

614

Appendix

Table of Contents

617	A Broader Impact	15
618	B Additional Related Works	15
619	B.1 Related RL formulations	15
620	B.2 Related literature of robustness in RL	16
621	B.3 Related literature of spurious correlation in RL	16
622	C Theoretical Analyses	17
623	C.1 Proof of Theorem 1	17
624	C.2 Proof of Theorem 2	19
625	C.3 Auxiliary results of SC-MDPs and RSC-MDPs	24
626	D Experiment Details	24
627	D.1 Architecture of the structural causal model	24
628	D.2 Environments	25
629	D.3 Computation resources	27
630	D.4 Hyperparameters	27
631	D.5 Discovered Causal Graph in SCM	27

635

A Broader Impact

636 Incorporating causality into reinforcement learning methods increases the interpretability of artificial
637 intelligence, which helps humans understand the underlying mechanism of algorithms and check
638 the source of failures. However, the learned causal transition model may contain human-readable
639 private information about the environment, which could raise privacy issues. To mitigate this potential
640 negative societal impact, the causal transition model needs to be encrypted and only accessible to
641 algorithms and trustworthy users.

642

B Additional Related Works

643 In this section, besides the most related formulation, robust RL introduced in Sec 3.3, we also
644 introduce some other related RL problem formulations partially shown in Figure 3. Then, we limit
645 our discussion to mainly two lines of work that are related to ours: (1) promoting robustness in RL;
646 (2) concerning the spurious correlation issues in RL.

647

B.1 Related RL formulations

648 **Robustness to noisy state: POMDPs and SA-MDPs.** State-noisy MDPs refer to the RL problem
649 that the agent can only access and choose the action based on a noisy observation rather than the true
650 state at each step, including two existing types of problems: Partially observable MDPs (POMDPs)
651 and state-adversarial MDPs (SA-MDPs), shown in Figure 3(b). In particular, at each step t , in
652 POMDPs, the observation o_t is generated by a fixed probability transition $\mathcal{O}(\cdot | s_t)$ (we refer to the
653 case that o_t only depends on the state s_t but not action); for state-adversarial MDPs, the observation
654 is an adversary $\nu(s_t)$ against and thus determined by the conducted policy, leading to the worst
655 performance by perturbing the state in a small set around itself. To against the state perturbation, both
656 POMDPs, and SA-MDPs are indeed robust to the noisy observation, or called agent-observed state,
657 but not the real state that transitions to the environment and next steps. In contrast, our RSC-MDPs

658 propose the robustness to the real state shift that will directly transition to the next state in the
659 environment, involving additional challenges induced by the appearance of out-of-distribution states.

660 **Robustness to unobserved confounder: MDPUC and confounded MDPs.** To address the mislead-
661 ing spurious correlations hidden in components of RL, people formulate RL problems as MDPs with
662 some additional components – unobserved confounders. In particular, the Markov decision process
663 with unobserved confounders (MDPUC) [35] serves as a general framework to concern all types of
664 possible spurious correlations in RL problems – at each step, the state, action, and reward are all
665 possibly influenced by some unobserved confounder, shown in Figure 2(d); confounded MDPs [19]
666 mainly concerns the misleading correlation between the current action and the next state, illustrated
667 in Figure 3(e). The proposed state-confounded MDPs (SC-MDPs) can be seen as a specified type of
668 MDPUC that focus on breaking the spurious correlation between different parts of the state space
669 itself (different from confounded MDPs which consider the correlation between action and next
670 state), motivated by various real-world applications in self-driving and control tasks. In addition, the
671 proposed formulation is more flexible and can work in both online and offline RL settings.

672 **Contextual MDPs (CMDPs).** A contextual MDP (CMDP) [36] is basically a set of standard MDPs
673 sharing the same state and action space but specified by different contexts within a context space.
674 In particular, the transition kernel, reward, and action of a CMDP are all determined by a (possibly
675 unknown) fixed context. The proposed robust state-confounded MDPs (RSC-MDPs) are similar
676 to CMDPs if we cast the unobserved confounder as the context in CMDPs, while different in two
677 aspects: (1) In a CMDP, the context is fixed throughout an episode, while the unobserved confounder
678 in RSC-MDPs can vary as $\{c_t\}_{1 \leq t \leq T}$; (2) In the online setting, the goal of CMDP is to beat the
679 optimal policy depending on the context, while RSC-MDPs seek to learn the optimal policy that does
680 not depend on the confounder $\{c_t\}_{1 \leq t \leq T}$.

681 B.2 Related literature of robustness in RL

682 **Robust RL (robust MDPs).** Concerning the robust issues in RL, a large portion of works focus on
683 robust RL with explicit uncertainty of the transition kernel, which is well-posed and a natural way
684 to consider the uncertainty of the environment [13, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46]. However,
685 to define the uncertainty set for the environment, most existing works use task structure-agnostic
686 and heuristic 'distance' such as KL divergence and total variation [14, 47, 48, 15, 49, 50, 51, 52]
687 to measure the shift between the training and test transition kernel, leading to a homogeneous
688 (almost structure-free) uncertainty set around the state space. In contrast, we consider a more general
689 uncertainty set that enables the robustness to a task-dependent heterogeneous uncertainty set shaped
690 by unobserved confounder and causal structure, in order to break the spurious correlation hidden in
691 different parts of the state space.

692 **Robustness in RL** Despite the remarkable success that standard RL has achieved, current RL
693 algorithms are still limited since the agent is vulnerable if the deployed environment is subject to
694 uncertainty and even structural changes. To address these challenges, a recent line of RL works
695 begins to concern robustness to the uncertainty or changes over different components of MDPs –
696 state, action, reward, and transition kernel, where a review [8] can be referred to. Besides robust
697 RL framework concerning the shift of the transition kernel and reward, to promote robustness in
698 RL, there exist various works [11, 12] that consider the robustness to action uncertainty, i.e., the
699 deployed action in the environment is distorted by an adversarial agent smoothly or circumstantially;
700 some works [9, 6, 10, 53, 54, 55] investigate the robustness to the state uncertainty including but not
701 limited to the introduced POMDPs and SA-MDPs in Appendix B.1 where the agent chooses the
702 action based on observation – the perturbed state determined by some restricted noise or adversarial
703 attack. The proposed RSC-MDPs can be regarded as addressing the state uncertainty since the
704 shift of the unobserved confounder leads to state perturbation. In contrast, RSC-MDPs consider
705 the out-of-distribution of the real state that will directly influence the subsequent transition in the
706 environment, but not the observation in POMDPs and SA-MDPs that will not directly influence the
707 environment.

708 B.3 Related literature of spurious correlation in RL

709 **Confounder in RL.** These works mainly focus on the confounder between action (treatment) and
710 state (effect), which is a long-standing problem that exists in the causal inference area. However,

711 we find that the confounder may cause problems from another perspective, where the confounder is
 712 built upon different dimensions of the state variable. Some people focus on the confounder between
 713 action and state, which is common in offline settings since the dataset is fixed and intervention is not
 714 allowed. But in the online setting, actions are controlled by an agent and intervention is available
 715 to eliminate spurious correlation. [56] reduces the spurious correlation between action and state in
 716 the offline setting. [57] deal with environment-irrelevant white noise; possible shift + causal [58].
 717 The confounder problem is usually easy to solve since agents can interact with the environment to do
 718 interventions. However, different from most existing settings, we find that even with the capability
 719 of intervention, the confounding between dimensions in states cannot be fully eliminated. Then the
 720 learned policy is heavily influenced if these confounder change during testing.

721 **Invariant Feature learning.** The problem of spurious correlation has attracted attention in the
 722 supervised learning area for a long time and many solutions are proposed to learn invariant features to
 723 eliminate spurious correlations. A general framework to remedy the ignorance of spurious correlation
 724 in empirical risk minimization (ERM) is invariant risk minimization (IRM) [59]. Other works tackle
 725 this problem with group distributional robustness [60], adversarial robustness [61], and contrastive
 726 learning [62]. These methods are also adapted to sequential settings. The idea of increasing the
 727 robustness of RL agents by training agents on multiple environments has been shown in previous
 728 works [63, 30, 30]. However, a shared assumption among these methods is that multiple environments
 729 with different values of confounder are accessible, which is not always true in the real world.

730 **Counterfactual Data Augmentation in RL.** One way to simulate multiple environments is data
 731 augmentation. However, most data augmentation works [24, 64, 25, 65, 66, 67, 68] apply image
 732 transformation to raw inputs, which requires strong domain knowledge for image manipulation and
 733 cannot be applied to other types of inputs. In RL, the dynamic model and reward model follow certain
 734 causal structures, which allow counterfactual generation of new transitions based on the collected
 735 samples. This line of work, named counterfactual data augmentation, is very close to this work.
 736 Deep generative models [69] and adversarial examples [70] are considered for the generation to
 737 improve sample efficiency in model-based RL. CoDA [71] and MocoDA [32] leverage the concept of
 738 locally factored dynamics to randomly stitch components from different trajectories. However, the
 739 assumption of local causality may be limited.

