# LEARNING ROBUST INTERVENTION REPRESENTA-TIONS WITH DELTA EMBEDDINGS

**Anonymous authors**Paper under double-blind review

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025 026 027

028

031

033

034

035

037

040

041

042

043

044

046

047

048

051

052

#### **ABSTRACT**

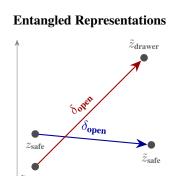
Causal representation learning has attracted significant research interest during the past few years, as a means for improving model generalization and robustness. Causal representations of interventional image pairs, have the property that only variables corresponding to scene elements affected by the intervention / action are changed between the start state and the end state. While most work in this area has focused on identifying and representing the variables of the scene under a causal model, fewer efforts have focused on representations of the interventions themselves. In this work, we show that an effective strategy for improving out of distribution (OOD) robustness is to focus on the representation of interventions in the latent space. Specifically, we propose that an intervention can be represented by a Causal Delta Embedding that is invariant to the visual scene and sparse in terms of the causal variables it affects. Leveraging this insight, we propose a method for learning causal representations from image pairs, without any additional supervision. Experiments in the Causal Triplet challenge demonstrate that Causal Delta Embeddings are highly effective in OOD settings, significantly exceeding baseline performance in both synthetic and real-world benchmarks.

# 1 Introduction

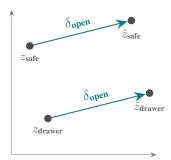
Understanding how the world changes in response to actions and external interventions is fundamental for artificial intelligence agents, especially those operating in dynamic environments. Although deep learning models are highly successful at capturing complex patterns from data, they often fail to generalize to new situations where the underlying data distribution changes, which is a critical limitation for real world deployment Hendrycks et al. (2021); Geirhos et al. (2020). To overcome this, agents must recover the underlying mechanisms that generate and transform data, enabling causal reasoning and robust generalization (Pearl, 2009).

This fundamental problem falls within the scope of Causal Representation Learning (CRL) (Schölkopf et al., 2021), which seeks to disentangle the causal variables of a system (Khemakhem et al., 2020). Despite its importance in practical applications such as robotics or healthcare (Gupta et al., 2024; Hellström, 2021; Sanchez et al., 2022), the challenge of learning disentangled and generalizable representations of the causal variables remains unresolved. Addressing this challenge requires accurate modelling of the underlying data generation process, a task which is guided by two fundamental assumptions within CRL. First, the Independent Causal Mechanisms (ICM) assumption, which posits that the distribution's generative process can be decomposed into autonomous and independent modules, each representing a distinct causal mechanism (Peters et al., 2017; Schölkopf et al., 2021). Second, the Sparse Mechanism Shift (SMS) assumption, which suggests that an intervention typically affects only a small, localized subset of these causal mechanisms (Schölkopf et al., 2021). Most existing methods focus on identifying these disentangled mechanisms from observations (Higgins et al., 2017; Khemakhem et al., 2020; Ahuja et al., 2022). Fewer methods have focused on learning generalizable representations of actions (interventions), which can be equally important in predicting the outcome of interventions, especially when faced with novel situations.

In this paper, we introduce Causal Delta Embedding (CDE), a novel framework for learning robust representations of interventions from image pairs. Using CDEs the intervention can be effectively isolated and represented as the vector difference between the latent representations of pre- and post-

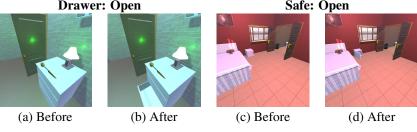


### **Consistent Representations**



(a) Baseline Model (ERM). Action representations depend on the object and scene features.

(b) Causal Delta Embeddings. The action representation  $\delta_{\rm open}$  is invariant to the object and scene context.



(c) Examples of intervention pairs from our dataset, showing pre- and post-intervention states for various actions and objects.

Figure 1: Visualizing Causal Delta Embeddings. Unlike a baseline model that produces entangled representations of action vectors (left), our model learns object invariant action representations (right), that generalize well to out of distribution samples. The model is trained on intervention pairs like those shown at the bottom.

intervention states if it satisfies the properties of (a) independence to causally irrelevant elements of the scene, in accordance to the ICM assumption (b) sparsity, in accordance to the SMS assumption and (c) object invariance, i.e., that the representation remains the same across objects. Using these properties as a guide, a learning strategy is proposed for learning CDEs from interventional image pairs.

We evaluate CDE on the Causal Triplet challenge (Liu et al., 2023), which encompasses 3 increasingly complex settings: single-object synthetic data, multi-object synthetic data and real world scenes from Epic Kitchens (Damen et al., 2022). Our experiments demonstrate that CDE establishes a new state of the art in OOD generalization for this challenge. Beyond quantitative performance, our analysis reveals that CDE learns a semantically structured intervention space, autonomously discovering anti-parallel relationships between opposing actions (e.g., open vs. close) without any explicit supervision.

Our main contributions are as follows:

- We introduce *Causal Delta Embedding (CDE)*, a novel approach for learning generalizable representations of interventions in a disentangled latent space.
- We propose a multi-objective loss function, designed to learn well separated, sparse and object invariant causal representations directly from visual data.
- We perform an extensive quantitative evaluation showing that our approach achieves stateof-the-art results in the Causal Triplet challenge.
- We show that our model discovers the semantic structure of the intervention space, including fundamental anti-parallel relationships between opposing actions, without any explicit supervision.

# 2 RELATED WORK

Causal Representation Learning Part of the research on CRL focuses on identifying latent causal variables from high dimensional observations (Khemakhem et al., 2020; Ahuja et al., 2022). These methods established identifiablity conditions for nonlinear ICA and demonstrated causal factor recovery under specific assumptions (Wendong et al., 2023; Monti et al., 2020). Another line of work focuses on causal disentanglement (Yang et al., 2021; Shen et al., 2020; Brehmer et al., 2022; Locatello et al., 2020a). These approaches often extend the Variational Autoencoder (VAE) framework (Kingma et al., 2013; Higgins et al., 2017). Object-centric learning methods have also been proposed in order to disentangle the visual scene into manipulable objects (Locatello et al., 2020b; Seitzer et al., 2022). More recent work has leveraged interventional data to improve causal disentanglement (Brehmer et al., 2022; Lippe et al., 2022; Squires et al., 2023; Lippe et al., 2023; Ahuja et al., 2022), showing that interventions provided crucial inductive biases for learning causal representations (Ahuja et al., 2023). While previous methods focus on identifying causal variables, our work instead models the interventional mechanisms, by learning generalizable embeddings that represent interventions in a way that remains invariant across different contexts.

