

---

# Counterfactual Data Augmentation using Locally Factored Dynamics

---

Silviu Pitis, Elliot Creager, Animesh Garg

Department of Computer Science, University of Toronto, Vector Institute  
{spitis, creager, garg}@cs.toronto.edu

## Abstract

Many dynamic processes, including common scenarios in robotic control and reinforcement learning (RL), involve a set of interacting subprocesses. Though the subprocesses are not independent, their interactions are often sparse, and the dynamics at any given time step can often be decomposed into *locally independent* causal mechanisms. Such local causal structures can be leveraged to improve the sample efficiency of sequence prediction and off-policy reinforcement learning. We formalize this by introducing local causal models (LCMs), which are induced from a global causal model by conditioning on a subset of the state space. We propose an approach to inferring these structures given an object-oriented state representation, as well as a novel algorithm for Counterfactual Data Augmentation (CoDA). CoDA uses local structures and an experience replay to generate counterfactual experiences that are causally valid in the global model. We find that CoDA significantly improves the performance of RL agents in locally factored tasks, including the batch-constrained and goal-conditioned settings.<sup>1</sup>

## 1 Introduction

High-dimensional dynamical systems are often composed of simple subprocesses that affect one another through sparse interaction. If the subprocesses *never* interacted, an agent could realize significant gains in sample efficiency by globally factoring the dynamics and modeling each subprocess independently [28, 29]. In most cases, however, the subprocesses do *eventually* interact and so the prevailing approach is to model the entire process using a monolithic, unfactored model. In this paper, we take advantage of the observation that *locally—during the time between their interactions—the subprocesses are causally independent*. By locally factoring dynamic processes in this way, we are able to capture the benefits of factorization even when their subprocesses interact on the global scale.

Consider a game of billiards, where each ball can be viewed as a separate physical subprocess. Predicting the opening break is difficult because all balls are mechanically coupled by their initial placement. Indeed, a dynamics model with dense coupling amongst balls may seem sensible when considering the expected outcomes over the course of the game, as each ball has a non-zero chance of colliding with the others. But at any given timestep, interactions between balls are usually sparse.

One way to take advantage of sparse interactions between otherwise disentangled entities is to use a structured state representation together with a graph neural network or other message passing transition model that captures the local interactions [26, 39]. When it is tractable to do so, such architectures can be used to model the world dynamics directly, producing transferable, task-agnostic models. In many cases, however, the underlying processes are difficult to model precisely, and model-free [46, 87] or task-oriented model-based [18, 63] approaches are less biased and exhibit superior performance. In this paper we argue that *knowledge of whether or not local interactions*

---

<sup>1</sup>Code available at <https://github.com/spitis/mrl>

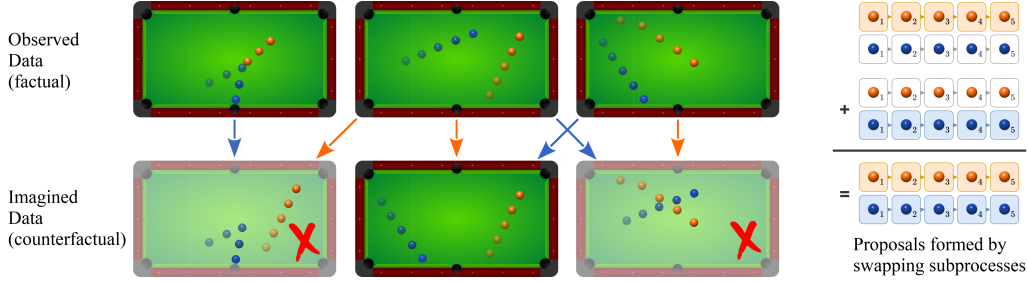


Figure 1: **Counterfactual Data Augmentation (CoDA)**. Given 3 factual samples, knowledge of the local causal structure lets us mix and match factored subprocesses to form counterfactual samples. The first proposal is rejected because one of its factual sources (the blue ball) is not locally factored. The third proposal is rejected because it is not itself factored. The second proposal is accepted, and can be used as additional training data for a reinforcement learning agent.

*occur is useful in and of itself*, and can be used to generate causally-valid counterfactual data even in absence of a forward dynamics model. In fact, if two trajectories have the same local factorization in their transition dynamics, then under mild conditions we can produce new counterfactually plausible data using our proposed **Counterfactual Data Augmentation (CoDA)** technique, wherein factorized subspaces of observed trajectory pairs are swapped (Figure 1). This lets us sample from a counterfactual data distribution by stitching together subsamples from observed transitions. Since CoDA acts only on the agent’s training data, it is compatible with any agent architecture (including unfactored ones).

In the remainder of this paper, we formalize this data augmentation strategy and discuss how it can improve performance of model-free RL agents in locally factored tasks. Our main contributions are:

1. We define local causal models (LCMs), which are induced from a global model by conditioning on a subset of the state space, and show how local structure can simplify counterfactual reasoning.
2. We introduce CoDA as a generalized data augmentation strategy that is able to leverage local factorizations to manufacture unseen, yet causally valid, samples of the environment dynamics. We show that goal relabeling [36, 1] and visual augmentation [2, 46] are instances of CoDA that use global independence relations and we propose a locally conditioned variant of CoDA that swaps independent subprocesses to form counterfactual experiences (Figure 1).
3. Using an attention-based method for discovering local causal structure in a disentangled state space, we show that our CoDA algorithm significantly improves the sample efficiency in standard, batch-constrained, and goal-conditioned reinforcement learning settings.

## 2 Local Causality in MDPs

### 2.1 Preliminaries and Problem Setup

The basic model for decision making in a controlled dynamic process is a Markov Decision Process (MDP), described by tuple  $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$  consisting of the state space, action space, transition function, reward function, and discount factor, respectively [72, 83]. Note that MDPs generalize uncontrolled Markov processes (set  $\mathcal{A} = \emptyset$ ), so that our work applies also to sequential prediction. We denote individual states and actions using lowercase  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ , and variables using the uppercase  $S$  and  $A$  (e.g.,  $s \in \text{range}(S) \subseteq \mathcal{S}$ ). A policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  defines a probability distribution over the agent’s actions at each state, and an agent is typically tasked with learning a parameterized policy  $\pi_\theta$  that maximizes value  $\mathbb{E}_{P, \pi} \sum_t \gamma^t R(s_t, a_t)$ .

In most non-trivial cases, the state  $s \in \mathcal{S}$  can be described as an object hierarchy together with global context. For instance, this decomposition will emerge naturally in any simulated process or game that is defined using a high-level programming language (e.g., the commonly used Atari [7] or Minecraft [35] simulators). In this paper we consider MDPs with a single, known top-level decomposition of the state space  $\mathcal{S} = \mathcal{S}^1 \oplus \mathcal{S}^2 \oplus \dots \oplus \mathcal{S}^n$  for fixed  $n$ , leaving extensions to hierarchical decomposition and

multiple representations [32, 17], dynamic factor count  $n$  [92], and (learned) latent representations [13] to future work. The action space might be similarly decomposed:  $\mathcal{A} = \mathcal{A}^1 \oplus \mathcal{A}^2 \oplus \dots \oplus \mathcal{A}^m$ .

Given such state and action decompositions, we may model time slice  $(t, t+1)$  using a structural causal model (SCM)  $\mathcal{M}_t = \langle V_t, U_t, \mathcal{F} \rangle$  ([65], Ch. 7) with directed acyclic graph (DAG)  $\mathcal{G}$ , where:

- $V_t = \{V_{t+1}^i\}_{i=0}^{2n+m} = \{S_t^1 \dots S_t^n, A_t^1 \dots A_t^m, S_{t+1}^1 \dots S_{t+1}^n\}$  are the nodes (variables) of  $\mathcal{G}$ .
- $U_t = \{U_{t+1}^i\}_{i=0}^{2n+m}$  is a set of noise variables, one for each  $V^i$ , determined by the initial state, past actions, and environment stochasticity. We assume that noise variables at time  $t+1$  are independent from other noise variables:  $U_{t+1}^i \perp\!\!\!\perp U_{t+1}^j \forall i, j$ . The instance  $u = (u^1, u^2, \dots, u^{2n+m})$  of  $U_t$  denotes an individual realization of the noise variables.
- $\mathcal{F} = \{f^i\}_{i=0}^{2n+m}$  is a set of functions (“structural equations”) that map from  $U_{t+1}^i \times \text{Pa}(V_{t+1}^i)$  to  $V_{t+1}^i$ , where  $\text{Pa}(V_{t+1}^i) \subset V_t \setminus V_{t+1}^i$  are the parents of  $V_{t+1}^i$  in  $\mathcal{G}$ ; hence each  $f^i$  is associated with the set of incoming edges to node  $V_{t+1}^i$  in  $\mathcal{G}$ ; see, e.g., Figure 2 (center).

Note that while  $V_t$ ,  $U_t$ , and  $\mathcal{M}_t$  are indexed by  $t$  (their distributions change over time), the structural equations  $f^i \in \mathcal{F}$  and causal graph  $\mathcal{G}$  represent the global transition function  $P$  and apply at all times  $t$ . To reduce clutter, we drop the subscript  $t$  on  $V$ ,  $U$ , and  $\mathcal{M}$  when no confusion can arise.

Critically, we require the set of edges in  $\mathcal{G}$  (and thus the number of inputs to each  $f^i$ ) to be *structurally minimal* ([67], Remark 6.6).

**Assumption (Structural Minimality).**  $V^j \in \text{Pa}(V^i)$  if and only if there exists some  $\{u^i, v^{-ij}\}$  with  $u^i \in \text{range}(U^i)$ ,  $v^{-ij} \in \text{range}(V \setminus \{V^i, V^j\})$  and pair  $(v_1^j, v_2^j)$  with  $v_1^j, v_2^j \in \text{range}(V^j)$  such that  $v_1^i = f^i(\{u^i, v^{-ij}, v_1^j\}) \neq f^i(\{u^i, v^{-ij}, v_2^j\}) = v_2^i$ .

Intuitively, structural minimality says that  $V^j$  is a parent of  $V^i$  if and only if setting the value of  $V^j$  can have a nonzero *direct* effect<sup>2</sup> on the child  $V^i$  through the structural equation  $f^i$ . The structurally minimal representation is unique [67].

Given structural minimality, we can think of edges in  $\mathcal{G}$  as representing global causal dependence. The probability distribution of  $S_{t+1}^i$  is fully specified by its parents  $\text{Pa}(S_{t+1}^i)$  together with its noise variable  $U_i$ ; that is, we have  $P(S_{t+1}^i | S_t, A_t) = P(S_{t+1}^i | \text{Pa}(S_{t+1}^i))$  so that  $S_{t+1}^i \perp\!\!\!\perp V^j | \text{Pa}(S_{t+1}^i)$  for all nodes  $V^j \notin \text{Pa}(S_{t+1}^i)$ . We call an MDP with this structure a *factored MDP* [37]. When edges in  $\mathcal{G}$  are sparse, factored MDPs admit more efficient solutions than unfactored MDPs [28].

## 2.2 Local Causal Models (LCMs)

**Limitations of Global Models** Unfortunately, even if states and actions can be cleanly decomposed into several nodes, in most practical scenarios the DAG  $\mathcal{G}$  is fully connected (or nearly so): since the  $f^i$  apply globally, so too does structural minimality, and edge  $(S_k^i, S_{k+1}^j)$  at time  $k$  is present so long as there is a single instance—at any time  $t$ , no matter how unlikely—in which  $S_t^i$  influences  $S_{t+1}^j$ . In the words of Andrew Gelman, “*there are (almost) no true zeros*” [24]. As a result, the factorized causal model  $\mathcal{M}_t$ , based on globally factorized dynamics, rarely offers an advantage over a simpler causal model that treats states and actions as monolithic entities (e.g., [12]).

**LCMs** Our key insight is that for each pair of nodes  $(V_t^i, S_{t+1}^j)$  with  $V_t^i \in \text{Pa}(S_{t+1}^j)$  in  $\mathcal{G}$ , there often exists a large subspace  $\mathcal{L}^{(j \perp i)} \subset \mathcal{S} \times \mathcal{A}$  for which  $S_{t+1}^j \perp\!\!\!\perp V_t^i | \text{Pa}(S_{t+1}^j) \setminus V_t^i, (s_t, a_t) \in \mathcal{L}^{(j \perp i)}$ . For example, in case of a two-armed robot (Figure 2), there is a large subspace of states in which the two arms are too far apart to influence each other physically. Thus, if we restrict our attention to  $(s_t, a_t) \in \mathcal{L}^{(j \perp i)}$ , we can consider a *local* causal model  $\mathcal{M}_t^{\mathcal{L}^{(j \perp i)}}$  whose local DAG  $\mathcal{G}^{\mathcal{L}^{(j \perp i)}}$  is strictly sparser than the global DAG  $\mathcal{G}$ , as the structural minimality assumption applied to  $\mathcal{G}^{\mathcal{L}^{(j \perp i)}}$  implies that there is no edge from  $V_t^i$  to  $S_{t+1}^j$ . More generally, for any subspace  $\mathcal{L} \subseteq \mathcal{S} \times \mathcal{A}$ , we can induce the Local Causal Model (LCM)  $\mathcal{M}_t^{\mathcal{L}} = \langle V_t^{\mathcal{L}}, U_t^{\mathcal{L}}, \mathcal{F}^{\mathcal{L}} \rangle$  with DAG  $\mathcal{G}^{\mathcal{L}}$  from the global model  $\mathcal{M}_t$  as:

<sup>2</sup>Thus parentage does describe knock-on effects, e.g.  $V_1$  on  $V_3$  in the Markov chain  $V_1 \rightarrow V_2 \rightarrow V_3$ .

