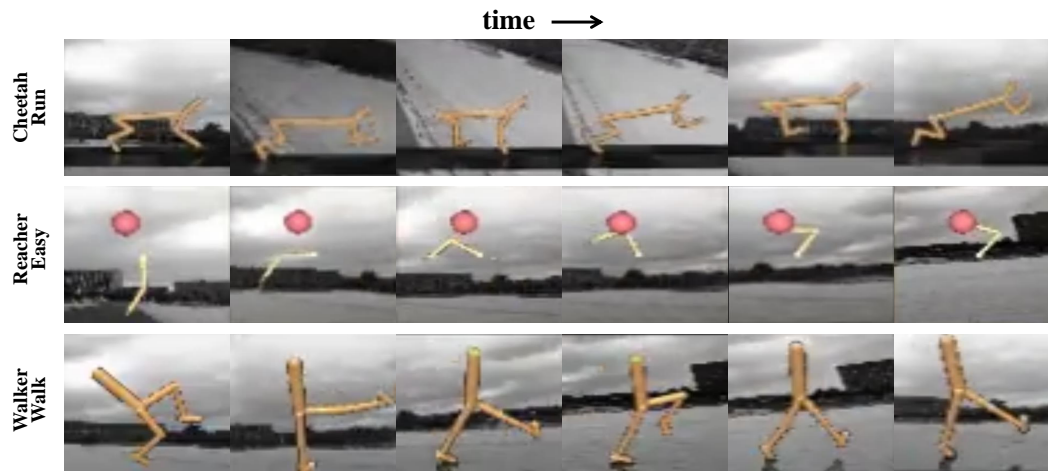


Figure 18: Visualization trajectories in 3 pixel-based locomotion tasks of DMControl environment with video backgrounds as distractors.

D.2.4 EFFECTIVENESS IN PIXEL-BASED TASKS

To further validate the effectiveness of our proposed framework in pixel-based environments, we evaluated **CIP** on three DMControl pixel-based tasks. We leverage IFactor for latent state processing and differentiation of uncontrollable state features to execute counterfactual data augmentation, alongside maximizing the mutual information between future states and actions for empowerment.

Figure 6 presents the learning curves, while Figure 18 shows the visualization trajectories. The proposed framework exhibits enhanced policy learning performance and effectively mitigates interference from background video, facilitating efficient locomotion. These findings reinforce the effectiveness and extensibility of our causal information prioritization framework, highlighting its potential to improve learning in complex, pixel-based environments.

D.3 PROPERTY ANALYSIS

D.3.1 ANALYSIS FOR REPLACING COUNTERFACTUAL DATA AUGMENTATION

In **CIP**, we exploit the causal relationship between states and rewards to perform counterfactual data augmentation on irrelevant state features, thus prioritizing critical state information. We compare this approach with an alternative method: masking irrelevant state features to achieve state abstraction for subsequent causal action empowerment and policy learning. To evaluate the efficacy of both approaches, we conduct experiments with **CIP** with counterfactual data augmentation (**CIP** w/i Cda) and **CIP** with causally-informed states (**CIP** w/i Cs) across three distinct environments.

Figure 19 illustrates comparative results for four manipulation skill learning tasks in the Meta-World environment. Both **CIP** variants achieve 100% task success rates with high sample efficiency, validating their effectiveness. Notably, **CIP** w/i Cda exhibits superior learning efficiency compared to

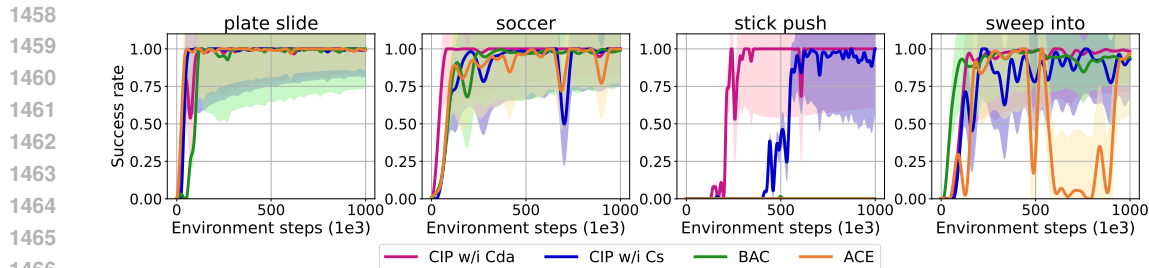


Figure 19: Experimental results in 4 manipulation skill learning tasks of Meta-World environment. w/i stands for with.

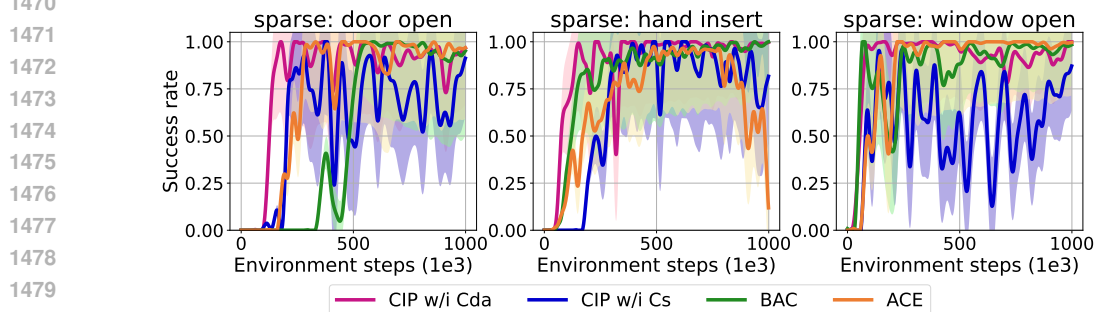


Figure 20: Experimental results in 3 manipulation skill learning tasks of Meta-World environment with sparse reward settings.