740 **Domain Randomization.** If we are allowed to control the data generation process, e.g., the underlying
 741 mechanism of the simulator, we can apply the golden rule in causality – Randomized Controlled
 742 Trial (RCT). The well-known technic, domain randomization [72], exactly follows the idea of RCT,
 743 which randomly perturb the internal state of the experiment in simulators. Later literature follows this
 744 direction and develops variants including randomization guided by downstream tasks in the target
 745 domain [73, 74], randomization to match real-world distributions [75, 76], and randomization to
 746 minimize data divergence [77]. However, it is usually impossible to randomly manipulate internal
 747 states in most situations in the real world. In addition, determining which variables to randomize is
 748 even harder given so many factors in complex systems.

749 **Discovering Spurious Correlations** Detecting spurious correlations helps models remove features
 750 that are harmful to generalization. Usually, domain knowledge is required to find such correlations [78,
 751 79, 80]. However, when prior knowledge is accessible, technics such as clustering can also be used to
 752 reveal spurious attributes [35, 81, 82]. When human inspection is available, recent works [83, 84, 85]
 753 also use explainability techniques to find spurious correlations. Another area for discovery is concept-
 754 level and interactive debugging [86, 87], which leverage concepts or human feedback to perform
 755 debugging.

756 C Theoretical Analyses

757 C.1 Proof of Theorem I

758 The proof follows the pipeline of proving the existence of the optimal policy for standard MDPs but
 759 tailored for RSC-MDPs since the additional components confounder C_s and the infimum operator.
 760 To begin with, recall that the goal is to find a policy $\tilde{\pi} = \{\tilde{\pi}_t\}_{1 \leq t \leq T}$ such that:

$$\tilde{V}_t^{\tilde{\pi}, \sigma}(s) = \tilde{V}_t^{*, \sigma}(s) := \sup_{\pi \in \Pi} \tilde{V}_t^{\pi, \sigma}(s) \quad \text{and} \quad \tilde{Q}_t^{\tilde{\pi}, \sigma}(s, a) = \tilde{Q}_t^{*, \sigma}(s, a) := \sup_{\pi \in \Pi} \tilde{Q}_t^{\pi, \sigma}(s, a). \quad (8)$$

761 Towards this, we start from the first claim in equation [8](#). Before proceeding, we let $\{S_t, A_t, R_t, C_t\}$
762 denote the random variables at time step t for all $1 \leq t \leq T$. Then due to the Markov properties, we
763 know that conditioned on current state s_t , the future state, action, and reward are all independent from
764 the previous $s_1, a_1, r_1, c_1, \dots, s_{t-1}, a_{t-1}, r_{t-1}, c_{t-1}$. For convenience, we introduce the following
765 notation:

$$\forall 1 \leq t \leq T : P_{+t} := \{P_k\}_{t \leq k \leq T} \quad \text{and} \quad \mathcal{U}^\sigma(P_{+t}^c) := \{\mathcal{U}^\sigma(P_k^c)\}_{t \leq k \leq T} \quad (9)$$

766 to represent the collection of variables from time step t to the end of the episode, and choose $\tilde{\pi}$ to
767 obey

$$\forall 1 \leq t \leq T : \pi_t(s) := \operatorname{argmax}_{a \in \mathcal{A}} \mathbb{E} \left[r_t(s, a) + \inf_{P_t \in \mathcal{U}^\sigma(P_{t,s,a}^c)} \mathbb{E}_{c_t \sim P_t} \left[\tilde{V}_{t+1}^{*,\sigma}(s_{t+1}) \right] \right] \quad (10)$$

768 With the above preparation in mind, for any $(t, s) \in \{1, 2, \dots, T\} \times \mathcal{S}$, one has

$$\begin{aligned} & \tilde{V}_t^{*,\sigma}(s) \\ & \stackrel{(i)}{=} \sup_{\pi \in \Pi} \inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} \tilde{V}_t^{\pi, P}(s) \stackrel{(ii)}{=} \sup_{\pi \in \Pi} \inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} \mathbb{E}_{\pi, P_{+t}} \left[\sum_{k=t}^T r_k(s_k, a_k) \right] \\ & \stackrel{(iii)}{=} \sup_{\pi \in \Pi} \inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} \mathbb{E}_{\pi_t} \left[r_t(s, a_t) \right. \\ & \quad \left. + \mathbb{E}_{c_t \sim P_t} \mathbb{E} \left[\sum_{k=t+1}^T r_k(s_k, a_k) \mid \pi, P_{+(t+1)}, (S_t, A_t, R_t, C_t, S_{t+1}) = (s, a_t, r_t, c_t, s_{t+1}) \right] \right] \\ & = \sup_{\pi \in \Pi} \mathbb{E}_{\pi_t} \left[\inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} r_t(s, a_t) + \inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} \mathbb{E}_{c_t \sim P_t} \right. \\ & \quad \left. \mathbb{E} \left[\sum_{k=t+1}^T r_k(s_k, a_k) \mid \pi, P_{+(t+1)}, (S_t, A_t, R_t, C_t, S_{t+1}) = (s, a_t, r_t, c_t, s_{t+1}) \right] \right] \end{aligned}$$

769 where (i) and (ii) holds by the definitions in equation [5](#) and equation [3](#) respectively, and (iii) follows
770 from expressing the term of interest by moving one step ahead and \mathbb{E}_{π_t} is taken with respect to
771 $a_t \sim \pi_t(\cdot \mid S_1 = s_1, A_1 = a_1, \dots, S_t = s)$, and the last equality arises from we can exchange the
772 operators \mathbb{E}_{π_t} and $\inf_{P \in \mathcal{U}^\sigma(P^c)}$ since they are independent.

773 To continue, we observe that the above equation can be rewritten and controlled as follows:

$$\begin{aligned} & \tilde{V}_t^{*,\sigma}(s) \\ & = \sup_{\pi \in \Pi} \mathbb{E}_{\pi_t} \left[\inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} r_t(s, a_t) + \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P_t} \inf_{P_{+(t+1)} \in \mathcal{U}^\sigma(P_{+(t+1)}^c)} \right. \\ & \quad \left. \mathbb{E} \left[\sum_{k=t+1}^T r_k(s_k, a_k) \mid \pi', P_{+(t+1)}, (S_t, A_t, R_t, C_t, S_{t+1}) = (s, a_t, r_t, c_t, s_{t+1}) \right] \right] \\ & \leq \sup_{\pi \in \Pi} \mathbb{E}_{\pi_t} \left[\inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} r_t(s, a_t) + \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P_t} \sup_{\pi' \in \Pi} \inf_{P_{+(t+1)} \in \mathcal{U}^\sigma(P_{+(t+1)}^c)} \right. \\ & \quad \left. \mathbb{E} \left[\sum_{k=t+1}^T r_k(s_k, a_k) \mid \pi', P_{+(t+1)}, (S_t, A_t, R_t, C_t, S_{t+1}) = (s, a_t, r_t, c_t, s_{t+1}) \right] \right] \\ & \stackrel{(i)}{=} \sup_{\pi \in \Pi} \mathbb{E}_{\pi_t} \left[r_t(s, a_t) + \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P_t} \left[\sup_{\pi' \in \Pi} \inf_{P_{+(t+1)} \in \mathcal{U}^\sigma(P_{+(t+1)}^c)} \mathbb{E}_{\pi', P_{+(t+1)}} \left[\sum_{k=t+1}^T r_k(s_k, a_k) \right] \right] \right] \\ & = \sup_{\pi \in \Pi} \mathbb{E}_{\pi_t} \left[r_t(s, a_t) + \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P_t} \left[\tilde{V}_{t+1}^{*,\sigma}(s_{t+1}) \right] \right] \end{aligned}$$

$$\begin{aligned}
&= \sup_{a_t \in \mathcal{A}} \mathbb{E}_{a_t} \left[r_t(s, a_t) + \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P_t} \left[\tilde{V}_{t+1}^{*,\sigma}(s_{t+1}) \right] \right] \\
&= \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E} \left[r_t(s, a_t) + \mathbb{E}_{c_t \sim P_t} \left[\tilde{V}_{t+1}^{*,\sigma}(s_{t+1}) \right] \mid a_t = \tilde{\pi}_t(s) \right], \tag{11}
\end{aligned}$$

774 where (i) holds by the Markov decision such that the rewards $\{r_k(s_k, a_k)\}_{t+1 \leq k \leq T}$ conditioned
775 on determined $(S_t, A_t, R_t, C_t, S_{t+1})$ or S_{t+1} are the same, and the last equality follows from the
776 definition of $\tilde{\pi}$ in equation [10](#) and the exchangeability of $\inf_{P_t \in \mathcal{U}^\sigma(P_t^c)}$ and $\mathbb{E}_{a_t}[\cdot]$.

777 Applying equation [11](#) recursively for $t+1, \dots, T$, we arrive at

$$\begin{aligned}
\tilde{V}_t^{*,\sigma}(s) &\leq \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E} \left[r_t(s, a_t) + \mathbb{E}_{c_t \sim P_t} \left[\tilde{V}_{t+1}^{*,\sigma}(s_{t+1}) \right] \mid a_t = \tilde{\pi}_t(s) \right] \\
&\leq \inf_{P_t \in \mathcal{U}^\sigma(P_t^c)} \inf_{P_{t+1} \in \mathcal{U}^\sigma(P_{t+1}^c)} \mathbb{E} \left[r_t(s, a_t) + \right. \\
&\quad \left. \mathbb{E}_{c_t \sim P_t} \left[r_{t+1}(s_{t+1}, a_{t+1}) + \mathbb{E}_{c_{t+1} \sim P_{t+1}} \left[\tilde{V}_{t+2}^{*,\sigma}(s_{t+1}) \right] \right] \mid (a_t, a_{t+1}) = (\tilde{\pi}_t(s), \tilde{\pi}_{t+1}(s_{t+1})) \right] \\
&\leq \dots \leq \inf_{P_{+t} \in \mathcal{U}^\sigma(P_{+t}^c)} \mathbb{E}_{\pi, P} \left[\sum_{k=t}^T r_k(s_k, a_k) \right] = \tilde{V}_t^{\tilde{\pi}, \sigma}(s). \tag{12}
\end{aligned}$$

778 where (i) holds by the Markov properties of the rewards.