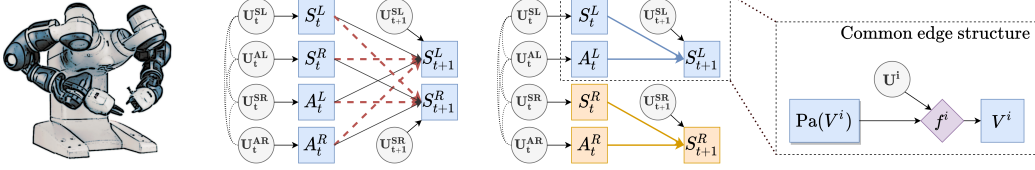


Figure 2: A two-armed robot (**left**) might be modeled as an MDP whose state and action spaces decompose into left and right subspaces:  $\mathcal{S} = \mathcal{S}^L \oplus \mathcal{S}^R, \mathcal{A} = \mathcal{A}^L \oplus \mathcal{A}^R$ . Because the arms can touch, the global causal model (**center left**) between time steps is fully connected, even though left-to-right and right-to-left connections (dashed red edges) are rarely active. By restricting our attention to the subspace of states in which left and right dynamics are independent we get a local causal model (**center right**) with two components that can be considered separately for training and inference.

- $V_t^{\mathcal{L}} = \{V_{t+1}^{\mathcal{L},i}\}_{i=0}^{2n+m}$ , where  $P(V_{t+1}^{\mathcal{L},i}) = P(V_{t+1}^i | (s_t, a_t) \in \mathcal{L})$ .
- $U_t^{\mathcal{L}} = \{U_{t+1}^{\mathcal{L},i}\}_{i=0}^{2n+m}$ , where  $P(U_{t+1}^{\mathcal{L},i}) = P(U_{t+1}^i | (s_t, a_t) \in \mathcal{L})$ .
- $\mathcal{F}^{\mathcal{L}} = \{f^{\mathcal{L},i}\}_{i=0}^{2n+m}$ , where  $f^{\mathcal{L},i} = f^i|_{\mathcal{L}}$  ( $f^i$  with range of input variables restricted to  $\mathcal{L}$ ). Due to structural minimality, the signature of  $f^{\mathcal{L},i}$  may shrink (as the range of the relevant variables is now restricted to  $\mathcal{L}$ ), and corresponding edges in  $\mathcal{G}$  will not be present in  $\mathcal{G}^{\mathcal{L}}$ .<sup>3</sup>

In case of the two-armed robot, conditioning on the arms being far apart simplifies the global DAG to a local DAG with two connected components (Figure 2). This can make counterfactual reasoning considerably more efficient: given a factual situation in which the robot’s arms are far apart, we can carry out separate counterfactual reasoning about each arm.

**Leveraging LCMs** To see the efficiency therein, consider a general case with global causal model  $\mathcal{M}$ . To answer the counterfactual question, “*what might the transition at time  $t$  have looked like if component  $S_t^i$  had value  $x$  instead of value  $y$ ?*”, we would ordinarily apply Pearl’s do-calculus to  $\mathcal{M}$  to obtain submodel  $\mathcal{M}_{\text{do}(S_t^i=x)} = \langle V, U, \mathcal{F}_x \rangle$ , where  $\mathcal{F}_x = \mathcal{F} \setminus f^i \cup \{S_t^i = x\}$  and incoming edges to  $S_t^i$  are removed from  $\mathcal{G}_{\text{do}(S_t^i=x)}$  [65]. The component distributions at time  $t + 1$  can be computed by reevaluating each function  $f^j$  that depends on  $S_t^i$ . When  $S_t^i$  has many children (as is often the case in the global  $\mathcal{G}$ ), this requires one to estimate outcomes for many structural equations  $\{f^j | V^j \in \text{Children}(V_t^i)\}$ . But if both the original value of  $S_t$  (with  $S_t^i = y$ ) and its new value (with  $S_t^i = x$ ) are in the set  $\mathcal{L}$ , the intervention is “within the bounds” of local model  $\mathcal{M}^{\mathcal{L}}$  and we can instead work directly with local submodel  $\mathcal{M}_{\text{do}(S_t^i=x)}^{\mathcal{L}}$  (defined accordingly). The validity of this follows from the definitions: since  $f^{\mathcal{L},j} = f^j|_{\mathcal{L}}$  for all of  $S_t^i$ ’s children, the nodes  $V_t^k$  for  $k \neq i$  at time  $t$  are held fixed, and the noise variables at time  $t + 1$  are unaffected, the distribution at time  $t + 1$  is the same under both models. When  $S_t^i$  has fewer children in  $\mathcal{M}^{\mathcal{L}}$  than in  $\mathcal{M}$ , this reduces the number of structural equations that need to be considered.

### 3 Counterfactual Data Augmentation

We hypothesize that local causal models will have several applications, and potentially lead to improved agent designs, algorithms, and interpretability. In this paper we focus on improving off-policy learning in RL by exploiting causal independence in local models for **Counterfactual Data Augmentation (CoDA)**. CoDA augments real data by making counterfactual modifications to a subset of the causal factors at time  $t$ , leaving the rest of the factors untouched. Following the logic outlined in the Subsection 2.2, this can be understood as manufacturing “fake” data samples using the counterfactual model  $\mathcal{M}_{\text{do}(S_t^{i \dots j} = x)}^{\mathcal{L}}$ , where we modify the causal factors  $S_t^{i \dots j}$  and resample their children. While this is always possible using a model-based approach if we have good models of the structural equations, it is particularly nice when the causal mechanisms are independent, as we can do counterfactual reasoning directly by reusing subsamples from observed trajectories.

<sup>3</sup>As a trivial example, if  $f^i$  is a function of binary variable  $V^j$ , and  $\mathcal{L} = \{(s, a) | V^j = 0\}$ , then  $f^{\mathcal{L},i}$  is not a function of  $V^j$  (which is now a constant), and there is no longer an edge from  $V^j$  to  $V^i$  in  $\mathcal{G}^{\mathcal{L}}$ .



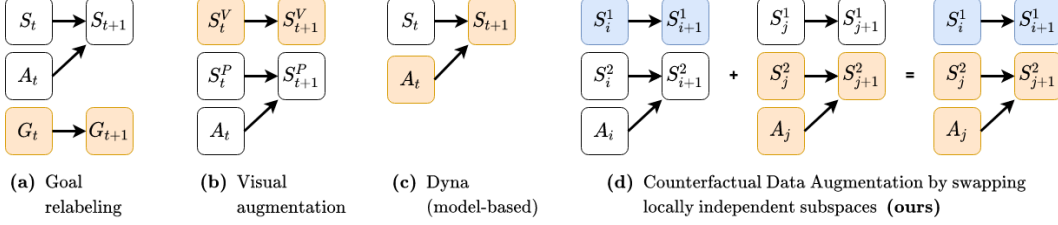


Figure 3: *Four instances of CoDA; orange nodes are relabeled, noise variables omitted for clarity. **First:** Goal relabeling [36], including HER [1], augments transitions with counterfactual goals. **Second:** Visual feature augmentation [2, 46] uses domain knowledge to change visual features  $S_t^V$  (such as textures, lighting, and camera positions) that the designer knows do not impact the physical state  $S_{t+1}^P$ . **Third:** Dyna [82], including MBPO [34], augments real states with new actions and resamples the next state using a learned dynamics model. **Fourth (ours):** Given two transitions that share local causal structures, we propose to swap connected components to form new transitions.*

**Definition.** The causal mechanisms represented by subgraphs  $\mathcal{G}_i, \mathcal{G}_j \subset \mathcal{G}$  are **independent** when  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are disconnected in  $\mathcal{G}$ .

When  $\mathcal{G}$  is divisible into two (or more) connected components, we can think of each subgraph as an independent causal mechanism that can be reasoned about separately.

Existing data augmentation techniques can be interpreted as specific instances of CoDA (Figure 3). For example, goal relabeling [36], as used in Hindsight Experience Replay (HER) [1] and Q-learning for Reward Machines [33], exploits the independence of the goal dynamics  $G_t \mapsto G_{t+1}$  (identity map) and the next state dynamics  $S_t \times A_t \mapsto S_{t+1}$  in order to relabel the goal variable  $G_t$  with a counterfactual goal. While the goal relabeling is done model-free, we typically assume knowledge of the goal-based reward mechanism  $G_t \times S_t \times A_t \times S_{t+1} \mapsto R_{t+1}$  to relabel the reward, ultimately mixing model-free and model-based reasoning. Similarly, visual feature augmentation, as used in reinforcement learning from pixels [46, 41] and sim-to-real transfer [2], exploits the independence of the physical dynamics  $S_t^P \times A_t \mapsto S_{t+1}^P$  and visual feature dynamics  $S_t^V \mapsto S_{t+1}^V$  such as textures and camera position, assumed to be static ( $S_{t+1}^V = S_t^V$ ), to counterfactually augment visual features. Both goal relabeling and visual data augmentation rely on *global* independence relationships.

We propose a novel form of Counterfactual Data Augmentation that uses knowledge of *local* independence relationships. In particular, we observe that whenever an environment transition is within the bounds of some local model  $\mathcal{M}^{\mathcal{L}}$  whose graph  $\mathcal{G}^{\mathcal{L}}$  has the locally independent causal mechanism  $\mathcal{G}_i$  as a disconnected subgraph (note:  $\mathcal{G}_i$  itself need not be connected), that transition contains an unbiased sample from  $\mathcal{G}_i$ . Thus, given two transitions in  $\mathcal{L}$ , we may mix and match the samples of  $\mathcal{G}_i$  to generate counterfactual data, so long as the resulting transitions are themselves in  $\mathcal{L}$ .

**Remark 3.1.** How much data can we generate using our CoDA algorithm? If we have  $n$  independent samples from subspace  $\mathcal{L}$  whose graph  $\mathcal{G}^{\mathcal{L}}$  has  $m$  connected components, we have  $n$  choices for each of the  $m$  components, for a total of  $n^m$  CoDA samples—an **exponential** increase in data! One might term this the “blessing of independent subspaces.”

**Remark 3.2.** Our discussion has been at the level of a single transition (time slice  $(t, t + 1)$ ), which is consistent with the form of data that RL agents typically consume. But we could also use CoDA to mix and match locally independent components over several time steps (see, e.g., Figure 1).

**Remark 3.3.** As is typical, counterfactual reasoning changes the data distribution. While off-policy agents are typically robust to distributional shift, future work might explore different ways to control or prioritize the counterfactual data distribution [77, 43]. We note, however, that certain prioritization schemes may introduce selection bias [31], effectively entangling otherwise independent causal mechanisms (e.g., HER’s “future” strategy [1] may introduce “hindsight bias” [45, 78]).

**Remark 3.4.** The global independence relations relied upon by goal relabeling and image augmentation are incredibly general, as evidenced by their wide applicability. We posit that certain local independence relations are similarly general. For example, the physical independence of objects separated by space (the billiards balls of Figure 1, the two-armed robot of Figure 2, and the environments used in Section 4), and the independence between an agent’s actions and the truth of (but not belief about) certain facts the agent is ignorant of (e.g., an opponent’s true beliefs).

---

**Algorithm 1** Mask-based Counterfactual Data Augmentation (CoDA)

---

<b>function</b> CODA(transition $t_1$ , transition $t_2$ ): $s_1, a_1, s_1' \leftarrow t_1$ $s_2, a_2, s_2' \leftarrow t_2$ $m_1, m_2 \leftarrow \text{MASK}(s_1, a_1), \text{MASK}(s_2, a_2)$ $D_1 \leftarrow \text{COMPONENTS}(m_1)$ $D_2 \leftarrow \text{COMPONENTS}(m_2)$ $d \leftarrow \text{random sample from } (D_1 \cap D_2)$ $\tilde{s}, \tilde{a}, \tilde{s}' \leftarrow \text{copy}(s_1, a_1, s_1')$ $\tilde{s}[d], \tilde{a}[d], \tilde{s}'[d] \leftarrow s_2[d], a_2[d], s_2'[d]$ $\tilde{D} \leftarrow \text{COMPONENTS}(\text{MASK}(\tilde{s}, \tilde{a}))$ <b>return</b> $(\tilde{s}, \tilde{a}, \tilde{s}')$ <b>if</b> $d \in \tilde{D}$ <b>else</b> $\emptyset$	<b>function</b> MASK(state $s$ , action $a$ ): Returns $(n+m) \times (n)$ matrix indicating if the $n$ next state components (columns) locally depend on the $n$ state and $m$ action components (rows).  <b>function</b> COMPONENTS(mask $m$ ): Using the mask as the adjacency matrix for $\mathcal{G}^{\mathcal{L}}$ (with dummy columns for next action), finds the set of connected components $C = \{C_j\}$ , and returns the set of independent components $D = \{\mathcal{G}_i = \bigcup_k C_k^i \mid C^i \subset \text{powerset}(C)\}$ .
---	---

---

**Implementing CoDA** We implement CoDA, as outlined above and visualized in Figure 3(d), as a function of two factual transitions and a mask function  $M(s_t, a_t) : \mathcal{S} \times \mathcal{A} \rightarrow \{0, 1\}^{(n+m) \times n}$  that represents the adjacency matrix of the sparsest local causal graph  $\mathcal{G}^{\mathcal{L}}$  such that  $\mathcal{L}$  is a neighborhood of  $(s_t, a_t)$ .<sup>4</sup> We apply  $M$  to each transition to obtain local masks  $m_1$  and  $m_2$ , compute their connected components, and swap independent components  $\mathcal{G}_i$  and  $\mathcal{G}_j$  (mutually disjoint and collectively exhaustive groups of connected components) between the transitions to produce a counterfactual proposal. We then apply  $M$  to the counterfactual  $(\tilde{s}_t, \tilde{a}_t)$  to validate the proposal—if the counterfactual mask  $\tilde{m}$  shares the same graph partitions as  $m_1$  and  $m_2$ , we accept the proposal as a CoDA sample. See Algorithm 1.