Visual Action Recognition and OOD Generalization Traditional action recognition methods rely on spatiotemporal patterns and achieve strong performance under IID conditions (Carreira & Zisserman, 2017; Feichtenhofer et al., 2019; Arnab et al., 2021). However, these correlation-based approaches struggle with distribution shifts (Geirhos et al., 2020) and often are associated with spurious correlations (Wang & Jordan, 2024). Recent work has explored domain adaptation (Chen et al., 2019; Munro & Damen, 2020) and causal approaches (Magliacane et al., 2018; Wang et al., 2023) for robust action understanding. Another category of methods uses large Vision Language Action (VLA) models (Kim et al., 2024; Zitkovich et al., 2023; Ma et al., 2024) to enable agents to perform actions in challenging environments. These models typically depend on large-scale pretraining on diverse data, yet generalization to unseen tasks remains an open challenge (Sapkota et al., 2025). Unlike these approaches, our method learns *causal* representations of interventions, in the sense that they satisfy properties resulting from the CRL assumptions. These representations are shown to generalize to novel object-action combinations without the need for finetuning.

Contrastive Learning and Sparse Representations Contrastive learning has proven effective for learning meaningful representations by contrasting similar and dissimilar examples (Chen et al., 2020; Khosla et al., 2020). However, existing methods contrast individual samples, rather than relationships between samples. The principle of sparse mechanism shifts (Schölkopf et al., 2021; Peters et al., 2017), suggests that interventions affect only a small subset of the causal mechanisms. Sparsity in causal representations has been explored in various works (Pfister & Peters, 2022; Xu et al., 2024) and has shown to improve disentanglement. However, combining sparsity with adversarial training (Liu et al., 2023) often leads to poor OOD performance, since other confounders might still be present in the scene, motivating our approach for stricter assumptions.

### 3 Problem Formulation

The central challenge this paper addresses is the development of a CRL framework that can robustly infer actions / interventions from high-dimensional observations, particularly under distribution shifts.

We formalize this challenge within the framework of the Structural Causal Model presented by Liu et al. (2023) (Figure 2). Let us consider a set of causal variables  $Z \in \mathcal{Z} \subset \mathbb{R}^l$ , representing the state of the underlying data generating mechanisms. These variables have dependencies that are defined through a set of structural equations:

$$Z_i := f_i(\operatorname{pa}(Z_i), \epsilon_i), \quad i = 1, \dots, l$$

where  $\operatorname{pa}(Z_i)$  denotes the set of causal parents of variable  $Z_i$ , and the  $\epsilon_i$  are mutually independent stochastic noise terms representing unmodeled factors. The high-dimensional visual observation  $x \in \mathcal{X} \subset \mathbb{R}^d$  is rendered from these latent variables via a complex, non-invertible generative function  $g: \mathcal{Z} \to \mathcal{X}$ , such that x = g(Z).

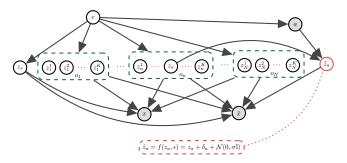


Figure 2: Causal Graph for a pair of observations  $(x, \tilde{x})$  before and after an action a, proposed by Liu et al. (2023). The data generating process is described by a set of latent factors, including global scene level factors  $z_s$  and local object level factors  $z_n^k$ , which are dependent due to confounders c. The action is assumed to influence only a few object level causal factors  $z_a$  in the scene and the effect of that influence is captured by  $\tilde{z}_a$ . The red dashed line indicates the structural equation assumed by our CDE approach.

We assume the latent space Z can be partitioned into scene-level variables  $Z_s$  (e.g., illumination, camera pose) and a set of object-level variables  $Z_o = \{Z_{n,k}\}_{n=1,k=1}^{N,K}$ , corresponding to the k-th property of the n-th object. An agent performs an action  $a \in \mathcal{A}$ , which performs an intervention on the system. This intervention transforms the pre-intervention state Z into a post-intervention state  $\tilde{Z}$ . Unobserved confounders (c) create spurious correlations and a training-testing distribution mismatch  $P_{\text{train}}(Z,a) \neq P_{\text{test}}(Z,a)$ . Following the *Independent Causal Mechanisms* principle (Schölkopf et al., 2012; Peters et al., 2017), we assume the true causal mechanism  $P(\tilde{Z}_a|Z_a,a)$  is invariant to this shift. Therefore, a robust model must learn this invariant mechanism instead of non-stationary correlations.

We investigate two challenging types of OOD shifts:

- Compositional Shifts: Training and test sets share the same object classes,  $O_{\text{train}} = O_{\text{test}}$ , but disjoint sets of object-action pairs.  $(A_{\text{train}} \times O_{\text{train}}) \cap (A_{\text{test}} \times O_{\text{test}}) = \emptyset$ .
- Systematic Shifts: The training and test sets of object classes are disjoint,  $O_{\text{train}} \cap O_{\text{test}} = \emptyset$ .

**Objective** Given a dataset of paired observations  $\mathcal{D} = \{(x, \tilde{x}, a)_j\}_{j=1}^M$ , where x and  $\tilde{x}$  are the pre- and post-intervention images respectively, and a is the corresponding action label, our objective is to learn a function  $\mathcal{F}: \mathcal{X} \times \mathcal{X} \to \mathcal{A}$ . This function must predict the action a by learning a representation that isolates the invariant causal signature of the intervention, thereby achieving high performance on OOD test data characterized by the compositional and systematic shifts defined above.

### 4 Causal Delta Embeddings

Consider an Encoder,  $\phi: \mathcal{X} \to \mathcal{Z}$  that maps a high-dimensional observation  $x \in \mathcal{X}$  to a point in the latent space  $\mathcal{Z}$ . A Delta Embedding is defined as follows.

**Definition 1 (Delta Embedding)** Given a pair of observations  $(x, \tilde{x})$  corresponding to the state of the world before and after an intervention  $a \in A$ , the Delta Embedding,  $\delta_a$ , is the vector difference

$$\delta_a := \phi(\tilde{x}) - \phi(x)$$

If the encoder is faithful to the data generating process illustrated by the model of Figure 2 and assuming identical noise across observations then for the Delta Embedding we have

$$\delta_a = \begin{bmatrix} 0 & \cdots & \tilde{z}_a - z_a & \cdots & 0 \end{bmatrix}^T \tag{1}$$

where  $z_a$  is the dimension (or subset of dimensions) of object n that is affected by action a. From equation 1 we observe the following properties of  $\delta_a$ .