CIP w/i Cs, underscoring the value of our counterfactual data augmentation approach in enhancing training data without additional environmental interactions. In three sparse reward setting tasks (Figure 20), **CIP w/i Cda** demonstrates superior policy performance. Further experiments across four locomotion environmental tasks corroborate these findings, consistently favoring the counterfactual data augmentation approach. These comprehensive experimental results strongly support the effectiveness and significance of incorporating counterfactual data augmentation in **CIP**, highlighting its potential to enhance reinforcement learning across diverse task domains.

D.3.2 EXTENSIVE ABLATION STUDY

Robot arm manipulation The ablation study results in the Meta-World and Adroit Hand environments are presented in Figure 22. The findings indicate that **CIP** without counterfactual data augmentation exhibits reduced learning efficiency and is unable to successfully complete tasks such as pick-and-place. This underscores the importance of incorporating counterfactual data augmentation, which prioritizes causal state information, to enhance learning efficiency by mitigating the influence of irrelevant state information and preventing policy divergence.

Furthermore, **CIP** without causal action empowerment demonstrates a significant decline in policy performance across robot arm manipulation tasks. In complex scenarios, such as Adroit Hand door opening and assembly, it fails to learn effective strategies for task completion. This outcome further corroborates the efficacy of the proposed causal action empowerment mechanism, as prioritizing causally informed actions facilitates more efficient exploration of the environment, ultimately enabling successful policy learning.

Sparse reward settings Figure 22 presents the results of the ablation study conducted across three sparse reward setting tasks. These findings underscore the substantial influence of causal action empowerment on the efficacy of policy learning, demonstrating its critical role in enhancing performance in challenging environments. Additionally, the incorporation of counterfactual data augmentation proves effective in mitigating the need for additional environmental interactions, thereby significantly improving sample efficiency. This approach not only facilitates more rapid learning but

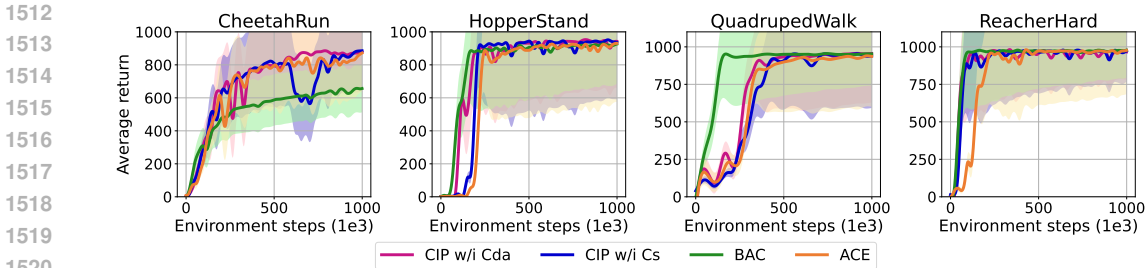


Figure 21: Experimental results in 4 locomotion tasks of DMControl environment.

also ensures that the agent can effectively navigate sparse reward scenarios by focusing on the most relevant causal information.

Locomotion We further conducted ablation experiments on locomotion tasks. The experimental results in the MuJoCo environment are shown in Figure 23, where it is evident that the performance of **CIP** without causal action empowerment declines significantly. Similarly, **CIP** without counterfactual data augmentation also exhibits reduced learning efficiency. Notably, in the 11 DMControl tasks, the decline in performance for **CIP** without causal action empowerment is particularly pronounced.

These experimental results further validate the effectiveness of our proposed method, which systematically analyzes the causal relationships between states, actions, and rewards. This analysis enables the execution of counterfactual data augmentation to avoid interference from irrelevant factors while prioritizing important state information. Subsequently, by leveraging the causal relationships between actions and rewards, we reweight actions to prioritize causally informed actions, thereby enhancing the agent’s controllability and overall learning efficacy.

D.3.3 HYPERPARAMETER ANALYSIS

We conduct a detailed analysis of the hyperparameters associated with the causal update interval (I) and sample size within the **CIP** framework. The experimental results for four distinct tasks are illustrated in Figure 25. Across all tasks, **CIP** demonstrates optimal performance with a causal update interval of $I = 2$ and a sample size of 10,000.

Our findings suggest that while a reduction in the causal update interval can lead to improved performance, it may also result in heightened computational costs. Additionally, we observe that higher update frequencies and increased sample sizes introduce greater instability, which significantly raises computational demands. This analysis underscores the importance of carefully balancing hyperparameter settings to optimize both performance and efficiency within the **CIP**.

Furthermore, we analyze the performance under different settings of the temperature factor α proposed in Eq. 9. The results across 3 tasks are shown in Figure 26. Our analysis reveals that **CIP** demonstrates robust performance across different values of α in manipulation tasks, while showing some instability in locomotion tasks when α is either too small or too large. Moreover, we observe that setting α to 0.2 yields optimal performance across all tasks, which motivated our choice of $\alpha = 0.2$ for all experiments.

Finally, we analyze the performance under different settings of the batch size and hidden size. The results across 3 tasks are shown in Figure 27. Our experimental results demonstrate that **CIP** exhibits robust performance across various parameter settings in coffee push and sparse hand insert tasks, while maintaining strong performance in hopper stand task. Based on these experimental results, we configure the hyperparameters as follows: for manipulation tasks, we set the batch size to 512 and hidden size to 1024, while for locomotion tasks, we use a batch size of 256 and hidden size of 256. All other hyperparameters remain constant across all tasks, as detailed in Table 3.