Note that masks  $m_1, m_2$  and  $\tilde{m}$  correspond to different neighborhoods  $\mathcal{L}_1, \mathcal{L}_2$  and  $\tilde{\mathcal{L}}$ , so it is not clear that we are “within the bounds” of any model  $\mathcal{M}^{\mathcal{L}}$  as was required in Subsection 2.2 for valid counterfactual reasoning. To correct this discrepancy we use the following proposition and additionally require the causal mechanisms (subgraphs) for independent components  $\mathcal{G}_i$  and  $\mathcal{G}_j$  to share structural equations in each local neighborhood:  $f^{\mathcal{L}_1, i} = f^{\mathcal{L}_2, i} = f^{\tilde{\mathcal{L}}, i}$  and  $f^{\mathcal{L}_1, j} = f^{\mathcal{L}_2, j} = f^{\tilde{\mathcal{L}}, j}$ .<sup>5</sup> This makes our reasoning valid in the local subspace  $\mathcal{L}^* = \mathcal{L}_1 \cup \mathcal{L}_2 \cup \tilde{\mathcal{L}}$ . See Appendix A for proof.

**Proposition 1.** *The causal mechanisms represented by  $\mathcal{G}_i, \mathcal{G}_j \subset \mathcal{G}$  are independent in  $\mathcal{G}^{\mathcal{L}_1 \cup \mathcal{L}_2}$  if and only if  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are independent in both  $\mathcal{G}^{\mathcal{L}_1}$  and  $\mathcal{G}^{\mathcal{L}_2}$ , and  $f^{\mathcal{L}_1, i} = f^{\mathcal{L}_2, i}, f^{\mathcal{L}_1, j} = f^{\mathcal{L}_2, j}$ .*

Since CoDA only modifies data within local subspaces, this biases the resulting replay buffer to have more factorized transitions. In our experiments below, we specify the ratio of observed-to-counterfactual data heuristically to control this selection bias, but find that off-policy agents are reasonably robust to large proportions of CoDA-sampled trajectories. We leave a full characterization of the selection bias in CoDA to future studies, noting that knowledge of graph topology was shown to be useful in mitigating selection bias for causal effect estimation [5, 6].

**Inferring local factorization** While the ground truth mask function  $M$  may be available in rare cases as part of a simulator, the general case either requires a domain expert to specify an approximate causal model (as in goal relabeling and visual data augmentation) or requires the agent to learn the local factorization from data. Given how common independence due to physical separation of objects is, the former option will often be available. In the latter case, we note that the same data could also be used to learn a forward model. Thus, there is an implicit assumption in the latter case that learning the local factorization is easier than modeling the dynamics. We think this assumption is rather mild, as an accurate forward dynamics model would subsume the factorization, and we provide some empirical evidence of its validity in Section 4.

Learning the local factorization is similar to conditional causal structure discovery [81, 75, 67], conditioned on neighborhood  $\mathcal{L}$  of  $(s_t, a_t)$ , except that the same structural equations must be applied globally (if the structural equations were conditioned on  $\mathcal{L}$ , Proposition 1 would fail). As there are many algorithms for general structure discovery [81, 75], and the arrow of time simplifies the inquiry

---

<sup>4</sup>If Jacobian  $\partial P / \partial x$  exists at  $x = (s_t, a_t)$ , the ground truth  $M(s_t, a_t)$  equals  $|(\partial P / \partial x)^T| > 0$ .

<sup>5</sup>To see why this is not trivially true, imagine there are two rooms, one of which is icy. In either room the ground conditions are locally independent of movement dynamics, but not so if we consider their union.

[27, 67], there may be many ways to approach this problem. For now, we consider a generalization of the global network mask approach used by MADE [25] (eq. 10) for autoregressive distribution modeling and GraN-DAG [44] (eq. 6) for causal discovery, which additionally conditions the mask on the current state and action.

This approach computes a locally conditioned network mask  $M(s_t, a_t)$  by taking the matrix product of locally conditioned layer masks:  $M(s_t, a_t) = \Pi_{\ell=1}^L M_\ell(s_t, a_t)$ . This mask can be understood as an upper bound on the network’s absolute Jacobian (see Appendix C). Again, there may be several models allow one to compute conditional layer masks. We tested two such models: a mixture of MLP experts and a single-head set transformer architecture [85, 47]. Each is described in more detail in Appendix C. Both are *trained* to model forward dynamics using an L2 prediction loss and induce a sparse network mask either via a sparsity penalty (in case of the mixture of experts model) or via a sparse attention mechanism (in case of the set transformer). In preliminary experiments (Appendix C) we found that the set transformer performed better, and proceed to use it in our main experiments (Section 4). The set transformer uses the attention mask at each layer as the layer mask for that layer, so that the network mask is simply the product of the attention masks. Though trained to model forward dynamics, the CoDA models are used by the agent to *infer* local factorization rather than to directly sample future states as is typical in model-based RL. We found this produced reasonable results in the tested domains (below). See Appendix C for details. Future work should consider other approaches to inferring local structure such as graph neural networks [39, 13].

## 4 Experiments

Our experiments evaluate CoDA in the online, batch, and goal-conditioned settings, in each case finding that CoDA significantly improves agent performance as compared to non-CoDA baselines. Since CoDA only modifies an agent’s training data, we expect these improvements to extend to other off-policy task settings in which the state space can be accurately disentangled. Below we outline our experimental design and results, deferring specific details and additional results to Appendix B.

**Standard online RL** We extend *Spriteworld* [89] to construct a “bouncing ball” environment (right), that consists of multiple objects (sprites) that move and collide within a confined 2D canvas. We use tasks of varying difficulty, where the agent must navigate  $N \in \{1, 2, 3, 4\}$  of 4 sprites to their fixed target positions. The agent receives reward of  $1/N$  for each of the  $N$  sprites placed; e.g., the hardest task (Place 4) gives  $1/4$  reward for each of 4 sprites placed. For each task, we use CoDA to expand the replay buffer of a TD3 agent [22] by about 8 times. We compare CoDA with a ground truth masking function (available via the *Spriteworld* environment) and learned masking function to the base TD3 agent, as well as a Dyna agent that generates additional training data by sampling from a model. For fair comparison, we use the same transformer used for CoDA masks for Dyna, which we pretrain using approximately 42,000 samples from a random policy. As in HER, we assume access to the ground truth reward function to relabel the rewards. The results in Figure 4 show that both variants of CoDA significantly improve sample complexity over the baseline. By contrast, the Dyna agent suffers from model bias, even though it uses the same model as CoDA.

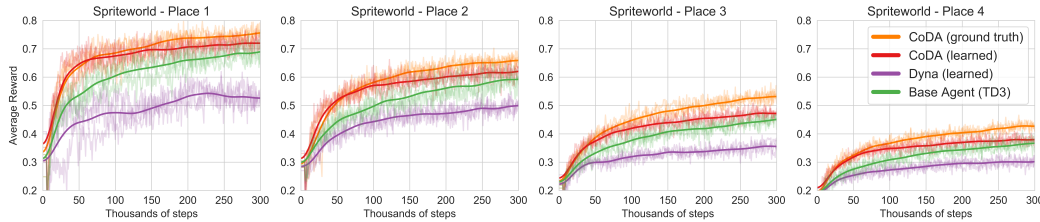
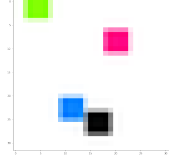


Figure 4: **Standard online RL** (3 seeds): CoDA with the ground truth mask always performs the best, validating our basic idea. CoDA with a pretrained model also offers a significant early boost in sample efficiency and maintains its lead over the base TD3 agent throughout training. Using the same model to generate data directly (a la Dyna [82]) performs poorly, suggesting significant model bias.

$ \mathcal{D} $ (1000s)	Real data 1R	MBPO 1R:1M	Ratio of Real:CoDA [:MBPO] data (ours)			
			1R:1C	1R:3C	1R:5C	1R:3C:1M
25	13.2 $\pm$ 0.7	18.5 $\pm$ 1.5	43.8 $\pm$ 2.8	40.9 $\pm$ 2.5	38.4 $\pm$ 4.9	<b>46.8 <math>\pm</math> 3.1</b>
50	22.8 $\pm$ 3.0	36.6 $\pm$ 4.3	66.6 $\pm$ 3.8	64.4 $\pm$ 3.1	62.5 $\pm$ 3.5	<b>70.4 <math>\pm</math> 3.8</b>
75	43.2 $\pm$ 4.9	46.0 $\pm$ 4.7	73.4 $\pm$ 2.8	<b>76.7 <math>\pm</math> 2.6</b>	75.0 $\pm$ 3.4	74.6 $\pm$ 3.2
100	63.0 $\pm$ 3.1	66.4 $\pm$ 4.9	77.8 $\pm$ 2.0	<b>82.7 <math>\pm</math> 1.5</b>	76.6 $\pm$ 3.0	73.7 $\pm$ 2.9
150	77.4 $\pm$ 1.2	72.6 $\pm$ 5.6	82.2 $\pm$ 1.8	<b>85.8 <math>\pm</math> 1.4</b>	84.2 $\pm$ 1.0	79.7 $\pm$ 3.6
250	78.2 $\pm$ 2.7	77.9 $\pm$ 2.4	85.0 $\pm$ 2.9	<b>87.8 <math>\pm</math> 1.8</b>	87.0 $\pm$ 1.0	78.3 $\pm$ 4.9

Table 1: **Batch RL** (10 seeds): Mean success ( $\pm$  standard error, estimated using 1000 bootstrap resamples) on Pong environment. CoDA with learned masking function more than doubles the effective data size, resulting in a 3x performance boost at smaller data sizes. Note that a 1R:5C Real:CoDA ratio performs slightly worse than a 1R:3C ratio due to distributional shift (Remark 3.3).

**Batch RL** A natural setting for CoDA is batch-constrained RL, where an agent has access to an existing transition-level dataset, but cannot collect more data via exploration [21, 48]. This makes any additional, high quality data invaluable. Another reason why CoDA is attractive in this setting is that there is no *a priori* reason to prefer the given batch data distribution to a counterfactual one. For this experiment we use a continuous control Pong environment based on RoboschoolPong [40]. The agent must hit the ball past the opponent, receiving reward of +1 when the ball is behind the opponent’s paddle, -1 when the ball is behind the agent’s paddle, and 0 otherwise. Since our transformer model performed poorly when used as a dynamics model, our Dyna baseline for batch RL adopts a state-of-the-art architecture [34] that employs a 7-model ensemble (MBPO). We collect datasets of up to 250,000 samples from an pretrained policy with added noise. For each dataset, we train both mask and reward functions (and in case of MBPO, the dynamics model) on the provided data and use them to generate different amounts of counterfactual data. We also consider combining CoDA with MBPO, by first expanding the dataset with MBPO and then applying CoDA to the result. We train the same TD3 agent on the expanded datasets in batch mode for 500,000 optimization steps. The results in Table 1 show that with only 3 state factors (two paddles and ball), applying CoDA is approximately equivalent to doubling the amount of real data.

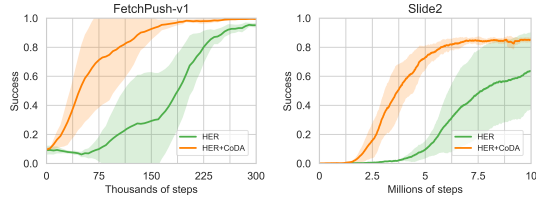
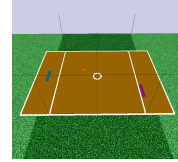
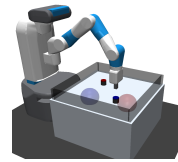


Fig. 5: **Goal-conditioned RL** (5 seeds): In FetchPush and the challenging Slide2 environment, a HER agent whose dataset has been enlarged with CoDA approximately doubles the sample efficiency of the base HER agent.

**Goal-conditioned RL** As HER [1] is an instance of prioritized CoDA that greatly improves sample efficiency in sparse-reward tasks, can our unprioritized CoDA algorithm further improve HER agents? We use HER to relabel goals on real data only, relying on random CoDA-style goal relabeling for CoDA data. After finding that CoDA obtains state-of-the-art results in FetchPush-v1 [71], we show that CoDA also accelerates learning in a novel and significantly more challenging Slide2 environment, where the agent must slide two pucks onto their targets (Figure 5). For this experiment, we specified a heuristic mask using domain knowledge (“objects are disentangled if more than 10cm apart”) that worked in both FetchPush and Slide2 despite different dynamics.



## 5 Related Work

Factored MDPs [28, 29, 90] consider MDPs where state variables are only influenced by a fixed subset of “parent” variables at the previous timestep. The notion of “context specific independence” (CSI), which was used to compactly represent single factors of a Bayes net [10] or MDP [9] for efficient inference and model storage,<sup>6</sup> is closely related to the local factorizations we study in this paper. CSI can be understood as going one step beyond CoDA, exploiting not only knowledge of

<sup>6</sup>Note that these methods were proposed to efficiently encode conditional probability tables at the graph nodes, requiring that all variables considered be discrete; CoDA on the other hand works in continuous tasks.

the local factorization, but also the structural equations at play in the local factorization; CSI could be leveraged for model-based RL approaches where faithful models of factored dynamics can be realized. Object-oriented and relational approaches to RL and prediction [15, 26, 39, 50, 93, 94] represent the dynamics as a set of interacting entities. Factored actions and policies have been used to formulate dimension-wise policy gradient baselines in standard and multi-agent settings [19, 54, 91].