- 1. *Independence*. Under the model of Figure 2, an action's representation is independent of the causally irrelevant elements of the scene, i.e. scene properties and objects not affected by *a*.
- 2. Sparsity. If the assumption of Sparse Mechanism Shifts Schölkopf et al. (2021); Liu et al. (2023) holds, then the action a will affect only a few underlying causal factors of the system, and the representation of the change,  $\delta_a$ , will be sparse.

To generalize to novel compositions of actions and objects, the action representation must satisfy additional properties. Specifically, even if the *Independence* and *Sparsity* properties are satisfied, if the action a affects different objects in a different way, a learning system will not be able to predict how the action will modify the representations of unseen objects, or even seen objects but without any examples of these objects with a in the training set.

We therefore introduce an additional requirement on the action representation, namely that it remains similar when applied to different objects, e.g., that the representation of action open is fundamentally the same, regardless of whether it is a door or a box that is being opened. We therefore introduce an additional property:

3. *Invariance*. The action representation  $\delta_a$  should not vary across different objects. One way to formalize this is through the variance of the delta embeddings across samples, i.e.,

$$Var_{x \sim P(X)}[\delta_a(x)] \approx \mathbf{0} \tag{2}$$

**Definition 2 (Causal Delta Embedding)** A Causal Delta Embedding (CDE) is a Delta Embedding that satisfies the properties of Independence, Sparsity and Invariance.

In terms of the SCM of Figure 2, Causal Delta Embeddings can be implemented by defining the structural equation of  $\tilde{z}_a$  as

$$\tilde{z}_a = f(z_a, \epsilon) = z_a + \delta_a + \mathcal{N}(0, \sigma \mathbb{I}) \tag{3}$$

where  $\sigma$  is small. The following section uses this definition to develop a strategy for learning Causal Delta Embeddings.

# 5 APPROACH

- 5.1 THE GLOBAL CAUSAL DELTA EMBEDDING MODEL
- 5.1.1 MODEL ARCHITECTURE

We first introduce a *global* model, i.e., a model that learns a single causal representation from the entire image. The model consists of three main components, as illustrated in Figure 3 (A).

The Encoder  $(\phi)$ : The encoder is responsible for mapping an input image x into the  $\mathcal{Z}$ . It is composed of two sub-modules. (i) **A Pre-trained Vision Backbone:** We use a powerful Vision Transformer (ViT) (Dosovitskiy et al., 2020), specifically one pre-trained with the DINO self-supervision algorithm (Caron et al., 2021). The backbone processes the input image and outputs a high-dimensional feature vector. We use the output corresponding to the '[CLS]' token as the global image representation. (ii) **A Causal Projector:** The feature vector from the backbone is then passed through a small multi-layer perceptron (MLP). This projector's role is to transform the general-purpose features into an l-dimensional representation satisfying the Causal Delta Embedding properties.

Delta Computation and the action classifier h: We compute the CDE according to Definition 1. Given the latent vectors for the pre-intervention image  $(z = \phi(x))$  and post-intervention image  $(\tilde{z} = \phi(\tilde{x}))$ , the delta is calculated via simple, element-wise subtraction,  $\delta = \tilde{z} - z$ . This vector is the sole input to a final classification head, which is an MLP that takes the l-dimensional delta and outputs logits for the different action classes in  $\mathcal{A}$ .

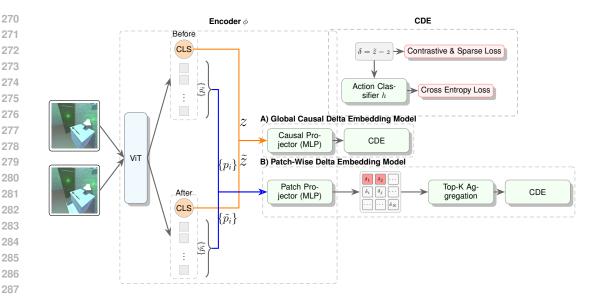


Figure 3: Model architecture. Model A (top) computes a global causal delta from CLS tokens. Model B (bottom) computes patch-wise deltas, aggregated to a causal delta. Both feed into a common action classifier.

### 5.1.2 IMPLEMENTATION OF THE LEARNING OBJECTIVE

To learn CDEs that satisfy the properties outlined in the Section 4, we combine three loss functions. (i) **Cross-Entropy Loss:** The primary objective is to ensure the delta embedding is useful for the downstream task. We use a standard Cross-Entropy loss,  $\mathcal{L}_{CE}$  between the predicted action logits  $h(\delta_i)$  and the one-hot ground-truth action label  $a_i$ . (ii) **Supervised Contrastive Loss:** To learn embeddings that are clustered together for the same action (Property 3, Invariance), we use the Supervised Contrastive Loss,  $\mathcal{L}_{contrast}$  (Khosla et al., 2020). For a batch of B delta embeddings, the loss for each embedding  $\delta_i$  (the "anchor") encourages it to be closer to other embeddings of the same class ("positives") than to all other embeddings in the batch.

$$\mathcal{L}_{\text{contrast}} = \sum_{i=1}^{B} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\text{sim}(\delta_i, \delta_p)/\tau)}{\sum_{j \neq i} \exp(\text{sim}(\delta_i, \delta_j)/\tau)}$$
(4)

where P(i) is the set of all positive samples for anchor i in the batch,  $sim(\cdot, \cdot)$  denotes the cosine similarity, and  $\tau$  is a scalar temperature hyperparameter. This loss component is also consistent with the structural equation 3. Finally, we introduce a (iii) **Sparsity Regularizer:** To encourage a minimal representation in line with the sparse mechanism shift hypothesis (Property 2, Sparsity), we apply an  $\ell_1$  regularization penalty. This loss penalizes the sum of the absolute values of the embedding dimensions, promoting solutions where most dimensions are zero.

$$\mathcal{L}_{\text{sparsity}} = \frac{1}{B} \sum_{i=1}^{B} \|\delta_i\|_1 = \frac{1}{B} \sum_{i=1}^{B} \sum_{k=1}^{l} |\delta_{i,k}|$$
 (5)

The final training objective is a weighted sum of these three components:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}} + \alpha_{\text{contrast}} \mathcal{L}_{\text{contrast}} + \alpha_{\text{sparsity}} \mathcal{L}_{\text{sparsity}}$$
 (6)

where  $\alpha_{\text{contrast}}$  and  $\alpha_{\text{sparsity}}$  are scalar hyperparameters that balance the influence of each loss component.