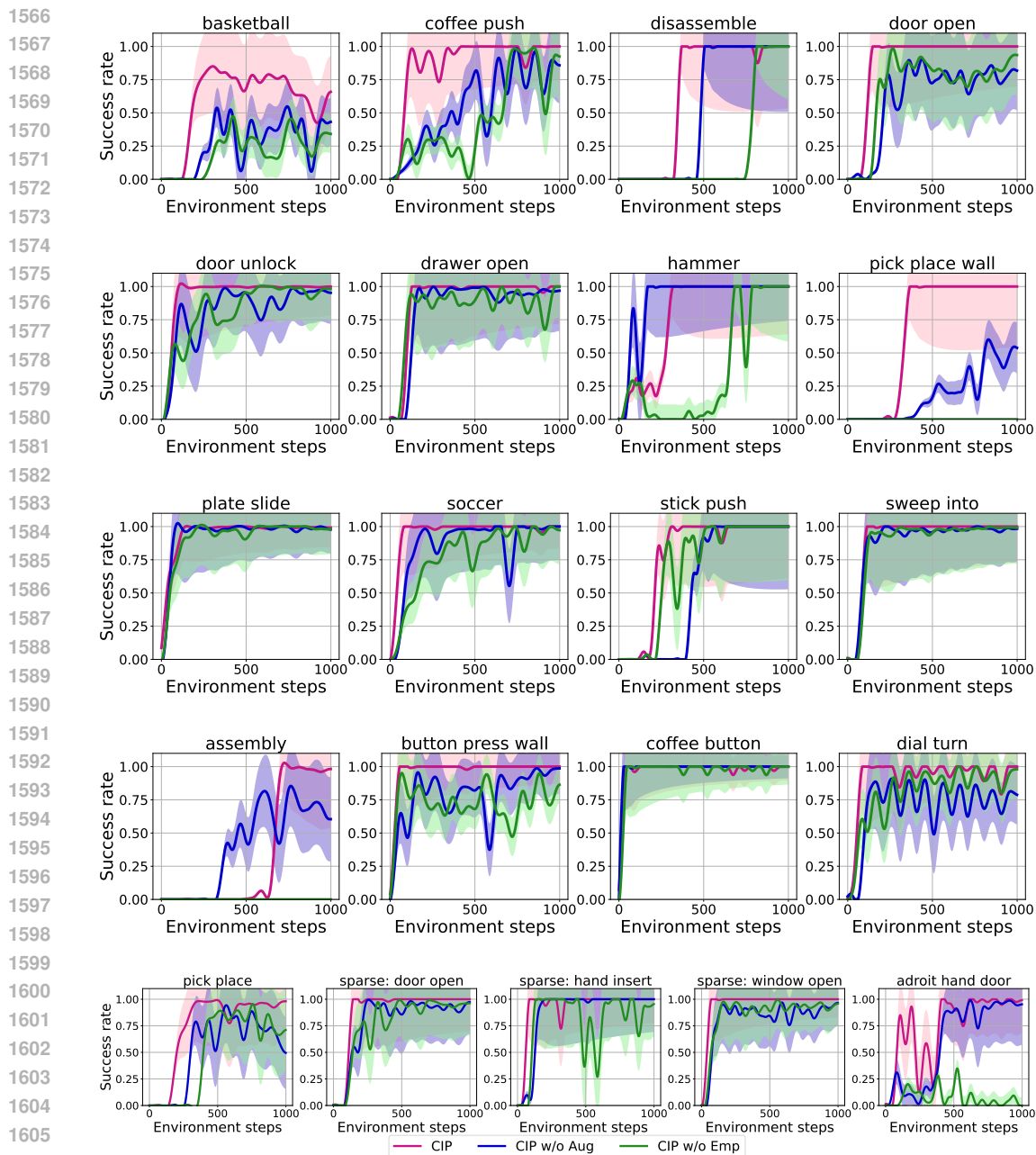


Figure 22: Ablation results across 21 manipulation skill learning tasks in Meta-World including sparse reward settings and adroit hand.

D.3.4 COMPUTATION COST ANALYSIS

We analyze the computational cost of the proposed framework. The computation time for all methods across 36 tasks is shown in Figure 28. Our experimental results demonstrate that CIP achieves its performance improvements with minimal additional computational burden - specifically less than 10% increase compared to SAC, less than 5% increase compared to ACE, and actually requiring less computation time than BAC. All experiments were conducted on the same computing platform with the same computational resources detailed in Appendix F. These experimental results verify that our proposed method achieves performance improvements without incurring significant additional computational costs.

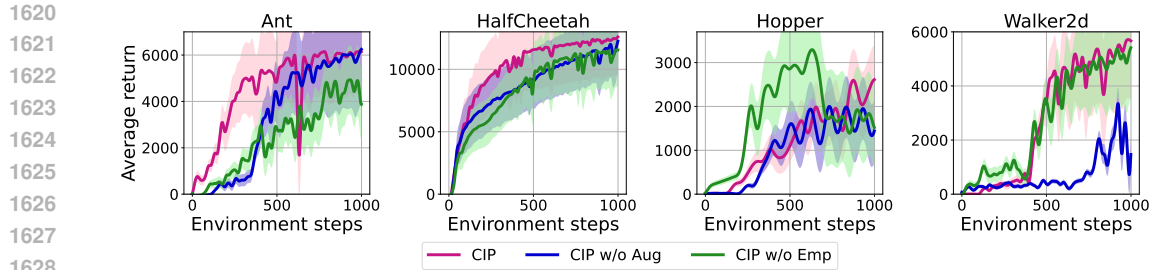


Figure 23: Ablation results across 4 locomotion tasks in MuJoCo environment.

