

---

# *HiMaCon: Discovering Hierarchical Manipulation Concepts from Unlabeled Multi-Modal Data*

---

RuiZhe Liu<sup>1</sup> Pei Zhou<sup>1</sup> Qian Luo<sup>1,4</sup> Li Sun<sup>1</sup>  
 Jun Cen<sup>3</sup> Yibing Song<sup>3</sup> Yanchao Yang<sup>1,2</sup>

<sup>1</sup>HKU Musketeers Foundation Institute of Data Science, The University of Hong Kong

<sup>2</sup>Department of Electrical and Electronic Engineering, The University of Hong Kong

<sup>3</sup>DAMO Academy, Alibaba Group <sup>4</sup>Transcengram

{zrllrz360, pezhou, qianluo, sunlids}@connect.hku.hk

{cenjun.cen, songyibing.syb}@alibaba-inc.com, yanchaoy@hku.hk

## Abstract

Effective generalization in robotic manipulation requires representations that capture invariant patterns of interaction across environments and tasks. We present a self-supervised framework for learning hierarchical manipulation concepts that encode these invariant patterns through cross-modal sensory correlations and multi-level temporal abstractions without requiring human annotation. Our approach combines a cross-modal correlation network that identifies persistent patterns across sensory modalities with a multi-horizon predictor that organizes representations hierarchically across temporal scales. Manipulation concepts learned through this dual structure enable policies to focus on transferable relational patterns while maintaining awareness of both immediate actions and longer-term goals. Empirical evaluation across simulated benchmarks and real-world deployments demonstrates significant performance improvements with our concept-enhanced policies. Analysis reveals that the learned concepts resemble human-interpretable manipulation primitives despite receiving no semantic supervision. This work advances both the understanding of representation learning for manipulation and provides a practical approach to enhancing robotic performance in complex scenarios. Code is available at: <https://github.com/zrllrz/HiMaCon>.

## 1 Introduction

Robot manipulation in diverse, unstructured environments remains a fundamental challenge. Despite advances in policy learning and architectures [4, 10, 19, 24], current approaches often fail when encountering unexpected variations or novel scenarios. As illustrated in Fig. 1, a policy trained to place cups into containers may succeed in familiar settings but fail when encountering unexpected barriers—revealing a critical generalization gap limiting real-world deployment.

We propose that addressing this challenge requires learning transferable *manipulation concepts*—hierarchical abstractions capturing fundamental manipulation patterns. These concepts connect low-level actions to high-level goals, enabling robust generalization. For example, the concept of “placing an object inside a container” encompasses invariant relational patterns that persist whether the container has barriers or not, allowing adaptation while maintaining core manipulation strategy.

To acquire these manipulation concepts, we propose a self-supervised framework that learns hierarchical latent representations without requiring labor-intensive human annotations [13, 28, 40]. Our approach operates through two complementary mechanisms: 1) *Cross-modal correlation learning* captures invariant patterns across different sensory modalities (vision, proprioception), enabling generalization across visual variations while preserving functional relationships. When placing ob-

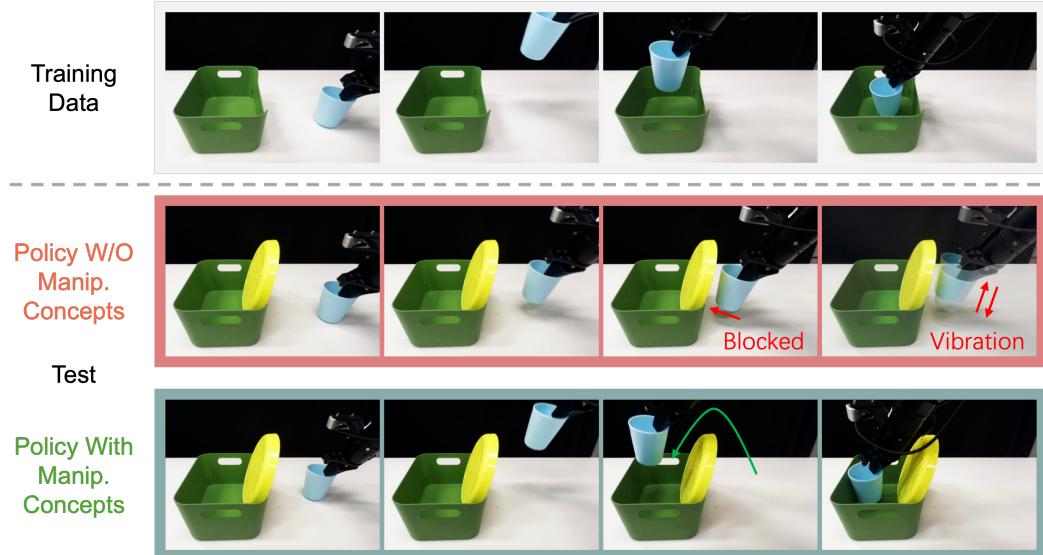


Figure 1: **Manipulation concepts enhance generalization.** Top: Training data with cups and containers without barriers. Middle: Without manipulation concepts, policies fail when encountering barriers. Bottom: With our concept enhancement, policies adapt accordingly.

jects in containers, these correlations encode the relationship between visual perception of container boundaries and proprioceptive feedback during placement, regardless of container appearance. 2) *Multi-horizon sub-goal organization* structures concepts hierarchically across temporal scales, from immediate actions (e.g., “align gripper with object”) to extended sequences (e.g., “transport object to container”). This hierarchical representation enables policies to simultaneously reason about immediate actions and longer-term goals, maintaining task coherence even when specific execution paths require adaptation.

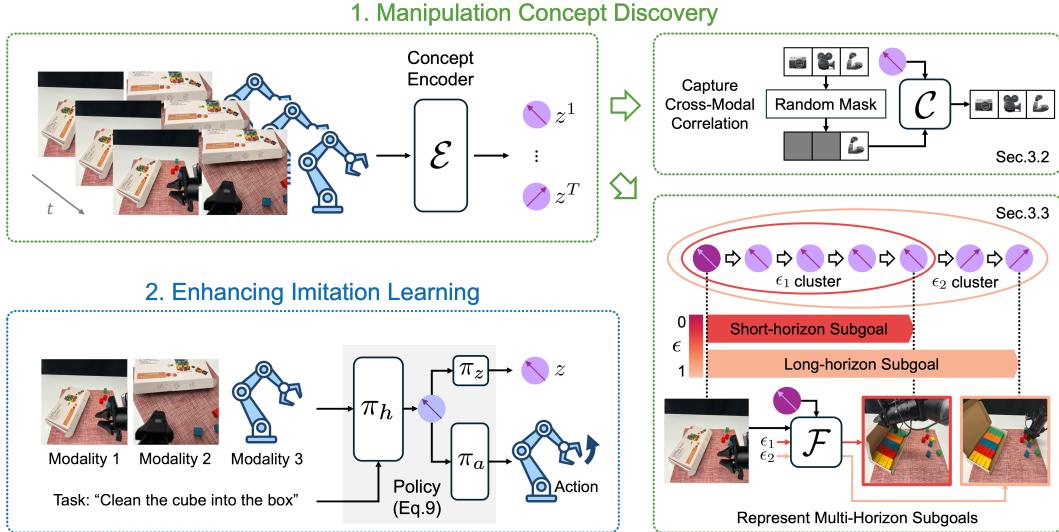
Our experiments across both simulated benchmark tasks and real-world robot deployments demonstrate that policies enhanced with these manipulation concepts consistently outperform conventional approaches, particularly in challenging scenarios requiring adaptation to novel objects, unexpected obstacles, and environmental variations (Fig. 1). The learned concepts form interpretable clusters that resemble meaningful manipulation primitives, providing insights into how robots perceive and reason about manipulation tasks.

In summary, our key contributions include: (1) a self-supervised framework that extracts structured hierarchical manipulation concepts from unlabeled multi-modal demonstrations, capturing both cross-modal correlations and multi-level temporal abstractions without human annotation; (2) an effective policy enhancement approach that integrates these concepts through joint prediction, maintaining compatibility with diverse policy architectures; and (3) comprehensive empirical evidence demonstrating significant performance improvements across diverse settings, with analyses revealing how learned concepts enable more robust generalization to novel environments.

## 2 Related Work

**Representation Learning in Robotics** Self-supervised representation learning has emerged as a powerful approach for extracting meaningful skills [29, 32, 48] from robotic data, avoiding the need for manual annotation in methods such as [13, 28, 37, 40]. Initial efforts explored single-modality approaches for vision-based [7, 11, 51, 65, 68] and proprioception-based [26, 39, 45, 52] representation learning. Recent work integrates multiple modalities, combining vision with language [22, 36, 42, 50, 64], proprioception with vision [47, 62, 66], even richer [6, 53, 69].

These approaches typically focus on cross-modal alignment but often overlook the structured temporal patterns inherent in manipulation tasks. Parallel developments in temporal representation learning have addressed this challenge through various approaches: time-contrastive learning [27, 34, 35, 42, 65], temporal masked auto-encoding [47], and explicit modeling of state transitions across different



**Figure 2: The proposed self-supervised manipulation concept discovery and policy enhancement.** *Stage 1:* The concept encoder ( $E$ ) processes multi-modal robot demonstrations to extract concept latents. These latents are refined through two objectives: (1) the Cross-Modal Correlation Network ( $C$ ) employs a mask-and-predict strategy to capture persistent patterns across sensing modalities (Sec. 3.2); (2) the Multi-Horizon Future Predictor ( $F$ ) enables concept latents to organize hierarchically into multi-horizon sub-goals based on coherence thresholds ( $\epsilon$ ) (Sec. 3.3). *Stage 2:* The learned concepts are integrated into policy learning through a backbone network ( $\pi_h$ ) with concept ( $\pi_z$ ) and action ( $\pi_a$ ) prediction heads, regularizing action generation with structured manipulation knowledge (Eq. 9).

timescales [25, 44, 54, 55]. Our work advances the field by simultaneously addressing both multi-modal integration and hierarchical temporal structures, creating representations that naturally align with manipulation sub-goals at varying time horizons while leveraging cross-modal correlational patterns that persist across different objects and contexts.

**Concept-Guided Robotic Policies** Concept-guided approaches enhance robotic policy performance by leveraging intermediate representations that bridge perception and action. These methods generally fall into two categories. First, two-module frameworks [3, 8, 29, 62, 68, 73] employ a dual-model architecture where one component extracts high-level task concepts while another generates the corresponding actions. While effective, these approaches often require specialized architectural designs that limit their applicability across different policy classes.

Second, in contrast, joint prediction approaches [18, 20, 56, 58, 67] integrate concept guidance by training policies to simultaneously predict both concepts and actions. This creates an implicit information pathway where concept understanding regularizes action generation. Our work adopts this more flexible approach, enabling seamless integration with diverse policy architectures while maintaining the interpretability benefits of explicit concept representations.

### 3 Method

We aim to encode robotic manipulation demonstrations into latent representations that capture task-induced patterns in multi-modal sensory-motor data. These representations should naturally cluster according to functional sub-goals, providing insights into manipulation objectives and enhancing policy learning. We term these clusters *manipulation concepts*—each representing action sequences targeting specific sub-goals—and call the learning process *manipulation concept discovery*.

Our self-supervised approach works without explicit sub-goal annotations, addressing the challenge of capturing meaningful manipulation patterns without labels. We design objective functions enforcing latent representations that reflect both temporal structure and cross-modal correlations. Our approach ensures: (1) integration of modality-specific features while encoding cross-modal correlations that persist across objects and contexts; (2) hierarchical organization of sub-processes representing sub-goals across temporal horizons and enabling action prediction guided by immediate and long-term

objectives. We validate these concepts through policy performance improvements and multiple analysis methods that demonstrate correspondence with meaningful manipulation primitives.

### 3.1 Problem Setup & Manipulation Concept Encoder

Given a dataset  $D = \{\tau_i\}_{i=1}^N$  of  $N$  manipulation trajectories, each  $\tau_i = \{(\mathbf{o}_i^t, a_i^t)\}_{t=1}^{T_i}$  contains observations  $\mathbf{o}_i^t$  and actions  $a_i^t$  at time  $t$ . For  $M$  modalities,  $\mathbf{o}_i^t = \{o_i^{1,t}, o_i^{2,t}, \dots, o_i^{M,t}\}$ , where  $o_i^{m,t}$  is the observation of the modality  $m$ . We denote  $\mathbf{o}_i^{S,t} = \{o_i^{m,t} \mid m \in S\}$  as observations of modalities in  $S \subseteq [M] = \{1, 2, \dots, M\}$ . We treat observations of both the same sensory modes (e.g., multiple views) and different modes as distinct modalities, as they are functionally different modalities in terms of complementary information.

The *Manipulation Concept Discovery* process assigns latent representation  $z_i^t \in \mathbb{R}^Z$  to each timestep  $t$  of the trajectory  $\tau_i$ , where  $z_i^t$  can be viewed as a noisy sampling of the underlying manipulation concept active at  $t$ . Since these representations cluster based on sub-goals, we refer to  $z_i^t$  as *manipulation concept latents* or simply *manipulation concepts*. We use continuous representations for differentiability and to avoid constraints of codebook-based discrete representations (e.g., finite capacity). To learn  $z_i^t$ , we introduce a manipulation concept encoder  $\mathcal{E}$  parameterized by  $\Theta_{\mathcal{E}}$ , which maps observation sequence  $\mathbf{o}_i = \{\mathbf{o}_i^t\}_{t=1}^{T_i}$  from trajectory  $\tau_i$  to concept sequence  $\mathbf{z}_i = \{z_i^t\}_{t=1}^{T_i}$ :

$$\mathbf{z}_i \leftarrow \mathcal{E}(\mathbf{o}_i; \Theta_{\mathcal{E}}) \quad (1)$$

We implement  $\mathcal{E}$  using a transformer to encode temporal dependencies (details in Sec. A.1). Next, we elaborate on the training strategies optimizing cross-modal and multi-horizon temporal correlation metrics (Sec. 3.2 and 3.3).

### 3.2 Capturing Multi-Modal Correlations

To enhance the utility of multi-modal information, we propose that manipulation concepts should capture *cross-modal correlations* rather than simply aggregating features from different modalities (e.g., concatenating multi-modal signals [9, 71]). Physiological evidence suggests that concept formation often occurs when correlations across sensory modalities are high [1, 15, 38, 57, 63]. These correlations remain consistent across scenarios involving the same concept, facilitating generalization. For instance, in container opening tasks, the correlated patterns between visual lid rotation, characteristic force feedback, and audio cues persist across different container types, enabling the transfer of the “opening” concept despite variations in object appearances.