Notice that although no loss component explicitly enforces Property 1, this property is directly satisfied by the use of the Delta Embedding and image pairs, where the observed changes are only due to a. This does not hold under the actionable counterfactual case (Liu et al., 2023), where additional scene variables may change across observations. If, however, these changes are not spuriously correlated with the action a in the data, then the use of  $\mathcal{L}_{CE}$  can still retrieve the CDE representation.

### 5.2 SPATIAL EXTENSION: THE PATCH-WISE MODEL

In complex scenes with multiple objects or significant background noise, an action may only affect a small, localized region of the image. A global embedding risks 'averaging out' this important local change, making it difficult to detect. To address this, we developed a patch-wise extension of our model.

# 5.2.1 ARCHITECTURE

The Patch-Wise model adapts the core architecture to operate on local regions, as shown in Figure 3 (B).

The architecture includes (i) **Patch-wise Feature Extraction:** We use a ViT backbone, but instead of taking the global '[CLS]' token, we retain the output feature vectors for each individual image patch. This gives us a sequence of patch features for both the before and after images, (ii) **Patch-wise Delta Computation:** A shared Causal Projector and the subtraction operation are applied independently to each corresponding pair of patch features. This yields a set of delta embeddings,  $\{\delta_p\}$ , one for each spatial patch location p. (iii) **The Aggregation Module** ( $\mathcal{G}$ ): Since the action classifier needs a single input vector, we must aggregate the set of patch-wise deltas. This module's task is to identify the region of change and produce a single, representative delta vector  $\bar{\delta}$ . We employed simple **Top-K Aggregation:** This strategy is based on the assumption that the action's primary effect is localized to a few patches. We identify the k patches with the largest change by measuring the  $L_2$  norm of their delta vectors ( $\|\delta_p\|_2$ ). The final delta  $\bar{\delta}$  is the average of these top k patch deltas.

The same loss function ( $\mathcal{L}_{total}$ ) is then applied to the aggregated delta vector  $\bar{\delta}$ .

### 6 EXPERIMENTS

This section evaluates the effectiveness of our CDE framework. We first describe our experimental setup, then present the main quantitative results demonstrating CDE's effectiveness in OOD generalization, followed by qualitative and ablation analyses that provide deeper insights into its learned representations and design choices.

### 6.1 Experimental Setup

Our evaluation is conducted on the Causal Triplet benchmark (Liu et al., 2023), specifically designed for intervention-centric causal representation learning. This benchmark features three distinct settings of increasing visual complexity: single-object synthetic scenes, multi-object synthetic scenes (both from ProcTHOR (Deitke et al., 2022)), and challenging real-world scenes from Epic-Kitchens (Damen et al., 2022). In all settings models are trained on pairs of pre- and post-intervention images with action labels and are evaluated for their ability to infer the action. Further details on the datasets and data filtering procedures are provided in the Appendix.

We follow the Causal Triplet protocol, evaluating models on both IID and OOD test sets. The OOD splits test two forms of generalization: Compositional Distribution Shifts, where the model encounters unseen combinations of actions and objects (e.g., open (drawer) when only open (door) and close (drawer) where seen during training); and Systematic Distribution Shifts, where generalization to entirely unseen object classes is required. Visualizations of these distribution shifts are available in the Appendix. All reported quantitative results are mean accuracies and standard deviations average over 3 random seeds. We set  $\alpha_{\rm contrast}=2.0$  and  $\alpha_{\rm sparsity}=1.0$  for all experiments (see the Appendix for more details).

We compare our two proposed models against the baselines from the Causal Triplet paper (Liu et al., 2023), including vanilla ResNets (He et al., 2016), methods incorporating causal regularization (ICM, SMS), and object-centric approaches (Slot Attention (Locatello et al., 2020b), GroupViT (Xu et al., 2022)).

Table 1: Single-object ProcTHOR results. Our Global Delta Embedding model significantly improves OOD generalization under both compositional and systematic shifts. (R: ResNet-18, V: Vit-Small)

Method	IID Acc.	OOD Comp.	OOD Syst.	Gap Syst. (↓)
Vanilla-R	<b>0.96</b> ±0.01	$0.36_{\pm0.13}$	$0.48_{\pm 0.08}$	0.48
Vanilla-V	$0.95_{\pm 0.01}$	$0.34_{\pm 0.27}$	$0.47_{\pm 0.11}$	0.48
ICM-R	$0.95_{\pm 0.01}$	$0.41_{\pm 0.15}$	$0.50_{\pm 0.09}$	0.45
ICM-V	$0.95_{\pm 0.01}$	$0.38_{\pm 0.26}$	$0.49_{\pm 0.01}$	0.46
SMS-R	$0.96_{\pm 0.01}$	$0.47_{\pm 0.18}^{-}$	$0.54_{\pm 0.07}$	0.42
SMS-V	$0.95_{\pm 0.01}$	$0.34_{\pm 0.27}$	$0.39_{\pm 0.04}$	0.56
Ours <sub>(Global)</sub>	<b>0.96</b> <sub>± 0.01</sub>	<b>0.91</b> <sub>± 0.02</sub>	<b>0.73</b> <sub>± 0.02</sub>	0.18

Table 2: Results across multi-object ProcTHOR and Epic-Kitchens (systematic shift).

Dataset	Method	IID Acc.	OOD Acc.	Gap
ProcTHOR	ResNet	$0.83_{\pm 0.01}$	$0.30_{\pm 0.08}$	0.53
	Oracle-mask	$0.90_{\pm 0.01}$	$0.42_{\pm 0.06}$	0.48
	Slot-avg	$0.49_{\pm 0.01}$	$0.15_{\pm 0.01}$	0.34
	Slot-dense	$0.51_{\pm 0.01}$	$0.19_{\pm 0.03}$	0.32
	Slot-match	$0.66_{\pm 0.01}$	$0.21_{\pm 0.01}$	0.45
	$Ours_{(Patch\text{-}wise)}$	$0.92_{\pm 0.01}$	$\textbf{0.45}_{\pm0.02}$	0.47
Epic-Kitchens	ResNet	$0.42_{\pm 0.03}$	$0.17_{\pm 0.03}$	0.25
•	CLIP	$0.45_{\pm 0.02}$	$0.24_{\pm 0.02}$	0.21
	Group-avg	$0.47_{\pm 0.03}$	$0.24_{\pm 0.03}$	0.23
	Group-dense	$0.50_{\pm 0.04}$	$0.26_{\pm 0.03}$	0.24
	Group-token	$0.52_{\pm 0.03}$	$0.27_{\pm 0.03}$	0.25
	Ours <sub>(Patch-wise)</sub>	$0.55 \pm 0.02$	<b>0.33</b> <sub>±0.02</sub>	0.22

### 6.2 Main Quantitative Results

Our CDE framework consistently delivers substantial improvements in OOD accuracy across all evaluation settings, establishing a new state of the art. For single-object scenes, our global CDE model cuts the generalization gap from 0.56 to 0.18 while matching IID accuracy (Table 1). In challenging multi-object and real-world settings (Table 2), our Patch-Wise model outperforms all baselines, including oracle methods that use ground-truth segmentations masks.