779 Observing from equation [12](#) that

$$\forall s \in \mathcal{S} : \quad \tilde{V}_t^{*,\sigma}(s) \leq \tilde{V}_t^{\tilde{\pi}, \sigma}(s) \leq \sup_{\pi \in \Pi} \tilde{V}_t^{\pi, \sigma}(s) = \tilde{V}_t^{*,\sigma}(s), \tag{13}$$

780 which directly verifies the first assertion in equation [8](#) $\tilde{V}_t^{\tilde{\pi}, \sigma}(s) = \tilde{V}_t^{*,\sigma}(s)$ for all $s \in \mathcal{S}$. The second
781 assertion in equation [8](#) can be achieved analogously. Until now, we verify that there exists at least a
782 policy $\tilde{\pi}$ that obeys equation [8](#), which we refer it as an optimal policy since its value is equal to or
783 larger than any other non-stationary and stochastic policies over all states $s \in \mathcal{S}$.

784 C.2 Proof of Theorem [2](#)

785 **Constructing a hard instance of the standard MDP.** In this section, we consider the following
786 standard MDP instance $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, P^0, T, r\}$, where $\mathcal{S} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ is the state
787 space consisting of four elements in dimension $n = 2$, and $\mathcal{A} = \{0, 1\}$ is the action space with only
788 two options. The transition kernel $P^0 = \{P_t^0\}_{1 \leq t \leq T}$ at different time steps $1 \leq t \leq T$ is defined as

$$P_1^0(s' | s, a) = \begin{cases} \mathbb{1}(s' = [0, 0])\mathbb{1}(a = 0) + \mathbb{1}(s' = [0, 1])\mathbb{1}(a = 1) & \text{if } (s, a) = ([0, 0], a) \\ \mathbb{1}(s' = s) & \text{otherwise} \end{cases}, \tag{14}$$

789 and

$$P_t^0(s' | s, a) = \mathbb{1}(s' = s), \quad \forall (t, s, a) \in \{2, 3, \dots, T\} \times \mathcal{S} \times \mathcal{A}. \tag{15}$$

790 Note that this transition kernel P^0 ensures the next state transitioned from the state $[0, 0]$ is either
791 $[0, 0]$ or $[0, 1]$. The reward function is specified as follows: for all time steps $1 \leq t \leq T$,

$$r_t(s, a) = \begin{cases} 1 & \text{if } s = [0, 0] \text{ or } s = [1, 1] \\ 0 & \text{otherwise} \end{cases}. \tag{16}$$

792 **The equivalence to one SC-MDP.** Then, we shall show that the constructed standard MDP \mathcal{M}
793 can be equivalently represented by one SC-MDP $\mathcal{M}_{\text{sc}} = \{\mathcal{S}, \mathcal{A}, T, r, \mathcal{C}, \{\mathcal{P}_t^i\}, P^c\}$ with $\mathcal{C} := [0, 1]$,
794 which yields the sequential observations $\{s_t, a_t, r_t\}_{1 \leq t \leq T}$ induced by any policy and any initial
795 state distribution in two processes are identical. To specify, $\mathcal{S}, \mathcal{A}, T, r$ are kept the same as \mathcal{M} .
796 Here, $\{\mathcal{P}_t^i\}$ shall be specified in a while, which determines the transition to each dimension of the
797 next state conditioned on the current state, action, and confounder for all time steps, i.e., $s_{t+1}^i \sim$
798 $\mathbb{E}_{c_t \sim P_t^c} [\mathcal{P}_t^i(\cdot | s_t, a_t, c_t)]$ for any i -th dimension of the state ($i \in \{1, 2\}$) and all timestep $1 \leq t \leq T$.

799 For convenience, we denote $\mathcal{P}_t := [\mathcal{P}_t^1, \mathcal{P}_t^2] \in \Delta(\mathcal{S})$ as the transition kernel towards the next state,
800 namely, $s_{t+1} \sim \mathbb{E}_{c_t \sim P_t^c} [\mathcal{P}_t(\cdot | s_t, a_t, c_t)]$.

801 To ensure the marginalized transition probability from any state-action pair (s_t, a_t) to the next state
802 s_{t+1} in \mathcal{M}_{sc} aligns with the one in the MDP \mathcal{M} , we set

$$P_t^c(c) = \mathbb{1}(c = 0), \quad \forall 1 \leq t \leq T. \quad (17)$$

803 In addition, before introducing the transition kernel $\{\mathcal{P}_t^i\}$ of the SC-MDP \mathcal{M}_{sc} , we introduce an
804 auxiliary transition kernel $P^{\text{sc}} = \{P_t^{\text{sc}}\}$ as follows:

$$P_1^{\text{sc}}(s' | s, a) = \begin{cases} \mathbb{1}(s' = [1, 0])\mathbb{1}(a = 0) + \mathbb{1}(s' = [1, 1])\mathbb{1}(a = 1) & \text{if } (s, a) = ([0, 0], 0) \\ \mathbb{1}(s' = s) & \text{otherwise} \end{cases}, \quad (18)$$

805 and

$$P_t^{\text{sc}}(s' | s, a) = \mathbb{1}(s' = s), \quad \forall (t, s, a) \in \{2, 3, \dots, T\} \times \mathcal{S} \times \mathcal{A}. \quad (19)$$

806 It can be observed that P^{sc} is similar to P^0 except for the transition in the state $[0, 0]$.

807 Armed with this transition kernel P^{sc} , the $\{\mathcal{P}_t^i\}$ of the SC-MDP \mathcal{M}_{sc} is set to obey

$$\mathcal{P}_1(s' | s, a, c) = \begin{cases} (1 - c)P_1^0(s' | s, a) + cP_1^{\text{sc}}(s' | s, a) & \text{if } (s, a) = ([0, 0], a) \\ \mathbb{1}(s' = s) & \text{otherwise} \end{cases}, \quad (20)$$

808 and

$$\mathcal{P}_t(s' | s, a, c) = \mathbb{1}(s' = s), \quad \forall (t, s, a, c) \in \{2, 3, \dots, T\} \times \mathcal{S} \times \mathcal{A} \times \mathcal{C}. \quad (21)$$

809 With the above preparation, we are ready to verify that the marginalized transition from the current
810 state and action to the next state in the SC-MDP \mathcal{M}_{sc} is identical to the one in MDP \mathcal{M} : for all
811 $(t, s_t, a_t, s_{t+1}) \in \{1, 2, \dots, T\} \times \mathcal{S} \times \mathcal{A} \times \mathcal{S}$:

$$\mathbb{P}(s_{t+1} | s_t, a_t) = \mathbb{E}_{c_t \sim P_t^c} [\mathcal{P}_t(s_{t+1} | s_t, a_t, c_t)] = \mathcal{P}_t(s_{t+1} | s_t, a_t, 0) = P^0(s_{t+1} | s_t, a_t) \quad (22)$$

812 where the second equality holds by the definition of P^c in equation 17, and the last equality holds by
813 the definitions of P^0 (see equation 14 and equation 15) and \mathcal{P} (see equation 20 and equation 21).

814 In summary, we verified that the standard MDP $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, P^0, T, r\}$ is equal to the above specified
815 SC-MDP \mathcal{M}_{sc} .

816 **Defining the corresponding RMDP and RSC-MDP.** Equipped with the equivalent MDP \mathcal{M}
817 and SC-MDP \mathcal{M}_{sc} , people could consider the robust variants of them respectively — a RMDP
818 $\mathcal{M}_{\text{rob}} = \{\mathcal{S}, \mathcal{A}, \mathcal{U}^{\sigma_1}(P^0), T, r\}$ with the uncertainty level σ_1 , and the proposed RSC-MDP
819 $\mathcal{M}_{\text{sc-rob}} = \{\mathcal{S}, \mathcal{A}, T, r, \mathcal{C}, \{\mathcal{P}_t^i\}, \mathcal{U}^{\sigma_2}(P^c)\}$ with the uncertainty level σ_2 .