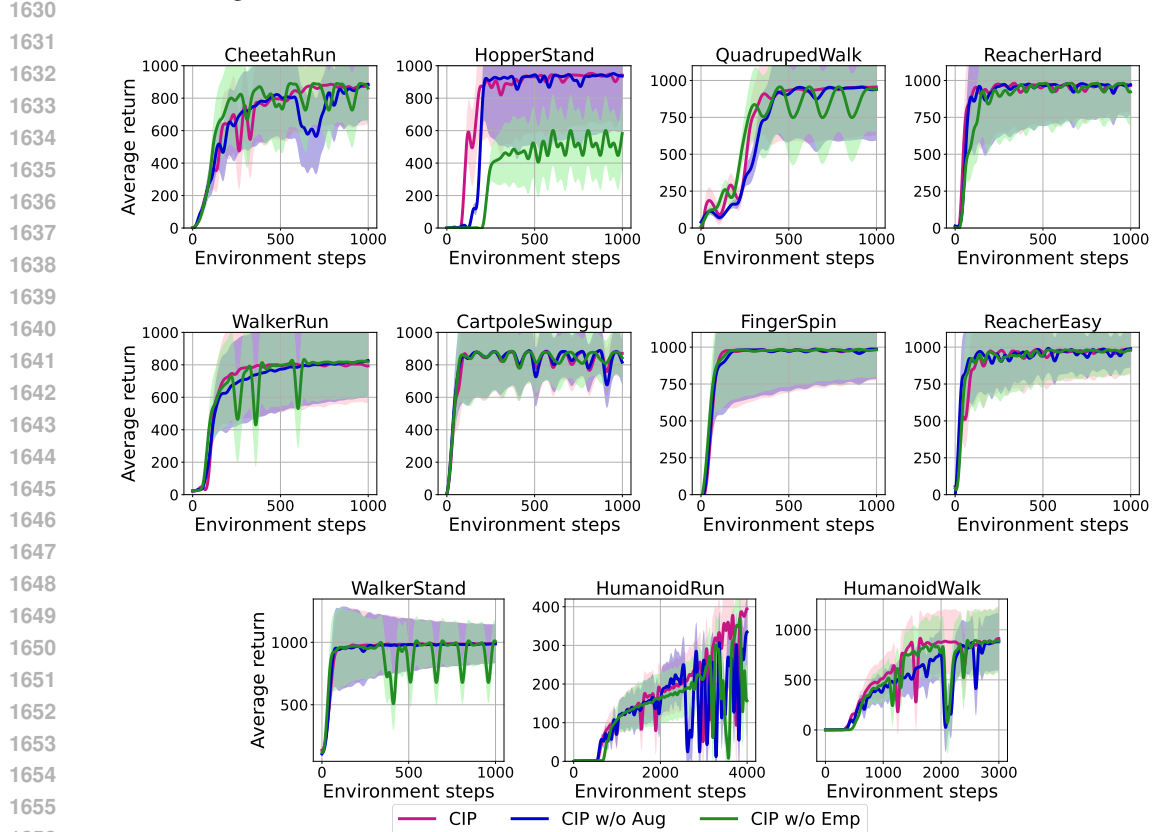


Figure 24: Ablation results across 11 locomotion tasks in DMControl environment.

1657
1658
1659
1660
1661

1662 D.3.5 STATISTICAL PERFORMANCE ANALYSIS

1663 To further validate the statistical significance of the performance, we select 3 statistical metrics (Agarwal et al., 2021) - IQM, Mean, and Median - for analysis across 8 locomotion tasks. The results are shown in Figure 29 and 30. Our findings indicate that **CIP** achieves notably superior performance across all tasks, with the sole exception of the ant task where it performs slightly below BAC.

1668

1669 D.3.6 GENERALIZATION ANALYSIS

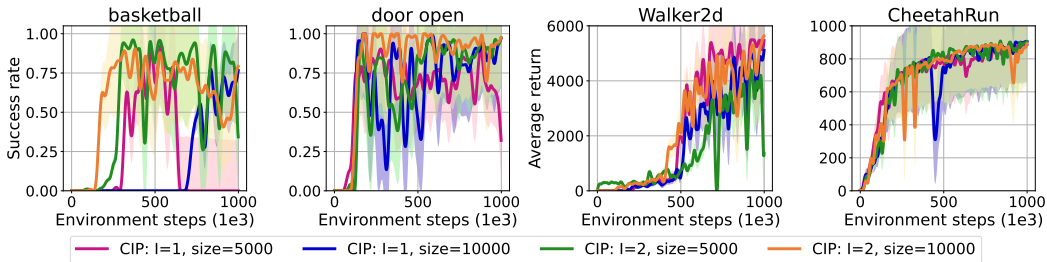
1670 We conduct multi-task experiments in the Meta-World environment (Yu et al., 2020) to validate the generalizability. We establish MT1 and MT10 tasks for generalization validation:

1671 **Multi-Task 1 (MT1):** Learning one multi-task policy that generalizes to 5 tasks belonging to the same environment. MT1 uses single Meta-World environments, with the training “tasks” corresponding to