To learn manipulation concepts that capture cross-modal correlations, we propose maximizing mutual information—a metric capable of modeling diverse correlations—between observations from different modalities, conditioned on the associated manipulation concept. Specifically, we maximize the conditional mutual information over bipartitions of modality observations:

$$\max_{\mathbf{Z}} \sum_{S \subsetneq [M], S \neq \emptyset} I(\mathbf{O}_S : \mathbf{O}_{[M] \setminus S} \mid \mathbf{Z}), \quad (2)$$

where  $\mathbf{O}_S$  are observations from a subset of modalities,  $\mathbf{O}_{[M] \setminus S}$  are observations from remaining modalities, and  $\mathbf{Z}$  is the manipulation concept. We implement Eq. 2 using a computationally efficient self-supervised *mask-and-predict* approach that stochastically samples bipartitions during training. This ensures scalability despite exponentially increasing bipartition numbers while integrating cross-modal correlation learning with multi-modal information compression.

Specifically, a Cross-Modal Correlation Network  $\mathcal{C}$  (CMCN) with parameters  $\Theta_c$  reconstructs full-modality observations from partial observations guided by manipulation concepts. During training, we mask observations from a random subset  $S$  of modalities and reconstruct all observations  $\mathbf{o}_i^t$  using the unmasked subset  $\mathbf{o}_i^{[M] \setminus S, t}$  and concept  $z_i^t$ :

$$\mathcal{L}_{mm}(t, \tau_i) = \mathbb{E}_S \left\| \mathcal{C} \left( \mathbf{o}_i^{[M] \setminus S, t}, z_i^t; \Theta_c \right) - \mathbf{o}_i^t \right\|, \quad (3)$$

where  $S \sim U(2^{[M]} \setminus \{\emptyset\})$  is a uniformly sampled non-empty subset of modality indices. By predicting full observations from partial inputs, we maximize the conditional mutual information in Eq. 2, forcing manipulation concepts  $z_i^t$  to capture cross-modal correlations. Additionally, when all modalities are masked, reconstruction solely from  $z_i^t$  ensures these representations compress and preserve essential multi-modal information (please see Sec. A.1 for more details).

### 3.3 Representing Multi-Horizon Sub-Goals

To complete tasks with hierarchical structures, manipulation concepts must encode multi-horizon sub-goal information. Physiological evidence shows human actions are hierarchically organized [17, 41], with coarse-grained goals defining overall tasks and fine-grained goals informing immediate actions. These multi-horizon sub-goals link ultimate goals with low-level actions, enabling smooth transitions while enhancing robustness.

We aim to make manipulation concepts organized to encode sub-processes across multiple temporal horizons without explicit annotations. Since concepts cluster by sub-goals, hierarchical sub-goals can emerge from these clusters at varying temporal scales. We propose that the temporal extent of a sub-process is determined by concept latent coherence within clusters, yielding a natural spectrum from short-horizon to long-horizon sub-goals. Specifically, given manipulation concept latents  $\mathbf{z}_i = \{z_i^t\}_{t=1}^{T_i}$  from trajectory  $\tau_i$ , we quantify their similarities using spherical distance:  $\text{dist}(z, u) = \frac{1}{\pi} \arccos \left\langle \frac{z}{\|z\|_2}, \frac{u}{\|u\|_2} \right\rangle$ . Concepts belong to the same sub-process if their distance falls below a coherence threshold  $\epsilon \in [0, 1]$ . More explicitly, sub-processes are derived as:

$$\begin{aligned} h(\mathbf{z}_i; \epsilon) &= \{[g_k, g_{k+1}] \mid k = 1, 2, \dots, K(\mathbf{z}_i; \epsilon)\}, \\ \text{where } g_1 &= 1, \quad g_{k+1} = \max_g \{g \mid g \in (g_k, T_i + 1] \cap \mathbb{N}^+ \wedge \forall t, t' \in [g_k, g), \text{dist}(z_i^t, z_i^{t'}) < \epsilon\}, \end{aligned} \quad (4)$$

where  $K(\mathbf{z}_i; \epsilon)$  is the number of clusters determined by  $\epsilon$ , and increasing  $\epsilon$  yields sub-processes spanning from short-horizon to long-horizon. Please see Alg. 1 for more details.

Furthermore, we propose learning objectives to ensure multi-horizon sub-processes from Eq. 4 align with meaningful sub-goal completion processes. Specifically, the manipulation concept guiding each sub-process should be informative about the state achieved upon sub-task completion [5, 33, 72]. For all coherence thresholds  $\epsilon$ , current observation  $\mathbf{O}$  and its associated concept  $\mathbf{Z}$  should be informative of the terminal observation  $\mathbf{O}^{\text{goal}(\epsilon)}$ , characterized by minimizing the following conditional entropy:

$$\forall \epsilon, \min_{\mathbf{Z}} \mathbb{H}(\mathbf{O}^{\text{goal}(\epsilon)} \mid \mathbf{O}, \mathbf{Z}), \quad (5)$$

To implement Eq. 5, we train a Multi-Horizon Future Predictor  $\mathcal{F}$  (MHFP) to hallucinate terminal observations of different sub-processes. For time step  $t$  in trajectory  $\tau_i$ , the terminal observation is determined by the ending time step of the interval containing  $t$ :

$$g(t; \mathbf{z}_i, \epsilon) = \min\{T_i, g_{k+1}\}, \text{ where } t \in [g_k, g_{k+1}) \in h(\mathbf{z}_i; \epsilon), \quad (6)$$

During training, the network  $\mathcal{F}$ , parameterized by  $\Theta_f$ , predicts this terminal observation based on current observation  $\mathbf{o}_i^t$ , manipulation concept  $z_i^t$ , and coherence threshold  $\epsilon$ :

$$\mathcal{L}_{\text{mh}}(t, \tau_i) = \mathbb{E}_\epsilon \left\| \mathcal{F}(\mathbf{o}_i^t, z_i^t, \epsilon; \Theta_f) - \mathbf{o}_i^{g(t; \mathbf{z}_i, \epsilon)} \right\|, \quad (7)$$

where  $\epsilon \sim U([0, 1])$  is sampled uniformly per iteration to improve efficiency by avoiding training over all  $\epsilon$  values. This training process iteratively improves both latents and sub-process derivation: we compute manipulation concepts using the encoder (Eq. 1), determine sub-process boundaries, then update all networks, including  $\mathcal{F}$  and the concept encoder. This improves future observation prediction and concept latents, which in turn refines sub-process derivation. By minimizing Eq. 7,  $z_i^t$  is ensured to encode multi-horizon sub-goal information, indicating hierarchical transitions to terminal states under various  $\epsilon$  while adjusting sub-processes by shaping concept latents for terminal state predictability. More details can be found in Sec. A.1.

**Final Objective for Manipulation Concept Discovery.** We jointly optimize the multi-modal correlation objective (Eq. 3) and multi-horizon sub-goal prediction objective (Eq. 7) to ensure manipulation concepts generated by the encoder  $\mathcal{E}$  (Eq. 1) satisfy both key properties:

$$\mathcal{L}_z(t, \tau_i) = \lambda_{\text{mm}} \mathcal{L}_{\text{mm}}(t, \tau_i) + \lambda_{\text{mh}} \mathcal{L}_{\text{mh}}(t, \tau_i), \quad (8)$$

where  $\lambda_{\text{mm}}, \lambda_{\text{mh}} > 0$  balance the two loss terms.

### 3.4 Enhancing Imitation Learning with Manipulation Concepts

After learning manipulation concepts through our self-supervised framework, we address how these concepts enhance policy learning. Unlike previous approaches that learn task-specific policies

directly from demonstrations [12, 28], we propose to leverage the learned manipulation concepts as an informative representation that bridges low-level actions and high-level goals.

Specifically, with manipulation concepts  $\mathbf{z}_i$  generated by encoder  $\mathcal{E}$ , we augment imitation learning by training policies to predict both ground-truth actions and corresponding concepts [21, 67, 70]. This approach uses concept prediction as a regularization that guides the policy to encode conceptual understanding alongside action planning:

$$\begin{aligned} h_i^t &= \pi_h(\mathbf{o}_i^t, \ell_i; \Theta_\pi^h), \quad \hat{z}_i^t = \pi_z(h_i^t; \Theta_\pi^z), \quad \hat{a}_i^t = \pi_a(h_i^t; \Theta_\pi^a), \\ \mathcal{L}_\pi(t, \tau_i, \ell_i) &= \|\hat{a}_i^t - a_i^t\| + \lambda_{mc} \|\hat{z}_i^t - z_i^t\|. \end{aligned} \quad (9)$$

The policy consists of: (1) A backbone  $\pi_h$  processing task descriptions  $\ell_i$  and observations  $\mathbf{o}_i^t$  to produce a shared representation  $h_i^t$ ; (2) A concept predictor  $\pi_z$  mapping  $h_i^t$  to predicted concepts  $\hat{z}_i^t$ ; and (3) An action decoder  $\pi_a$  mapping  $h_i^t$  to predicted actions  $\hat{a}_i^t$ . This joint objective enforces the policy to leverage concept information encoded within  $h_i^t$  while predicting actions. Even though concepts are learned task-agnostically for generalization, the policy receives task descriptions in a multi-task setting, serving as a mechanism to learn the reuse of concepts. The learning objective balances action and concept prediction using  $\lambda_{mc} > 0$ . More details are provided in Sec. A.2.

## 4 Experiments

We evaluate our manipulation concept discovery approach through experiments addressing four key questions: (1) Do learned concepts enhance policy performance on tasks used for concept discovery, validating our strategies for encoding cross-modal correlations (Sec. 3.2) and multi-horizon sub-goals (Sec. 3.3)? (2) Can concepts learned from one task set transfer effectively to different tasks sharing underlying manipulation patterns? (3) Does our concept discovery mechanism generalize to novel tasks with decreased overlap in manipulation patterns? (4) What interpretable properties emerge in the learned concepts that explain their effectiveness for robotic manipulation? Through these investigations, we demonstrate both the immediate benefits of our approach for imitation learning and its broader applicability for transfer learning and generalization in manipulation tasks.

### 4.1 Experimental Setup

**Dataset and Environment** Sec. 4.2 and 4.3 conduct experiments using the **LIBERO** benchmark [30], a comprehensive platform for robotic learning built on Robosuite [75]. We utilize three distinct task sets:

- **LIBERO-90**: A diverse collection of 90 manipulation tasks serving as our primary training domain for concept discovery and initial policy learning.
- **LIBERO-LONG**: 10 novel long-horizon tasks, each composed of two LIBERO-90 tasks in sequence, designed to evaluate transfer to more complex task structures.
- **LIBERO-GOAL**: 10 tasks in an entirely novel environment unseen during concept discovery, used to evaluate the generalization of learned concepts to unfamiliar contexts.

Each task includes a natural language description and 50 expert demonstrations. For multi-modal observations, we use: *Agentview vision*: 128×128 RGB third-person camera capturing the entire environment; *Eye-in-hand vision*: 128×128 RGB gripper-mounted camera; *Proprioceptive state*: 9D vector encoding gripper position, rotation, and physical states.

**Manipulation Concept Discovery Methods** We compare our approach with several state-of-the-art concept discovery baselines (implementation details in Sec. A.3):

- **InfoCon** [31]: A VQ-VAE type of method for single-hierarchy concept discovery.
- **XSkill** [65]: Contrastive learning for manipulation skill extraction from demonstration videos.
- **DecisionNCE** [27]: Learns reward-relevant representations from demonstrations with language annotations, evaluated in two variants: using task instructions (DecisionNCE-task) and using elementary action labels (DecisionNCE-motion).
- **RPT** [47]: Temporally and modality-masked autoencoder for multi-modal sequence modeling.
- **All**: A simplified variant of our approach that predicts all modalities from concepts without modeling cross-modal correlations.

Table 1: **Evaluation of manipulation concept discovery methods across different task settings.** Success rates (%) of ACT and Diffusion Policy (DP) models when enhanced with manipulation concepts from various discovery methods. All concept encoders were trained only on LIBERO-90, and evaluated on: original tasks (**L90-90**), novel long-horizon compositions (**L90-L**), and entirely new environments (**L90-G**). Values in parentheses show standard deviations across 4 seeds. **Bold** and underlined values indicate best and second-best results.

<b>L90-90</b>		InfoCon	XSkill	RPT	All	Next	CLIP	DINOv2	DecisionNCE task	NCE motion	Plain	<b>Ours</b>
ACT	66.5 (0.8)	<u>73.4</u> (0.8)	68.8 (0.8)	64.1 (2.0)	68.0 (0.4)	63.8 (0.5)	71.9 (0.3)	69.0 (0.1)	66.8 (0.8)	46.6 (1.9)	46.6 (0.8)	<b>74.8</b>
	DP	78.2 (0.6)	<u>87.7</u> (0.6)	84.3 (0.1)	81.5 (0.5)	82.6 (0.1)	80.7 (0.9)	79.4 (0.1)	75.7 (0.8)	82.7 (0.6)	75.1 (0.6)	<b>89.6</b>
<b>L90-L</b>		InfoCon	XSkill	RPT	All	Next	CLIP	DINOv2	DecisionNCE task	NCE motion	Plain	<b>Ours</b>
ACT	55.5 (0.9)	55.0 (1.0)	<u>59.0</u> (1.0)	55.5 (0.9)	55.0 (1.0)	51.0 (1.0)	55.0 (1.0)	53.0 (1.0)	49.3 (0.9)	54.0 (0.9)	54.0 (1.0)	<b>63.0</b>
	DP	75.0 (1.0)	73.0 (1.0)	61.3 (0.9)	79.3 (0.9)	<u>83.0</u> (1.0)	67.0 (1.0)	63.0 (1.0)	58.7 (0.9)	52.7 (0.9)	34.1 (1.1)	<b>89.0</b>
<b>L90-G</b>		InfoCon	XSkill	RPT	All	Next	CLIP	DINOv2	DecisionNCE task	NCE motion	Plain	<b>Ours</b>
ACT	67.0 (1.0)	77.0 (1.0)	75.0 (1.0)	69.0 (1.0)	71.0 (1.0)	77.0 (1.0)	<u>77.3</u> (0.9)	70.0 (0.9)	75.0 (0.5)	57.0 (1.0)	57.0 (1.0)	<b>81.0</b>
	DP	92.7 (0.9)	<u>93.0</u> (1.0)	91.5 (0.9)	91.0 (1.0)	91.3 (0.9)	92.0 (0.9)	91.0 (0.7)	92.0 (0.8)	93.0 (1.0)	90.7 (0.9)	<b>95.7</b>

- **Next:** Predicts adjacent time-step observations, a common approach adopted in [7, 68].
- **CLIP** [46]: Language-aligned visual features from a pretrained foundation model.
- **DINOv2** [43]: Self-supervised visual representations without temporal modeling.
- **Plain:** Standard imitation learning without manipulation concepts.