A growing body of work applies causal reasoning the RL setting to improve sample efficiency, interpretability and learn better representations [55, 57, 74]. Particularly relevant is the work by Buesing et al. [12], which improves sample efficiency by using a causal model to sample counterfactual trajectories, thereby reducing variance of off-policy gradient estimates in a guided policy search framework. These counterfactuals use coarse-grained representations at the trajectory level, while our approach uses factored representations within a single transition. Batch RL [48, 21, 58] and more generally off-policy RL [88, 60] are counterfactual by nature, and are particularly important when it is costly or dangerous to obtain on-policy data [84]. The use of counterfactual goals to accelerate learning goal-conditioned RL [36, 76, 71] is what inspired our local CoDA algorithm.

Data augmentation is also widely used in supervised learning, and is considered a required best practice in high dimensional problems [42, 46, 66]. Heuristics for data augmentation often encode a causal invariance statement with respect to certain perturbations on the inputs. Thus model performance on counterfactual/augmented data can be seen as a measure of *robustness*. Assuring that models perform robustly in this sense is relevant to applications where *fairness* is a concern, as counterfactuals can be used to achieve robust performance and debias data [23, 64, 8].

## 6 Conclusion

In this paper we proposed a local causal model (LCM) framework that captures the benefits of decomposition in settings where the global causal model is densely connected. We used our framework to design a local Counterfactual Data Augmentation (CoDA) algorithm that expands available training data with counterfactual samples by stitching together locally independent subsamples from the environment. Empirically, we showed that CoDA can more than double the sample efficiency and final performance of reinforcement learning agents in locally factored environments.

There are several interesting avenues for future work. First, the sizable gap between ground truth and learned CoDA in our Spriteworld results suggest there is room for improvement in our approach to learning the masking function. Second, we have applied CoDA in a random, unprioritized fashion, but past work [1, 77] suggests there is significant benefit to prioritization. Third, we have applied CoDA in a way that might be considered model-free, insofar as we reuse subsamples from the environment dynamics rather than generating samples using our models of the causal mechanisms. However, our LCM formalism allows for mixing model-free and model-based reasoning, which could further improve sample efficiency. Fourth, we used fully-observable, disentangled state spaces with a fixed top-level decomposition, but ultimately we would like to deploy CoDA in partially observable, entangled settings (e.g. RL from pixels) with multiple dynamic, multi-level decompositions [32, 17, 92]. Unsupervised learning of factorized latent representations is an active area of research [16, 39, 52, 53, 89], and it would be interesting to combine these methods with CoDA. Finally, it would be interesting to explore applications of our LCM framework to other areas such as interpretability [56, 59], exploration [61, 86], and off-policy evaluation [84].

## Acknowledgments and Disclosure of Funding

We thank Jimmy Ba, Harris Chan, Seyed Kamyar Seyed Ghasemipour, James Lucas, David Madras, Yuhuai Wu and Lunjun Zhang for helpful comments and discussions. We also thank the anonymous reviewers for their feedback, which improved the final manuscript. SP is supported by an NSERC PGS-D award. EC is a student researcher at Google Brain in Toronto. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute (<https://vectorinstitute.ai/partners>).

## Broader Impact

**Considerations related to counterfactual reasoning** CoDA transforms the observational data distribution into a counterfactual one. This incurs several risks and benefits, listed below.

- Modern machine learning models and reinforcement learning agents often generalize poorly under distributional shift (sometimes called “covariate shift”) [49, 14]. CoDA has the potential to produce out-of-distribution data, that would never show up in the observational distribution. Thus, care should be taken when applying agent modules that have been trained only on observational data to counterfactual data, as their performance could decline sharply, thereby creating safety risks. We anticipate that work on uncertainty will be essential to controlling the risks associated with distributional shift [80].
- The fact that CoDA creates distributional shift can also provide benefits in the form of distributional robustness and fairness. Training on out-of-distribution data can make models more robust [79], increasing their trustworthiness and practical applicability. As noted in Section 5, counterfactual reasoning can be leveraged in areas where *fairness* is a concern. For example, [64] propose to reduce gender bias in natural language processing by generating counterfactual sentences with swapped gender pronouns. We anticipate that a version of CoDA that prioritizes fairness concerns in a similar manner could be applied in the reinforcement learning and computer vision contexts.
- In reinforcement learning specifically, a shift in data distribution requires either (1) the use of an off-policy algorithm, or (2) high variance off-policy corrections for on-policy algorithms. In the former case, it should be noted that in case of function approximation, even off-policy algorithms may be negatively affected by large shifts in their training distribution [83, 20]. More work is needed to quantify these effects and their implications for agent performance.

**Improving RL in batch-constrained settings** In many settings, such as medicine and education, obtaining large quantities of observational data using a random policy is prohibitively expensive and/or unethical [84, 60]. As such, agents that can efficiently learn effective policies from batch-constrained data are needed, as are accurate ways to estimate agent performance from off-policy data [58, 62, 73]. We see CoDA as complementary to both goals, as subspace swapping is a powerful tool to produce large quantities of counterfactual data given a modest observational dataset. However, subspace swapping alone may be insufficient to generate plausible “exploratory” data for evaluating and learning new policies. For example, medical records from certain demographic groups may be unavailable or improperly collected/labeled. It is conceivable that a CoDA with a suitable prioritization scheme could compensate for such sample bias, but applied work in batch-constrained domains that characterizes the effect of sample bias on CoDA should nevertheless be carried out.

**General considerations related to artificial agency** To the extent that CoDA is a general technique for improving the ability of artificial agents to achieve their goals, it inherits the potential risks and benefits associated with empowered artificial agency, including but not limited to:

- (a) the pursuit of misguided or dangerous goals, whether due to mispecification by a benevolent principal, the self-serving motives of its principals, or interference by malicious parties or other deviations from proper intents,
- (b) the unsafe and improper pursuit of goals due to poor modeling or representation, resource constraints and lack of capacity, constraint mispecification, partial observability, or inadequate encoding and understanding of human values, and
- (c) improvements to capital processes and automation of human labor, which could improve economic efficiency and raise the overall social welfare, but also run the risk of increased inequality, workforce displacement, and technological unemployment.

The risk associated with points (a) and (b) may be exacerbated in case of CoDA due to the risks associated with counterfactual reasoning outlined above: to the extent that CoDA is done with a poorly fit local factorization model, or with a local factorization model that does not generalize well to the counterfactual distribution, this could cause the agent to pursue poorly formulated counterfactual (imagined) goals, or create causally invalid data that hurts agent performance.

## References

- [1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in neural information processing systems*, pages 5048–5058, 2017.
- [2] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [5] Elias Bareinboim and Judea Pearl. Controlling selection bias in causal inference. In *Artificial Intelligence and Statistics*, pages 100–108, 2012.
- [6] Elias Bareinboim, Jin Tian, and Judea Pearl. Recovering from selection bias in causal and statistical inference. In *AAAI*, pages 2410–2416. Citeseer, 2014.
- [7] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- [8] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357, 2016.
- [9] Craig Boutilier, Richard Dearden, Moises Goldszmidt, et al. Exploiting structure in policy construction. In *IJCAI*, volume 14, pages 1104–1113, 1995.
- [10] Craig Boutilier, Nir Friedman, Moises Goldszmidt, and Daphne Koller. Context-specific independence in bayesian networks. *arXiv preprint arXiv:1302.3562*, 2013.
- [11] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [12] Lars Buesing, Theophane Weber, Yori Zwols, Sebastien Racaniere, Arthur Guez, Jean-Baptiste Lespiau, and Nicolas Heess. Woulda, coulda, shoulda: Counterfactually-guided policy search. *International Conference on Learning Representations*, 2019.
- [13] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.
- [14] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. *arXiv preprint arXiv:1812.02341*, 2018.
- [15] Carlos Diuk, Andre Cohen, and Michael L Littman. An object-oriented representation for efficient reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pages 240–247, 2008.
- [16] Aysegul Dundar, Kevin J Shih, Animesh Garg, Robert Pottorf, Andrew Tao, and Bryan Catanzaro. Unsupervised disentanglement of pose, appearance and background from images and videos. *arXiv preprint arXiv:2001.09518*, 2020.
- [17] Babak Esmaeili, Hao Wu, Sarthak Jain, Aican Bozkurt, N Siddharth, Brooks Paige, Dana H Brooks, Jennifer Dy, and Jan-Willem Meent. Structured disentangled representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2525–2534, 2019.
- [18] Amir-massoud Farahmand, Andre Barreto, and Daniel Nikovski. Value-aware loss function for model-based reinforcement learning. In *Artificial Intelligence and Statistics*, pages 1486–1494, 2017.
- [19] Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [20] Justin Fu, Aviral Kumar, Matthew Soh, and Sergey Levine. Diagnosing bottlenecks in deep q-learning algorithms. *arXiv preprint arXiv:1902.10250*, 2019.



- [21] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. *arXiv preprint arXiv:1812.02900*, 2018.
- [22] Scott Fujimoto, Herke Van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- [23] Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 219–226, 2019.
- [24] Andrew Gelman. Causality and statistical learning, 2011.
- [25] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In *International Conference on Machine Learning*, pages 881–889, 2015.
- [26] Anirudh Goyal, Alex Lamb, Jordan Hoffmann, Shagun Sodhani, Sergey Levine, Yoshua Bengio, and Bernhard Schölkopf. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*, 2019.
- [27] Clive WJ Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society*, pages 424–438, 1969.
- [28] Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored mdps. *Journal of Artificial Intelligence Research*, 19:399–468, 2003.
- [29] Assaf Hallak, François Schnitzler, Timothy Mann, and Shie Mannor. Off-policy model-based learning under unknown factored dynamics. In *International Conference on Machine Learning*, pages 711–719, 2015.
- [30] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [31] MA Hernán and JM Robins. Causal inference: What if. *Boca Raton: Chapman & Hill/CRC*, 2020.
- [32] Irina Higgins, Nicolas Sonnerat, Loic Matthey, Arka Pal, Christopher P Burgess, Matko Bošnjak, Murray Shanahan, Matthew Botvinick, Demis Hassabis, and Alexander Lerchner. SCAN: Learning hierarchical compositional visual concepts. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rkN2I1-RZ>.
- [33] Rodrigo Toro Icarte, Toryn Klassen, Richard Valenzano, and Sheila McIlraith. Using reward machines for high-level task specification and decomposition in reinforcement learning. In *International Conference on Machine Learning*, pages 2107–2116, 2018.
- [34] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems*, pages 12498–12509, 2019.
- [35] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The malmo platform for artificial intelligence experimentation. In *IJCAI*, pages 4246–4247, 2016.
- [36] Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, pages 1094–1099. Citeseer, 1993.
- [37] Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *IJCAI*, volume 16, pages 740–747, 1999.
- [38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [39] Thomas Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world models. *arXiv preprint arXiv:1911.12247*, 2019.
- [40] Oleg Klimov and John Schulman. Roboschool, 2017.
- [41] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.
- [42] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [43] Aviral Kumar, Abhishek Gupta, and Sergey Levine. Discor: Corrective feedback in reinforcement learning via distribution correction. *arXiv preprint arXiv:2003.07305*, 2020.

- [44] Sébastien Lachapelle, Philippe Brouillard, Tristan Deleu, and Simon Lacoste-Julien. Gradient-based neural dag learning. *arXiv preprint arXiv:1906.02226*, 2019.
- [45] Sameera Lanka and Tianfu Wu. Archer: Aggressive rewards to counter bias in hindsight experience replay. *arXiv preprint arXiv:1809.02070*, 2018.
- [46] Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *arXiv preprint arXiv:2004.14990*, 2020.
- [47] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam R Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. *arXiv preprint arXiv:1810.00825*, 2018.
- [48] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [49] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [50] Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. Towards practical multi-object manipulation using relational reinforcement learning. *arXiv preprint arXiv:1912.11032*, 2019.
- [51] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [52] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. *arXiv preprint arXiv:1811.12359*, 2018.
- [53] Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. *arXiv preprint arXiv:2002.02886*, 2020.
- [54] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in neural information processing systems*, pages 6379–6390, 2017.
- [55] Chaochao Lu, Bernhard Schölkopf, and José Miguel Hernández-Lobato. Deconfounding reinforcement learning in observational settings. *arXiv preprint arXiv:1812.10576*, 2018.
- [56] Daoming Lyu, Fangkai Yang, Bo Liu, and Steven Gustafson. Sdrl: Interpretable and data-efficient deep reinforcement learning leveraging symbolic planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2970–2977, 2019.
- [57] Prashan Madumal, Tim Miller, Liz Sonenberg, and Frank Vetere. Explainable reinforcement learning through a causal lens. *arXiv preprint arXiv:1905.10958*, 2019.
- [58] Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *AAMAS*, pages 1077–1084, 2014.
- [59] Alexander Mott, Daniel Zoran, Mike Chrzanowski, Daan Wierstra, and Danilo Jimenez Rezende. Towards interpretable reinforcement learning using attention augmented agents. In *Advances in Neural Information Processing Systems*, pages 12329–12338, 2019.
- [60] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient off-policy reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 1054–1062, 2016.
- [61] Ashvin V Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforcement learning with imagined goals. In *Advances in Neural Information Processing Systems*, pages 9191–9200, 2018.
- [62] Michael Oberst and David Sontag. Counterfactual off-policy evaluation with gumbel-max structural causal models. *arXiv preprint arXiv:1905.05824*, 2019.
- [63] Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. In *Advances in Neural Information Processing Systems*, pages 6118–6128, 2017.
- [64] Ji Ho Park, Jamin Shin, and Pascale Fung. Reducing gender bias in abusive language detection. *arXiv preprint arXiv:1808.07231*, 2018.