1672

1673

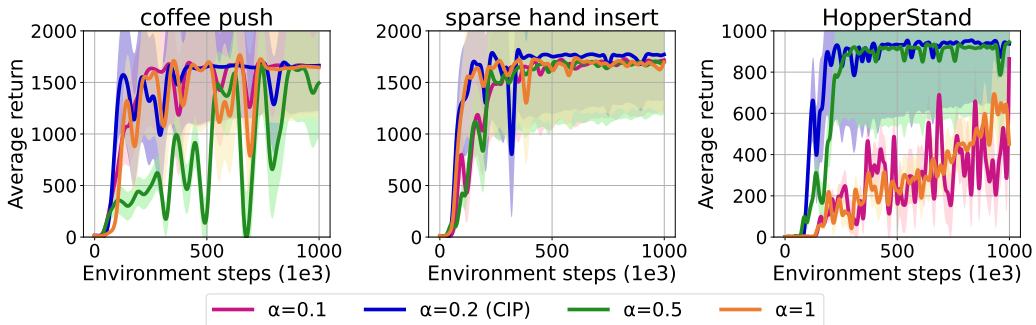
1674
1675
1676
1677
1678
1679
1680
1681
1682



1683
1684
1685

Figure 25: Hyperparameter study. Learning curves of **CIP** with different hyperparameter settings. The shaded regions are the standard deviation of each policy.

1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696



1697

Figure 26: Hyperparameter analysis of temperature factor α across 3 task.

1698
1699
1700
1701
1702
1703

5 random initial object and goal positions. The goal positions are provided in the observation and are a fixed set, as to focus on the ability of algorithms in acquiring a distinct skill across multiple goals, rather than generalization and robustness.

1704
1705
1706
1707
1708
1709
1710

Multi-Task 10 (MT10): This task involves learning a single multi-task policy that generalizes to 50 tasks across 10 training environments, totaling 500 training tasks. A crucial step towards rapid adaptation to distinctly new tasks is the ability to train a single policy capable of solving multiple distinct training tasks. The multi-task evaluation in Meta-World tests the ability to learn multiple tasks simultaneously, without accounting for generalization to new tasks. The MT10 evaluation encompasses 10 environments: reach, push, pick and place, open door, open drawer, close drawer, press button top-down, insert peg side, open window, and close window.

1711
1712
1713
1714
1715
1716
1717
1718

We adapt our proposed **CIP** to multi-task learning by incorporating a one-hot task ID as input, comparing **MT-CIP** with MT-SAC. The results in Figure 31 show that **MT-CIP** outperforms MT-SAC in both MT1 (soccer) and MT10 tasks, achieving average success rates above 50% and 40% respectively. Notably, **MT-CIP** exhibits strong performance in specific MT10 tasks like drawer close and window open. The superior performance of **MT-CIP** stems from its effective learning of causal information, enabling robust task transfer across diverse domains. While these results are promising, future work will focus on causal state abstraction for enhanced generalization and sample efficiency. All experiments were conducted under the same hyperparameter settings, and the implementation will be made publicly available.

1719
1720
1721

D.3.7 CAUSAL DISCOVERY ANALYSIS

1722
1723
1724
1725
1726
1727

In **CIP**, we use the linear causal discovery method DirectLiNAM for causal structure learning. To explore alternative approaches, we compare it with two other causal discovery methods: score-based GES (Chickering, 2002) and constraint-based PC (Spirtes et al., 2001). The experimental results in Figure 32 across three tasks demonstrate that our chosen DirectLiNAM method exhibits superior performance compared to both alternatives. During experimentation, we also observe that both GES and PC methods incur significant computational overhead and frequently encounter memory constraints. In contrast, our proposed method **CIP**, which is fundamentally reward-guided, efficiently

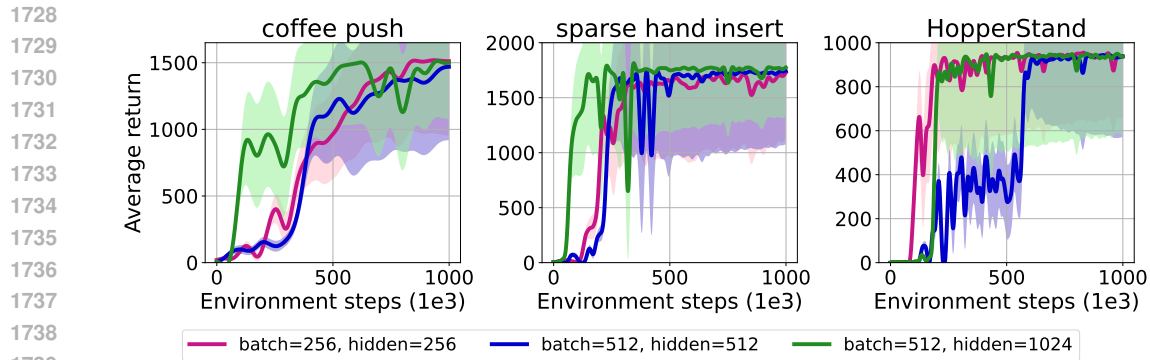


Figure 27: Hyperparameter analysis of batch size and hidden size across 3 task.