### 6.3 ACTION RELATIONSHIPS IN CAUSAL DELTA SPACE

To study the semantic structure of the learned delta space, we investigated whether the model could discover fundamental relationships between actions on its own. We computed the pairwise cosine similarity between all learned action representations. The result is visualized in the appendix (Figure 7). The analysis reveals that the model has learned a perfect *anti-parallel relationship* for opposite actions. The cosine similarity between the representations for open and close, for dirty and clean, as well as for turn on and turn off, is -1.0. This demonstrates that our framework not only separates the action concepts but also discovers opposing relationships between them, organizing the representations in a meaningful way. A similar pattern is observed in the more challenging real-world dataset where the model learns the anti-parallel representations for the open and close action pair, as well as for the fold and stretch pair (see Figure 8 in the Appendix for details).

In summary, the combination of strong predictive properties and consistent semantic structure demonstrates that our CDE framework learns meaningful representations of interventions. For further geometric analysis of the delta space, including UMAP projections and k-NN classifier performance, refer to the Appendix.

Table 3: Ablation study of our method's components on the ViT-Small model. Results are for the single-object systematic shift setting, showing the impact on OOD accuracy when each core component is removed.

<b>Model Configuration</b>	<b>IID Acc.</b> (%)	<b>OOD Acc.</b> (%)	
Full Model	0.95	0.73	
Ablations			
w/o Sparsity Loss	0.96	0.70	
w/o Contrastive Loss	0.95	0.60	
Baseline (CE Loss only)	0.94	0.59	

#### 6.4 ABLATION STUDY

To understand the contribution of each component of our CDE framework, we also conducted a series of ablation studies, by analyzing the impact of each major loss component on the performance of our primary model with a ViT-Small backbone. Table 3 presents the results, comparing our full model against versions where each loss component is removed, and a baseline trained only with standard CE loss.

The results demonstrate the effectiveness of our approach. Our full model achieves an OOD accuracy of 73.0%, a +14 point improvement over the baseline trained solely with a CE objective, validating that explicitly structuring the representation space is critical for generalization. Removing the supervised contrastive loss component causes a 13-point drop in OOD accuracy. Removing the sparsity loss term causes another 3-point drop. Please refer to the Appendix for further ablation experiments.

### 7 CONCLUSION

This paper introduces the *Causal Delta Embedding (CDE)* framework, a simple yet effective approach to interventional causal representation learning. By explicitly modeling interventions as delta vectors in a structured latent space, CDE inherently satisfies the properties of independence, sparsity and invariance, leading to improved generalization. Our empirical validation on the Causal Triplet challenge demonstrates that CDE achieves state-of-the-art OOD generalization, outperforming prior methods across synthetic and real world datasets. Beyond quantitative gains, we show that CDE learns semantically meaningful representations without supervision, where opposing actions have anti-parallel representations. Despite the promising results, we acknowledge that limitations remain for real-world data, since both IID and OOD accuracies are still low for real world deployment, and also the use of universal delta embeddings for each action limits its ability to capture context-dependent visual transformations of actions. Future research directions include dynamically identifying modified regions of the input states through attention mechanisms, extending the framework to video streams for modeling causal dynamics in temporal sequences, and investigating compositional properties of delta embeddings to enable multi-step interventions and generalization to novel action sequences.

**Reproducibility statement:** The previous sections have outlined the main building blocks of the proposed method, as well as the approach followed in the experiments, with the Appendix providing additional information and results. The code to reproduce the experiments is attached as supplementary material (without any identifiable information of authors) and will be made publicly available upon acceptance. Finally, all experiments were carried out by strictly following the Causal Triplet benchmark (Liu et al., 2023) evaluation protocols, which relies on the publicly available ProcTHOR (Deitke et al., 2022) and Epic-Kitchens (Damen et al., 2022) datasets.

### REFERENCES

Kartik Ahuja, Jason S Hartford, and Yoshua Bengio. Weakly supervised representation learning with sparse perturbations. *Advances in Neural Information Processing Systems*, 35:15516–15528,

486 2022.

- Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal representation learning. In *International conference on machine learning*, pp. 372–407. PMLR, 2023.
- Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6836–6846, 2021.
- Johann Brehmer, Pim De Haan, Phillip Lippe, and Taco S Cohen. Weakly supervised causal representation learning. *Advances in Neural Information Processing Systems*, 35:38319–38331, 2022.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6299–6308, 2017.
- Min-Hung Chen, Zsolt Kira, Ghassan AlRegib, Jaekwon Yoo, Ruxin Chen, and Jian Zheng. Temporal attentive alignment for large-scale video domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6321–6330, 2019.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PmLR, 2020.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision*, pp. 1–23, 2022.
- Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Kiana Ehsani, Jordi Salvador, Winson Han, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Procthor: Large-scale embodied ai using procedural generation. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6202–6211, 2019.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Tarun Gupta, Wenbo Gong, Chao Ma, Nick Pawlowski, Agrin Hilmkil, Meyer Scetbon, Marc Rigter, Ade Famoti, Ashley Juan Llorens, Jianfeng Gao, et al. The essential role of causality in foundation world models for embodied ai. *arXiv preprint arXiv:2402.06665*, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Thomas Hellström. The relevance of causation in robotics: A review, categorization, and analysis. *Paladyn, Journal of Behavioral Robotics*, 12(1):238–255, 2021.

- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8340–8349, 2021.
  - Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2017.
  - Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvarinen. Variational autoencoders and nonlinear ica: A unifying framework. In *International conference on artificial intelligence and statistics*, pp. 2207–2217. PMLR, 2020.
  - Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
  - Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
  - Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.
  - Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Stratis Gavves. Citris: Causal identifiability from temporal intervened sequences. In *International Conference on Machine Learning*, pp. 13557–13603. PMLR, 2022.
  - Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M Asano, Taco Cohen, and Efstratios Gavves. Biscuit: Causal representation learning from binary interactions. In *Uncertainty in Artificial Intelligence*, pp. 1263–1273. PMLR, 2023.
  - Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pp. 38–55. Springer, 2024.
  - Yuejiang Liu, Alexandre Alahi, Chris Russell, Max Horn, Dominik Zietlow, Bernhard Schölkogpf, and Francesco Locatello. Causal triplet: An open challenge for intervention-centric causal representation learning. In *Conference on Causal Learning and Reasoning*, pp. 553–573. PMLR, 2023.
  - Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. In *International conference on machine learning*, pp. 6348–6359. PMLR, 2020a.
  - Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in neural information processing systems*, 33:11525–11538, 2020b.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
  - Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. A survey on vision-language-action models for embodied ai. *arXiv preprint arXiv:2405.14093*, 2024.
  - Sara Magliacane, Thijs Van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions. *Advances in neural information processing systems*, 31, 2018.
  - Ricardo Pio Monti, Kun Zhang, and Aapo Hyvärinen. Causal discovery with general non-linear relationships using non-linear ica. In *Uncertainty in artificial intelligence*, pp. 186–195. PMLR, 2020.

- Jonathan Munro and Dima Damen. Multi-modal domain adaptation for fine-grained action recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 122–132, 2020.
  - Judea Pearl. Causality. Cambridge university press, 2009.

- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- Niklas Pfister and Jonas Peters. Identifiability of sparse causal effects using instrumental variables. In *Uncertainty in Artificial Intelligence*, pp. 1613–1622. PMLR, 2022.
- Pedro Sanchez, Jeremy P Voisey, Tian Xia, Hannah I Watson, Alison Q O'Neil, and Sotirios A Tsaftaris. Causal machine learning for healthcare and precision medicine. *Royal Society Open Science*, 9(8):220638, 2022.
- Ranjan Sapkota, Yang Cao, Konstantinos I Roumeliotis, and Manoj Karkee. Vision-language-action models: Concepts, progress, applications and challenges. *arXiv preprint arXiv:2505.04769*, 2025.
- Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij. On causal and anticausal learning. *arXiv preprint arXiv:1206.6471*, 2012.
- Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.
- Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, et al. Bridging the gap to real-world object-centric learning. *arXiv preprint arXiv:2209.14860*, 2022.
- Xinwei Shen, Furui Liu, Hanze Dong, Qing Lian, Zhitang Chen, and Tong Zhang. Disentangled generative causal representation learning, 2020.
- Chandler Squires, Anna Seigal, Salil S Bhate, and Caroline Uhler. Linear causal disentanglement via interventions. In *International conference on machine learning*, pp. 32540–32560. PMLR, 2023.
- Shanshan Wang, Yiyang Chen, Zhenwei He, Xun Yang, Mengzhu Wang, Quanzeng You, and Xingyi Zhang. Disentangled representation learning with causality for unsupervised domain adaptation. In *Proceedings of the 31st ACM international conference on multimedia*, pp. 2918–2926, 2023.
- Yixin Wang and Michael I Jordan. Desiderata for representation learning: A causal perspective. *Journal of Machine Learning Research*, 25(275):1–65, 2024.
- Liang Wendong, Armin Kekić, Julius von Kügelgen, Simon Buchholz, Michel Besserve, Luigi Gresele, and Bernhard Schölkopf. Causal component analysis. Advances in Neural Information Processing Systems, 36:32481–32520, 2023.
- Danru Xu, Dingling Yao, Sébastien Lachapelle, Perouz Taslakian, Julius Von Kügelgen, Francesco Locatello, and Sara Magliacane. A sparsity principle for partially observable causal representation learning. *arXiv preprint arXiv:2403.08335*, 2024.
- Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18134–18144, 2022.
- Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye Hao, and Jun Wang. Causalvae: Disentangled representation learning via neural structural causal models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9593–9602, 2021.
- Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pp. 2165–2183. PMLR, 2023.

# A DATASET DETAILS

This section provides further details on the datasets used in our evaluation.

**ProcTHOR** The ProcTHOR dataset Deitke et al. (2022) provides synthetic indoor scenes. For our single-object scenes, each scene contains one manipulated object, ensuring a clear focus on the intervention. In multi-object scenes, multiple objects are present, increasing the visual complexity of the scene, although only one object is again manipulated. We follow the dataset generation and filtering procedures as described in Liu et al. (2023) to ensure consistency with the Causal Triplet benchmark.

**Epic-Kitchens** The Epic-Kitchens dataset Damen et al. (2022) comprises real world egocentric videos of diverse kitchen activities. From this, we extract pre- and post-intervention image pairs. Unlike synthetic environments, Epic-Kitchens introduces significant real world challenges such as camera motion, varying lighting conditions, occlusions and dynamic backgrounds, making the task of isolating interventions particularly challenging. To ensure dataset quality, a two-stage filtering process using Grounding DINO Liu et al. (2024) for zero-shot object detection is applied. For each extracted pair, the pipeline verifies that the target object appears clearly in both frames with a detection confidence above a set threshold t=0.45. This automated filtering removes cases with poor object visibility or excessive motion blur.

### A.1 VISUALIZATIONS OF OOD SHIFTS

Figures 4, 5 and 6 visually illustrate the compositional and systematic distribution shifts utilized in the Causal Triplet benchmark.

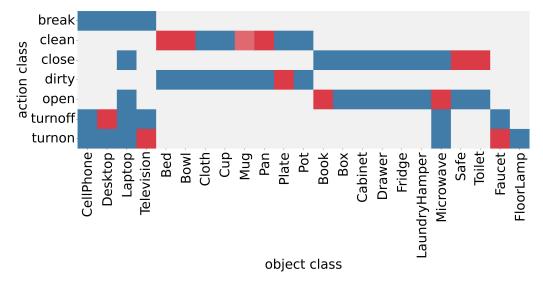


Figure 4: Compositional Distribution Shift in the ProcThor dataset. Blue boxes indicate IID data, while red boxes indicate novel OOD action-object combinations.

### B GEOMETRIC ANALYSIS OF CAUSAL DELTA EMBEDDINGS

This section provides additional analysis of the geometric properties of the learned Causal Delta Embeddings, complementing the insights presented along with the experimental results.