820 In this section, without loss of generality, we consider total deviation as the ‘distance’ function ρ for
821 the uncertainty sets of both RMDP \mathcal{M}_{rob} and RSC-MDP $\mathcal{M}_{\text{sc-rob}}$, i.e., for any probability vectors
822 $P', P \in \Delta(\mathcal{C})$ (or $P', P \in \Delta(\mathcal{S})$), $\rho(P', P) := \frac{1}{2} \|P' - P\|_1$. Consequently, for any uncertainty set
823 $\sigma \in [0, 1]$, the uncertainty set $\mathcal{U}^{\sigma_1}(P^0)$ of the RMDP (see equation 1) and $\mathcal{U}^{\sigma_2}(P^c)$ of the RSC-MDP
824 $\mathcal{M}_{\text{sc-rob}}$ (see equation 4) are defined as follows:

$$\begin{aligned} \mathcal{U}^\sigma(P^0) &:= \otimes \mathcal{U}^\sigma(P_{t,s,a}^0), & \mathcal{U}^\sigma(P_{t,s,a}^0) &:= \left\{ P_{t,s,a} \in \Delta(\mathcal{S}) : \frac{1}{2} \|P_{t,s,a} - P_{t,s,a}^0\|_1 \leq \sigma \right\}, \\ \mathcal{U}^\sigma(P^c) &:= \otimes \mathcal{U}^\sigma(P_t^c), & \mathcal{U}^\sigma(P_t^c) &:= \left\{ P \in \Delta(\mathcal{C}) : \frac{1}{2} \|P - P_t^c\|_1 \leq \sigma \right\}. \end{aligned} \quad (23)$$

825 To continue, the proof is established by specifying the robust optimal policy $\pi_{\text{RMDP}}^{*,\sigma_1}$ associated with
826 \mathcal{M}_{rob} and $\pi_{\text{RSC}}^{*,\sigma_2}$ associated with $\mathcal{M}_{\text{sc-rob}}$ and then compare their performance on RSC-MDP with
827 some initial state distribution.

828 **The performance comparisons between $\pi_{\text{RMDP}}^{*,\sigma_1}$ of RMDP \mathcal{M}_{rob} and $\pi_{\text{RSC}}^{*,\sigma_2}$ of RSC-MDP $\mathcal{M}_{\text{sc-rob}}$.**

829 To begin, we introduce the following lemma which specifies the robust optimal policy $\pi_{\text{RMDP}}^{*,\sigma_1}$
830 associated with the RMDP \mathcal{M}_{rob} .

831 **Lemma 1.** For any $\sigma_1 \in (0, 1]$, the robust optimal policy and its corresponding robust SC-value
 832 functions satisfy

$$\pi_{\text{RM DP}}^{*,\sigma_1}(0 | s) = 1, \quad \text{for } s \in \mathcal{S}. \quad (24a)$$

833 In addition, we characterize the robust SC-value functions of the RSC-MDP $\mathcal{M}_{\text{sc-rob}}$ associated with
 834 any policy, combined with the robust optimal policy $\pi_{\text{RSC}}^{*,\sigma_2}$ of $\mathcal{M}_{\text{sc-rob}}$ — the optimal robust SC-value
 835 functions, shown in the following lemma.

836 **Lemma 2.** Consider any $\sigma_2 \in (\frac{3}{4}, 1]$ and the RSC-MDP $\mathcal{M}_{\text{sc-rob}} = \{\mathcal{S}, \mathcal{A}, T, r, \mathcal{C}, \{\mathcal{P}_t^i\}, \mathcal{U}^{\sigma_2}(P^c)\}$.
 837 For any policy π , the corresponding robust SC-value functions satisfy

$$\tilde{V}_1^{\pi,\sigma_2}([0, 0]) = 1 + (T - 1) \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[\pi_1(0 | [0, 0])(1 - c_1) + \pi_1(1 | [0, 0])c_1 \right]. \quad (25a)$$

838 In addition, the optimal robust SC-value function and the robust optimal policy $\pi_{\text{RSC}}^{*,\sigma_2}$ of the RMDP
 839 $\mathcal{M}_{\text{sc-rob}}$ obeys:

$$\tilde{V}_1^{\pi_{\text{RSC}}^{*,\sigma_2},\sigma_2}([0, 0]) = \tilde{V}_1^{*,\sigma_2}([0, 0]) = 1 + \frac{T - 1}{2}. \quad (26)$$

840 Applying Lemma 2 with policy $\pi = \pi_{\text{RM DP}}^{*,\sigma_1}$ in Lemma 1 one has

$$\tilde{V}_1^{\pi_{\text{RM DP}}^{*,\sigma_1},\sigma_2}([0, 0]) = 1 + (T - 1) \inf_{P \in \mathcal{U}_2^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[1 - c_1 \right] \leq 1 + \frac{T - 1}{4}, \quad (27)$$

841 where the last inequality holds by the probability distribution P obeying $P_1(0) = \frac{1}{4}$ and $P_1(1) = \frac{3}{4}$
 842 is inside the uncertainty set $\mathcal{U}_2^\sigma(P_1^c)$.

843 Finally, putting equation 27 and equation 26 together, we complete the proof by showing that with
 844 the initial state distribution ϕ define as $\rho(s_1 = [0, 0]) = 1$, we arrive at

$$\tilde{V}_1^{\pi_{\text{RSC}}^{*,\sigma_2},\sigma_2}(\phi) - \tilde{V}_1^{\pi_{\text{RM DP}}^{*,\sigma_1},\sigma_2}(\phi) = \tilde{V}_1^{*,\sigma_2}(\phi) - \tilde{V}_1^{\pi_{\text{RM DP}}^{*,\sigma_1},\sigma_2}(\phi) \geq \frac{T - 1}{4} \approx \frac{T}{4}. \quad (28)$$

845 C.2.1 Proof of Lemma 1

846 **Specifying the minimum of the robust value functions in different states.** For any uncertainty set
 847 $\sigma_1 \in (0, 1]$, we first characterize the robust value function of any policy π over different states. To
 848 start, we denote the minimum of the robust value function over states at each time step t as below:

$$V_{\min,t}^{\pi,\sigma_1} := \min_{s \in \mathcal{S}} V_t^{\pi,\sigma_1}(s) \geq 0, \quad (29)$$

849 where the last inequality holds by that the reward function defined in equation 16 is always non-
 850 negative. Obviously, there exists at least one state $s_{\min,t}^\pi$ that satisfies $V_t^{\pi,\sigma_1}(s_{\min,t}^\pi) = V_{\min,t}^{\pi,\sigma_1}$.

851 With this in mind, we shall verify that for any policy π ,

$$\forall 1 \leq t \leq T : \quad V_t^{\pi,\sigma_1}([0, 1]) = V_t^{\pi,\sigma_1}([1, 0]) = 0. \quad (30)$$

852 To achieve this, we will use a recursive argument. First, the base case can be verified since when
 853 $t + 1 = T + 1$, the value functions are all zeros at $T + 1$ step, i.e., $V_{t+1}^{\pi,\sigma_1}(s) = V_{T+1}^{\pi,\sigma_1}(s) = 0$ for all
 854 $s \in \mathcal{S}$. Then, the goal is to verify the following fact

$$V_t^{\pi,\sigma_1}([0, 1]) = V_t^{\pi,\sigma_1}([1, 0]) = 0 \quad (31)$$

855 with the assumption that $V_{t+1}^{\pi,\sigma_1}(s) = 0$ for any state $s = \{[0, 1], [1, 0]\}$. It is easily observed that for
 856 any policy π , the robust value function when state $s = \{[0, 1], [1, 0]\}$ at any time step t obeys

$$\begin{aligned} 0 \leq V_t^{\pi,\sigma_1}(s) &= \mathbb{E}_{a \sim \pi_t(\cdot | s)} \left[r_t(s, a) + \inf_{P \in \mathcal{U}^{\sigma_1}(P_{t,s,a}^0)} P V_{t+1}^{\pi,\sigma_1} \right] \\ &\stackrel{(i)}{=} 0 + (1 - \sigma_1) V_{t+1}^{\pi,\sigma_1}(s) + \sigma_1 V_{\min,t+1}^{\pi,\sigma_1} \stackrel{(ii)}{=} 0 + \sigma_1 V_{\min,t+1}^{\pi,\sigma_1} \end{aligned}$$

$$\leq 0 + \sigma_1 V_{t+1}^{\pi, \sigma_1}(s) = 0 \quad (32)$$

857 where (i) holds by $r_t(s, a) = 0$ for all $s = \{[0, 1], [1, 0]\}$, the fact $P_t^0(s | s, a) = 1$ (see equation 14
858 and equation 15), and the definition of the uncertainty set $\mathcal{U}^{\sigma_1}(P^0)$ in equation 23, (ii) follows from
859 the recursive assumption $V_{t+1}^{\pi, \sigma_1}(s) = 0$ for any state $s = \{[0, 1], [1, 0]\}$, and the last equality holds
860 by $V_{\min, t+1}^{\pi, \sigma_1} \leq V_{t+1}^{\pi, \sigma_1}(s)$ (see equation 29). Until now, we complete the proof for equation 31 and
861 then verify equation 30.

862 Note that equation 30 directly leads to

$$\forall 1 \leq t \leq T : \quad V_{\min, t}^{\pi, \sigma_1} = 0. \quad (33)$$

863 **Considering the robust value function at state $[0, 0]$.** Armed with above facts, we are now ready to
864 derive the robust value function for the state $[0, 0]$.

865 When $2 \leq t \leq T$, one has

$$\begin{aligned} V_t^{\pi, \sigma_1}([0, 0]) &= \mathbb{E}_{a \sim \pi_t(\cdot | [0, 0])} \left[r_t([0, 0], a) + \inf_{P \in \mathcal{U}^{\sigma_1}(P_{t, [0, 0], a})} PV_{t+1}^{\pi, \sigma_1} \right] \\ &\stackrel{(i)}{=} 1 + \left[(1 - \sigma_1) V_{t+1}^{\pi, \sigma_1}([0, 0]) + \sigma_1 V_{\min, t+1}^{\pi, \sigma_1} \right] \\ &= 1 + (1 - \sigma_1) V_{t+1}^{\pi, \sigma_1}([0, 0]) \end{aligned} \quad (34)$$

866 where (i) holds by $r_t([0, 0], a) = 1$ for all $a \in \{0, 1\}$ and the definition of P^0 (see equation 14 and
867 equation 15), and the last equality arises from equation 33.