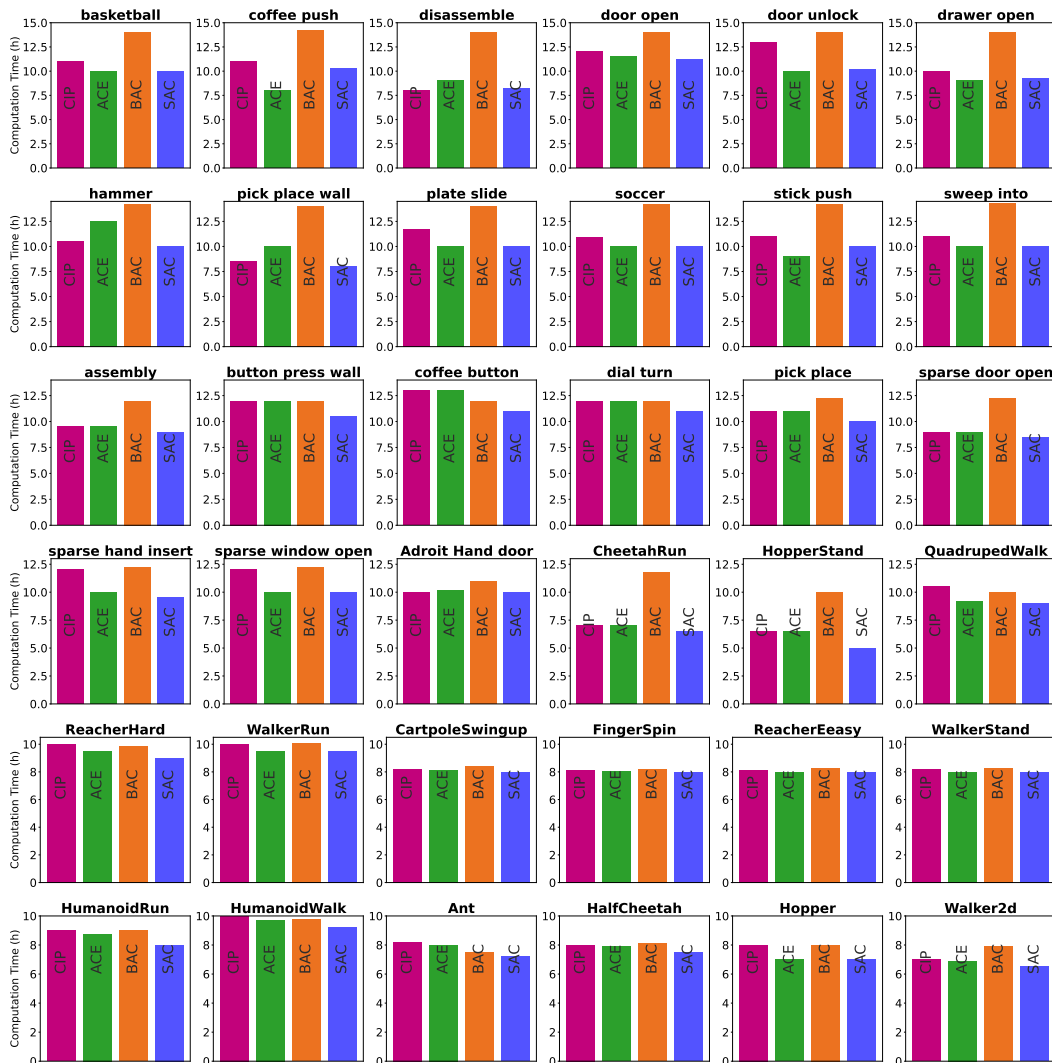


Figure 28: Computation time in 36 tasks.

discovers causal relationships between dimensional factors in states and actions with respect to rewards. This approach better aligns with the requirements of policy learning while maintaining minimal computational costs.