- [65] Judea Pearl. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.
- [66] Luis Perez and Jason Wang. The effectiveness of data augmentation in image classification using deep learning. *arXiv preprint arXiv:1712.04621*, 2017.
- [67] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.
- [68] Silviu Pitis, Harris Chan, and Jimmy Ba. Protoge: Prototype goal encodings for multi-goal reinforcement learning. *The 4th Multidisciplinary Conference on Reinforcement Learning and Decision Making (RLDM2019)*, 2019.
- [69] Silviu Pitis, Harris Chan, and Stephen Zhao. mrl: modular rl. <https://github.com/spitis/mrl>, 2020.
- [70] Silviu Pitis, Harris Chan, Stephen Zhao, Bradly Stadie, and Jimmy Ba. Maximum entropy gain exploration for long horizon multi-goal reinforcement learning. In *Proceedings of the Thirty-seventh International Conference on Machine Learning*, 2020.
- [71] Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell, Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, et al. Multi-goal reinforcement learning: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018.
- [72] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [73] Aniruddh Raghu, Omer Gottesman, Yao Liu, Matthieu Komorowski, Aldo Faisal, Finale Doshi-Velez, and Emma Brunskill. Behaviour policy estimation in off-policy policy evaluation: Calibration matters. *arXiv preprint arXiv:1807.01066*, 2018.
- [74] Danilo J Rezende, Ivo Danihelka, George Papamakarios, Nan Rosemary Ke, Ray Jiang, Theophane Weber, Karol Gregor, Hamza Merzic, Fabio Viola, Jane Wang, et al. Causally correct partial models for reinforcement learning. *arXiv preprint arXiv:2002.02836*, 2020.
- [75] Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D Mahecha, Jordi Muñoz-Marí, et al. Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13, 2019.
- [76] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International conference on machine learning*, pages 1312–1320, 2015.
- [77] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- [78] Yannick Schroecker and Charles Isbell. Universal value density estimation for imitation learning and goal-conditioned reinforcement learning. *arXiv preprint arXiv:2002.06473*, 2020.
- [79] Aman Sinha, Hongseok Namkoong, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.
- [80] Jasper Snoek, Yaniv Ovadia, Emily Fertig, Balaji Lakshminarayanan, Sebastian Nowozin, D Sculley, Joshua Dillon, Jie Ren, and Zachary Nado. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*, pages 13969–13980, 2019.
- [81] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- [82] Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.
- [83] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [84] Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 2139–2148, 2016.
- [85] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

- [86] Maya Zhe Wang and Benjamin Y Hayden. Monkeys are curious about counterfactual outcomes. *Cognition*, 189:1–10, 2019.
- [87] Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- [88] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [89] Nicholas Watters, Loic Matthey, Matko Bosnjak, Christopher P Burgess, and Alexander Lerchner. Cobra: Data-efficient model-based rl through unsupervised object discovery and curiosity-driven exploration. *arXiv preprint arXiv:1905.09275*, 2019.
- [90] Théophane Weber, Nicolas Heess, Lars Buesing, and David Silver. Credit assignment techniques in stochastic computation graphs. *arXiv preprint arXiv:1901.01761*, 2019.
- [91] Cathy Wu, Aravind Rajeswaran, Yan Duan, Vikash Kumar, Alexandre M Bayen, Sham Kakade, Igor Mordatch, and Pieter Abbeel. Variance reduction for policy gradient with action-dependent factorized baselines. *arXiv preprint arXiv:1803.07246*, 2018.
- [92] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in neural information processing systems*, pages 3391–3401, 2017.
- [93] Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Relational deep reinforcement learning. *arXiv preprint arXiv:1806.01830*, 2018.
- [94] Guangxiang Zhu, Zhiao Huang, and Chongjie Zhang. Object-oriented dynamics predictor. In *Advances in Neural Information Processing Systems*, pages 9804–9815, 2018.

## A Proof of Proposition 1

**Lemma 1.** *If  $V^j \in \text{Pa}^{\mathcal{L}}(V^i)$  in DAG  $\mathcal{G}^{\mathcal{L}}$  of (local) causal model  $\mathcal{M}^{\mathcal{L}}$ , and  $\mathcal{L} \subset \mathcal{X}$ , then  $V^j \in \text{Pa}^{\mathcal{X}}(V^i)$  in DAG  $\mathcal{G}^{\mathcal{X}}$  corresponding to causal model  $\mathcal{M}^{\mathcal{X}}$ .*

*Proof.* By minimality, there exist  $\{u^i, v^{-j}, v_1^j\}$  and  $\{u^i, v^{-j}, v_2^j\}$  with  $v^{-j} \in \text{Pa}^{\mathcal{L}}(V^i) \setminus V^j$  for which  $f^i(\{u^i, v^{-j}, v_1^j\}) \neq f^i(\{u^i, v^{-j}, v_2^j\})$ . Expand  $\{v^{-j}, v_1^j\}$  and  $\{v^{-j}, v_2^j\}$  to  $(s_1, a_1), (s_2, a_2) \in \mathcal{L}$  (with any values of other components). But  $\mathcal{L} \subset \mathcal{X}$ , so  $(s_1, a_1), (s_2, a_2) \in \mathcal{X}$  and it follows from minimality in  $\mathcal{X}$  that  $V^j \in \text{Pa}^{\mathcal{X}}(V^i)$ .  $\square$

**Corollary 1.** *If  $\mathcal{L} \subset \mathcal{X}$ ,  $\mathcal{G}^{\mathcal{L}}$  is sparser (has fewer edges) than  $\mathcal{G}^{\mathcal{X}}$ .*

**Proposition 1.** *The causal mechanisms represented by  $\mathcal{G}_i, \mathcal{G}_j \subset \mathcal{G}$  are independent in  $\mathcal{G}^{\mathcal{L}_1 \cup \mathcal{L}_2}$  if and only if  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are independent in both  $\mathcal{G}^{\mathcal{L}_1}$  and  $\mathcal{G}^{\mathcal{L}_2}$ , and  $f^{\mathcal{L}_1, i} = f^{\mathcal{L}_2, i}, f^{\mathcal{L}_1, j} = f^{\mathcal{L}_2, j}$ .*

*Proof.* ( $\Rightarrow$ ) If  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are independent in  $\mathcal{G}^{\mathcal{L}_1 \cup \mathcal{L}_2}$ , independence in  $\mathcal{G}^{\mathcal{L}_1}$  and  $\mathcal{G}^{\mathcal{L}_2}$  follows from Corollary 1. That  $f^{\mathcal{L}_1, i} = f^{\mathcal{L}_2, i}$  (and  $f^{\mathcal{L}_1, j} = f^{\mathcal{L}_2, j}$ ), on their shared domain, follows since each is a restriction of the same function  $f^{\mathcal{L}_1 \cup \mathcal{L}_2, i}$  (or  $f^{\mathcal{L}_1 \cup \mathcal{L}_2, j}$ ).

( $\Leftarrow$ ) Suppose  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are independent in  $\mathcal{G}^{\mathcal{L}_1}$  and  $\mathcal{G}^{\mathcal{L}_2}$  but not  $\mathcal{G}^{\mathcal{L}_1 \cup \mathcal{L}_2}$ . By the definition of independence applied to  $\mathcal{G}^{\mathcal{L}_1 \cup \mathcal{L}_2}$ , we have that, without loss of generality, there is a  $V_i \in \mathcal{G}_i, V_j \in \mathcal{G}_j$  with  $V_j \in \text{Pa}^{\mathcal{L}_1 \cup \mathcal{L}_2}(V_i)$ . Then, from the definition of minimality, it follows that there exist  $(s_1, a_1), (s_2, a_2) \in \mathcal{L}_1 \cup \mathcal{L}_2$  that differ only in the value of  $V_j$ , and  $u_i \in \text{range}(U_i)$  for which  $f^i(s_1, a_1, u_i) \neq f^i(s_2, a_2, u_i)$ .

Clearly, if  $(s_1, a_1)$  and  $(s_2, a_2)$  are both in  $\mathcal{L}_1$  (or  $\mathcal{L}_2$ ), there will be an edge from  $V_j$  to  $V_i$  in  $\mathcal{G}^{\mathcal{L}_1}$  (or  $\mathcal{G}^{\mathcal{L}_2}$ ) and the claim follows by contradiction. Thus, the only interesting case is when, without loss of generality,  $(s_1, a_1) \in \mathcal{L}_1$  and  $(s_2, a_2) \in \mathcal{L}_2$ . The key observation is that  $(s_1, a_1)$  and  $(s_2, a_2)$  differ only in the value of node  $V_j \notin \mathcal{G}_i$ : since  $\mathcal{G}_i$  is an independent causal mechanism in both  $\mathcal{G}^{\mathcal{L}_1}$  and  $\mathcal{G}^{\mathcal{L}_2}$  and the parents of  $V_i$  take on the same values in each, we have that  $f^i(s_1, a_1, u_i) = f^{i, \mathcal{L}_1}(s_1, a_1, u_i) = f^{i, \mathcal{L}_2}(s_2, a_2, u_i) = f^i(s_2, a_2, u_i)$  and the claim follows by contradiction.  $\square$

## B Additional Experiment Details

This section provides training details for the experiments discussed in Section 4 as well as some additional results. Code is available at <https://github.com/spitis/mrl> [69].

### B.1 Online RL

Here we detail the procedure used in the Online RL experiments in the Spriteworld environment.

**Implementation** We work with the original TD3 codebase, architecture, and hyperparameters (except batch size; see below), and focus our efforts solely on modifying the agent’s training distribution.

**Environment** We extend the base Spriteworld framework [89] as follows:

1. We add basic collisions, wherein a sprites velocity is reversed if it would overlap with another sprite on the next frame. This induces a local factorization and prevents a global factorization from being applicable.
2. We use a new continuous action space that consists of a 2-dimensional “click” on the canvas. The click changes the velocity of the closest sprite according to distance and angle between the click and the sprite’s center.
3. We add a disentangled state renderer that returns the position and velocity of each sprite (a total of 16-dimensions in tasks with four sprites).
4. We add a mask renderer that returns the ground truth masking function. This allows us to evaluate our masking function in Appendix C.

5. We add a suite of partial and sparse reward tasks that we use for experiments. The tasks vary in difficulty along the following two dimensions

- **Number of shapes:** The agent must place  $N \in \{1, 2, 3, 4\}$  of the 4 sprites at their respective target (one of four quadrants).
- **Sparsity of reward signal:** The agent receives reward either *Partial*,  $\frac{1}{N}$  reward for each of the  $N$  sprites placed near its target (up to a maximum reward of 1); or *Sparse*, reward of 1 is only observed when *all*  $N$  sprites reach their target.

For example, hardest task—“sparse reward, place 4” in the results—requires placing 4 shapes with sparse reward. Episodes run for 50 timesteps, and there is no termination signal (i.e., the environment is continuous from the agent’s perspective). We also tried a variant with dense rewards, but found that the baseline agent was able to easily solve this task and no benefit was conferred by data augmentation.

These extensions will included with the release of our code.

**Data augmentation** Every 1000 environment steps, we sample 2000 pairs of random transitions from the agent’s replay buffer, and apply CoDA to produce a maximum of 5 unique CoDA samples per pair. We apply two forms of CoDA, using (1) an oracle / ground truth mask function that we back out of the simulator, and (2) the mask of a pre-trained transformer model (see Section C). The mask function was trained using approximately 42,000 samples from a random policy (5/6 of 50,000, with the rest of the data used for validation). CoDA samples are added to a second CoDA replay buffer. For purposes of this experiment both buffers have effectively infinite capacity (they are never filled). During training, the agent’s batches are sampled proportionally from the real and CoDA replay buffers (this means that approximately 7/8 of the data that the agent trains on is counterfactual).

**Baselines** In addition to the base TD3 agent and CoDA, we also use the transformer model that is used as a mask function to generate data by performing forward rollouts with a random policy, as in Dyna [82]. So that this baseline produces approximately the same number of samples as CoDA with the learned mask, we roll the model out for 5 steps from 1500 random samples from the replay buffer, again every 1000 environment steps. This produces 7500 model-based samples for every 1000 environment transitions.

As an additional model-based Dyna baseline, we also fit a plain feedforward neural network with 2 hidden layers of 256 ReLU units as a feedforward model. We found that this forward model performed much worse than our transformer-based dyna baseline. The results are shown in Figure 6.

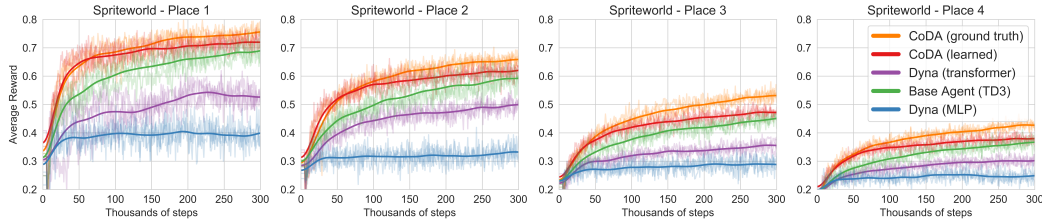


Figure 6: Average reward in Partial reward setting with MLP-based Dyna baseline.

Use of the transformer mask function for CoDA requires setting the threshold value  $\tau$ , which we do by monitoring accuracy and F1 scores for sparsity prediction (as discussed in Appendix C) on validation data, ultimately using the value  $\tau = 0.05$ .

**Batch Size** Since CoDA samples are plentiful we increase the agent’s batch size from 256 to 1000 to allow it to train on more environment samples in the same number of training steps. We found that this slightly improved the performance of the base TD3 agent. An increase in batch size also allows the agent to see more of its own on-policy data in the face of many off-policy CoDA samples.