868 Applying equation 34 recursively for $t, t+1, \dots, T$ yields that

$$V_t^{\pi, \sigma_1}([0, 0]) = \sum_{k=t}^T (1 - \sigma_1)^{k-t} \geq 1. \quad (35)$$

869 At the first step, the robust value function obeys:

$$\begin{aligned} V_1^{\pi, \sigma_1}([0, 0]) &= \mathbb{E}_{a \sim \pi_1(\cdot | [0, 0])} \left[r_t([0, 0], a) + \inf_{P \in \mathcal{U}^{\sigma_1}(P_{1, [0, 0], a})} PV_2^{\pi, \sigma_1} \right] \\ &\stackrel{(i)}{=} 1 + \pi_1(0 | [0, 0]) \inf_{P \in \mathcal{U}^{\sigma_1}(P_{1, [0, 0], 0})} PV_2^{\pi, \sigma_1} + \pi_1(1 | [0, 0]) \inf_{P \in \mathcal{U}^{\sigma_1}(P_{1, [0, 0], 1})} PV_2^{\pi, \sigma_1} \\ &\stackrel{(ii)}{=} 1 + \pi_1(0 | [0, 0]) \left[(1 - \sigma_1) V_2^{\pi, \sigma_1}([0, 0]) + \sigma_1 V_{\min, 2}^{\pi, \sigma_1} \right] \\ &\quad + \pi_1(1 | [0, 0]) \left[(1 - \sigma_1) V_2^{\pi, \sigma_1}([0, 1]) + \sigma_1 V_{\min, 2}^{\pi, \sigma_1} \right] \\ &= 1 + \pi_1(0 | [0, 0]) (1 - \sigma_1) V_2^{\pi, \sigma_1}([0, 0]) \end{aligned} \quad (36)$$

870 where (i) holds by $r_t([0, 0], a) = 1$ for all $a \in \{0, 1\}$, (ii) follows from the definition of P^0 (see
871 equation 14 and equation 15), and the last equality arises from equation 30 and equation 33.

872 **The optimal policy $\pi_{\text{RM DP}}^{*, \sigma_1}$.** Observing that the positive value of $V_2^{\pi, \sigma_1}([0, 0])$ verified in equa-
873 tion 35, as $V_1^{\pi, \sigma_1}([0, 0])$ is increasing monotonically as $\pi_1(0 | [0, 0])$ is larger, we directly have that
874 $\pi_{\text{RM DP}}^{*, \sigma_1}(0 | [0, 0]) = 1$.

875 Considering that the action does not influence the state transition for all other states $s \neq [0, 0]$,
876 without loss of generality, we choose the robust optimal policy to obey

$$\forall s \in \mathcal{S} : \quad \pi_{\text{RM DP}}^{*, \sigma_1}(0 | s) = 1. \quad (37)$$

877 C.2.2 Proof of Lemma 2

878 To begin with, for any uncertainty level $\sigma_2 \in (\frac{1}{2}, 1]$ and any policy $\pi = \{\pi_t\}$, we consider the robust
879 SC-value function $\tilde{V}_1^{\pi, \sigma_2}$ of the RSC-MDP $\mathcal{M}_{\text{sc-rob}}$.

880 **Deriving $\tilde{V}_t^{\pi, \sigma_2}$ for $2 \leq t \leq T$.** Towards this, for any $2 \leq t \leq T$ and $s \in \mathcal{S}$, one has

$$\tilde{V}_t^{\pi, \sigma_2}(s) \stackrel{(i)}{=} \inf_{P \in \mathcal{U}^{\sigma_2}(P^c)} \tilde{V}_t^{\pi, P}(s) = \inf_{P \in \mathcal{U}^{\sigma_2}(P_t^c)} \mathbb{E}_{a \sim \pi_t(s)} \left[\tilde{Q}_t^{\pi, P}(s, a) \right]$$

$$\begin{aligned}
&\stackrel{\text{(ii)}}{=} \inf_{P \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{a \sim \pi_t(s)} \left[r_t(s, a) + \mathbb{E}_{c_t \sim P} \left[\mathcal{P}_{t,s,a,c_t} \tilde{V}_{t+1}^{\pi,\sigma} \right] \right] \\
&\stackrel{\text{(iii)}}{=} r_t(s, a) + \inf_{P \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P} \left[\mathcal{P}_{t,s,a,c_t} \tilde{V}_{t+1}^{\pi,\sigma} \right] \\
&= r_t(s, a) + \tilde{V}_{t+1}^{\pi,\sigma}(s), \tag{38}
\end{aligned}$$

881 where (i) holds by the definition in equation 5, (ii) follows from the *state-confounded* Bellman
882 consistency equation in equation 47, (iii) follows from that the reward function r and \mathcal{P}_t are all
883 independent from the action (see equation 16, equation 17 and equation 21), and the last inequality
884 holds by $\mathcal{P}_t(s' | s, a, c) = \mathbf{1}(s' = s)$ is independent from c_t (see equation 21).

885 Applying the above fact recursively for $t, t+1, \dots, T$ leads to that for any $s \in \mathcal{S}$,

$$\begin{aligned}
\tilde{V}_t^{\pi,\sigma^2}(s) &= r_t(s, a_t) + \tilde{V}_{t+1}^{\pi,\sigma}(s) = r_t(s, a) + r_{t+1}(s, a_{t+1}) + \tilde{V}_{t+2}^{\pi,\sigma}(s) \\
&= \dots = r_t(s, a_t) + \sum_{k=t+1}^T r_k(s_k, a_k), \tag{39}
\end{aligned}$$

886 which directly yields

$$\tilde{V}_2^{\pi,\sigma^2}([0, 0]) = \tilde{V}_2^{\pi,\sigma^2}([1, 1]) = T - 1 \quad \text{and} \quad \tilde{V}_2^{\pi,\sigma^2}([0, 1]) = \tilde{V}_2^{\pi,\sigma^2}([1, 0]) = 0. \tag{40}$$

887 **Characterizing $\tilde{V}_1^{\pi,\sigma^2}([0, 0])$ for any policy π .** In this section, we are especially interested in the
888 value of $\tilde{V}_1^{\pi,\sigma^2}$ on the state $[0, 0]$. To proceed, one has

$$\begin{aligned}
\tilde{V}_1^{\pi,\sigma^2}([0, 0]) &\stackrel{\text{(i)}}{=} \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \tilde{V}_1^{\pi,P}([0, 0]) = \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{a \sim \pi_1([0, 0])} \left[\tilde{Q}_1^{\pi,P}([0, 0], a) \right] \\
&\stackrel{\text{(ii)}}{=} \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{a \sim \pi_1([0, 0])} \left[r_1([0, 0], a) + \mathbb{E}_{c_1 \sim P} \left[\mathcal{P}_{1,[0,0],a,c_1} \tilde{V}_2^{\pi,\sigma} \right] \right] \\
&\stackrel{\text{(iii)}}{=} 1 + \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[\left(\pi_1(0 | [0, 0]) \mathcal{P}_{1,[0,0],0,c_1} + \pi_1(1 | [0, 0]) \mathcal{P}_{1,[0,0],1,c_1} \right) \tilde{V}_2^{\pi,\sigma} \right] \\
&\stackrel{\text{(iv)}}{=} 1 + \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[\pi_1(0 | [0, 0]) \left((1 - c_1) P_{1,[0,0],0}^0 + c_1 P_{1,[0,0],0}^{\text{sc}} \right) \tilde{V}_2^{\pi,\sigma} \right. \\
&\quad \left. + \pi_1(1 | [0, 0]) \left((1 - c_1) P_{1,[0,0],1}^0 + c_1 P_{1,[0,0],1}^{\text{sc}} \right) \tilde{V}_2^{\pi,\sigma} \right] \\
&\stackrel{\text{(v)}}{=} 1 + \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[\pi_1(0 | [0, 0]) \left((1 - c_1) \tilde{V}_2^{\pi,\sigma}([0, 0]) + c_1 \tilde{V}_2^{\pi,\sigma}([1, 0]) \right) \right. \\
&\quad \left. + \pi_1(1 | [0, 0]) \left((1 - c_1) \tilde{V}_2^{\pi,\sigma}([0, 1]) + c_1 \tilde{V}_2^{\pi,\sigma}([1, 1]) \right) \right] \\
&= 1 + (T - 1) \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[\pi_1(0 | [0, 0]) (1 - c_1) + \pi_1(1 | [0, 0]) c_1 \right] \\
&= 1 + (T - 1) \pi_1(0 | [0, 0]) + (T - 1) \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} \left[c_1 (1 - 2\pi_1(0 | [0, 0])) \right], \tag{41}
\end{aligned}$$

889 where (i) holds by the definition in equation 5, (ii) follows from the *state-confounded* Bellman
890 consistency equation in equation 47, (iii) follows from $r_1([0, 0], a) = 1$ for all $a \in \{0, 1\}$ which
891 is independent from c_t . (iv) arises from the definition of \mathcal{P} in equation 20, (v) can be verified by
892 plugging in the definitions from equation 14 and equation 18, and the penultimate equality holds by
893 equation 40.