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

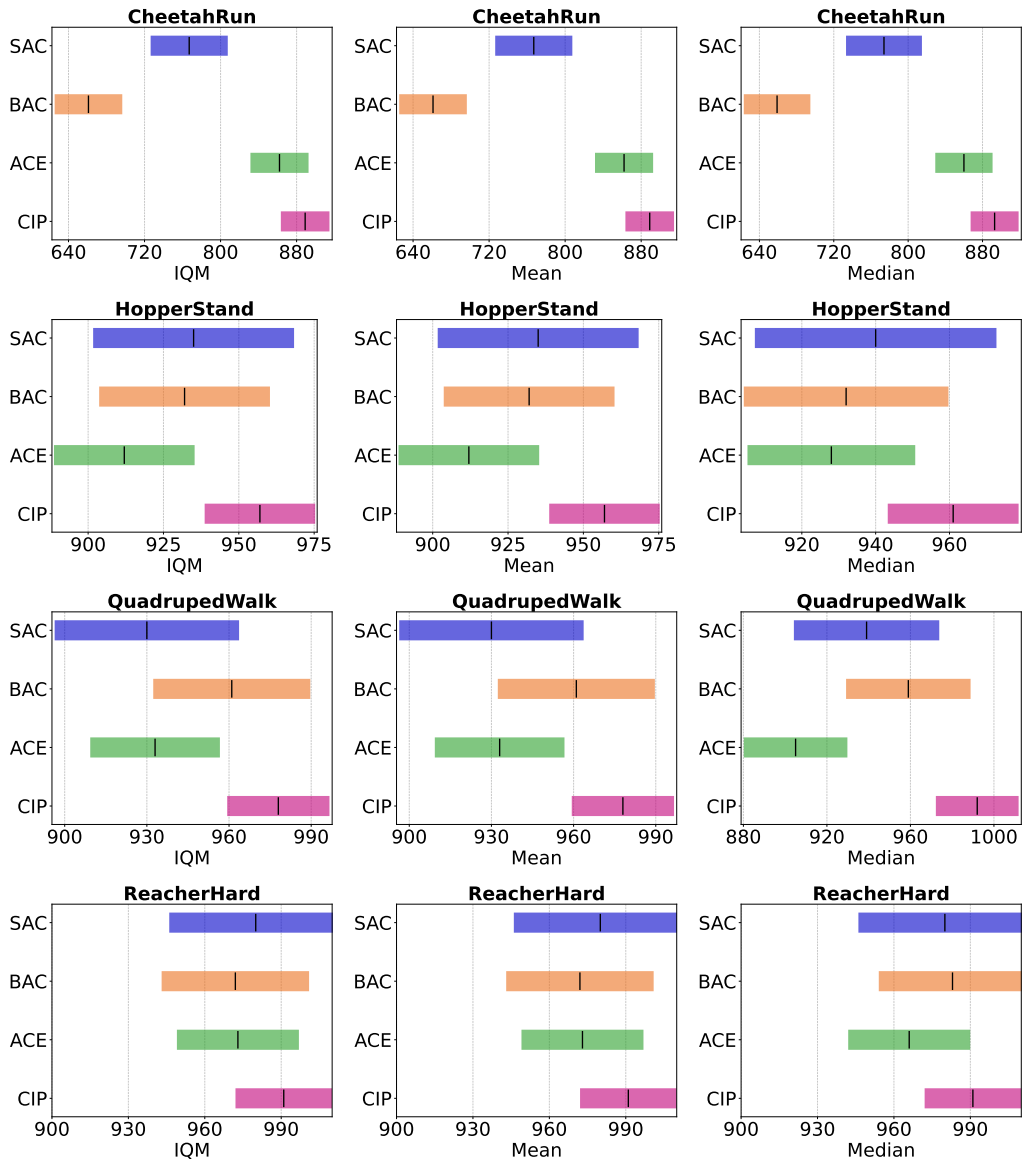


Figure 29: statistical metrics of IQM, Mean, and Median (higher values are better) on 4 DMControl tasks.

1836

1837

1838

1839

1840

1841

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

1866

1867

1868

1869

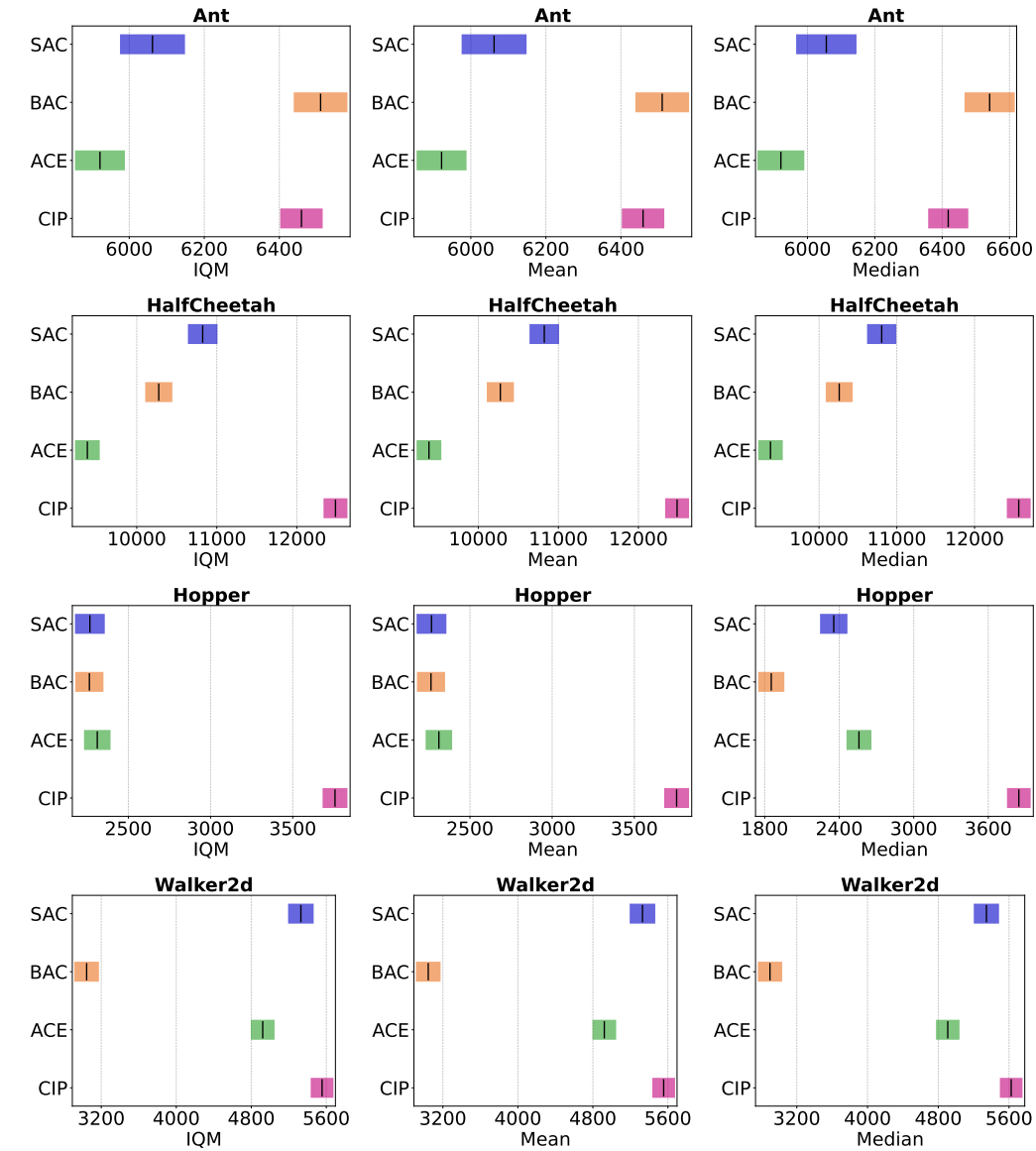
1870

1871

1872

1873

1874



1875 Figure 30: statistical metrics of IQM, Mean, and Median (higher values are better) on 4 MuJoCo
1876 tasks.

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

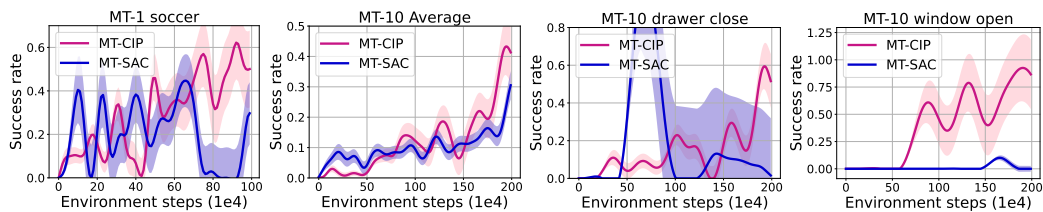


Figure 31: Generalization results in MT1 and MT10 tasks.