**Policies for Concept-Enhanced Imitation Learning** To evaluate the effectiveness of our discovered manipulation concepts, we integrate them into two established imitation learning frameworks using the joint prediction approach described in Sec. 3.4:

- **ACT** [71]: A transformer-based conditional variational autoencoder that predicts action chunks.
- **Diffusion Policy (DP)** [19]: A 1D convolutional UNet that generates actions through denoising.

For both policy architectures, we add the concept prediction head ( $\pi_z$  in Eq. 9) to predict manipulation concepts from the shared concept-aware representations. Implementation details appear in Sec. A.2. All experiments are reported with results aggregated across 4 random seeds.

## 4.2 Evaluating Policy Performance with Learned Manipulation Concepts

- **Performance on Original Training Tasks** We first evaluate our concept discovery method on the same tasks used for concept training. As shown in the **L90-90** results (Tab. 1), our approach consistently outperforms all baselines with both policy architectures. The performance gap between our method and *Next/InfoCon* demonstrates the importance of multi-hierarchical sub-goal modeling, while improvements over *All* highlight the value of explicitly capturing cross-modal correlations. Our method also surpasses *DecisionNCE* variants despite not requiring language supervision, validating the effectiveness of our self-supervised objectives.
- **Transfer to Long-Horizon Tasks** To evaluate concept transferability to more complex compositions, we apply concept encoders trained on LIBERO-90 directly to LIBERO-LONG demonstrations featuring novel long-horizon tasks. The **L90-L** results show our method maintains its performance advantage in this challenging transfer setting. This demonstrates that our approach learns manipulation concepts that effectively decompose hierarchical tasks, enabling policies to better handle novel complex task compositions requiring sequential execution of multiple sub-goals.

Table 2: **Impact of modality combinations on concept discovery performance.** Success rates (%) of ACT and DP policies using manipulation concepts discovered with different input modality combinations. All models were trained and evaluated on LIBERO-90, with specific modalities excluded (marked with “–”). A: agentview vision, H: eye-in-hand vision, P: proprioceptive state.

	<b>Ours</b>	– H P	A – P	A H –	– – P	– H –	A – –
ACT	74.8±0.8	70.5±1.8	71.3±0.3	70.1±1.2	67.5±0.8	68.7±0.6	69.4±0.4
DP	89.6±0.6	85.8±0.2	85.6±0.3	84.3±0.5	84.8±0.1	83.7±0.1	85.3±0.5

- **Generalization to Novel Environments** We further test generalization by applying concept encoders trained on LIBERO-90 *directly* to LIBERO-GOAL demonstrations featuring unseen environments and tasks. The *L90-G* results show our method continues to outperform all baselines in this challenging scenario. This indicates our approach discovers fundamental manipulation primitives that transfer effectively across environmental variations.
- **Impact of Multi-Modal Observations** Our ablation study (Tab. 2) shows that performance consistently improves as more modalities are incorporated. The most significant drops occur when removing proprioceptive information, highlighting its importance for grounding visual observations with physical interaction states and confirming the value of our cross-modal correlation approach.

### 4.3 Analyzing Manipulation Concept Properties

**Enhanced Cross-Modal Correlation** To verify our Cross-Modal Correlation Network’s effectiveness (Sec. 3.2), we measure mutual information between modalities conditioned on concept latents (Sec. A.4). Tab. 3 shows that our approach achieves higher conditional mutual information than the **All** baseline. This confirms that our mask-and-predict strategy enables the concept encoder to capture persistent cross-modal patterns that generalize across different objects and contexts, providing a robust representational basis for policies.

**Alignment with Semantic Sub-Goals** We evaluate whether our concepts align with human-understandable semantics by grouping latents from different demonstrations based on human-identified sub-goals and computing similarities between these groupings:

$$\langle C_i, C_j \rangle = \frac{1}{|C_i||C_j|} \sum_{z_i \in C_i} \sum_{z_j \in C_j} \left\langle \frac{z_i}{\|z_i\|_2}, \frac{z_j}{\|z_j\|_2} \right\rangle, \quad (10)$$

where  $C_i$ ,  $C_j$  represent human-identified sub-goal categories, and  $z_i$ ,  $z_j$  are latents within each category (details in Sec. C.2). As shown in Fig. 4, similarity matrices consistently show the highest values along the diagonal, demonstrating that our approach discovers concepts that exhibit clustering patterns corresponding to meaningful manipulation primitives.

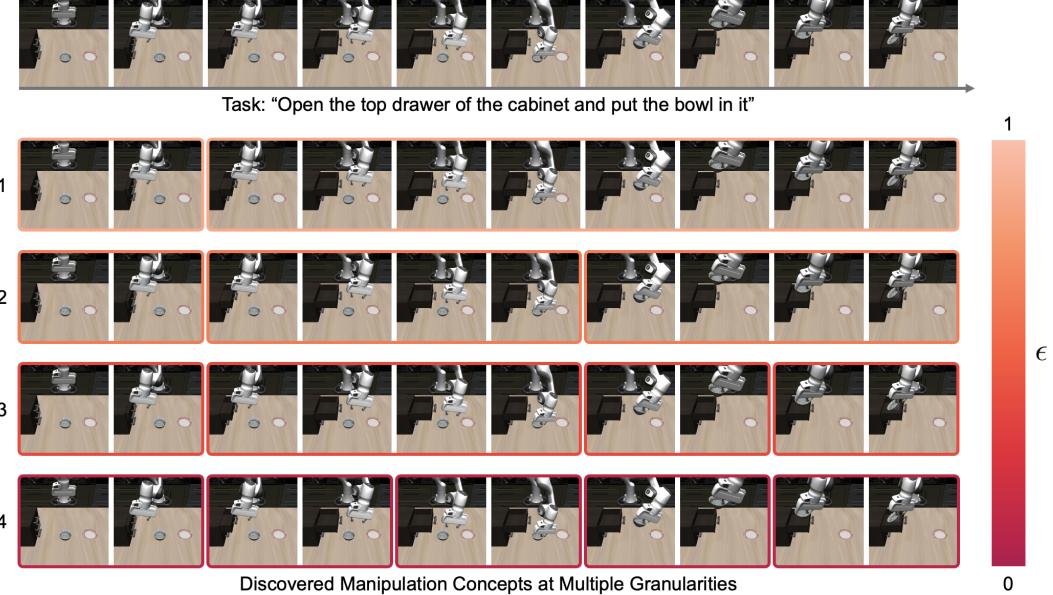
**Multi-Level Hierarchical Structure** Varying the coherence threshold  $\epsilon$  in Eq. 4 reveals the hierarchical organization of our learned concepts. Fig. 3 (and Sec. C.5) shows larger  $\epsilon$  values identify coarse-grained phases, while smaller values capture fine-grained actions. This emergent hierarchy enables policies to simultaneously reason about immediate actions and longer-term goals without explicit hierarchical supervision, contributing to improved performance on complex sequential tasks that require coordinated execution across multiple temporal scales.

### 4.4 Real-World Validation

**Real-World Generalization Study** To study generalization capabilities, we deploy concept-enhanced policies on a Mobile ALOHA robot [16] in “cleaning cup” tasks (Fig. 5). Training data includes only simple container arrangements with consistent color pairings.

Table 3: **Conditional mutual information between modality pairs.** Values conditioned on concept latents from our method versus the **All** baseline that does *not* model cross-modal correlations. A: agentview, H: eye-in-hand vision, P: proprioception.

	<b>Ours</b>	<b>All</b>
$I(O_H : o_A   z)$	3.7999	2.0080
$I(o_P : o_A   z)$	4.8319	3.1312
$I(o_P : o_H   z)$	4.8255	3.1322



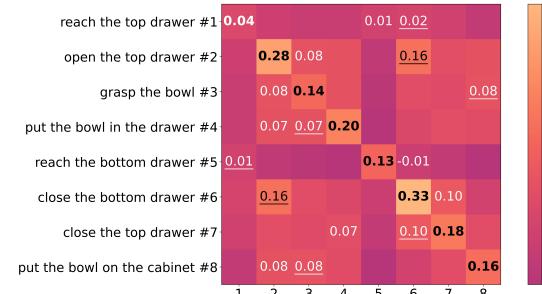
**Figure 3: Multi-granular task decomposition through concept latent clustering.** Visualization of sub-processes derived by clustering manipulation concept latents at different coherence thresholds ( $\epsilon$ ) for the task “open the top drawer and put the bowl in it.” Higher  $\epsilon$  values (top rows) produce coarser decompositions, while lower values (bottom rows) yield finer-grained segmentation. The emergent sub-processes naturally align with semantic task components, for example, the third segment in row 2 corresponds to “put bowl in drawer,” while the second segment in row 4 corresponds to “pull drawer open.” This demonstrates our method’s ability to discover hierarchical, human-interpretable task structures without explicit supervision.

We test on six increasingly challenging variations: **(1) Novel Placements:** Cups and containers in unseen arrangements; **(2) Color Composition:** Altered cup-container color pairings; **(3) Novel Objects:** Entirely unseen containers, cups, and plates; **(4) Obstacles:** Objects between the robot and the cups obstructing vision; **(5) Barriers:** Internal dividers within containers impeding placement; **(6) Grasping Together:** Two adjacent cups requiring simultaneous grasp.

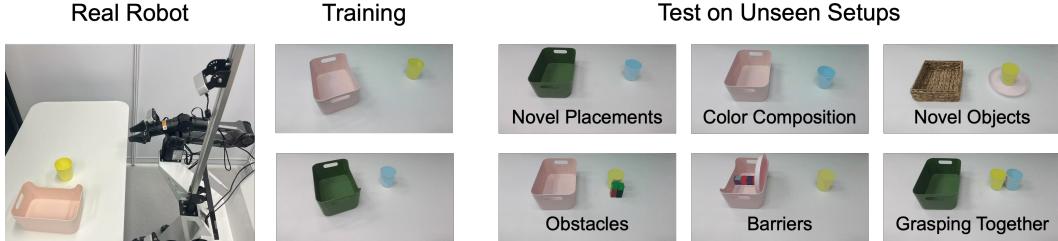
As shown in Tab. 4, policies enhanced with our manipulation concepts consistently outperform baselines across all scenarios, with advantages in challenging conditions. We suggest that the following two mechanisms behind learned manipulation concepts improve generalization:

**1. Relational focus:** Concept-enhanced policies prioritize transferable relational patterns (e.g., “object inside container”) over surface features. Our cross-modal correlation learning (Sec. 3.2) enables this capability by identifying patterns that remain invariant across modalities. This relational emphasis explains the stronger performance on scenarios that alter visual appearance while preserving task structure. For instance, while **Novel Placement** tests spatial variation alone, the other protocols introduce substantial visual perturbations (different colors, objects, or occlusions) that shift the appearance distribution. The consistent performance gains across these visually diverse scenarios (Tab. 4) suggest that the learned concepts successfully capture the underlying relational invariant—placing cups into containers—rather than memorizing superficial visual patterns.

**2. Hierarchical awareness:** Concept-enhanced policies exhibit more systematic failure recovery than baselines, suggesting better tracking of sub-goal completion. Baseline failures frequently exhibit



**Figure 4: Semantic alignment of learned concepts.** Cosine similarity between concept latents grouped by human-defined sub-goals. Diagonal patterns demonstrate that our approach discovers concepts that exhibit clustering patterns corresponding to meaningful manipulation primitives.



**Figure 5: Real-world generalization evaluation with Mobile ALOHA robot.** Left: Mobile ALOHA robot setup for cup cleaning tasks. Center: Training conditions with simple, consistent cup-container color pairings. Right: Six test variations with increasing complexity: novel placements, altered color combinations, unfamiliar objects, external obstacles, internal barriers, and simultaneous grasping of multiple cups. These variations test the policy’s ability to generalize beyond training conditions by systematically introducing new challenges.

premature task abandonment: the robot moves toward containers without having grasped objects, or hovers near placement locations without executing placement. In contrast, when concept-enhanced policies fail initial grasp attempts, they consistently retry grasping (typically 2-3 attempts) before proceeding, demonstrating recognition of incomplete sub-goals. Although these recoveries ultimately fail due to time limits or object displacement, they reveal structured task progression rather than blind action execution.

These mechanisms may enable manipulation concepts to promote policy generalization by encoding fundamental spatial and functional relationships that remain consistent across environmental variations. Details are provided in Sec. C.6.

**Multi-Horizon Goal Prediction Visualization** To visualize the temporal information encoded in our manipulation concepts, we examine outputs from our Multi-Horizon Goal Predictor (MHGP,  $\mathcal{F}$  in Eq. 7) using the BridgeDataV2 dataset [60].

Fig. 7 (Sec. C.7) shows predicted goal states when conditioned on the current observation, manipulation concept, and various coherence thresholds ( $\epsilon$ ).

The predictions capture essential task structures – such as anticipated arm trajectories and object interactions – rather than attempting pixel-perfect reconstructions. This abstraction of scene-specific details in favor of functional relationships is crucial for cross-environment generalization. Importantly, as  $\epsilon$  increases, the predictions correspond to states progressively further into the future, with smaller values showing immediate next steps and larger values revealing final goal states. This demonstrates that our learned concepts encode meaningful temporal structures at multiple time horizons, enabling policies to simultaneously reason about immediate actions and longer-term objectives. Details are provided in Sec. C.7.

## 5 Discussion

We demonstrate that self-supervised discovery of hierarchical manipulation concepts significantly enhances robot policy performance across original tasks, novel compositions, and entirely new environments. Three key strengths emerge: (1) our representations naturally resemble semantically meaningful manipulation primitives without requiring explicit labels, as evidenced by diagonal clustering in similarity matrices; (2) the concepts bridge low-level actions and high-level goals through hierarchical organization, enabling reasoning at multiple temporal scales; and (3) concept-enhanced policies focus on transferable relational patterns rather than superficial features, explaining their robust generalization to scenarios with substantial distribution shifts. These findings highlight the potential of learning manipulation concepts from unlabeled multi-modal demonstrations for creating more adaptable and interpretable robotic systems. Limitations are discussed in Sec. D.

**Table 4: Real-world generalization success rates (%)** for ACT policies with and without manipulation concepts (MC). Test conditions: Placements (novel layouts), Color (new pairings), Objects (unseen items), Obstacles (external barriers), Barrier (internal dividers), and Multi-grasp (two cups simultaneously).