894 **Characterizing the optimal robust SC-value functions.**

895 To further consider equation 41, we recall the fact that $\mathcal{U}^\sigma(P_1^c) = \{P \in \Delta(\mathcal{C}) : \frac{1}{2} \|P - P_1^c\|_1 \leq \sigma_2\}$.

896 Observing from equation 41 that for any fixed $\pi_1(0 | [0, 0])$, $c_1(1 - 2\pi_1(0 | [0, 0]))$ is monotonously
 897 increasing with c_1 when $1 - 2\pi_1(0 | [0, 0]) \geq 0$ and decreasing with c_1 otherwise, it is easily
 898 verified that the solution of

$$f(\pi_1(0 | [0, 0])) := (T - 1) \inf_{P \in \mathcal{U}^\sigma(P_1^c)} \mathbb{E}_{c_1 \sim P} [c_1(1 - 2\pi_1(0 | [0, 0]))] \quad (42)$$

899 satisfies

$$f(\pi_1(0 | [0, 0])) = \begin{cases} 0 & \text{if } \pi_1(0 | [0, 0]) \geq \frac{1}{2} \\ (T - 1)\sigma_2(1 - 2\pi_1(0 | [0, 0])) & \text{otherwise} \end{cases}. \quad (43)$$

900 And note that the value of $\tilde{V}_1^{\pi, \sigma_2}([0, 0])$ only depends on $\pi_1(\cdot | [0, 0])$ which can be represent by
 901 $\pi_1(0 | [0, 0])$. Plugging in equation 43 into equation 41, we have that when $\pi_1(0 | [0, 0]) \geq \frac{1}{2}$,

$$\begin{aligned} \max_{\pi} \tilde{V}_1^{\pi, \sigma_2}([0, 0]) &= \max_{\pi_1(0 | [0, 0]) \geq \frac{1}{2}} 1 + (T - 1)\pi_1(0 | [0, 0]) + (T - 1)\sigma_2(1 - 2\pi_1(0 | [0, 0])) \\ &= 1 + (T - 1)\sigma_2 + (T - 1) \max_{\pi_1(0 | [0, 0]) \geq \frac{1}{2}} (1 - 2\sigma_2)\pi_1(0 | [0, 0]) \\ &= 1 + \frac{T - 1}{2}, \end{aligned} \quad (44)$$

902 where the last equality holds by $\sigma_2 > \frac{1}{2}$ and letting $\pi_1(0 | [0, 0]) = \frac{1}{2}$. Similarly, when $\pi_1(0 | [0, 0]) <$
 903 $\frac{1}{2}$,

$$\max_{\pi} \tilde{V}_1^{\pi, \sigma_2}([0, 0]) = \max_{\pi_1(0 | [0, 0]) < \frac{1}{2}} 1 + (T - 1)\pi_1(0 | [0, 0]) < 1 + \frac{T - 1}{2}. \quad (45)$$

904 Consequently, we complete the proof by concluding that

$$\tilde{V}_1^{\pi_{\text{RSC}}^*, \sigma_2}([0, 0]) = \tilde{V}_1^{*, \sigma_2}([0, 0]) = \max_{\pi} \tilde{V}_1^{\pi, \sigma_2}([0, 0]) = 1 + \frac{T - 1}{2}. \quad (46)$$

905 C.3 Auxiliary results of SC-MDPs and RSC-MDPs

906 **Facts about SC-MDPs.** For any state-confounded MDPs (SC-MDPs) $\mathcal{M}_{\text{SC}} = \{\mathcal{S}, \mathcal{A}, T, r,$
 907 $\mathcal{C}, \{\mathcal{P}_t^i\}, P^c\}$, denoting the optimal policy as π^* and the corresponding optimal SC-value func-
 908 tion as \tilde{V} , any policy π satisfies the corresponding *state-confounded* Bellman consistency equation as
 909 below:

$$\tilde{Q}_t^{\pi, P^c}(s, a) = r_t(s, a) + \mathbb{E}_{c_t \sim P_t^c} [\mathcal{P}_{t, s, a, c_t} \tilde{V}_{t+1}^{\pi, \sigma}], \quad (47)$$

910 where $\mathcal{P}_{t, s, a, c_t} \in \mathbb{R}^{1 \times S}$ such that $\mathcal{P}_{t, s, a, c_t}(s') := \mathcal{P}_t(s' | s, a, c_t)$ for $s' \in \mathcal{S}$.

911 **Facts about RSC-MDPs.** It is easily verified that for any RSC-MDP $\mathcal{M}_{\text{sc-rob}} = \{\mathcal{S}, \mathcal{A}, T, r,$
 912 $\mathcal{C}, \{\mathcal{P}_t^i\}, \mathcal{U}^{\sigma_2}(P^c)\}$, any policy π and the optimal policy π^* satisfy the corresponding *robust state-*
 913 *confounded* Bellman consistency equation and Bellman optimality equation shown below, respec-
 914 tively:

$$\begin{aligned} \tilde{Q}_t^{\pi, \sigma}(s, a) &= r_t(s, a) + \inf_{P \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P} [\mathcal{P}_{t, s, a, c_t} \tilde{V}_{t+1}^{\pi, \sigma}], \\ \tilde{Q}_t^{*, \sigma}(s, a) &= r_t(s, a) + \inf_{P \in \mathcal{U}^\sigma(P_t^c)} \mathbb{E}_{c_t \sim P} [\mathcal{P}_{t, s, a, c_t} \tilde{V}_{t+1}^{*, \sigma}], \end{aligned} \quad (48)$$

915 where $\mathcal{P}_{t, s, a, c_t} \in \mathbb{R}^{1 \times S}$ such that $\mathcal{P}_{t, s, a, c_t}(s') := \mathcal{P}_t(s' | s, a, c_t)$ for $s' \in \mathcal{S}$, and $\tilde{V}_{t+1}^{*, \sigma}(s) =$
 916 $\max_a \tilde{Q}_{t+1}^{*, \sigma}(s, a)$.

917 D Experiment Details

918 D.1 Architecture of the structural causal model

919 We plot the architecture of the structural causal model we used in our method in Figure 6. In normal
 920 neural networks, the input is treated as a whole to pass through linear layers or convolution layers.

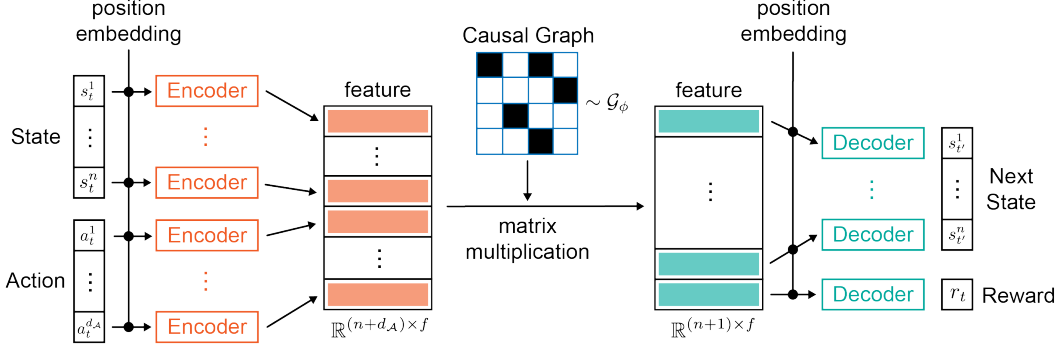


Figure 6: Model architecture of the structural causal model. Encoder, Decoder, position embedding, and Causal Graph are learnable during the training stage.

921 However, this structure blends all information in the input, making the causal graph useless to separate
 922 cause and effect. Thus, in our model, we design an encoder that is shared across all dimensions of the
 923 input. Since different dimensions could have exactly the same values, we add a learnable position
 924 embedding to the input of the encoder. In summary, the input dimension of the encoder is $1 + d_{pos}$,
 925 where d_{pos} is the dimension of the position embedding.

926 After the encoder, we obtain a set of independent features for each dimension of the input. We
 927 now multiply the features with a learnable binary causal graph \mathcal{G} . The element (i, j) of the graph
 928 is sampled from a Gumbel-Softmax distribution with parameter $\phi_{i,j}$ to ensure the loss function is
 929 differentiable w.r.t ϕ .

930 The multiplication of the causal graph and the input feature creates a linear combination of the input
 931 feature with respect to the causal graph. The obtained features are then passed through a decoder
 932 to predict the next state and reward. Again, the decoder is shared across all dimensions to avoid
 933 information leaking between dimensions. Position embedding is included in the input to the decoder
 934 and the output dimension of the decoder is 1.

935 D.2 Environments

936 We design four self-driving tasks in the Carla simulator [22] and four manipulation tasks in the
 937 Robosuite platform [23]. All of these realistic tasks contain strong spurious correlations that are
 938 explicit to humans. We provide detailed descriptions of all these environments in the following.

939 **Brightness.** The nominal environments are shown in the 1th column of Figure 7 where the brightness
 940 and the traffic density are correlated. When the ego vehicle drives in the daytime, there are many
 941 surrounding vehicles (first row). When the ego vehicle drives in the evening, there is no surrounding
 942 vehicle (second row). The shifted environment swaps the brightness and traffic density in the nominal
 943 environment, i.e., many surrounding vehicles in the evening and no surrounding vehicles in the
 944 daytime.