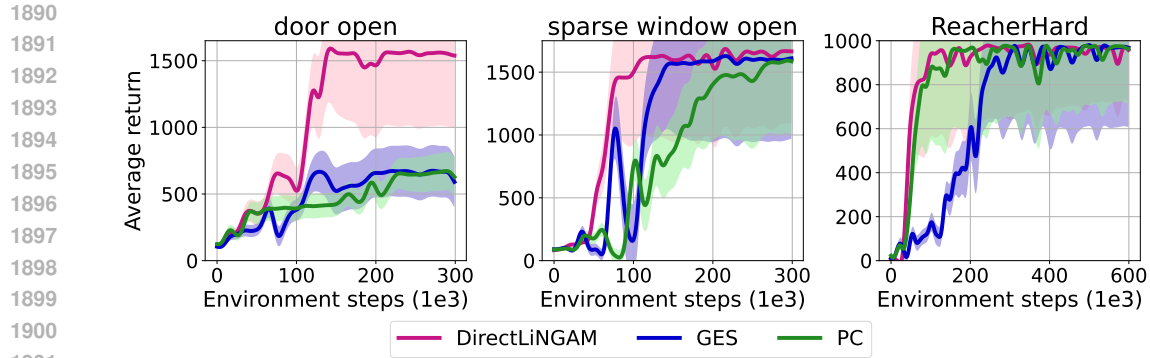


Figure 32: Compared performance with 2 different causal discovery methods across 3 task.

E DETAILS ON THE PROPOSED FRAMEWORK

Algorithm 1 lists the full pipeline of **CIP** below.

F EXPERIMENTAL PLATFORMS AND LICENSES

F.1 EXPERIMENTAL PLATFORMS

All experiments of this approach are implemented on 2 Intel(R) Xeon(R) Gold 6430 and 2 NVIDIA Tesla A800 GPUs.

F.2 LICENSES

In our code, we have utilized the following libraries, each covered by its respective license agreements:

- PyTorch (BSD 3-Clause "New" or "Revised" License)
- Numpy (BSD 3-Clause "New" or "Revised" License)
- Tensorflow (Apache License 2.0)
- Meta-World (MIT License)
- MuJoCo (Apache License 2.0)
- Deep Mind Control (Apache License 2.0)
- Adroit Hand (Creative Commons License 3.0)

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

Algorithm 1 Causal information prioritization for efficient RL

Input: Q network Q_{π_c} , policy network π_c , inverse dynamics model ϕ_c with Q network Q_{ϕ_c} , replay buffer \mathcal{D} , local causal buffer \mathcal{D}_c , causal update interval I , causal matrix $M^{a \rightarrow s}$ and $M^{a \rightarrow r}$.

for each environment step t **do**
 Collect data with π_θ from real environment
 Add to replay buffer \mathcal{D} and local buffer \mathcal{D}_c
end for

Step 1: Counterfactual data augmentation

if every I environment step **then**
 Sample transitions \mathcal{D}_s from local buffer \mathcal{D}_c
 Learn causal mask matrix $M^{a \rightarrow r}$ with $\{(s, a, r, s')\}^{|\mathcal{D}_s|}$ for causal state prioritization
 Compute uncontrollable set \mathcal{U}_s followed by Eq. 4
 Sample $(s, a, r, s') \in \mathcal{D}_s$
for $s^i \in \mathcal{U}_s$ **do**
 Sample $(\hat{s}, \hat{a}, \hat{r}, \hat{s}') \sim \mathcal{D}_s$
if state $\hat{s}^i \in \mathcal{U}_s$ **then**
 Construct a counterfactual transition $(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$ by swapping (s^i, s'^i) with (\hat{s}^i, \hat{s}'^i)
 Add $(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$ to local buffer \mathcal{D}_c
end if
end for
end if

Step 2: Causal weighted matrix learning

if every I environment step **then**
 Sample transitions \mathcal{D}_a from local buffer \mathcal{D}_c
 Learn causal weighted matrix $M^{a \rightarrow r}$ with $\{(s, a, r, s')\}^{|\mathcal{D}_a|}$ for causal action prioritization
end if

Step 3: Policy optimization with causal action empowerment

for each gradient step **do**
 Sample N transitions (s, a, r, s') from \mathcal{D}
 Compute causal action empowerment followed by Eq. 8.
 Calculate the target Q_{ϕ_c} value
 Update Q_{ϕ_c} by $\min_{\phi_c} (\mathcal{T}Q_{\phi_c} - Q_{\phi_c})^2$
 Update ϕ_c by $\max(Q_{\phi_c}(s, a))$
 Calculate the target Q_{π_c} value
 Update Q_{π_c} by $\min_{\pi_c} (\mathcal{T}_c Q_{\pi_c} - Q_{\pi_c})^2$
 Update π_c by $\max_c (Q_{\pi_c}(s, a) + \mathcal{E}_{\pi_c}(s))$
end for