	Place	Color	Obj.	Obst.	Barf.	Multi
w/o MC	53.3	46.7	40.0	20.0	0.0	0.0
w/ MC	<b>73.3</b>	<b>60.0</b>	<b>53.3</b>	<b>33.3</b>	<b>20.0</b>	<b>13.3</b>

## Acknowledgments and Disclosure of Funding

This work is supported by the Early Career Scheme of the Research Grants Council (RGC) grant # 27207224, the HKU-100 Award, a donation from the Musketeers Foundation, in part by the JC STEM Lab of Robotics for Soft Materials funded by The Hong Kong Jockey Club Charities Trust, and DAMO Academy through the Alibaba Innovative Research Program.

## References

- [1] Nikolai Axmacher, Florian Mormann, Guillen Fernández, Christian E Elger, and Juergen Fell. Memory formation by neuronal synchronization. *Brain research reviews*, 52(1):170–182, 2006.
- [2] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference on machine learning*, pages 531–540. PMLR, 2018.
- [3] Suneel Belkhale, Tianli Ding, Ted Xiao, Pierre Sermanet, Quon Vuong, Jonathan Tompson, Yevgen Chebotar, Debidatta Dwibedi, and Dorsa Sadigh. Rt-h: Action hierarchies using language. *arXiv preprint arXiv:2403.01823*, 2024.
- [4] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- [5] Kevin Black, Mitsuhiro Nakamoto, Pranav Atreya, Homer Rich Walke, Chelsea Finn, Aviral Kumar, and Sergey Levine. Zero-shot robotic manipulation with pre-trained image-editing diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [6] Rogerio Bonatti, Sai Vemprala, Shuang Ma, Felipe Frujeri, Shuhang Chen, and Ashish Kapoor. Pact: Perception-action causal transformer for autoregressive robotics pre-training. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3621–3627. IEEE, 2023.
- [7] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- [8] Qingwen Bu, Hongyang Li, Li Chen, Jisong Cai, Jia Zeng, Heming Cui, Maoqing Yao, and Yu Qiao. Towards synergistic, generalized, and efficient dual-system for robotic manipulation. *arXiv preprint arXiv:2410.08001*, 2024.
- [9] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.
- [10] Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram Madukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyun Kim,

Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open x-embodiment: Robotic learning datasets and rt-x models, 2024.

- [11] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at datasets for visuo-motor pre-training. In Conference on Robot Learning, pages 1183–1198. PMLR, 2023.
- [12] Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with large language models for object rearrangement. In 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 2086–2092. IEEE, 2023.
- [13] Ben Eisner, Harry Zhang, and David Held. Flowbot3d: Learning 3d articulation flow to manipulate articulated objects. arXiv preprint arXiv:2205.04382, 2022.
- [14] Haoquan Fang, Markus Grotz, Wilbert Pumacay, Yi Ru Wang, Dieter Fox, Ranjay Krishna, and Jiafei Duan. Sam2act: Integrating visual foundation model with a memory architecture for robotic manipulation. In 7th Robot Learning Workshop: Towards Robots with Human-Level Abilities.
- [15] Juergen Fell and Nikolai Axmacher. The role of phase synchronization in memory processes. Nature reviews neuroscience, 12(2):105–118, 2011.
- [16] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. arXiv preprint arXiv:2401.02117, 2024.
- [17] Scott T Grafton and Antonia F de C Hamilton. Evidence for a distributed hierarchy of action representation in the brain. Human movement science, 26(4):590–616, 2007.
- [18] Yanjiang Guo, Yucheng Hu, Jianke Zhang, Yen-Jen Wang, Xiaoyu Chen, Chaochao Lu, and Jianyu Chen. Prediction with action: Visual policy learning via joint denoising process. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.

- [19] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- [20] Zhiwei Jia, Vineet Thumuluri, Fangchen Liu, Linghao Chen, Zhiao Huang, and Hao Su. Chain-of-thought predictive control. In *Forty-first International Conference on Machine Learning*, 2024.
- [21] Zhiwei Jia, Vineet Thumuluri, Fangchen Liu, Linghao Chen, Zhiao Huang, and Hao Su. Chain-of-thought predictive control. In *International Conference on Machine Learning*, pages 21768–21790. PMLR, 2024.
- [22] Siddharth Karamcheti, Suraj Nair, Annie S Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang. Language-driven representation learning for robotics. *arXiv preprint arXiv:2302.12766*, 2023.
- [23] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025.
- [24] 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.
- [25] Thomas Kipf, Yujia Li, Hanjun Dai, Vinicius Zambaldi, Alvaro Sanchez-Gonzalez, Edward Grefenstette, Pushmeet Kohli, and Peter Battaglia. Compile: Compositional imitation learning and execution. In *International Conference on Machine Learning*, pages 3418–3428. PMLR, 2019.
- [26] Seungjae Lee, Yibin Wang, Haritheja Etukuru, H Jin Kim, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Behavior generation with latent actions. *arXiv preprint arXiv:2403.03181*, 2024.
- [27] Jianxiong Li, Jinliang Zheng, Yinan Zheng, Liyuan Mao, Xiao Hu, Sijie Cheng, Haoyi Niu, Jihao Liu, Yu Liu, Jingjing Liu, et al. Decisionncc: Embodied multimodal representations via implicit preference learning. In *International Conference on Machine Learning*, pages 29461–29488. PMLR, 2024.
- [28] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. *arXiv preprint arXiv:2312.16217*, 2023.
- [29] Zhixuan Liang, Yao Mu, Hengbo Ma, Masayoshi Tomizuka, Mingyu Ding, and Ping Luo. Skilldiffuser: Interpretable hierarchical planning via skill abstractions in diffusion-based task execution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16467–16476, 2024.
- [30] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 44776–44791. Curran Associates, Inc., 2023.
- [31] Ruizhe Liu, Qian Luo, and Yanchao Yang. Infocon: Concept discovery with generative and discriminative informativeness. In *The Twelfth International Conference on Learning Representations*, 2024.
- [32] Yuyao Liu, Jiayuan Mao, Joshua B Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. One-shot manipulation strategy learning by making contact analogies. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 15387–15393. IEEE, 2025.
- [33] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In *Conference on robot learning*, pages 1113–1132. PMLR, 2020.

- [34] Yecheng Jason Ma, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. In *International Conference on Machine Learning*, pages 23301–23320. PMLR, 2023.
- [35] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. *arXiv preprint arXiv:2210.00030*, 2022.
- [36] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? *Advances in Neural Information Processing Systems*, 36:655–677, 2023.
- [37] Jiayuan Mao, Tomás Lozano-Pérez, Joshua B Tenenbaum, and Leslie Pack Kaelbling. Learning reusable manipulation strategies. In *Conference on Robot Learning*, pages 1467–1483. PMLR, 2023.
- [38] Lucia Melloni, Carlos Molina, Marcela Pena, David Torres, Wolf Singer, and Eugenio Rodriguez. Synchronization of neural activity across cortical areas correlates with conscious perception. *Journal of neuroscience*, 27(11):2858–2865, 2007.
- [39] Atharva Mete, Haotian Xue, Albert Wilcox, Yongxin Chen, and Animesh Garg. Quest: Self-supervised skill abstractions for learning continuous control. *arXiv preprint arXiv:2407.15840*, 2024.
- [40] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 909–918, 2019.
- [41] John D Murray, Alberto Bernacchia, David J Freedman, Ranulfo Romo, Jonathan D Wallis, Xinying Cai, Camillo Padoa-Schioppa, Tatiana Pasternak, Hyojung Seo, Daeyeol Lee, et al. A hierarchy of intrinsic timescales across primate cortex. *Nature neuroscience*, 17(12):1661–1663, 2014.
- [42] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- [43] Maxime Oquab, Timothée Darzet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [44] Karl Pertsch, Oleh Rybkin, Jingyun Yang, Shenghao Zhou, Konstantinos Derpanis, Kostas Daniilidis, Joseph Lim, and Andrew Jaegle. Keyframing the future: Keyframe discovery for visual prediction and planning. In *Learning for Dynamics and Control*, pages 969–979. PMLR, 2020.
- [45] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [47] Ilija Radosavovic, Baifeng Shi, Letian Fu, Ken Goldberg, Trevor Darrell, and Jitendra Malik. Robot learning with sensorimotor pre-training. In *Conference on Robot Learning*, pages 683–693. PMLR, 2023.
- [48] Seungeun Rho, Laura Smith, Tianyu Li, Sergey Levine, Xue Bin Peng, and Sehoon Ha. Language guided skill discovery. In *The Thirteenth International Conference on Learning Representations*.

- [49] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [50] Dominik Schmidt and Minqi Jiang. Learning to act without actions. In *The Twelfth International Conference on Learning Representations*, 2024.
- [51] Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter Abbeel. Masked world models for visual control. In *Conference on Robot Learning*, pages 1332–1344. PMLR, 2023.
- [52] Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning  $k$  modes with one stone. *Advances in neural information processing systems*, 35:22955–22968, 2022.
- [53] Rutav Shah, Roberto Martín-Martín, and Yuke Zhu. Mutex: Learning unified policies from multimodal task specifications. *arXiv preprint arXiv:2309.14320*, 2023.
- [54] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020.
- [55] Shashank Sharma, Vinay Namboodiri, and Janina Hoffmann. Multi-resolution skill discovery for hierarchical reinforcement learning. In *NeurIPS 2023 Workshop on Goal-Conditioned Reinforcement Learning*, 2023.
- [56] Lucy Xiaoyang Shi, Archit Sharma, Tony Z Zhao, and Chelsea Finn. Waypoint-based imitation learning for robotic manipulation. In *Conference on Robot Learning*, pages 2195–2209. PMLR, 2023.
- [57] Wolf Singer. Consciousness and neuronal synchronization. *The neurology of consciousness*, pages 43–52, 2011.
- [58] Yang Tian, Sizhe Yang, Jia Zeng, Ping Wang, Dahua Lin, Hao Dong, and Jiangmiao Pang. Predictive inverse dynamics models are scalable learners for robotic manipulation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [59] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [60] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023.
- [61] Weikang Wan, Yifeng Zhu, Rutav Shah, and Yuke Zhu. Lotus: Continual imitation learning for robot manipulation through unsupervised skill discovery. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 537–544. IEEE, 2024.
- [62] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [63] Thilo Womelsdorf and Pascal Fries. The role of neuronal synchronization in selective attention. *Current opinion in neurobiology*, 17(2):154–160, 2007.
- [64] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li, Minghuan Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training for visual robot manipulation. *arXiv preprint arXiv:2312.13139*, 2023.
- [65] Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and Shuran Song. Xskill: Cross embodiment skill discovery. In *Conference on Robot Learning*, pages 3536–3555. PMLR, 2023.
- [66] Jonathan Yang, Dorsa Sadigh, and Chelsea Finn. Polybot: Training one policy across robots while embracing variability. *arXiv preprint arXiv:2307.03719*, 2023.

- [67] Mengjiao Sherry Yang, Dale Schuurmans, Pieter Abbeel, and Ofir Nachum. Chain of thought imitation with procedure cloning. *Advances in Neural Information Processing Systems*, 35:36366–36381, 2022.
- [68] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo, Jianwei Yang, Baolin Peng, Ajay Mandlekar, Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretraining from videos. *arXiv preprint arXiv:2410.11758*, 2024.
- [69] Jia Zeng, Qingwen Bu, Bangjun Wang, Wenke Xia, Li Chen, Hao Dong, Haoming Song, Dong Wang, Di Hu, Ping Luo, et al. Learning manipulation by predicting interaction. *arXiv preprint arXiv:2406.00439*, 2024.
- [70] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. *arXiv preprint arXiv:2503.22020*, 2025.
- [71] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023.
- [72] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A 3d vision-language-action generative world model. In *Forty-first International Conference on Machine Learning*, 2024.
- [73] Pei Zhou, Ruizhe Liu, Qian Luo, Fan Wang, Yibing Song, and Yanchao Yang. Autocgp: Closed-loop concept-guided policies from unlabeled demonstrations. In *The Thirteenth International Conference on Learning Representations*.
- [74] Yifeng Zhu, Peter Stone, and Yuke Zhu. Bottom-up skill discovery from unsegmented demonstrations for long-horizon robot manipulation. *IEEE Robotics and Automation Letters*, 7(2):4126–4133, 2022.
- [75] Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, Yifeng Zhu, and Kevin Lin. robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv preprint arXiv:2009.12293*, 2020.

## A Implementation details

### A.1 Manipulation concept discovery (Ours)

This section details the neural network architectures and training procedures employed in our manipulation concepts discovery framework (Sec. 3) as implemented on the LIBERO benchmark.