**Additional results in sparse reward task variants** In addition to the partial reward tasks described in Section 4, we also tested CoDA in four sparse reward tasks of varying difficulty. These are the same as the partial reward tasks, except that a sparse reward of 1 is granted only when *all*  $N$  sprites

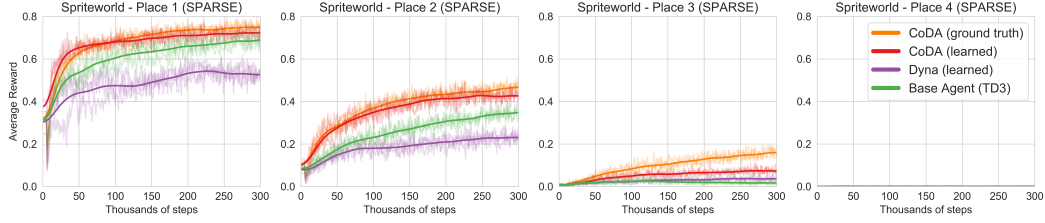


Figure 7: **Average reward on Sparse reward bouncing ball environment.** As in the partial reward case (Figure 4), we observe that CoDA agents outperform the other agents (except in Place 4, where no agent achieves any reward).

are in their target locations. While these tasks were much harder (and perhaps impossible in the case of Place 4 due to moving sprites), as shown in Figure 7, the CoDA agents maintain a clear advantage.

**Results plot** The results plot shows the 3 seeds in reduced opacity together with their smoothed mean.

## B.2 Batch RL

Here we detail the procedure used in the Batch RL experiments in the Pong environment.

**Implementation** This experiment works with a different codebase than the Spriteworld experiment, in order to simplify the use of CoDA in a Batch RL setting. Our experiment first builds the agent’s dataset (consisting of real data, dyna data and/or CoDA data), then instantiates a TD3 agent by filling its replay buffer with the dataset. The replay buffer is always expanded to include the entire enlarged dataset (for the 5x CoDA ratio at 250,000 data size this means the buffer has 1.5E6 experiences). The agent is run for 500,000 optimization steps.

**Hyperparameters** We used similar hyperparameters to the original TD3 codebase, with the following differences:

- We use a discount factor of  $\gamma = 0.98$  instead of 0.99.
- Since Pong is a sparse reward task with  $\gamma = 0.98$ , we clip critic targets to  $(-50, 50)$ .
- We use networks of size (128, 128) instead of (256, 256).
- As for our Spriteworld experiment, we use a batch size of 1000.

**Environment** We base our Pong environment on RoboSchoolPong-v1 [40]. The original environment allowed the ball to teleport back to the center after one of the players scored, offered a small dense reward signal for hitting the ball, and included a stray “timeout” feature in the agent’s state representation. We fix the environment so that the ball does not teleport, and instead have the episode reset every 150 steps, and also 10 steps after either player scores. The environment is treated as continuous and never returns a done signal that is not also accompanied by a `TimeLimit.truncated` indicator [11]. We change the reward to be strictly sparse, with reward of  $\pm 1$  given when the ball is behind one of the players’ paddles. Finally, we drop the stray “timeout” feature, so that the state space is 12-dimensional, where each set of 4 dimensions is the x-position, y-position, x-velocity, and y-velocity of the corresponding object (2 paddles and one ball).

The opponent in the environment uses a small pretrained policy that was included in the original environment [40] (the effectiveness of this policy was unaffected by the removal of the timeout feature). To collect the batch dataset, we use the same policy as an “expert” policy, but replace 50% of its actions by random actions.

**Training the CoDA model** Without access to a ground truth mask, we needed to train a masking function  $C$  to identify local disentanglement. We also forewent the ground truth reward, instead training our own reward classifier. In each case we used the batch dataset given, and so we trained different models for each random seed. For our masking model, we stacked two single-head



transformer blocks (without positional encodings) and used the product of their attention masks as the mask. Each block consists of query  $Q$ , key  $K$ , and value  $V$  networks that each have 3-layers of 256 neurons each, with the attention computed as usual [85, 47]. The transformer is trained to minimize the L2 error of the next state prediction given the current state and action as inputs. The input is disentangled, and so has shape (batch\_size, num\_components, num\_features). In each row (component representation) of each sample, features corresponding to other components are set to zero. The transformer is trained for 2000 steps with a batch size of 256, Adam optimizer [38] with learning rate of  $3e-4$  and weight decay of  $1e-5$ . For our reward function we use a fully-connected neural network with 1 hidden layer of 128 units. The reward network accepts an  $(s, a, s')$  tuple as input (not disentangled) and outputs a softmax over the possible reward values of  $[-1, 0, 1]$ . It is trained for 2000 steps with a batch size of 512, Adam optimizer [38] with learning rate of  $1e-3$  and weight decay of  $1e-4$ . All hyperparameters were rather arbitrary (we used the default setting, or in case it did not work, the first setting that gave reasonable results, as was determined by inspection). To ensure that our model and reward functions are trained appropriately (i.e., do not diverge) for each seed, we confirm that the average loss of the CoDA model is below 0.005 at the training and that the average loss of the reward model is below 0.1, which values were found by inspection of a prototype run. These conditions were met by all seeds.

When used to produce masks, we chose a threshold of  $\tau = 0.02$  by inspection, which seemed to produce reasonable results. A more principled approach would do cross-validation on the available data, as we did for Spriteworld (Appendix C).

**Tested configurations** Table 1 reports a subset of our tested configurations. We report our results in full below. We considered the following configurations:

1. **Real data only.**
2. **CoDA + real data:** after training the CoDA model, we expand the base dataset by either 2, 3, 4 or 6 times.
3. **Dyna (using CoDA model):** after training the CoDA model, we use it as a forward dynamics model instead of for CoDA; we use 1-step rollouts with random actions from random states in the given dataset to expand the dataset by 2x. We also tried with 5-step rollouts, but found that this further hurt performance (not shown). Note that Dyna results use only 5 seeds.
4. **Dyna (using MBPO model):** as the CoDA model exhibits significant model bias when used as a forward dynamics model, we replicate the state-of-the-art model-based architecture used by MBPO [34] and use it as a forward dynamics model for Dyna; we experimented with 1-step and 5-step rollouts with random actions from random states in the given dataset to expand the dataset by 2x. This time we found that the 5-step rollouts do better, which we attribute to the lower model bias together with the ability to create a more diverse dataset (1-step not shown). The MBPO model is described below. We use the same reward model as CoDA to relabel rewards for the MBPO model, which only predicts next state.
5. **MBPO + CoDA:** as MBPO improved performance over the baseline (real data only) at lower dataset sizes, we considered using MBPO together with CoDA. We use the base dataset to train the MBPO, CoDA, and reward models, as described above. We then use the MBPO model to expand the base dataset by 2x, as described above. We then use the CoDA model to expand *the expanded dataset* by 3x the original dataset size. Thus the final dataset is 5x as large as the original dataset (1 real : 1 MBPO : 3 CoDA).

All configurations alter only the training dataset, and the same agent architecture/hyperparameters (reported above) are used in each case.

**MBPO model** Since using the CoDA model for Dyna harms rather than helps, we consider using a stronger, state-of-the-art model-based approach. In particular, we adopt the model used by Model-Based Policy Optimization [34]. This model consists of a size 7 ensemble of neural networks, each with 4 layers of 200 neurons. We use ReLU activations, Adam optimizer [38] with weight decay of  $5e-5$ , and have each network output a the mean and (log) diagonal covariance of a multi-variate Gaussian. We train the networks with a maximum likelihood loss. To sample from the model, we choose an ensemble member uniformly at random and sample from its output distribution, as in [34].

**Full results** See Table 2 for results.

$ \mathcal{D} $ (1000s)	Real data 1R	Dyna 1R:1M	MBPO 1R:1M	Ratio of Real:CoDA [:MBPO] data (ours)				
				1R:1C	1R:2C	1R:3C	1R:5C	1R:3C:1M
25	13.2 $\pm$ 0.7	12.2 $\pm$ 1.4	18.5 $\pm$ 1.5	43.8 $\pm$ 2.7	41.6 $\pm$ 3.7	40.9 $\pm$ 2.5	38.4 $\pm$ 4.9	<b>46.8</b> $\pm$ 3.2
50	22.8 $\pm$ 3.2	17.4 $\pm$ 3.0	36.6 $\pm$ 4.0	66.6 $\pm$ 3.7	58.4 $\pm$ 4.7	64.4 $\pm$ 3.1	62.5 $\pm$ 3.5	<b>70.4</b> $\pm$ 3.7
75	43.2 $\pm$ 4.7	23.3 $\pm$ 4.7	46.0 $\pm$ 4.5	73.4 $\pm$ 2.8	71.6 $\pm$ 3.8	<b>76.7</b> $\pm$ 2.6	75.0 $\pm$ 3.3	74.6 $\pm$ 3.3
100	63.0 $\pm$ 3.0	25.9 $\pm$ 7.4	66.4 $\pm$ 5.2	77.8 $\pm$ 2.0	77.4 $\pm$ 1.8	<b>82.7</b> $\pm$ 1.5	76.6 $\pm$ 3.0	73.7 $\pm$ 2.9
150	77.4 $\pm$ 1.3	34.1 $\pm$ 1.5	72.6 $\pm$ 5.6	82.2 $\pm$ 1.7	84.0 $\pm$ 1.2	<b>85.8</b> $\pm$ 1.4	84.2 $\pm$ 1.0	79.7 $\pm$ 3.5
250	78.2 $\pm$ 2.7	44.2 $\pm$ 4.0	77.9 $\pm$ 2.3	85.0 $\pm$ 2.8	<b>88.1</b> $\pm$ 1.0	87.8 $\pm$ 1.7	87.0 $\pm$ 1.0	78.3 $\pm$ 5.0

Table 2: Extended Batch RL results. Mean success ( $\pm$  standard error, estimated using 1000 bootstrap resamples) on Pong environment. All results average over 10 seeds, except Dyna, which uses 5.

### B.3 Goal-conditioned RL

Here we detail the procedure used in the Goal-conditioned RL experiments on the Fetchpush-v1 and Slide2 environments.

**Implementation** This experiment uses the same codebase as our Batch RL, which provides state-of-the-art baseline HER agents and will be released with the paper.

**Hyperparameters** For Fetchpush-v1 we use the default hyperparameters from the codebase (see [70] for details on how they were selected), which outperform the original HER agents of [1, 71] and follow-up works. We do not tune the CoDA agent (but see additional CoDA hyperparameters below). They are as follows:

- Off-policy algorithm: DDPG [51]
- Hindsight relabeling strategy: `futureactual_2_2` [68], using exclusively `future` [1] relabeling for the first 25,000 steps
- Optimizer: Adam [38] with default hyperparameters
- Batch size: 2000
- Optimization frequency: 1 optimization step every 2 environment steps after the 5000th environment step
- Target network updates: update every 10 optimization steps with a Polyak averaging coefficient of 0.05
- Discount factor: 0.98
- Action l2 regularization: 0.01
- Networks: 3x512 layer-normalized [4] hidden layers with ReLU activations
- Target clipping: (-50, 0)
- Action noise: 0.1 Gaussian noise
- Epsilon exploration [71]: 0.2, with an initial 100% exploration period of 10,000 steps
- Observation normalization: yes
- Buffer size: 1M

On Slide2 we tried to tune the baseline hyperparameters somewhat, but note that this is a fairly long experiment (10M timesteps) and so only a few settings were tested due to constraints. In particular, we considered the following modifications:

- Expanding the replay buffer to 2M (effective)
- Reducing the batch size to 1000 (effective)\* (used for results)
- Using the `future_4` strategy (agent fails to learn in 10M steps)
- Reducing optimization step frequency to 1 step every 4 environment steps (about the same performance)

We tried similar adjustments to our CoDA agent, but found the default hyperparameters (used for results) performed well. We found that the CoDA agent outperforms the base HER agent on all tested settings.

For CoDA, we used the following additional hyperparameters:

- CoDA buffer size: 3M
- Make CoDA data every: 250 environment steps
- Number of source pairs from replay buffer used to make CoDA data: 2000
- Number of CoDA samples per source pair: 2
- Maximum ratio of CoDA:Real data to train on: 3:1

**Environment** On `FetchPush-v1` the standard state features include the relative position of the object and gripper, which entangles the two. While this could be dealt with by dynamic relabeling (as used for HER’s reward), we simply drop the corresponding features from the state.

`Slide2` has two pucks that slide on a table and bounce off of a solid railing. Observations are 40-dimensional (including the 6-dimensional goal), and actions are 4-dimensional. Initial positions and goal positions are sampled randomly on the table. During training, the agent gets a sparse reward of 0 (otherwise -1) if *both* pucks are within 5cm of their ordered target. At test time we count success as having both picks within 7.5cm of the target on the last step of the episode. Episodes last 75 steps and there is no done signal (this is intended as a continuous task).

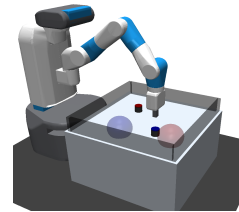


Figure 8: The `Slide2` environment.

**CoDA Heuristic** For these experiments we use a hand-coded heuristic designed with domain knowledge. In particular, we assert that the action is always entangled with the gripper, and that gripper/action and objects (pucks or blocks) are disentangled whenever they are more than 10cm apart. This encodes independence due to physical separation, which we hypothesize is a very generally heuristic that humans implicitly rely on all the time. The pucks have a radius of 2.5cm and height of 4cm, and the blocks are 5cm x 5cm x 5cm, so this heuristic is quite generous / suboptimal. Despite being suboptimal, it demonstrates the ease with which domain knowledge can be injected via the CoDA mask: we need only a high precision (low false positive rate) heuristic—the recall is not as important. It is likely that an agent could learn a better mask that also takes into account velocity.