945 **Behavior.** The nominal environments are shown in the 2nd column of Figure 7, where the other
 946 vehicle has aggressive driving behavior. When the ego vehicle is in front of the other vehicle, the
 947 other vehicle always accelerates and overtakes the ego vehicle in the left lane. When the ego vehicle
 948 is behind the other vehicle, the other vehicle will always accelerate. In the shifted environment, the
 949 behavior of the other vehicle is conservative, i.e., the other vehicle always decelerates to block the
 950 ego vehicle.

951 **Crossing.** The nominal environments are shown in the 3rd column of Figure 7, where the pedes-
 952 trian follows the traffic rule and only cross the road when the traffic light is green. In the shifted
 953 environment, the pedestrian disobeys the traffic rule and crosses the road when the traffic light is red.

954 **CarType.** The nominal environments are shown in the 4th column of Figure 7, where the type of
 955 vehicle and the speed of the vehicle are correlated. When the vehicle is a truck, the speed is low and
 956 when the vehicle is a motorcycle, the speed is high. In the shifted environment, the truck drives very
 957 fast and the motorcycle drives very slow.

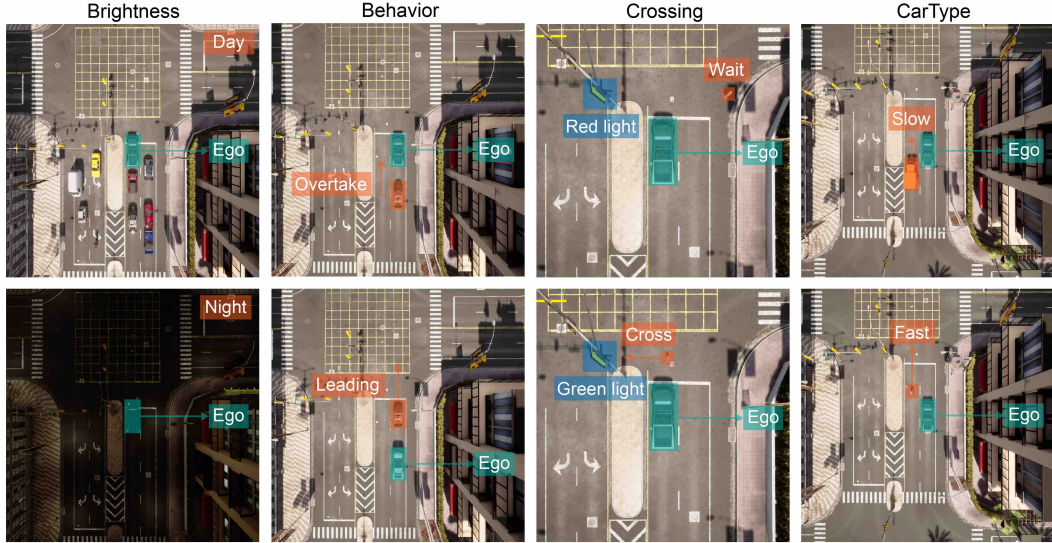


Figure 7: Illustration of tasks in the Carla simulator.

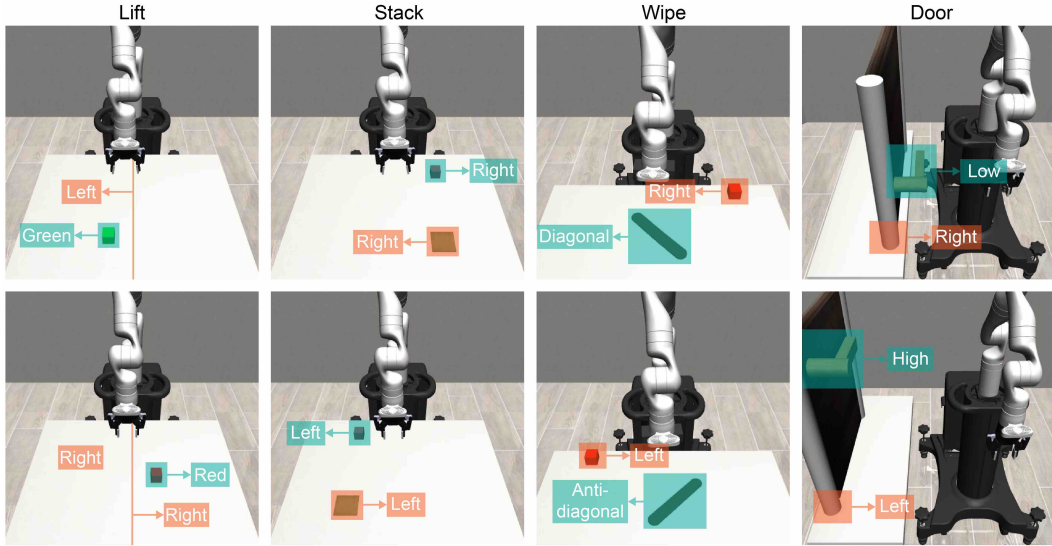


Figure 8: Illustration of tasks in the Robosuite simulator.

958 **Lift.** The nominal environments are shown in the 1th column of Figure 8, where the position of the
 959 cube and the color of the cube are correlated. When the cube is in the left part of the table, the color
 960 of the cube is green, when the cube is in the right part of the table, the color of the cube is red. The
 961 shifted environment swaps the color and position of the cube in the nominal environment, i.e., the
 962 cube is green when it is in the right part and the cube is red when it is in the left part.

963 **Stack.** The nominal environments are shown in the 2nd column of Figure 8, where the position of the
 964 red cube and green plate are correlated. When the cube is in the left part of the table, the plate is also
 965 in the left part; when the cube is in the right part of the table, the plate is also in the right part. In
 966 the shifted environment, the relative position of the cube and the plate changes, i.e., When the cube is in
 967 the left part of the table, the plate is in the right part; when the cube is in the right part of the table,
 968 the plate is in the left part.

969 **Wipe.** The nominal environments are shown in the 3rd column of Figure 8, where the shape of the
 970 dirty region is correlated to the position of the cube. When the dirty region is diagonal, the cube is
 971 on the right-hand side of the robot arm. When the dirty region is anti-diagonal, the cube is on the

972 left-hand side of the robot arm. In the shifted environment, the correlation changes, i.e., when the
 973 dirty region is diagonal, the cube is on the left-hand side of the robot arm; when the dirty region is
 974 anti-diagonal, the cube is on the right-hand side of the robot arm.

975 **Door.** The nominal environments are shown in the 4th column of Figure 8, where the height of the
 976 handle and the position of the door is correlated. When the door is closed to the robot arm, the handle
 977 is in a low position. When the door is far from the robot arm, the handle is in a high position. In
 978 the shifted environment, the correlation changes, i.e., when the door is closed to the robot arm, the
 979 handle is in a high position; when the door is far from the robot arm, the handle is in a low position.

980 D.3 Computation resources

981 Our algorithm is implemented on top of the Tianshou [88] package. All of our experiments are
 982 conducted on a machine with an Intel i9-9900K CPU@3.60GHz (16 core) CPU, an NVIDIA GeForce
 983 GTX 1080Ti GPU, and 64GB memory.

984 D.4 Hyperparameters

985 We summarize all hyper-parameters used in the Carla experiments (Table 5) and Robosuite experi-
 986 ments (Table 6). The source code of experiments will be released after double-blind reviewing.

Table 5: Hyper-parameters in Carla experiments

Parameters	Notation	Environment			
		Brightness	Behavior	Crossing	CarType
Horizon steps	T	100	100	100	100
State dimension	n	24	12	12	12
Action dimension	$d_{\mathcal{A}}$	2	2	2	2
Max training steps		1×10^5	1×10^5	5×10^5	5×10^5
Weight of $\ \mathcal{G}\ _p$	λ	0.1	-	-	-
norm of $\ \mathcal{G}\ _p$	p	0.1	-	-	-
Actor learning rate		3×10^{-4}	-	-	-
Critic learning rate		1×10^{-3}	-	-	-
Batch size		256	-	-	-
Discount factor	γ in SAC	0.99	-	-	-
Soft update weight	τ in SAC	0.005	-	-	-
Weight of entropy	α in SAC	0.1	-	-	-
Hidden layers		[256, 256, 256]	-	-	-
Returns estimation step		4	-	-	-
Buffer size		1×10^5	-	-	-
Steps per update		10	-	-	-

987 D.5 Discovered Causal Graph in SCM

988 To show the performance of our learned SCM, we plot the estimated causal graphs of all experiments
 989 in Figure 9, Figure 10, Figure 11, Figure 12, and Figure 13

Table 6: Hyper-parameters in Robosuite experiments

Parameters	Notation	Environment			
		Lift	Stack	Door	Wipe
Horizon steps	T	300	300	300	500
Control frequency (Hz)		20	20	20	20
State dimension	n	50	110	22	30
Action dimension	$d_{\mathcal{A}}$	4	4	8	7
Controller type		OSC position	OSC position	Joint velocity	Joint velocity
Max training steps		1×10^6	5×10^6	1×10^6	1×10^6
Weight of $\ \mathbf{G}\ _p$	λ	0.01	-	-	-
norm of $\ \mathbf{G}\ _p$	p	0.1	-	-	-
Actor learning rate		3×10^{-4}	-	-	-
Critic learning rate		1×10^{-3}	-	-	-
Batch size		128	-	-	-
Discount factor	γ in SAC	0.99	-	-	-
Soft update weight	τ in SAC	0.005	-	-	-
alpha learning rate	lr_{α} in SAC	3×10^{-4}	-	-	-
Hidden layers		[256, 256, 256]	-	-	-
Returns estimation step		4	-	-	-
Buffer size		1×10^6	-	-	-
Steps per update		10	-	-	-