**Manipulation Concept Encoder (Sec. 3.1)** The manipulation concepts encoder  $\mathcal{E}$  (Eq. 1) first encodes the multi-modal observations at each time step of the input demonstration into an encoded vector. It then utilizes a self-attention transformer to process the sequence of encoded vectors into a sequence of manipulation concepts. For the observation encoding process, our experiments on LIBERO incorporate two vision observations: **agent-view vision** and **eye-in-hand vision**. The original images are tensors of shape  $128 \times 128 \times 3$ . To enhance processing efficiency, we preprocess the images for each time step using the VAE encoder from stable diffusion [49], compressing each image into a tensor of shape  $16 \times 16 \times 4$ , which is then flattened into a 1024-dimensional vector. In addition to the two vision observations, we include a 9-dimensional robot state at each time step of each demonstration as the proprioceptive state observation. For these three observations at each time step, we employ three distinct 2-layer MLPs to process each observation into a feature vector of the hidden size (256) used by the subsequent transformer. The encoded features from these observations are then summed to form a 256-dimensional representation that encapsulates the sensing information from the three modalities.

$$\begin{aligned} h_{\text{av}} &= \text{MLP}_{\text{av}}(I_{\text{av-compress}}) & h_{\text{ev}} &= \text{MLP}_{\text{ev}}(I_{\text{ev-compress}}) & h_{\text{prop}} &= \text{MLP}_{\text{prop}}(s_{\text{prop}}) \\ h &= h_{\text{av}} + h_{\text{ev}} + h_{\text{prop}} \end{aligned} \quad (11)$$

Here,  $I_{\text{av-compress}}$  represents the 1024-dimensional compressed **agent-view vision**, and  $I_{\text{ev-compress}}$  represents the 1024-dimensional compressed **eye-in-hand vision**.  $s_{\text{prop}}$  denotes the 9-dimensional proprioceptive state observation. The output of the hidden layers from the three MLPs is 1024 dimensions. The  $h$  in Eq. 11 represents the encoded observation feature at each time step of a given demonstration  $\tau_i$ :  $(h_i^1, h_i^2, \dots, h_i^{T_i})$ . The next module in  $\mathcal{E}$  is a 12-layer self-attention (MHA in Eq. 12) transformer, enabling each time step to aggregate information from every other time step in the input sequence. In our implementation, we do not input the entire demonstration; instead, the transformer processes a fixed input sequence length of  $T_{\text{context}} = 60$ . A learnable temporal embedding, represented as a tensor of shape  $60 \times 256$ , is added to the input sequence to enhance temporal representation. The hidden feature dimension at each time step is 256, and each self-attention layer consists of 8 heads. Moreover, since spherical distance is utilized in Sec. 3.3, the output manipulation concepts are normalized to have a unit length with respect to the 2-norm:

$$(z_i^t, z_i^{t+1}, \dots, z_i^{t+T_{\text{context}}-1}) \leftarrow \text{Norm}_2 \left( [\text{MHA}]_{\times 12} \left( h_i^t, h_i^{t+1}, \dots, h_i^{t+T_{\text{context}}-1} \right) \right) \quad (12)$$

The output manipulation concept sequence in Eq. 1 represents the predicted manipulation concepts at time-steps  $t, t+1, \dots, (t+T_{\text{context}}-1)$  of the demonstration  $\tau_i$ . During training, demonstrations with lengths shorter than  $T_{\text{context}}$  are padded to  $T_{\text{context}}$  by repeating the observations from the last time-step at the end of each demonstration. During inference, when  $\mathcal{E}$  is used to label the demonstrations in the original dataset, the manipulation concepts at each time step are designed to incorporate information from as many future time steps as possible. This approach aims to better capture motion pattern dynamics, aligning with prior works that generate the dynamics at the current time step based on information derived from the dynamics spanning the current to future time steps [68]. Specifically:

- For each time-step  $t \leq T_i - T_{\text{context}}$ , the corresponding manipulation concepts are derived when the input to Eq. 12 starts from this time-step and spans a length of  $T_{\text{context}}$ :  $(h_i^t, h_i^{t+1}, \dots, h_i^{t+T_{\text{context}}-1})$ .
- For each time step  $t > T_i - T_{\text{context}}$ , the corresponding manipulation concepts are derived when the input to Eq. 12 begins at time step  $h_i^{T_i-T_{\text{context}}+1}$  and spans a length of  $T_{\text{context}}$ , ensuring that the final time step corresponds to the end of the demonstration:  $(h_i^{T_i-T_{\text{context}}+1}, h_i^{T_i-T_{\text{context}}+2}, \dots, h_i^{T_i})$ .
- If the original demonstration length is smaller than  $T_{\text{context}}$ , the manipulation concepts correspond to the input appended with repeated observations as described earlier.

However, we do not firmly believe this is the optimal approach for labeling manipulation concepts. Further exploration of inference-time strategy design is left for future work, as it is not a core focus of the manipulation concept discovery methodology presented.

**Learning Multi-Modal Features and Correlations (Sec. 3.2)** The Cross-Modal Correlation Network  $\mathcal{C}$  (Eq. 3) shares a similar structure with  $\mathcal{E}$  (Eq. 1). First, it includes four 2-layer MLP encoders, analogous to the three encoders in Eq. 11, with an additional encoder for processing the manipulation concepts. Each of these four MLPs outputs a hidden feature of dimension 1024, which is then reduced to a 256-dimensional encoded feature. These encoded features are summed to represent the combined information from the three observations and the manipulation concepts. Second, it incorporates a 4-layer self-attention transformer to process the sequence of features (with the same fixed length  $T_{\text{context}} = 60$ ) produced by the four MLPs. Following this, three 3-layer MLP decoders map the transformer’s output to the reconstructed observations at each time step. Unlike in Eq. 12, the transformer’s output does not require normalization. Each decoder MLP has hidden layers with a dimension of 1024. As described in Eq. 3, for the three observations—**agent-view camera vision**, **eye-in-hand camera vision**, and **proprioceptive state observation**—we randomly mask these modalities, ensuring that at least one modality is masked during each iteration. The  $2^3 - 1 = 7$  possible masking scenarios follow a uniform distribution, with each scenario appearing with a probability of  $\frac{1}{7}$ . For the sampled masks, all observations of the corresponding masked modalities in the input sequence are replaced with zero tensors. The loss is applied separately to the reconstruction of the three different observations. Specifically, L2 loss is applied to the two vision observations, while L1 loss is applied to the proprioceptive state observations. The loss weight  $\lambda_{\text{mm}}$  in Eq. 8 is set to 1.0.

**Learning Multi-Hierarchical Sub-goals (Sec. 3.3)** The Multi-Horizon Future Predictor  $\mathcal{F}$  (Eq. 7) shares a similar structure with  $\mathcal{C}$  (Eq. 3). The key differences are:

- $\mathcal{F}$  does not require a masking strategy.
- The transformer in  $\mathcal{F}$  is a 4-layer causal self-attention transformer. Causal attention is used because, in Eq. 7, the prediction is made from each current time step to certain future time steps. Therefore, for each time-step input in  $\mathcal{F}$ , access to information from subsequent time steps is restricted.
- To incorporate the granularity parameter  $\epsilon \in [0, 1]$ , we discretize the continuous range into 1000 uniform bins  $\{0.000, 0.001, \dots, 0.999\}$  and learn a corresponding VQ-VAE codebook [59] with 1000 entries, each represented as a 256-dimensional embedding vector. In each transformer block, the feed-forward layer receives the concatenation of the attention output and the embedding corresponding to the sampled  $\epsilon$  value.
- The output predictions correspond to the observations at the time steps determined by the rules described in Sec. 3.3 (Eqs. 4 and 6). Still, the loss is applied separately to the reconstruction of the three types of observations. Specifically, L2 loss is used for the two vision observations, while L1 loss is applied to the proprioceptive state observations. The loss weight  $\lambda_{\text{mh}}$  in Eq. 8 is set to 1.0.

**Training Details** We train the manipulation concept discovery process for 200,000 iterations with a batch size of 512. Each item in the batch is a segment of demonstration with a fixed length of  $T_{\text{context}} = 60$ . The training process uses the AdamW optimizer with a weight decay of 0.001 and momentum parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ . The base learning rate is set to 0.001. Initially, the model is trained with a 100-iteration warmup phase, during which the learning rate increases linearly from 0.0001 to 0.001. After the warmup, the model is trained for the remaining iterations using a cosine decay schedule, gradually reducing the learning rate back to 0.0001. This training setup is compatible with GPUs such as the GeForce RTX 3090 or 4090. However, we leverage the A800 GPU for improved efficiency, completing the training process in 1.5 days.

## A.2 Enhancing Imitation Learning

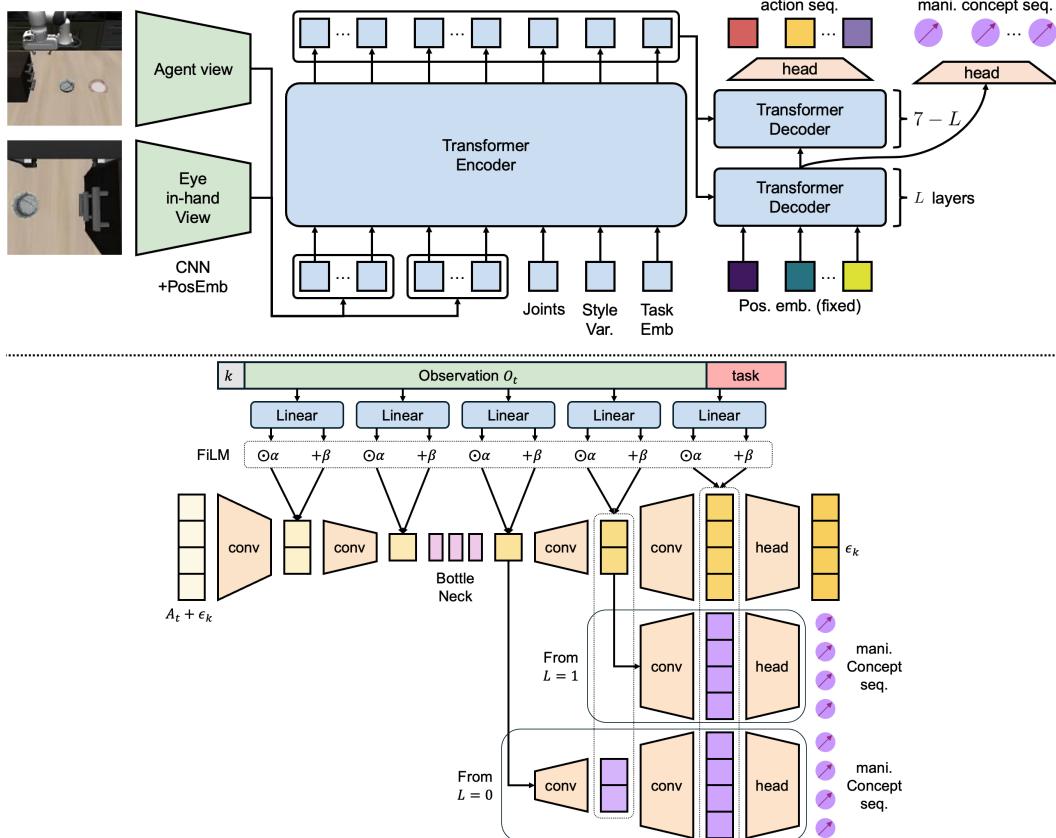


Figure 6: Upper: Enhanced ACT (decoder part); lower: Enhanced Diffusion Policy

This section introduces the neural network architectures and training details used in the Enhancing Imitation Learning process (Sec. 3.4). It focuses on modifying the original neural network policy to enable the prediction of manipulation concepts, thereby enhancing performance. Moreover, the implementation of the base policy follow [39].

**ACT** The pipeline (we focus on the CVAE decoder as it is the only modified component) is shown in the upper part of Fig. 6. Following [39], the transformer encoder in ACT’s CVAE decoder is modified to incorporate task embeddings provided by CLIP. The transformer decoder in ACT’s CVAE decoder is adapted to predict manipulation concepts. Specifically, the output of the  $L$ -th layer in the transformer decoder is processed by an additional decoding head, which is nearly identical to the one used for outputting action chunks, with the only modification being the output dimension. This decoding head outputs manipulation concept chunks corresponding to the same time steps as the action chunks, with its parameters adjusted to match the dimensionality of the manipulation concepts at each time step (256). Other training and testing settings follow [39]. Moreover, the transformer decoder in ACT’s CVAE decoder, as implemented by [39], consists of 7 layers. During our experiments, we tested various combinations of  $L$ -th layers to determine the optimal layer for processing by the manipulation concept decoding head. Our results indicate that  $L = 2$  provides slightly better performance than other configurations. We present the ablation study on  $L$  and the weight  $\lambda_{mc}$  in Eq. 9 for ACT on LIBERO-90 tasks, as shown in Tab. 5<sup>1</sup>. However, we believe this raises an interesting and challenging direction for future work: systematically investigating the rationale and insights behind the selection of  $L$ , even beyond the context of our setting.

Table 5: Ablation study on the intermediate layer outputs ( $L$ ) used as inputs to the manipulation concept decoder and the loss weight  $\lambda_{mc}$  in Eq. 9 for Enhancing Imitation Learning in ACT, evaluated on LIBERO-90 tasks.

ACT	$\lambda_{mc} = 1.0$	$\lambda_{mc} = 0.1$	$\lambda_{mc} = 0.01$	$\lambda_{mc} = 0.001$
$L = 2$	<b>74.8±0.8</b>	70.6±0.8	69.0±0.1	68.7±0.5
$L = 3$	70.0±0.4	69.9±0.2	68.8±1.0	68.7±0.6
$L = 4$	72.6±0.5	69.9±0.3	69.6±0.2	67.3±0.5

**Diffusion Policy** The pipeline is illustrated in the lower part of Fig. 6. Following [39], the convolution-based Diffusion Policy is modified to concatenate the noise level ( $k$  and corresponding  $\epsilon_k$ ) embedding, observation, and task embedding as the conditional input to the diffusion model network, using the FiLM strategy. We further introduce an additional manipulation concept decoding up-sampling module, nearly identical to the one used for outputting action chunks, with the only modification being the output dimension, to decode intermediate outputs from the corresponding up-sampling layer of the diffusion model. This decoding head can be configured to process intermediate outputs to predict manipulation concept chunks corresponding to the time steps of the predicted (noise of) action chunk outputs. The figure illustrates the cases for  $L = 0$  and  $L = 1$ . During our experiments, we tested various combinations of  $L$ -th layers to identify the optimal layer for processing by the manipulation concept decoding head. Our results suggest that  $L = 1$  achieves better performance than other configurations. We present the ablation study on  $L$  and the weight  $\lambda_{mc}$  in Eq. 9 for Diffusion Policy on LIBERO-90 tasks, as shown in Tab. 6. Similar to ACT, we believe this topic needs further systematic study to uncover deeper insights. Other training and testing settings follow [39].

Table 6: Ablation study on the intermediate layer outputs ( $L$ ) used as inputs to the manipulation concept decoder and the loss weight  $\lambda_{mc}$  in Eq. 9 for Enhancing Imitation Learning in Diffusion Policy, evaluated on LIBERO-90 tasks.

DP	$\lambda_{mc} = 1.0$	$\lambda_{mc} = 0.1$	$\lambda_{mc} = 0.01$	$\lambda_{mc} = 0.001$
$L = 0$	83.5±0.8	78.9±0.4	78.7±0.3	75.6±0.6
$L = 1$	80.0±0.4	<b>89.6±0.6</b>	82.0±0.2	79.9±0.1

<sup>1</sup>We provide sample rollouts (supplementary/rollout\_summary) and .gif logs of test-time ACT and DP performance (supplementary/rollout\_video\_samples.gif) in the supplementary materials.

Our future work includes a deeper study of modification strategies for various policies to adapt to the Enhancing Imitation Learning framework, following the methodology outlined in Sec. 3.4.