### B.1 ACTION REPRESENTATION RELATIONSHIPS LEARNED FROM REAL-WORLD DATASETS

Figure 7 illustrates the pairwise cosine similarities between the embeddings learned for all actions in the ProcTHOR dataset, while Figure 8 presents the same information for the more challeng-

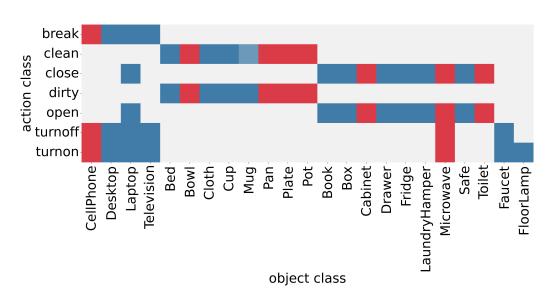


Figure 5: Systematic Distribution Shift in the ProcThor dataset. Blue boxes indicate IID data, while red boxes indicate novel OOD objects that the model has not encountered during training.

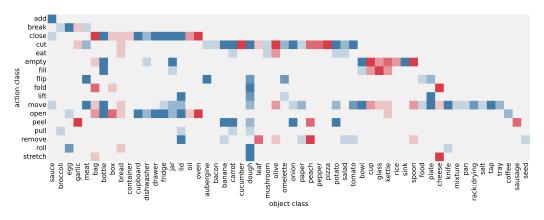


Figure 6: Systematic Distribution Shift in the EpicKitchens dataset. Blue boxes indicate IID data, while red boxes indicate novel OOD objects that the model has not encountered during training.

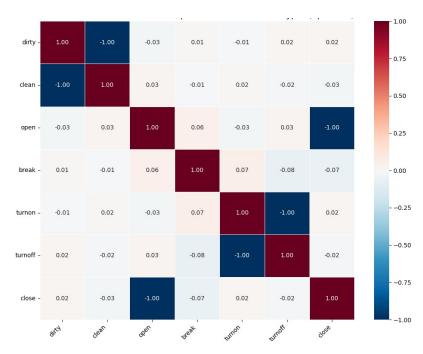


Figure 7: Heatmap of pairwise cosine similarities between all learned delta embeddings. The strong blue squares (similarity near -1.0) reveal a near-perfect anti-parallel relationship for opposite action pairs, which was discovered entirely from the data.

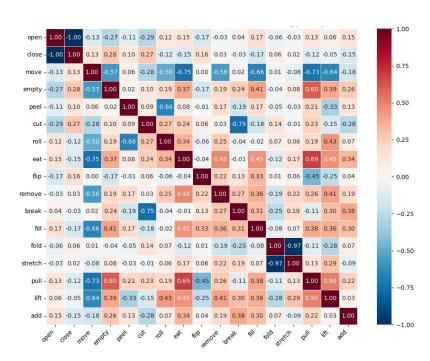


Figure 8: Heatmap of pairwise cosine similarities between all learned action prototypes for the EpicKitchens dataset.

ing real-world Epic Kitchens dataset. We observe that in both cases the learned relationships for opposing actions such as open and close as well as fold and stretch are antiparallel in the embedding space.

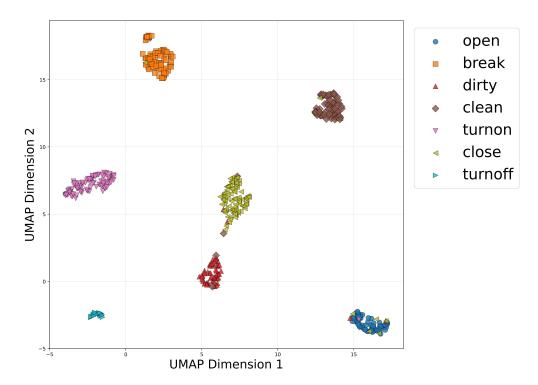


Figure 9: UMAP projection of individual delta embeddings from the IID test set. Embeddings are shaped by their ground-truth action. The plot reveals strong global separation between different action clusters.

### B.2 ANALYSIS OF LEARNED ACTION REPRESENTATIONS

To study the properties of the action representations resulting from our method, we first tested if the resulting delta embeddings could reliably predict the outcome of an intervention. To do this, for each sample in the OOD test set, we took the 'before' state embedding (z) and added the corresponding average action vector  $(\mu_{action})$  that was computed using the training set samples. We then measured the cosine similarity between this predicted 'after' state and the ground-truth 'after' state  $(\tilde{z})$ . Our framework showed remarkable predictive power, achieving an average cosine similarity of 0.98 in the single object systematic shift setting. This near perfect score confirms that the learned action prototypes function as true, generalizable transformation vectors, providing strong evidence that our model has learned the underlying mechanics of interventions.

#### **B.3** UMAP PROJECTION

Figures 9 and 11 present the UMAP projection of individual delta embeddings from the IID and OOD test set of the single-object environment respectively. The delta embeddings in the IID setting achieve a clear separation between each action, leading to a near perfect IID accuracy as was presented by our experiments. On the other hand, while strong intra-class cohesion is visible, the global separation between these action clusters is not always visually distinct in the OOD setting. This suggests that while representations remain locally coherent, action representations are not as clearly discriminated compared to the IID setting. It is worth mentioning, however, that the 2D projection may not fully capture the features of the high-dimensional latent space.

### B.3.1 k-NN CLASSIFIER PERFORMANCE

To quantitatively assess the quality of the local structure of the learned representations, we evaluate the performance of a simple, non-parametric k-Nearest Neighbors (k-NN) classifier directly on the Causal Delta Embeddings, where we set the number of neighbors k=5. A high k-NN accuracy

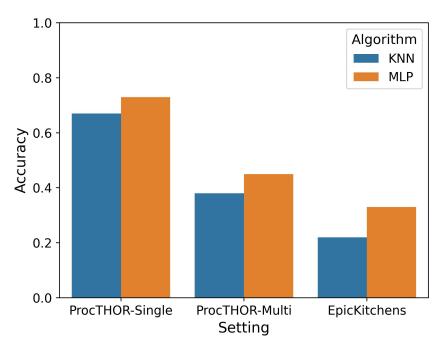


Figure 10: Comparison of classifier performance on OOD test sets. The k-NN classifier, is able to classify OOD samples if trained with Delta Embeddings.

indicates that the local neighborhoods are semantically meaningful and highly predictive of the action class. Figure 10 compares the k-NN classifier accuracy with the use of CDE with an MLP head (as in Figure 3) across all three benchmark settings. Although k-NN does not match the effectiveness of the MLP head, it still achieves comparable results in novel OOD samples, especially in the synthetic dataset. This provides an indication that the invariance property holds in the OOD case.