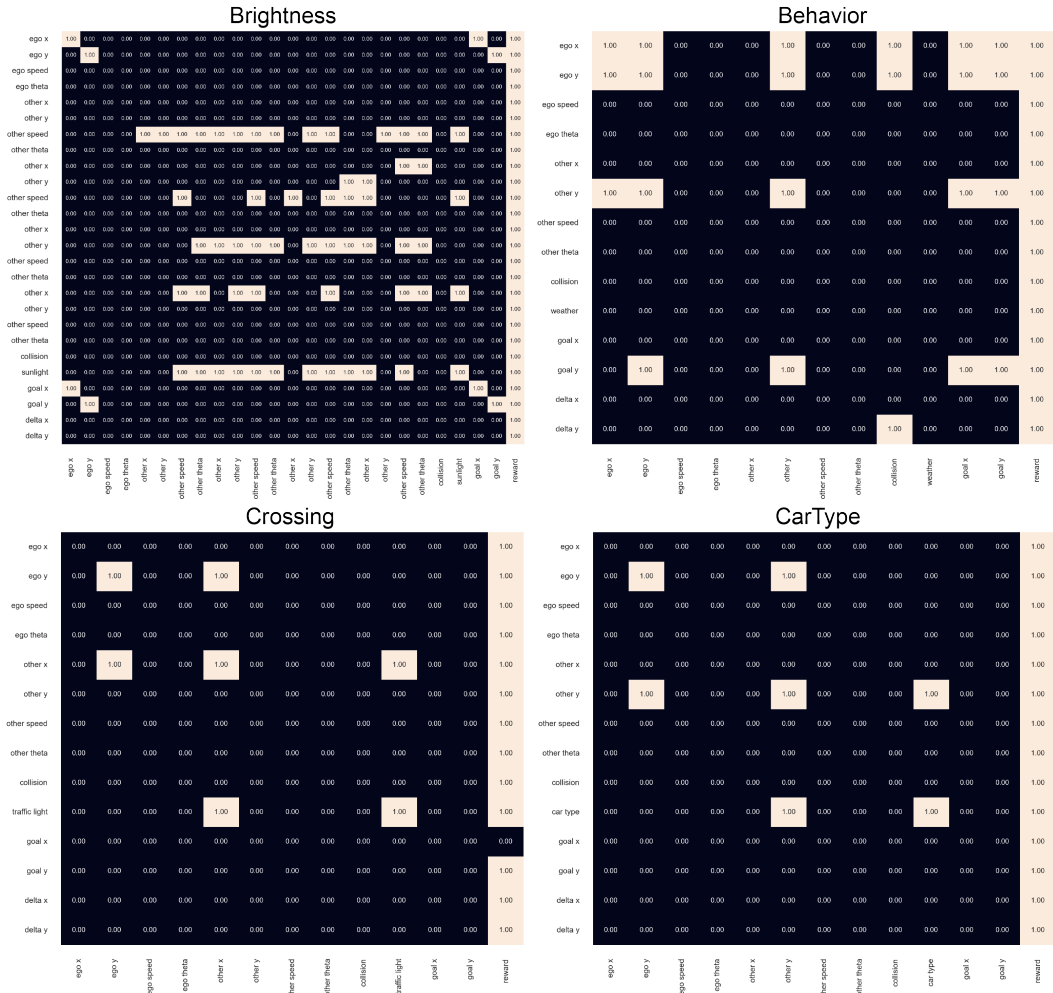


Figure 9: Estimated Causal Graphs of four tasks in Carla.

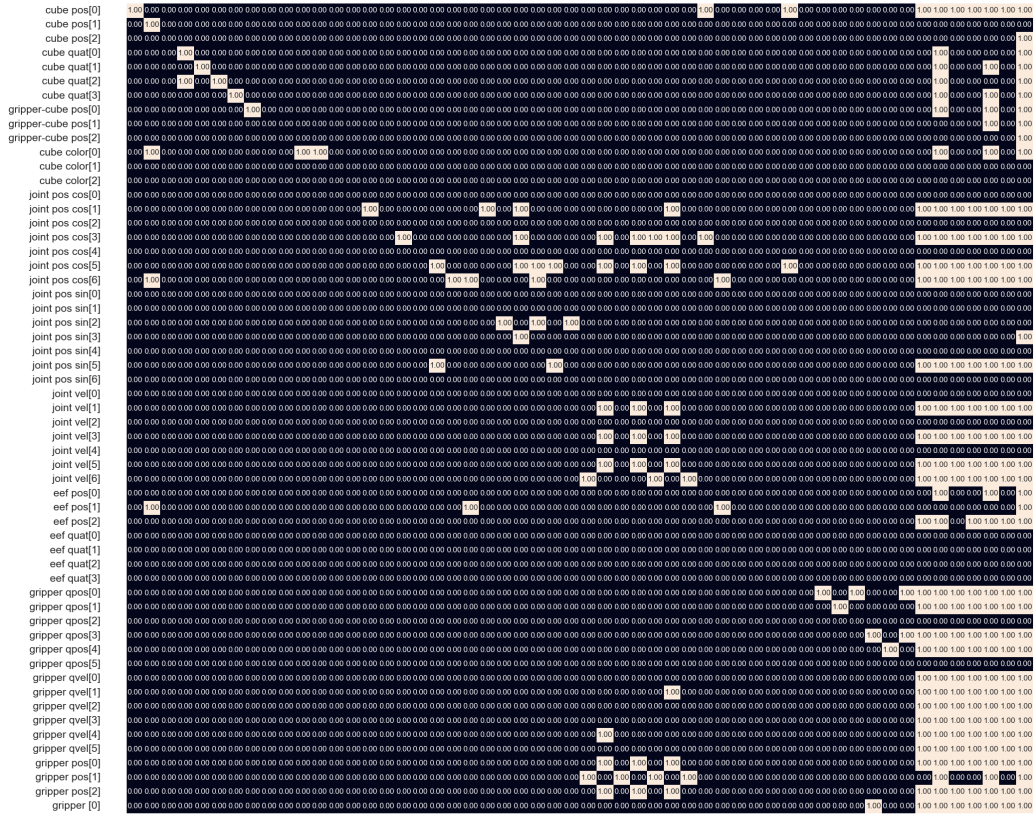


Figure 10: Estimated Causal Graphs of the Lift task in Robosuite.

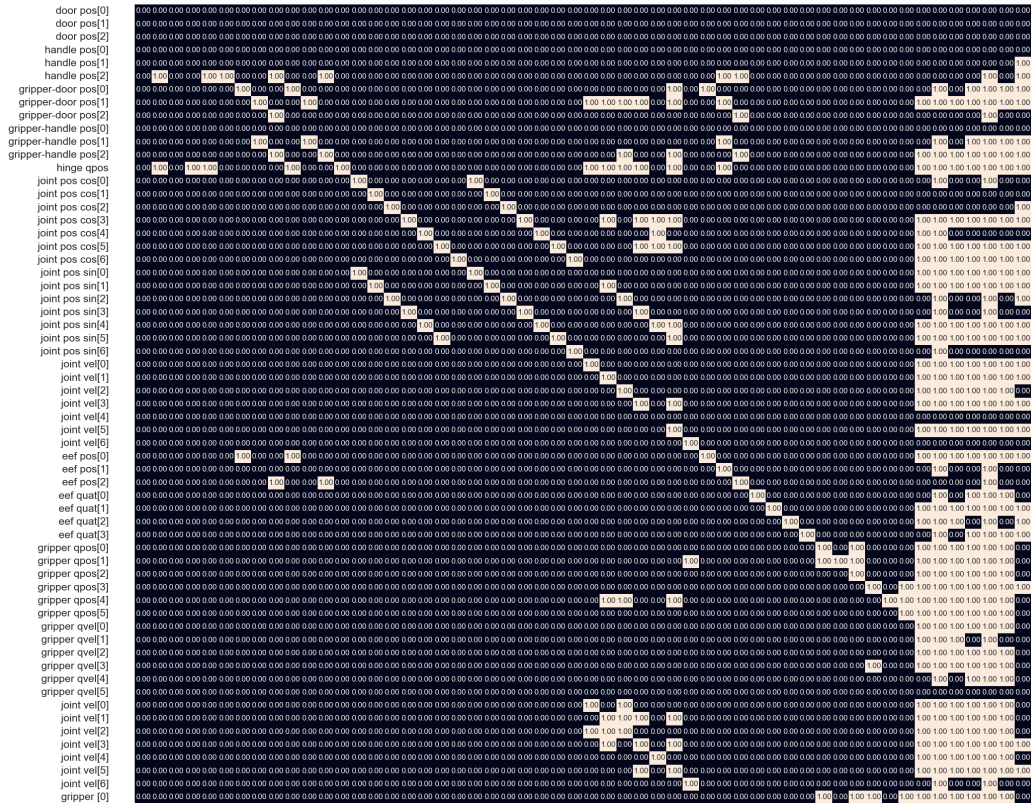


Figure 12: Estimated Causal Graphs of the Door task in Robosuite.



Figure 13: Estimated Causal Graphs of the Wipe task in Robosuite.