**Results plot** The plot shows mean  $\pm 1$  standard deviation of the smoothed data over 5 seeds.

## C Inferring Local Factorization

Here we present several approaches to inferring the local factorization of subspaces, a crucial subroutine of CoDA. We note that in many cases where domain knowledge is available, simple heuristics may suffice, e.g. in our Goal-conditioned RL experiments discussed in Section 4 where a simple distance-based indicator function in the state space was used. However, as such heuristics may not be universally available, the question of whether data-driven approaches can successfully infer local factorization is of general interest. We note that the performance of CoDA will improve alongside future improvements in this inference task (motivating future work in this area), and that inferring the local factorization in general is an easier task than learning the environment dynamics.

We begin by presenting two methods for inferring local factorization, derived from variants of a next-state prediction task, which we here refer to as SANDy for Sparse Attention Neural Dynamics. To verify the merits of SANDy in the local factorization inference task, we evaluate two SANDy variants in controlled settings where the ground truth factorization is known: first in a synthetic Markov process (MDP without actions), and second in `Spriteworld`. This label information is used only to evaluate performance, and not to train the SANDy parameters. The SANDy-Transformer model was used for the Online RL experiments presented in Section 4.

## C.1 Methods

We propose two Sparse Attention for Neural Dynamics (SANDy) models. In each case we seek to learn a function (or mask)  $M(s, a) \rightarrow \{0, 1\}^{(|S|+|A|) \times |S|}$  whose output represents the adjacency matrix of the local causal graph, conditioned on the state and action. We note that  $M(s, a)_{i,j} = 0$  alone is insufficient, in general, to determine the local subspace  $\mathcal{L} \subset \mathcal{S} \times \mathcal{A}$ , since there may be multiple disconnected subspaces  $\mathcal{L}_k$  with  $M_k(s, a)_{i,j} = 0$  whose union  $\bigcup_k \mathcal{L}_k$  has  $M_{\bigcup_k}(s, a)_{i,j} = 1$ . This can be resolved by our Proposition 1 if we also force the relevant structural equations to be the same. For now, we assume the mask determines the local subspace, and leave exploration of this possibility to future work. Empirically, we will see that our assumption is reasonable.

**SANDy-Mixture** The first model is a mixture-of-MLPs model with an attention mechanism that is computed from the current state. Each component of the mixture is a neural dynamics model with sparse local dependencies. For a given component, the key idea is to train a neural dynamics model to predict the next state  $h(s_t, a_t) \approx s_{t+1}$  and approximate the masking function by thresholding the (transpose of the) network Jacobian of  $h$ ,  $[\mathbf{J}(s, a)]_{i,j} = \frac{\delta}{\delta[s, a]_j} [h(s, a)]_i$ . Intuitively, we can think of  $\mathbf{J}$  as providing the first-order element-wise dependencies between the predicted next state and the network input. We then derive the local factorization by thresholding the absolute Jacobian

$$M_\tau(s, a) = \mathbf{1}(|\mathbf{J}(s, a)| > \tau), \quad (1)$$

where  $\mathbf{1}(\cdot)$  represents the indicator function and  $\tau$  is a threshold hyperparameter.

To estimate the network Jacobian, we note that for standard activation functions (sigmoid, tanh, relu), it can be bound from above by the matrix product of its weight matrices. To see this, let  $h_\theta$  be an  $L$ -layer MLP parameterized by  $\theta = (\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(L)})$  with activation  $\sigma$  with bounded derivative  $\sigma'(x) \leq 1$ , and note that for each layer  $\mathbf{h}^{(\ell)} = \sigma(\mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)})$  we have:

$$\left| \frac{d\mathbf{h}_j^{(\ell)}}{d\mathbf{h}_i^{(\ell-1)}} \right| \leq |\mathbf{W}_{ji}^{(\ell)}|.$$

Then, using the chain rule, triangle inequality, and the identity  $|ab| = |a||b|$ , we can compute:

$$\begin{aligned} \left| \frac{d\mathbf{h}_j^{(\ell)}}{d\mathbf{h}_i^{(\ell-2)}} \right| &\leq \left| \frac{d}{d\mathbf{h}_i^{(\ell-2)}} \mathbf{W}_{j\cdot}^{(\ell)} \mathbf{h}^{(\ell-1)} \right| = \left| \mathbf{W}_{j\cdot}^{(\ell)} \cdot \frac{d\mathbf{h}^{(\ell-1)}}{d\mathbf{h}_i^{(\ell-2)}} \right| \\ &\leq \sum_k \left| \mathbf{W}_{jk}^{(\ell)} \frac{d\mathbf{h}_k^{(\ell-1)}}{d\mathbf{h}_i^{(\ell-2)}} \right| \leq |\mathbf{W}_{j\cdot}^{(\ell)}| \cdot |\mathbf{W}_{\cdot i}^{(\ell-1)}|. \end{aligned}$$

Expanding this out to multiple layers, we see that  $|\mathbf{J}_\theta(s, a)| \leq \prod_{\ell \in [L]} |\mathbf{W}^{(\ell)}|$ , as desired. We use this upper bound to approximate the network Jacobian of an MLP by setting  $\hat{\mathbf{J}} = \prod_{\ell \in [L]} |\mathbf{W}^{(\ell)}|$ . A similar idea is used by [25, 44], among others, to control element-wise input-output dependencies.

We use this static approximation  $\hat{\mathbf{J}}$  to facilitate learning a sparse dynamic mask by specifying a mixture model  $h(s, a) = \sum_i \alpha_\phi^{(i)}(s, a) h_\theta^{(i)}(s, a)$  over the environment dynamics (with  $\sum_i \alpha_\phi^{(i)}(s, a) = 1 \forall (s, a)$ ), where each component  $h_\theta^{(i)}$  is an MLP as specified above with a sparsity prior on its Jacobian bound  $\hat{\mathbf{J}}^{(i)}$  to encourage sparse (i.e. well-factorized) local solutions. The dynamic mask is computed by first approximating the Jacobian,  $\hat{\mathbf{J}}(s, a) = \sum_i \alpha_\phi^{(i)}(s, a) \hat{\mathbf{J}}^{(i)}$ , then thresholding by  $\tau$  as in Equation 1.

Note that we assume the mixture components  $\alpha$  are a function of the current state and action alone; in other words the factorization (captured by the network Jacobian of each component) can be *locally inferred*. The next-state prediction is given by  $\hat{\mathbf{s}}_{t+1} = (h_\theta^{(1)}(\mathbf{s}_t, \mathbf{a}_t), \dots, h_\theta^{(N)}(\mathbf{s}_t, \mathbf{a}_t))^T \alpha_\phi(\mathbf{s}_t, \mathbf{a}_t)$ . To train the model, we optimize the objective:

$$\underset{\theta, \phi}{\text{minimize}} \frac{1}{|\mathcal{D}|} \left( \sum_{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \in \mathcal{D}} \|\mathbf{s}_{t+1} - h(\mathbf{s}_t, \mathbf{a}_t)\|_2^2 \right) + \lambda_1 S(\theta) + \lambda_2 R(\phi) + \lambda_3 \|(\theta, \phi)\|_2 \quad (2)$$

where  $S(\theta) = \frac{1}{K} \sum_i |\hat{\mathbf{J}}_\theta^{(i)}|_1$  puts an  $\ell_1$  prior on each mixture component to induce sparsity, and  $R(\phi)$  encourages high entropy in the attention probabilities

$$R(\phi) = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{s} \in \mathcal{D}} \sqrt{\frac{1}{N} \sum_{j \in [N]} [A_\phi(\mathbf{s}, \mathbf{a})]_j}.$$

We note that more sophisticated methods of computing Jacobians of neural networks—including architectural changes and optimization strategies [3]—have been proposed, and could in principle be used here as well.

**SANDy-Transformer** As an alternative model, we use a transformer-like architecture that applies self-attention between a set of (potentially multi-dimensional) inputs [85, 47]. Our architecture is composed of a stacked self-attention blocks. Each block accepts a set of inputs  $X = \{x_i \in \mathbb{R}^n\}$  and composes three functions of each input: query  $Q : \mathbb{R}^n \rightarrow \mathbb{R}^d$ , key  $K : \mathbb{R}^n \rightarrow \mathbb{R}^d$ , and value  $V : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . The block returns a set of outputs  $\{y_i \in \mathbb{R}^m\}$  of size  $|X|$ , each of which is computed as:  $y_i = A_i^T V(X)$ , where  $A_i = \text{softmax}(\sum_j (Q(x_i)_j K(x_j)_i))$  (note that  $V(X)$  is a matrix of size  $|X| \times m$ ). We approximate the block mask (approximate Jacobian) as  $A = [A_1, A_2, \dots, A_{|X|}] \in \mathbb{R}^{|X| \times |X|}$ , and the full network mask as the product of the block masks (as in the SANDy-Mixture). We used two-layer MLPs for each function  $K, Q, V$  in Spriteworld and three-layer MLPs in Pong.

To apply this architecture to our problem, we first embed each state and action component (single feature, or group of features) into  $\mathbb{R}^n$  to produce a set of inputs  $X$  and pass this through each stacked self-attention block to obtain a set of outputs  $Y$ . We then discard any output components that correspond to the action features to obtain the next state prediction and mask. The network  $h$  is trained to minimize mean squared error:

$$\underset{\theta}{\text{minimize}} \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \in \mathcal{D}} \|\mathbf{s}_{t+1} - h(\mathbf{s}_t, \mathbf{a}_t)\|_2^2. \quad (3)$$

Unlike the SANDy-Mixture, no sparsity regularizers are applied, as we found the sparsity induced by the softmax attention mechanism to be sufficient.

## C.2 Evaluation environments

**Synthetic Markov Processes** We investigate the capacity of the SANDy-Mixture to learn simple factorized transition dynamics under a globally factored Markov Process (STATIONARYMP). Unlike the general MDP setting, no agent/policy is considered. However the ability to train an unconstrained dynamics model to approximate factorized environment dynamics is an important subtask within our overall approach. Assuming a spherical Gaussian prior over the initial states  $\mathbf{s}^0 \in \mathbb{R}^9$ , the STATIONARYMP is entirely specified by transition distribution  $p(\mathbf{s}^{t+1}|\mathbf{s}^t)$ . We assume that the state transitions factorize (globally) into three parts. Denoting by  $\mathbf{s}_{n \dots m}^t$  the  $n$ -th through  $m$ -th dimensions of the time  $t$  state  $\mathbf{s}^t$ , we have:

$$p(\mathbf{s}_{1 \dots 9}^{t+1}|\mathbf{s}_{1 \dots 9}^t) = p(\mathbf{s}_{1 \dots 4}^{t+1}|\mathbf{s}_{1 \dots 4}^t)p(\mathbf{s}_{5 \dots 7}^{t+1}|\mathbf{s}_{5 \dots 7}^t)p(\mathbf{s}_{8,9}^{t+1}|\mathbf{s}_{8,9}^t).$$

In other words, we have a block-diagonal transition matrix comprising three blocks with sizes 4, 3, and 2, respectively. In our case, all transition factors are deterministic non-linear mappings, e.g.  $p(\mathbf{s}_{1 \dots 4}^{t+1}|\mathbf{s}_{1 \dots 4}^t) = \delta(g_{1 \dots 4}(\mathbf{s}_{1 \dots 4}^t))$ , with  $g_{1 \dots 4} : \mathbb{R}^4 \rightarrow \mathbb{R}^4$  is a randomly-initialized single-hidden-layer neural network with 32 hidden units and GELU nonlinearity [30]. In this case, we can alternatively express the deterministic dynamics via the transition function

$$\begin{aligned} \mathbf{s}^{t+1} &= (\mathbf{s}_{1 \dots 4}^{t+1}, \mathbf{s}_{5 \dots 7}^{t+1}, \mathbf{s}_{8,9}^{t+1}) \\ &= (g_{1 \dots 4}(\mathbf{s}_{1 \dots 4}^t), g_{5 \dots 7}(\mathbf{s}_{5 \dots 7}^t), g_{8,9}(\mathbf{s}_{8,9}^t)). \end{aligned} \quad (\text{STATIONARYMP})$$

We now turn to a more sophisticated MP with locally factored dynamics ( $\epsilon$ -NONSTATIONARYMP), and investigate whether the SANDy-Mixture can learn to recognize local disentanglement.

We refer to the component functions  $g_{1 \dots 4} : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ ,  $g_{5 \dots 7} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ , and  $g_{8,9} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  used in the STATIONARYMP as *local* transitions. We now introduce *global* interactions via the global transition functions,  $G_{1 \dots 4} : \mathbb{R}^4 \rightarrow \mathbb{R}^9$ ,  $G_{5 \dots 7} : \mathbb{R}^4 \rightarrow \mathbb{R}^9$ , and  $G_{8,9} : \mathbb{R}^4 \rightarrow \mathbb{R}^9$ , which respectively

map the local state factors onto the global state space<sup>7</sup>. This allows us to extend the STATIONARYMP by adding global state transitions to the local state transitions whenever the norm of the local state factors exceeds the value of a hyperparameter  $\epsilon$  (lower  $\epsilon$  indicates more global interaction). Denoting by  $\mathbb{1}(\cdot)$  the indicator function, we have

$$\begin{aligned} \mathbf{s}^{t+1} &= [(\mathbf{s}_1^{t+1}, \mathbf{s}_2^{t+1}, \mathbf{s}_3^{t+1}, \mathbf{s}_4^{t+1}), (\mathbf{s}_5^{t+1}, \mathbf{s}_6^{t+1}, \mathbf{s}_7^{t+1}), (\mathbf{s}_8^{t+1}, \mathbf{s}_9^{t+1})] \\ &= [g_{1\dots 4}(\mathbf{s}_{1\dots 4}^t), g_{5\dots 7}(\mathbf{s}_{5\dots 7}^t), g_{8,9}(\mathbf{s}_{8,9}^t)] \\ &\quad + G_{1\dots 4}(\mathbf{s}_{1\dots 4}^t) \mathbb{1}(\|\mathbf{s}_{1\dots 4}^t\|_2 > \epsilon) \\ &\quad + G_{5\dots 7}(\mathbf{s}_{5\dots 7}^t) \mathbb{1}(\|\mathbf{s}_{5\dots 7}^t\|_2 > \epsilon) \\ &\quad + G_{8,9}(\mathbf{s}_{8,9}^t) \mathbb{1}(\|\mathbf{s}_{8,9}^t\|_2 > \epsilon). \end{aligned} \quad (\epsilon\text{-NONSTATIONARYMP})$$

**Spriteworld** Since we extended Spriteworld with a ground truth mask renderer, we are able to directly evaluate our SANDy models in Spriteworld as well. See the main text and Appendix B for a description of the Spriteworld environment.