### A.3 Manipulation concept discovery (Baselines)

- **InfoCon.** Based on the design of InfoCon [31], All the size of hidden features output by transformers and concept features is 256. The state encoder (also process video clips consisting of concatenated, compressed vision observations and proprioceptive states, as outlined in Sec. A.1) uses a 12-layer transformer. The state reconstructor uses a 4-layer transformer. The goal-based policy uses a 4-layer transformer. The predictor for the generative goal uses a 4-layer transformer. For hyper-network used for discriminative goals, we use 2 hidden layers in the goal function. The number of concepts is fixed, the maximum number of 30 manipulation concepts for all the tasks. We employ the AdamW optimizer, coupled with a warm-up cosine annealing scheduler same as Sec. A.1. The weight decay is always  $1.0 \times 10^{-3}$ . We use a batch size of 512 during training. We train our model for 200,000 iterations with a base learning rate of  $1.0 \times 10^{-3}$  on a single A800 GPU within 1.5 days.
- **XSkill.** Following the design of XSkill [65], we implement its skill discovery framework on LIBERO-90, focusing exclusively on the “robot” embodiment and the Skill Discovery component from the XSkill pipeline. To ensure comparable model capacity and support multi-modality, our implementation employs a 12-layer Transformer as the temporal skill encoder. This encoder processes video clips consisting of concatenated, compressed vision observations and proprioceptive states, as outlined in Sec. A.1, along with a trainable token to predict skill representations, which are subsequently used for skill prototype prediction. To augment the concatenated video clips containing multi-modality information, Gaussian noise with  $\sigma = 1.0 \times 10^{-3}$  is applied. This unified augmentation approach accommodates the nature of proprioceptive states, as standard image augmentation techniques are not directly suitable for robotic proprioception. The training process employs a batch size of 512 and a learning rate of  $1.0 \times 10^{-3}$  for 200,000 iterations on a single A800 GPU within 1.5 days.
- **DecisionNCE** We fine-tune the DecisionNCE-T model (<https://github.com/2toinf/DecisionNCE>) on our dataset, as it outperforms DecisionNCE-P in our analysis of the experimental results in [27]. We use two types of language annotations: (1) the original task descriptions (Decision-task), and (2) detailed subtask labels derived by decomposing each task into meaningful subprocesses (Decision-motion). To construct the latter, we manually segment each demonstration based on changes in the robot’s proprioceptive state (e.g., movement direction, gripper open/close status). Segments corresponding to the same task are then assigned unified subtask labels across demonstrations, with remaining inconsistencies resolved through manual adjustment.
- **RPT.** We modify the original RPT design [47] to adapt it for our task of discovering manipulation concept latents in the LIBERO-90 setting. We employ a 16-layer self-attention transformer to process inputs consisting of 60 consecutive, interleaved agent-view and eye-in-hand vision frames. Vision inputs are compressed using a stable diffusion VAE encoder, similar to the method in Sec. A.1. The total sequence length processed by the transformer is  $60 \times 3 = 180$ . Each modality is mapped to a 256-dimensional embedding vector using an MLP, as defined in Eq. 11. The transformer’s output is then decoded to reconstruct the original inputs using a 3-layer MLP with 1024-dimensional hidden layers. We follow the masking strategy outlined in [47] to perform temporal MAE training for the transformer. To label manipulation concept latents using the trained transformer, we extract the intermediate output of the 12th layer when the input consists of the full observation without masking. Notice that we select the output at the proprioceptive state input positions of the transformer to represent the manipulation concept latent at each time-step. The labeling process follows the procedure introduced in Sec. A.1. For training, we use a batch size of 512 and a learning rate of  $1.0 \times 10^{-3}$ , running for 200,000 iterations on a single A800 GPU, which completes within 3 days.
- **All.** This is an ablation version of our manipulation concept discovery method, focusing on the design for capturing multi-modal correlations (Sec. 3.2). Specifically, this baseline replaces the loss in Eq. 3 with a loss that does not use partial masking but instead always masks all modalities:  $\mathcal{L}_{\text{all}}(t, \tau_i) = \|\mathcal{C}(\emptyset | z_i^t; \Theta_c) - \mathbf{o}_i^t\|$ . Based on our reasoning and the experiment results show in Tab. 3, we think this method may not be good at learning correlation between different modalities. Other settings follow Sec. A.1.

- **Next.** This is an ablation version of our manipulation concept discovery method, focusing on the design for representing multi-horizon subgoals (Sec. 3.3). Specifically, this baseline replaces the loss in Eq. 7 with a loss that always predicts the next adjacent time-step observation:  $\mathcal{L}_{\text{next}}(t, \tau_i) = \|\mathcal{F}(\mathbf{o}_i^t, z_i^t; \Theta_f) - \mathbf{o}_i^{t+1}\|$ . We observe that this setting is commonly used in recent works [7, 68], which learn representations based on adjacent time-step observations or observations separated by a fixed time horizon. We suggest that learning based on a fixed time horizon is conceptually similar to adjacent time-step settings, as the fixed time horizon can be interpreted as a unified time step. Our method differs by considering the temporal correlation across multiple variable horizons, which is also addressed by baseline methods like RPT. Other settings follow Sec. A.1.
- **CLIP.** To ensure compatibility with other baselines, which have an output dimension of 256, we select the ViT-B/32 CLIP model from the original source (<https://github.com/openai/CLIP>). This model outputs a 512-dimensional feature vector, the closest to 256 among the accessible CLIP models from this codebase when given an image.
- **DINOv2.** To match the output dimension of 256 used by other baselines, we select the dinov2-small DINOv2 model from the source at <https://huggingface.co/facebook/dinov2-small>. This model produces a 384-dimensional feature vector when given an image.

Note that DecisionNCE, CLIP, and DINOv2 baselines use only vision (and language) information for concept discovery. We preserve their original modality structure rather than adapting them to include proprioceptive states, as this would deviate from their pretraining foundations.

#### A.4 Mutual information estimation

The estimation of mutual information is based on MINE [2], which uses batchwise samples drawn from a joint distribution and employs a neural network to estimate the mutual information. To extend this approach for estimating conditional mutual information (CMI), we reformulate CMI by decomposing it into mutual information terms, as shown below:

$$\mathbb{I}(X : Y | Z) = \mathbb{I}(X : Y) + \mathbb{I}(XY : Z) - \mathbb{I}(X : Z) - \mathbb{I}(Y : Z), \quad (13)$$

where  $XY$  denotes the random variable sampled from the joint distribution of  $X$  and  $Y$  and is represented as the concatenation of their encoded vectors. The neural network in MINE has two layers, with the hidden layer size set to 1.5 times of the dimensions of the two random variables.

## B Pseudocode

Here we provide pseudocode for (i) Deriving subprocess from manipulation concept latents (Alg. 1). (ii) Manipulation concept disocvery training process of our method (Alg. 2).

## C More Study on Learned Manipulation Concepts

### C.1 Additional Experiments on Enhanced Imitation Learning

**Sampling Strategies** In this part, we focus on methodology for deriving hierarchical structures from learned representations (Sec. 3.3). While we adopt a threshold-based hierarchy derivation method (Eq. 4) as a proof of concept, we acknowledge that alternative derivation methodologies warrant further investigation (see Sec. D). For the threshold-based approach, we employ uniform sampling of the threshold  $\epsilon$  during training. This choice ensures full coverage of all possible hierarchical structures, as we do not know a priori which threshold values might be suboptimal. To validate this design choice, we conduct an ablation study comparing different sampling strategies for  $\epsilon$  in Eq. 7:

As shown in Tab. 7, uniform sampling currently achieves the best performance across both policy architectures. We hypothesize that while task-specific sampling strategies might excel on particular subsets, uniform sampling provides robust performance across diverse tasks due to its comprehensive coverage of the threshold space. Future work could explore adaptive sampling strategies tailored to specific task distributions.

Table 7: **Sampling Strategies Ablation.** We compare different sampling strategies for  $\epsilon$  in Eq. 7. Manipulation concepts are learned from LIBERO-90 and applied to policy learning on LIBERO-90. We report success rates (%).

Sampling Strategy	Description	ACT	DP
Uniform (Ours)	$\epsilon \sim \mathcal{U}(0, 1)$	74.8±0.8	89.6±0.6
Sparse	$\epsilon \sim \{0.1, 0.2, \dots, 1.0\}$	67.6±0.5	81.1±0.8
Biased	$\epsilon \sim \mathcal{U}(\frac{1}{3}, \frac{2}{3})$	65.6±0.7	78.7±0.4

**Learning Methodology Contribution** We conduct an ablation study to isolate the contributions of our two core learning methodologies: Capturing Multi-Modal Correlations (Sec. 3.2) and Representing Multi-Horizon Sub-Goals (Sec. 3.3). Tab. 8 compares three configurations: (1) *Cross-modal only*: learning with only cross-modal alignment objectives in Eq. 3, (2) *Multi-horizon only*: learning with multi-horizon sub-goal prediction in Eq. 7 but without cross-modal alignment, and (3) *Full method*: combining both cross-modal alignment and multi-horizon prediction.

Table 8: **Methodology Contribution Ablation.** We evaluate the contribution of each learning component by training manipulation concepts on LIBERO-90 and applying them to policy learning on LIBERO-90. We report success rates (%).

Method	ACT	DP
Cross-modal only	69.1±0.6	82.8±1.0
Multi-horizon only	71.6±0.4	80.5±0.5
Ours (Full method)	74.8±0.8	89.6±0.6

The results in Tab. 8 reveal that both components make substantial and complementary contributions to performance. We attribute this synergy to the distinct roles of each component: cross-modal alignment grounds the understanding of correlations across different modalities, while multi-horizon prediction captures hierarchical temporal structure. Together, they enable the learning of manipulation concepts that are both correlationally coherent and temporally structured, leading to more robust policy learning.

**Data Constraint Experiments** We evaluate whether manipulation concepts can help mitigate the challenges of imitation learning under limited data. Specifically, we vary the amount of data available for training both the manipulation concept encoder (Eq. 1) and the enhanced imitation learning framework (Sec. 3.4) to assess their impact on policy success rates. We conduct experiments on LIBERO-90 tasks using the diffusion policy. As shown in Tab. 9, incorporating manipulation concepts consistently improves policy performance compared to settings without them, even under restricted data conditions. This demonstrates that learning and leveraging manipulation concepts can make imitation learning more data-efficient and effective.

Table 9: **Performance under data constraints.** Success rates of diffusion policies with and without manipulation concept enhancement, evaluated on LIBERO-90 (*L90-90*). In each setting, the number of demonstrations per task available for training both the manipulation concept encoder and the policy is limited as indicated.

	50 demos/task	25 demos/task	10 demos/task
Ours	89.6 ± 0.6	77.6 ± 0.5	61.2 ± 1.1
Plain	75.1 ± 0.6	70.1 ± 0.3	59.1 ± 0.9

**Distance Metric** We conduct an ablation study comparing spherical distance and cosine distance  $\frac{1-\cos(\cdot)}{2}$  for  $\text{dist}(\cdot, \cdot)$  in Eq. 4. Tab. 10 reports the performance when concepts are learned and applied

to LIBERO-90 tasks. Further investigation into distance-threshold-based subprocess derivation methods represents a promising direction for future work.

Table 10: Ablation study on distance metrics for concept learning on LIBERO-90. Spherical distance consistently outperforms cosine distance across both baseline methods.

	Cosine Distance	Spherical Distance (Ours)
ACT	67.8±0.5	74.8±0.8
DP	82.0±0.4	89.6±0.6

**Sub-process Derivation** We conduct an ablation study comparing two approaches for constraining manipulation concept latents within each sub-process in Eq. 4. Our proposed method enforces proximity among all concept latents throughout the sub-process (“Sequential Constraint”), while the baseline only constrains the distance between the initial and final concept latents (“Endpoint Constraint”). We evaluate both approaches on LIBERO-90, where concept discovery and policy enhancement are performed. Tab. 11 reports the task success rates when integrating the learned manipulation concepts with different policy architectures.

Table 11: Ablation study on sub-process derivation constraints. We compare enforcing proximity among all manipulation concept latents within each sub-process (Sequential Constraint) versus constraining only the initial and final latents (Endpoint Constraint). Results show average success rates (%) with standard errors across LIBERO-90 tasks.

	Sequential Constraint	Endpoint Constraint
ACT	74.8±0.8	68.4±0.8
DP	89.6±0.6	79.8±0.5

**Future Prediction Strategy** Apart from the different sub-goal determination strategies we compared (**Next** and **InfoCon** in Sec. 4.1), we evaluate two additional future prediction strategies.

- **Next-n.** Unlike our sub-process derivation strategy (Eq. 4), this baseline encodes future observations at varying time horizons by randomly sampling a future timestep. Specifically:  $\mathcal{L}_{\text{next-n}}(t, \tau_i) = \mathbb{E}_{n \sim U\{1, 2, \dots, T_i - t\}} \|\mathcal{F}(\mathbf{o}_i^t, z_i^t, n; \Theta_f) - \mathbf{o}_i^{t+n}\|$ .
- **Next-random.** This strategy builds upon **Next-n** but differs in how future targets are selected. We first randomly segment training demonstrations into sub-processes for concept discovery. Then, for a state at time-step  $t$ , the prediction target is randomly selected from among the end-states of subsequent sub-processes. For example, if a demonstration is segmented into 5 sub-processes and time-step  $t$  is in the 2nd sub-process, the model will randomly predict one of the end-states from the 2nd, 3rd, 4th, or 5th sub-processes during concept discovery learning.

We evaluated diffusion policies enhanced by these strategies, with results presented in Tab. 12. The data demonstrates that our manipulation concepts yield better policy enhancement compared to the alternative strategies. This highlights the importance of carefully designing which future observations to predict and validates the effectiveness of our self-supervised sub-goal derivation and learning method. Specifically, the performance decrease observed with **Next-n** and **Next-random**, despite their consideration of multi-horizon futures, likely stems from the fact that not all future states effectively represent sub-goal completion. Intermediate movement states may be reached through multiple alternative trajectories that ultimately achieve the same sub-goal, thus providing limited information about the underlying task structure.

**Usage of Manipulation Concept Encoder** We investigate two strategies for leveraging the manipulation concept encoder from Eq. 1 in downstream policy learning. The encoder serves as an intermediate module that extracts manipulation concept representations from demonstrations. We compare the following approaches: (1) **Direct Conditioning**: The trained encoder directly processes current observations to generate manipulation concepts, which are then concatenated with observations as additional input features to the policy network. (2) **Joint Prediction (Ours)**: The policy

Table 12: **Comparison of Additional Future Prediction Strategies.** Success rates of diffusion policies enhanced with manipulation concepts discovered using our method versus two alternative future prediction strategies on the LIBERO-90 benchmark.

<b>L90-90</b>	<b>Ours</b>	<b>Next-n</b>	<b>Next-random</b>
DP	$89.6 \pm 0.6$	$83.0 \pm 0.3$	$82.8 \pm 0.4$

network is trained to jointly predict both future actions and future manipulation concepts from current observations, as described in Sec. 3.4. Tab. 13 presents the comparative results across two policy architectures.