# C ABLATION STUDIES

In order to understand the effectiveness of each component of our method, we conducted a series of ablation studies to evaluate the impact of different backbone architectures, the impact of loss hyperparameters  $\alpha$  and the impact of the hyperparameter k in the Top-K selection procedure for our Patch-Wise model. All the subsequent experiments ran on the single-object systematic shifts setting, except for the Top-K ablation study which ran on the multi-object systematic shifts setting.

#### C.1 IMPACT OF BACKBONE ARCHITECTURE

To isolate the contribution of our CDE framework from the choice of feature extractor, we conducted a controlled comparison between our ViT-Small backbone and the ResNet-18 backbone used in the original Causal Triplet benchmark.

# C.2 IMPACT OF LOSS HYPERPARAMETERS

In order to select values for  $\alpha_{contrast}$  and  $\alpha_{sparsity}$ , we conducted an ablation study comparing various values and combinations between them. Table 5 compares some of the combinations of the values that we experimented with. Selecting a larger value for  $\alpha_{contrast}$ , rather than  $\alpha_{sparsity}$ , helps the model learn better representations, thus achieving better OOD accuracy. We set  $\alpha_{contrast}=2.0$  and  $\alpha_{sparsity}=1.0$  in all our main experiments.

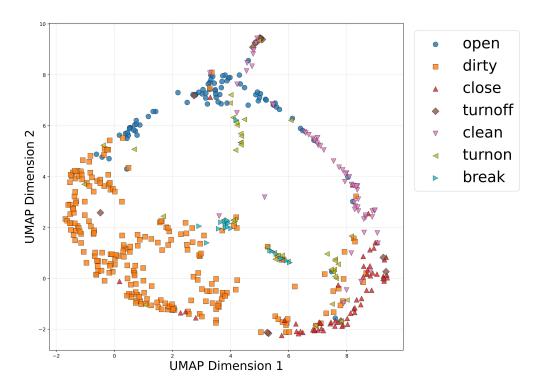


Figure 11: UMAP projection of individual delta embeddings from the OOD test set. Embeddings are shaped by their ground-truth action. The plot reveals strong local cohesion (points of the same shape cluster together) but shows a lack of clear global separation between the different action clusters.

Table 4 compares the OOD performance of the best configuration for each backbone against the benchmark's state-of-the-art ResNet-18 result. Our final ViT-based model significantly outperforms the best ResNet-based model, demonstrating that while the richer features from ViT enhance performance, the substantial gains are primarily driven by our proposed CDE learning framework.

Table 4: Comparison of OOD performance with different backbone architectures on the single-object systematic shift benchmark.

Backbone	Method	<b>OOD Acc.</b> (%)
ResNet-18 ResNet-18	Liu et al. (2023) Ours*	$0.54 \\ 0.45$
ViT-Small ViT-Small	Ours (CE Only) Ours (Full Model)	0.59 <b>0.73</b>

Best ResNet performance from our experiments was with CE + Con Loss.

Table 5: Comparison of various hyperparameter values for  $\alpha_{\text{contrast}}$  and  $\alpha_{\text{sparsity}}$  on the single-object systematic shift benchmark.

$\alpha_{\mathrm{contrast}}$	$\alpha_{ extsf{sparsity}}$	OOD Acc. (%)
0.0	0.0	$0.21_{\pm 0.02}$
0.1	1.0	$0.27_{\pm 0.11}$
1.0	0.1	$0.28_{\pm 0.04}$
0.5	0.5	$0.29_{\pm 0.07}$
1.0	2.0	$0.31_{\pm 0.07}$
2.0	1.0	$0.33_{\pm 0.07}$

# C.3 TOP-K SELECTION

In order to select the hyperparameter k in multi-object and real world data settings, we executed an ablation study to understand the sensitivity of our method to this parameter. As presented in Table 6, we can see that OOD accuracy increases as k increases too. This observation makes sense, since bigger objects (e.g. Fridge, Bed) would need more patches for their representations in order to be captured effectively. Thus, we set the value of k=4 across all our multi-object and real world experiments.

Table 6: Comparison of OOD performance with k values for the patch selection process in multi-object settings.

k	ProcTHOR	EpicKitchens
k = 1	$0.42_{\pm 0.07}$	$0.12_{\pm 0.03}$
k = 2	$0.45_{\pm 0.06}$	$0.13_{\pm 0.03}$
k = 3	$0.47_{\pm 0.04}$	$0.13_{\pm 0.02}$
k = 4	$0.48_{\pm 0.04}$	$0.15_{\pm 0.02}$

# D EXPERIMENTAL DETAILS

#### D.1 HYPERPARAMETERS

Table 7 summarizes the key hyperparameters used across all experiments. These values were selected based on ablation studies and remained consistent across different experimental settings unless otherwise noted.

#### D.2 EXECUTION ENVIRONMENT

All experiments were run on a NVIDIA A100 GPU with the Slurm Workload Manager. The code was implemented in Python, using the Pytorch library. Each run takes approximately one hour to complete for the ProcTHOR and two hours for the Epic-Kitchens dataset.

### D.3 IMAGE AUGMENTATIONS

We do not apply any augmentations to the images, since we do not want to modify the interventional nature of the pairs. Augmentation in this problem could harm our assumptions. For example, a rotation could affect Equation equation 1 and eliminate the faithfulness of the encoder. We leave it as future work to investigate whether augmentations can boost OOD performance under different assumptions. We only resize images to  $224 \times 224$  pixels and apply zero-mean normalization with unit variance.

Table 7: Summary of hyperparameters used across all experiments.

Parameter	Value
Learning Rate	$1 \times 10^{-4}$
Backbone LR	$1 \times 10^{-5}$
Batch Size	128
Epochs	50 (100 for Epic-Kitchens)
Weight Decay	0.05
$lpha_{ m contrast}$	2.0
$\alpha_{ m sparsity}$	1.0
Temperature $(\tau)$	0.07
Top-K(k)	4
Embedding Dim. (l)	256 (512 for Epic-Kitchens)
Input Resolution	$224 \times 224$

# D.4 OPTIMIZATION

We use a batch size of 128 and an AdamW Loshchilov & Hutter (2017) optimizer with a cosine annealing learning scheduler for 50 epochs. In the real world setting, we instead train for 100 epochs. The ViT feature extractor is not frozen but fine-tuned with a reduced learning rate of 10% of the network's base learning rate, which is set to  $1\times10^{-4}$ . The weight decay parameter is 0.05. All reported results include standard deviations computed over three independent runs with different random seeds.