### C.3 Results

In this Subsection we measure ability of the proposed SANDy algorithm (in its two variants) to correctly infer local factorization. At each transition we can query the environment for the ground-truth connectivity pattern of the local causal graph: given  $|S| + |A|$  dimensions of current state and action and  $|S|$  dimensions of next state, this corresponds an adjacency matrix  $\mathbf{Y} \in \{0, 1\}^{|S|+|A| \times |S|}$ . We note that accessing these evaluation labels—which are not used to train SANDy—requires a controlled synthetic environment like the ones we consider, and we leave design of an evaluation protocol suitable for real-world environments to future work.

We learn the SANDy network parameters using a training dataset of 40,000 transitions, with an additional validation dataset of 10,000 transitions used for early stopping and hyperparameter selection. We used the Adam optimizer with learning rate of 0.001 and default hyperparameters. In the  $\epsilon$ -NONSTATIONARYMP, we set  $\epsilon = 1.5$ , while in the Spriteworld setting we collect training trajectories by deploying a random-action agent in the environment, and randomly resetting the environment with 5% probability at every step to increase diversity of experiences. We then evaluate the SANDy models by computing local factorizations  $M_\tau(s, a) : |S| \times |A| \rightarrow \{0, 1\}^{|S|+|A| \times |S|}$  as a function of the threshold  $\tau$  for each transition in a held-out test dataset of 10,000 trajectories. We compute true and false positive rates for various values of  $\tau$  to produce ROC plots.

Figure 9 shows that while the SANDy-Mixture is sufficient to solve the simpler synthetic MP settings (with avg AUC of 0.96 and 0.91 for the stationary and non-stationary variants), it scales poorly to the Spriteworld environment. While the modest inductive bias of sparse local connections and high-entropy mixture components in SANDy-Mixture makes it widely applicable, we hypothesize that its sensitivity to hyperparameters makes it difficult to tune in complex settings. Fortunately, SANDy-Transformer, performs favorably in Spriteworld by incorporating a stronger inductive bias about the state subspace structure. Thus we use SANDy-Transformer to perform local factorization inference in the remaining experiments.

Figure 10 provides some qualitative intuition as to how the two variants of SANDy differ in their attention strategy in the Spriteworld environment.

## D Fitting dynamics models to Spriteworld

Since CoDA is a data augmentation strategy, it is reasonable to consider an alternative approach to augmenting the experience buffer: sampling from a dynamics model as in model-based RL. Here we present some qualitative results from our efforts in fitting dynamics models to the Spriteworld environment. We found that while dynamics models achieve a decent error in the next-state prediction task, they fail to produce a diverse set of trajectories when used as autoregressive samplers. In particular, the autoregressive sampling did not model collisions well and often produced trajectories

<sup>7</sup> Like the local transition functions, global transition functions are implemented as randomly-initialized single-hidden-layer neural networks.

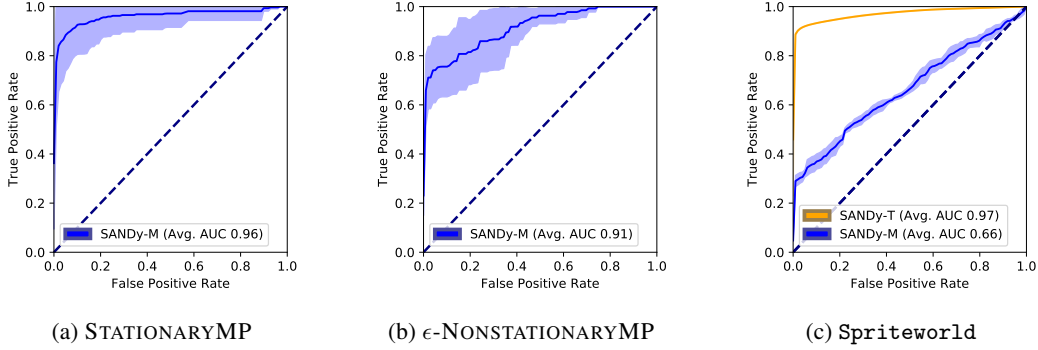


Figure 9: ROC plots for correct sparsity pattern prediction on the three environments. On held-out test transitions we derive the ground truth local connectivity per step—label information that is *not* used to train the attention model—and measure (over 5 runs; 1 std. dev. shaded) true and false positive rates while sweeping the mask threshold  $\tau$  over its allowable range. An accurate model generates an Area Under the Curve (AUC) close to 1. We observe that while SANDy-Mixture is sufficient for (nearly) solving the simpler synthetic MP environments, it underperforms in Spriteworld. SANDy-Transformer, which has a stronger inductive bias, is sufficient to (nearly) solve Spriteworld.

where sprites converged to fixed points in space after a short number of steps. Figure 11 shows trajectories sampled autoregressively from Linear, MLP, and LSTM-based dynamics models, alongside the ground truth trajectory. Note that all dynamics models were trained to minimize error in next-state prediction given the current state and action. In other words the LSTM auto-regressively predicts successive dimensions of the next state rather than modeling multiple time steps of the trajectory. Nevertheless the environment is truly Markov because instantaneous velocities are observed, so this information should be sufficient in theory to capture the environment dynamics.

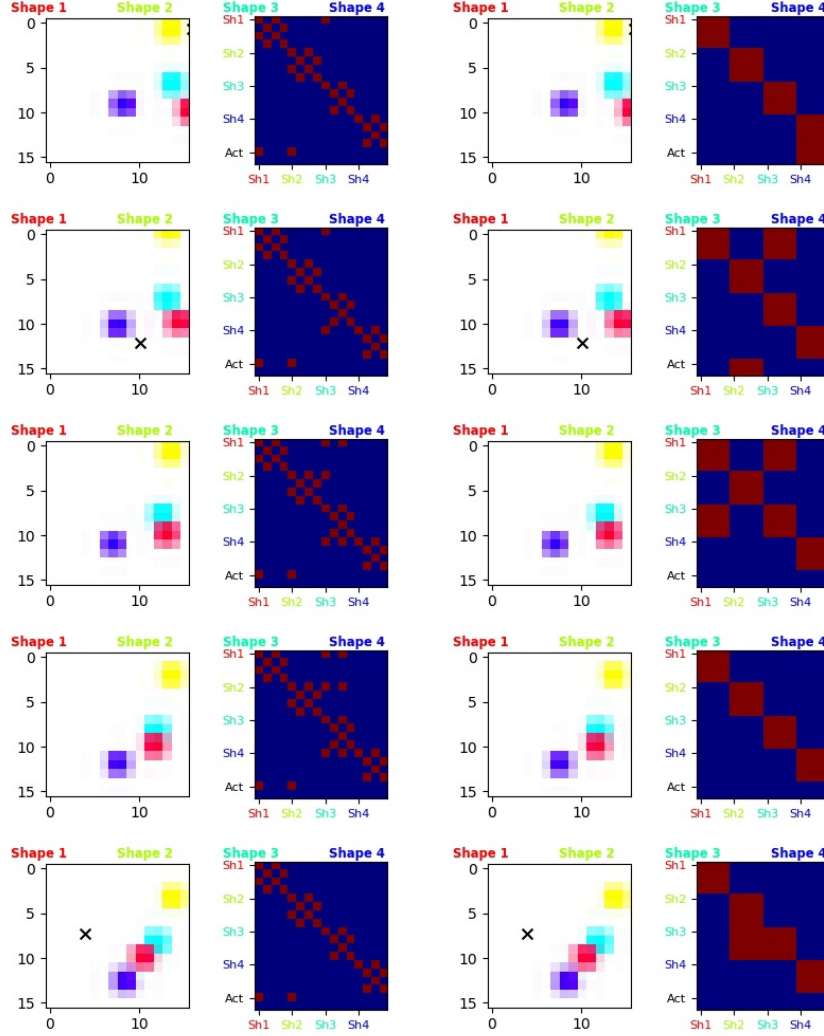
## E Sample efficient dynamics modeling with CoDA

If we had access to the ground truth local factorization, even for a few samples (e.g., we could have humans label them), how much more efficient would it be to train a dynamics model? In Figure 12, we sample 2000 transitions from a random policy in Spriteworld and use the data to train a forward dynamics model using MSE loss. The validation loss throughout training is plotted. Our baseline uses only the initial dataset, and quickly overfits the training set, showing increasing error after the initial few epochs. The same applies to a “random” CoDA strategy, that does CoDA using an identity attention mask ( $M(s, a) = I \forall (s, a)$ ) to randomly relabel the components. The random strategy does a bit better than the no CoDA strategy, since the randomness acts as a regularizer. Finally, we train a model using an additional 35,000 unique counterfactual CoDA transitions, and find that it significantly improves validation loss and prevents the model from overfitting. Note that we could have generated many more CoDA samples: from 2000 base transitions, if 80% of them do not involve collisions and there are 4 connected components in each, we could generate as many as  $1600^4$  (6.5 trillion!) counterfactual samples.

## F Compute Infrastructure

Experiments were run on a mix of local machines and a compute cluster, with a mix of GTX 1080 Ti, Titan XP, and Tesla P100 GPUs. This was solely to run jobs in parallel, and all experiments can be run locally (GPU optional for Spriteworld and Pong, but recommended for Fetch experiments).





(a) SANDy-Mixture

(b) SANDy-Transformer

Figure 10: Qualitative comparison of two attention mechanisms on the same Spriteworld trajectory. **SANDy-Mixture** (left) has a weaker inductive bias as it relies on sparsity regularization alone. Accordingly, it can learn a more compact subspace (e.g. grid patterns within a shape indicate that  $x$  and  $y$  coordinates move independently), but is less reliable in attending to collisions between shapes, and completely fails to attend to the action’s affect on shapes. **SANDy-Transformer** has a stronger inductive bias and can more reliably infer the local interaction pattern between the five subspaces (four shapes and one action).

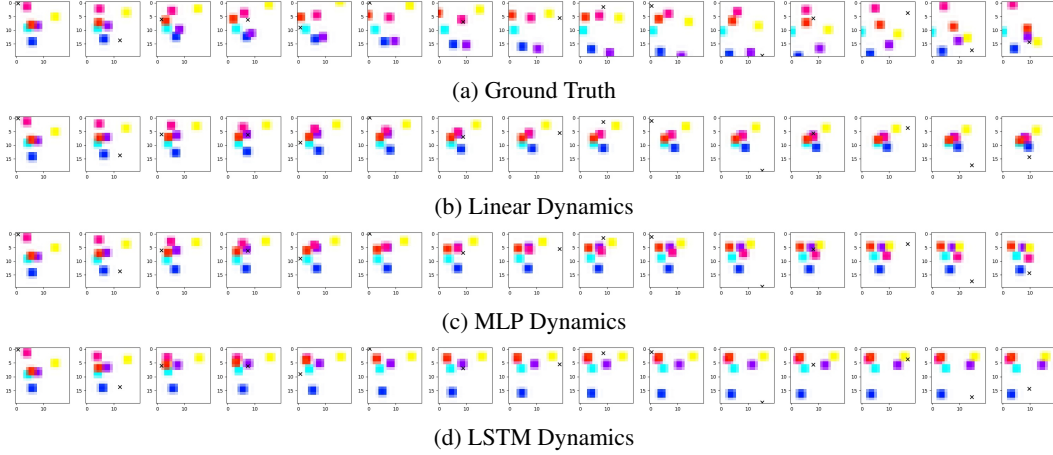


Figure 11: Auto-regressive model-based rollouts for a variety of dynamics models fit to Spriteworld. While the dynamics models were able to achieve relatively low error in the next-state prediction task, they fail to capture collisions and long-term dependencies in the sampled trajectories, and thus were omitted as baselines in the RL experiments. All trajectories share the same initial state.

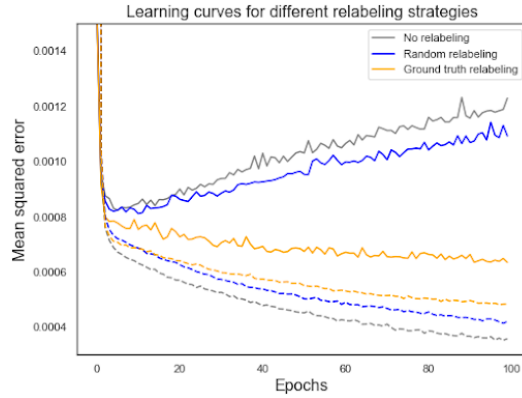


Figure 12: Learning curves for training a forward dynamics model using data from a random policy. Dotted lines indicate training performance, whereas solid lines indicate validation performance. We see that using the ground truth mask prevents overfitting and allows us to achieve much better validation performance.