Table 13: **Comparison of Manipulation Concept Usage Strategies.**

<b>Policy</b>	<b>Direct Conditioning</b>	<b>Joint Prediction (Ours)</b>
ACT	$71.1 \pm 0.4$	$74.8 \pm 0.8$
DP	$79.3 \pm 0.9$	$89.6 \pm 0.6$

The performance gap stems from a temporal alignment mismatch between concept representations and action predictions. In **Direct Conditioning**, the encoder extracts concepts from current or past observations, creating a temporal lag: the policy receives historical concept information when planning future actions. In contrast, Joint Prediction enforces temporal coherence by training the policy to predict future manipulation concepts alongside future actions, ensuring that the predicted concepts align temporally with the planned action sequence.

This temporal alignment is critical in multi-phase manipulation tasks. For example, consider a pick-and-place scenario: immediately after grasping an object, the current observation encodes grasping-related dynamics. However, to execute the subsequent placement action, the policy requires placement-relevant information. Joint Prediction learns to anticipate these future task-phase concepts, providing the policy with forward-looking contextual information. Direct Conditioning, by contrast, conditions the policy on backward-looking grasping concepts that offer limited guidance for placement planning.

While our results demonstrate the advantages of temporal alignment through joint prediction, we acknowledge that direct conditioning on historical concepts may benefit tasks requiring explicit long-horizon memory or reactive behaviors based on past states [14]. Future work will explore hybrid architectures that combine both strategies.

## C.2 Alignment with Semantic Sub-Goals

We evaluate whether the manipulation concept latents learned by our method resemble human-interpretable semantics. Specifically, we assess whether latents assigned to time steps of demonstrations (Sec. 3.1) exhibit higher pairwise similarity when those steps belong to sub-processes pursuing the same human-defined sub-goal.

To analyze the learned representations, we first group manipulation concept latents according to human-annotated sub-goals. For instance, in the task “open the top drawer”, latents from time steps where the robot reaches for the top drawer handle are categorized as “reach the top drawer”. Latents from other demonstrations and tasks involving identical processes (reaching the top drawer) are placed in the same category. We then quantify the similarity between two categories by calculating the average cosine similarity between their respective latents, as defined in Eq. 10.

Fig. 9 shows results from analyzing demonstrations from three tasks:

- Task #1: Open the top drawer of the cabinet and put the bowl in it;
- Task #2: Close the bottom drawer of the cabinet and open the top drawer;
- Task #3: Close the top drawer of the cabinet and put the black bowl on top.

We selected these tasks because they clearly demonstrate overlapping subgoals across different tasks (e.g., Task #1 and Task #2 both include “opening the top drawer”). This enables testing whether the

latents capture similar subgoal semantics across different tasks—an essential capability for cross-task learning efficiency (Sec. 1). Manipulation concept latents are grouped based on human-defined sub-goals, with similarities between category pairs visualized as heatmaps. Three heatmaps are presented, each using a different granularity of sub-goal annotation:

1. Top-1st heatmap: Omits task-specific distinctions, merging similar manipulation processes across tasks into the same category
2. Top-2nd heatmap: Further merges similar manipulation processes, disregarding distinctions like “top drawer” versus “bottom drawer”
3. Top-3rd heatmap: Consolidates manipulation processes further, treating actions like bowl transitions as the same concept regardless of context

In each heatmap, the entry at position  $(i, j)$  represents the average similarity ( $\times 10.0$ ) between categories  $i$  and  $j$ . For readability, only the top three similarity values in each row are displayed.

We emphasize that testing semantic capture at different “description granularity levels” is important because semantics naturally exist at multiple levels of abstraction, from highly specific details to broadly generalizable patterns. Finer-grained descriptions provide more precise details but limited generalization, while coarser-grained descriptions capture more general features applicable across diverse scenarios. For example, the general instruction “close the drawer” applies broadly to subprocesses in both Task #2 and Task #3, whereas the more specific “close the top drawer” incorporates spatial features that make it applicable in Task #3 but not in Task #2. Through this multi-granularity analysis, we evaluate whether our manipulation concept latents successfully capture both fine-grained semantics needed for specific scenarios and coarse-grained semantics that enable transfer across more scenarios.

What we observe is that the highest similarity values consistently appear along the diagonal in each heatmap in Fig. 9, so concept latents from the same category show higher similarity compared with different categories. This indicates that the learned latent clusters resemble clusters derived from human-interpretable sub-goal classifications, suggesting that our model captures meaningful semantic structure in the manipulation processes. Moreover, the patterns observed across the three heatmaps with different description granularities reveal that the latents encode semantics at multiple levels of abstraction. They capture both generalizable semantics applicable across tasks and scenes, while simultaneously preserving fine-grained scene-specific details.

Furthermore, Fig. 10(b) provides a t-SNE visualization of manipulation concept latents from all 90 tasks in LIBERO-90. For each task, latents ( $z_i^t$ ) were extracted at every time step of demonstrations. In the plot, latents are color-coded by their originating tasks. We observe that clusters often contain latents from diverse tasks, as indicated by the mixed colors in each cluster. This further supports our finding that the learned latents generalize across tasks and capture shared semantic structures.<sup>2</sup>

### C.3 Motion Study

We evaluate whether the learned manipulation concept latents capture the robot’s motion. Using Eq. 10, we calculate the average similarity ( $\times 100.0$ ) between movements based on manipulation concept latents corresponding to specific gripper actions. Specifically, we collect latents for the following movements from task demonstrations in LIBERO-90:

1. Forward-backward motion: Latents for time-steps where the robot moves forward, backward, or remains still along the forward-backward axis.
2. Left-right motion: Latents for time-steps where the robot moves left, right, or remains still along the left-right axis.
3. Up-down motion: Latents for time-steps where the robot moves up, down, or remains still along the up-down axis.
4. Gripper state: Latents for time-steps where the gripper opens or closes.

---

<sup>2</sup>It should be noted that t-SNE performs extreme dimensionality reduction, so these clusters may not perfectly reflect similarity in the high-dimensional space. This visualization should therefore be considered as supplementary evidence.

Movements with velocities below 20% of the maximum observed velocity are classified as “still”. Using these collected latents, we generate heatmaps (similar to Fig. 9) to visualize the average cosine similarity across different movement directions and gripper states (Fig. 10(a)).

The heatmaps reveal that the highest cosine similarity values often appear along the diagonal. This demonstrates that latents corresponding to the same motion patterns exhibit greater similarity to each other than to those from different motion patterns, indicating that the latents effectively capture different movement directions and gripper states. However, we observe that forward-backward motion is captured with lower accuracy compared to other dimensions. We hypothesize that incorporating additional 3D-informative modalities, such as depth maps, beyond the current proprioceptive states could improve the representation of motion along the forward-backward axis. We leave the exploration of such modality incorporation to future work.

#### C.4 Diversity & Discrimination Study

We analyze the diversity and discriminability of learned manipulation concepts by comparing concept latents from our method (Sec. 3) and the baselines in *Manipulation Concept Discovery Baselines* (Sec. 4.1). Specifically, we cluster latents from these methods and examine the number of clusters under varying granularities. The number of clusters reflects concept diversity: more clusters indicate a wider variety of concepts. Clustering granularity determines whether clusters are fine-grained (fine granularity) or general (coarse granularity). Additionally, small granularity perturbations test discriminability, as less discriminative latents lead to significant clustering changes under small granularity variations. For each method, We collect manipulation concept latents from 90 LIBERO-90 tasks (one demonstration per task) and use DBSCAN to cluster them while varying the density parameter  $\text{Eps}$ , which controls clustering granularity. Fig. 11 shows the cluster counts across different  $\text{Eps}$  values. From Fig. 11, our manipulation concept discovery method (**Ours**) demonstrates two key advantages: 1) At higher granularities ( $\text{Eps} > 0.2$ ), **Ours** maintains a **higher number of clusters**. 2) The **decline in cluster count is relatively smooth and gradual**, showing stability under small  $\text{Eps}$  changes. These results highlight the superior diversity and discriminability of our manipulation concept discovery method.

#### C.5 Multi-Level Hierarchical Structure

In Fig. 3, we present a visualization example of the **Multi-Level Hierarchical Structure** described in Sec. 4.3. Additional visualization results are available in the supplementary materials under the directory `supplementary/vis_multi_h`.

#### C.6 Real Robot Experiments Details

**Training Data.** As shown in Fig. 5, the training data for the “cleaning cup” task consists of demonstrations using mobile ALOHA [16] to place the cup from the table into the container. Each demonstration features a scene containing exactly one cup and one container. There are two pairings of color combinations: blue cups with green containers and yellow cups with pink containers. For each pairing, we collect 27 demonstrations with varied spatial arrangements.

**Evaluation Setting.** For evaluation, we test our model on six scenarios that introduce variations absent from the training data:

- **Novel Placements.** Objects maintain the same color pairings as in training but appear in previously unseen spatial arrangements.
- **Color Composition.** We rearrange color pairings (blue cups with pink containers and yellow cups with green containers) to test generalization across color combinations.
- **Novel Objects.** We introduce unseen objects, such as bamboo-woven containers, pink cups not present in training, or cups initially placed on plates rather than directly on the table.
- **Obstacles.** We position obstacles in front of cups to challenge visual perception.
- **Barriers.** We place a plate inside the container, requiring the robot to lift the cup high enough to clear this barrier when depositing it.
- **Grasping Together.** We position two cups adjacent to one another, requiring the robot to grasp both simultaneously at their contact point and deposit them together in the container.



**Figure 7: Multi-horizon goal prediction with learned manipulation concepts.** Visualization of future states predicted by our Multi-Horizon Goal Predictor (MHGP, Eq. 7) when conditioned on the current observation, a manipulation concept latent ( $z$ ), and varying coherence thresholds ( $\epsilon$ ). From left to right, as  $\epsilon$  increases from 0 to 1, predictions extend progressively further into the future, demonstrating how our manipulation concepts encode temporal abstraction at multiple horizons. Note that predictions capture essential functional relationships (robot-object interactions) rather than pixel-perfect reconstructions, facilitating generalization across environments.

**Manipulation Concept Discovery.** The model architecture and hyperparameter configuration for manipulation concept discovery follow the methodology described in Sec. A.1. Since the dataset is relatively small, we adapt smaller transformers: a 4-layer concept encoder ( $\mathcal{E}$ , Eq. 1), a 4-layer Cross-Modal Correlation Network ( $\mathcal{C}$ , Eq. 3), and a 4-layer Multi-Horizon Future Predictor ( $\mathcal{F}$ , Eq. 7). For data collected using mobile ALOHA [16], we incorporate the following modalities: three  $640 \times 480$  resolution cameras (left-gripper, right-gripper, and upper-gripper) and 42-dimensional proprioception states (comprising 14-dimensional joint torque, position, and velocity measurements). All image data undergoes preprocessing as detailed in Sec. A.1.

**Enhancing Imitation Learning.** Please refer to ACT section in Sec. A.2.

### C.7 Multi-Horizon Goal Prediction Visualization

We provide visualization results of the **Multi-Horizon Goal Prediction Visualization** (Sec. 4.4) in Fig. 7 and supplementary materials under the directory `supplementary/prediction`. Below are the details of the experiments:

**Dataset.** For our experiments, we utilized the BridgeDataV2 dataset [60]. Since multi-view data is not universally available across all demonstrations, we selected two specific modalities: the robot’s proprioceptive states (7DoF) and the third-person camera view. The camera images were preprocessed to  $128 \times 128$  resolution following the procedure outlined in Sec. A.2.

**Manipulation Concept Discovery.** We implemented the model architecture and hyperparameter configuration as detailed in Sec. A.1, adapting it specifically to operate with the two modalities described in the **Dataset** section above.

### C.8 Preliminary VLA Integration

We present a preliminary exploration of integrating manipulation concepts with vision-language-action models (VLAs). We build upon OpenVLA-OFT [23], which fine-tunes OpenVLA using pretrained parameters and a novel action adapter for downstream tasks. The action adapter processes hidden layer features from the original pretrained VLA model. Following this architecture, we introduce an additional “concept adapter” that implements the method described in Sec. 3.4, enabling the integration of manipulation concepts into the VLA.

To evaluate the data efficiency gains from manipulation concepts, we fine-tune the enhanced VLA on **50% of the training data** used for LIBERO-10 tasks in the original OpenVLA-OFT study [23]. We compare fine-tuning performance with and without manipulation concept integration. Fig. 8 presents the results, where the x-axis indicates training epochs and the y-axis shows success rates for checkpoints at each epoch. The solid lines labeled “best” represent the highest success rate achieved up to that epoch.

The results demonstrate that manipulation concepts improve data utilization. With only **half the training data**, the concept-enhanced approach consistently achieves higher success rates throughout training. Notably, the original OpenVLA-OFT achieved 94.5% success with the full dataset [23], while our concept-enhanced model with **half the data** reaches comparable performance levels, indicating substantially improved data efficiency.

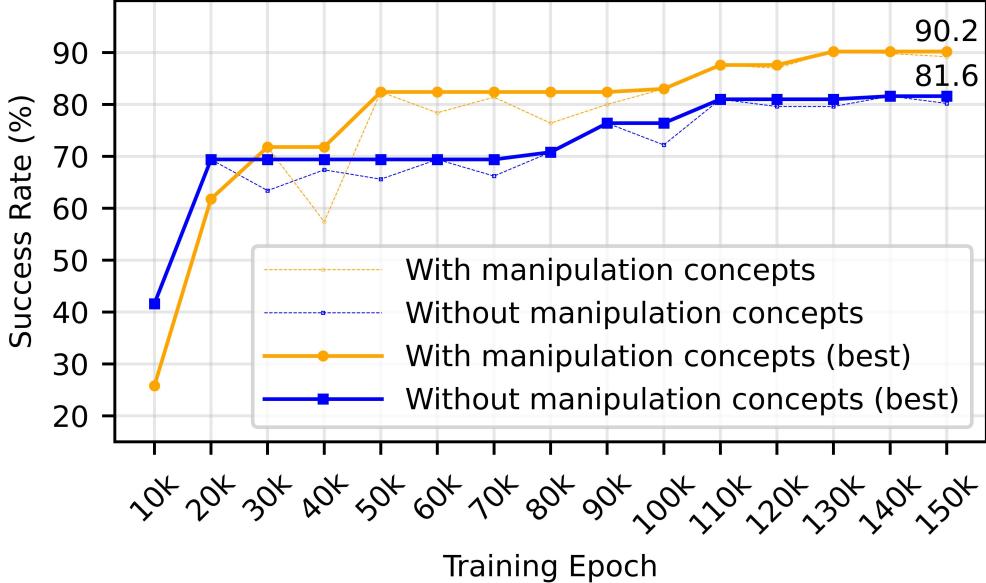


Figure 8: Data efficiency comparison on LIBERO-10 tasks with **50% training data**. Solid lines show best performance up to each epoch for models with and without manipulation concepts.

We hypothesize that this improvement stems from HiMaCon’s ability to capture manipulation dynamics at multiple abstraction levels. The learned concepts provide explicit intermediate representations that bridge high-level task instructions and low-level control actions, thereby reducing the learning burden on VLAs by supplying structured manipulation knowledge rather than requiring learning of complex sensorimotor patterns from scratch. Further investigation of this integration will be pursued in future work.

## D Limitations & Future works

**Further Exploration of multi-modality.** We propose enhancing robotic data collection with richer modalities and studying how these modalities can derive more effective manipulation concepts. While current robotics research primarily focuses on visual information, human manipulation relies on multiple sensory inputs, particularly tactile feedback to complement vision. This is especially crucial for robotic systems with limited tactile capabilities. Future work should investigate which modalities contribute most significantly to performance improvements and how to fully leverage their potential.

**Further Exploration of multi-horizon sub-goal.** Our work proposes methods to derive sub-processes for achieving sub-goals across multiple horizons, though several improvements remain possible. Current methods inadequately capture relationships between different values of  $\epsilon$  in Eq. 4, failing to reflect the natural tree structure of hierarchical sub-goals. Future research could explicitly derive tree structures [61, 74] where long-horizon sub-goals serve as parent nodes to short-horizon child nodes. Additionally, our cosine similarity approach for determining sub-goal correspondence could be refined with more sophisticated metrics.

**Scaling up.** Computational constraints have limited our exploration of how our method scales with larger datasets. We plan to leverage pretrained multi-modal foundation models, adopting structures inspired by [7] and extending pretraining beyond robotics data as in [68]. We also aim to further investigate whether our manipulation concepts can enhance advanced policies like Vision-Language-Action models [4, 19, 23, 24].

---

**Algorithm 1** Derive Subprocess  $h(\mathbf{z}_i; \epsilon)$ 


---

**Input:** manipulation concept vectors  $\mathbf{z}_i = \{z_i^t\}_{t=1}^{T_i}$ , coherence parameter  $\epsilon \in [0, 1]$ .  
**Initialize:**  $End = []$ ,  $g_b = 1$   
**while**  $g_b \leq T_i$  **do**  
     $g_e = g_b + 1$   
    **while** *true* **do**  
        **if**  $\exists u \in [g_b, g_e]$ , s.t.  $\text{dist}(z_i^u, z_i^{g_e}) \geq \epsilon$  **or**  $g_e > T_i$  **then**  
            **break**  
        **end if**  
         $g_e = g_e + 1$   
    **end while**  
     $End.append([g_b, g_e])$   
     $g_b = g_e$   
**end while**  
**Return**  $End$

---



---

**Algorithm 2** Manipulation Concept Discovery Training (one demonstration per batch)

---

**Input:** demonstrations  $\tau_i \in D$ , where  $\tau_i = \{(o_i^{1,t}, o_i^{2,t}, \dots, o_i^{M,t}, a_i^t)\}_{t=1}^{T_i}$   
**Initialize:** Manipulation concept assignment encoder  $\mathcal{E}(\cdot; \Theta_{\mathcal{E}})$   
**Initialize:** Modality Correlation Learner  $\mathcal{C}(\cdot; \Theta_c)$ , Subgoal Learner  $\mathcal{F}(\cdot; \Theta_f)$   
**while** *true* **do**  
    **for**  $\tau_i$  in  $D$  **do**  
         $(z_i^1, \dots, z_i^{T_i}) \leftarrow \mathcal{E}\left((o_i^{1,1}, \dots, o_i^{M,1}), (o_i^{1,2}, \dots, o_i^{M,2}), \dots, (o_i^{1,T_i}, \dots, o_i^{M,T_i}); \Theta_{\mathcal{E}}\right)$   
        **while** True **do**  
            Randomly generate a tuple  $(m_1, m_2, \dots, m_M)$ , where  $m_i \in \{0, 1\}$   
            **if**  $\sum_{i=1}^M m_i < M$  **then**  
                **break**  
            **end if**  
        **end while**  
         $(\hat{o}_i^{1,t}, \dots, \hat{o}_i^{M,t})_{t=1}^{T_i} \leftarrow \mathcal{C}\left((o_i^{1,t} \cdot m_1, o_i^{2,t} \cdot m_2, \dots, o_i^{M,t} \cdot m_M, z_i^t)_{t=1}^{T_i}; \Theta_c\right)$   
         $\mathcal{L}_{mm} = \sum_{t=1}^{T_i} \sum_{m=1}^M \|\hat{o}_i^{m,t} - o_i^{m,t}\|$   
         $\epsilon \sim U([0, 1])$   
         $End = h(z_i^1, \dots, z_i^{T_i}; \epsilon)$  {Alg. 1}  
        **for**  $t = 1$  **to**  $T_i$  **do**  
             $g_t = \min(\{g_e \mid [g_b, g_e] \in End, g_e > t\} \cup \{T_i\})$   
        **end for**  
         $(\bar{o}_i^{1,t}, \dots, \bar{o}_i^{M,t})_{t=1}^{T_i} \leftarrow \mathcal{F}\left((o_i^{1,t}, o_i^{2,t}, \dots, o_i^{M,t}, z_i^t, \epsilon)_{t=1}^{T_i}; \Theta_f\right)$   
         $\mathcal{L}_{mh} = \sum_{t=1}^{T_i} \sum_{m=1}^M \|\bar{o}_i^{m,t} - o_i^{m,t}\|$   
    **end for**  
**end while**

---

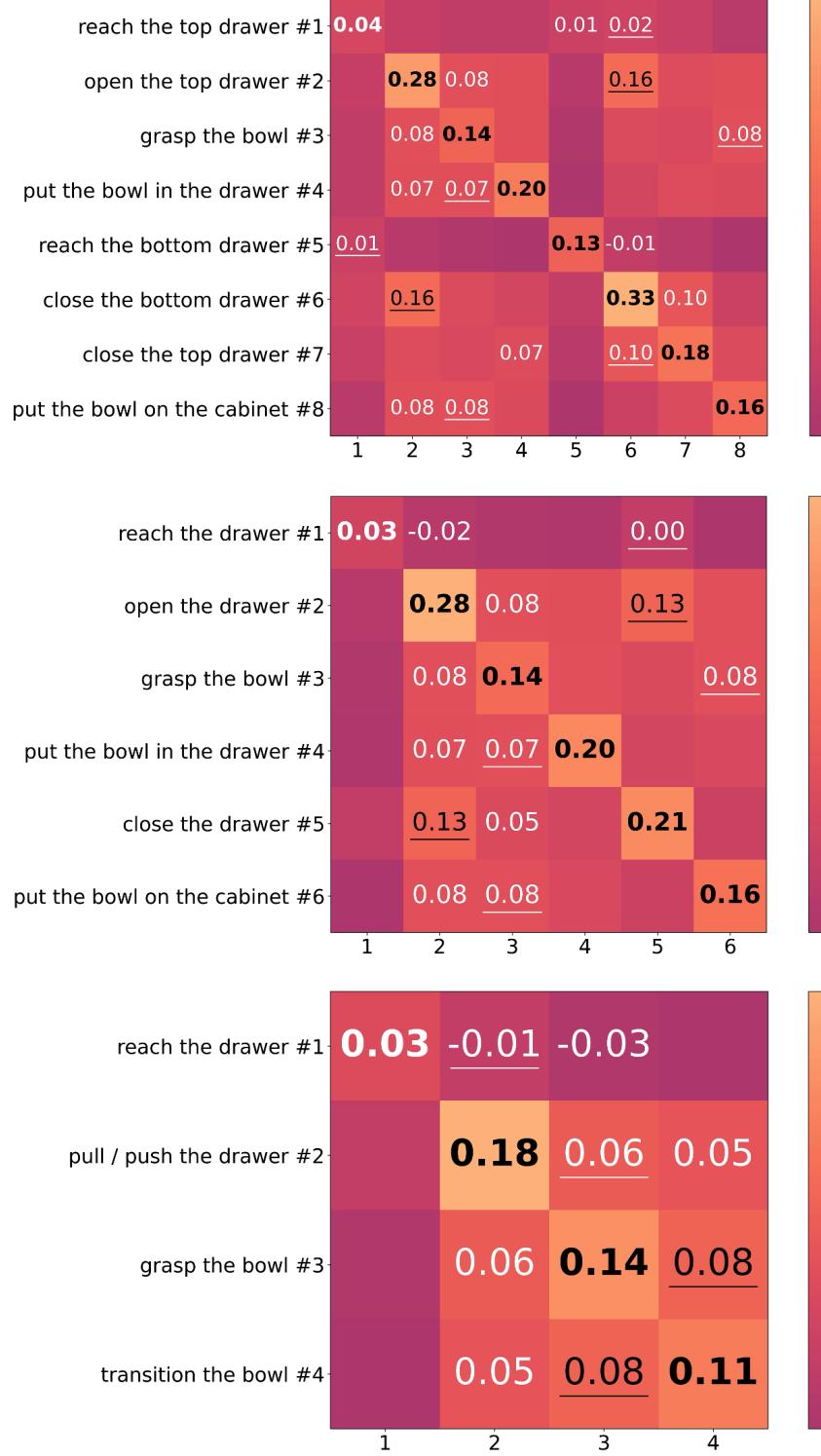
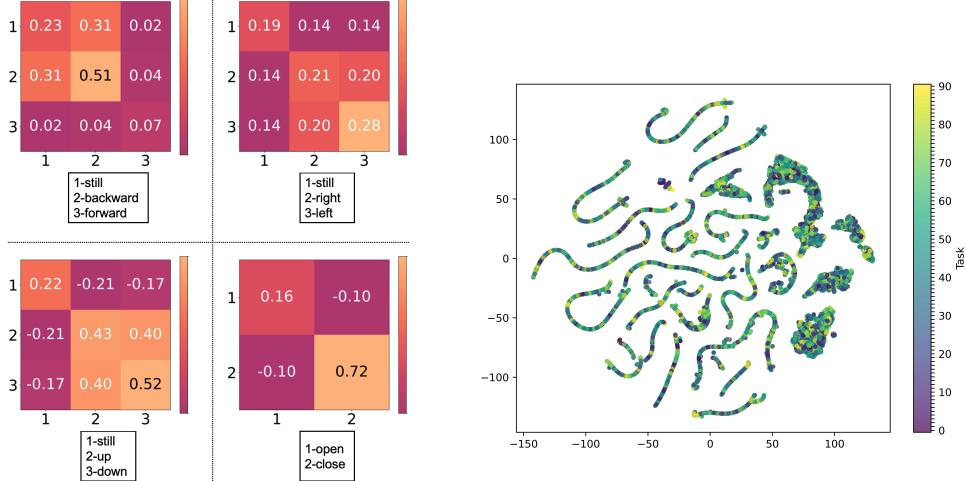


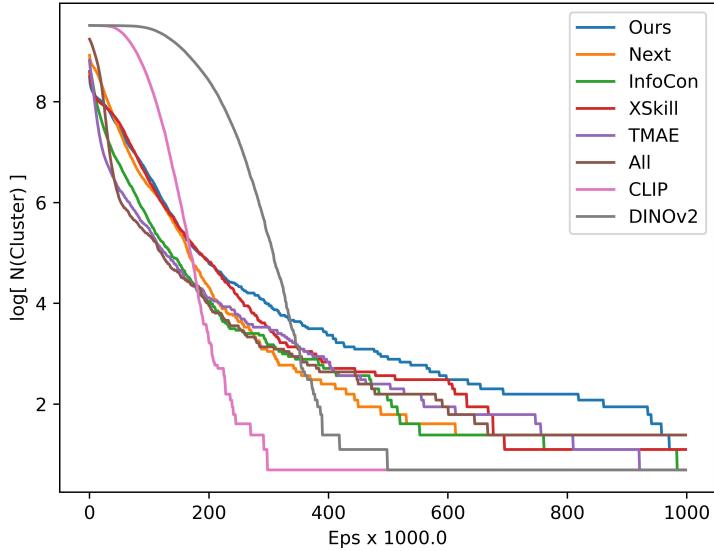
Figure 9: Average cosine similarity between pairs of sub-goal categories (defined by human semantics) computed using manipulation concept latents learned by our method (Sec.3). In each heatmap, the value at the  $i$ -th row and  $j$ -th column represents the average cosine similarity between latent vectors from the  $i$ -th and  $j$ -th categories. Three levels of labeling are provided across the heatmaps; please refer to Sec. C.2 for details.



(a) Average cosine similarity between pairs of movement categories (defined by human semantics) computed using manipulation concept latents learned by our method (Sec. 3).

(b) **t-SNE Clustering of Manipulation Concept Latents corresponding to tasks.** We perform t-SNE clustering on the manipulation concepts at each time step. These concepts are generated by our method (Sec. 3). Each sample is colored according to its task, representing one of 90 possible tasks as indicated by the colorbar.

Figure 10



**Figure 11: DBSCAN Clustering Analysis of Manipulation Concept Latents' Diversity and Discrimination.** Clustering is performed on manipulation concept latents generated by our method and the baseline methods described in **Manipulation Concept Discovery Baselines** (Sec. 4.1), across 90 tasks from the LIBERO-90 dataset. The figure shows the (log) number of clusters obtained using DBSCAN for clustering density  $\epsilon \in [0, 1]$ , with no points classified as noise.

## NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer **[Yes]**, **[No]**, or **[NA]**.
- **[NA]** means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

**The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "**[Yes]**" is generally preferable to "**[No]**", it is perfectly acceptable to answer "**[No]**" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "**[No]**" or "**[NA]**" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer **[Yes]** to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”.**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: **[Yes]**

Justification: The abstract gives a summary of our contribution on self-supervised learning of manipulation concepts.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: **[Yes]**

Justification: We discuss limitations including improvements to hierarchy derivation, further work on scaling up, and modality balance.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: We mainly make use of established theoretical frameworks (such as mutual information) for clarification and modeling of our method.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide details in the appendix and supplementary materials.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will open source all code and newly-created datasets upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Details are specified in the Experiments section and in the appendix and supplementary materials.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: For the policy success rates we currently include standard deviation.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Please refer to the details provided in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Current experiments and topics do not conflict with the Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Currently, the experiments are carried out in simulations and on robots in laboratories.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Currently, we have not encountered any safeguard issues.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have checked the sources we used.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We will release all new assets we created (code/models/datasets) upon acceptance.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We currently do not have crowdsourcing experiments.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We currently do not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

#### **16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA].

Justification: The core method development in this research does not involve LLMs.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.