

Discovering Robotic Interaction Modes with Discrete Representation Learning

Anonymous Author(s)

1 **Abstract:** Human actions manipulating articulated objects, such as opening and
2 closing a drawer, can be categorized into multiple modalities we define as interaction
3 modes. Traditional robot learning approaches lack discrete representations of
4 these modes, which are crucial for empirical sampling and grounding. In this paper,
5 we present ActAIM2, which learns a discrete representation of robot manipulation
6 interaction modes in a purely unsupervised fashion, without the use of expert labels
7 or simulator-based privileged information. Utilizing novel data collection methods
8 involving simulator rollouts, ActAIM2 consists of an interaction mode selector
9 and a low-level action predictor. The selector generates discrete representations
10 of potential interaction modes with self-supervision, while the predictor outputs
11 corresponding action trajectories. Our method is validated through its success
12 rate in manipulating articulated objects and its robustness in sampling meaningful
13 actions from the discrete representation. Extensive experiments demonstrate
14 ActAIM2’s effectiveness in enhancing manipulability and generalizability over
15 baselines and ablation studies. For videos and additional results, see our website:
16 <https://actaim2.github.io/>.

17 **Keywords:** robot manipulation, discrete representation learning, interaction
18 mode, self-supervised

1 Introduction

20 Humans exhibit an exceptional aptitude for manipulating articulated objects by utilizing prior knowledge
21 and learning through imitation. Generally, the outcome after manipulating the articulated objects
22 is categorical such as opening or closing the door. Motivated by this cognitive process, our study
23 seeks to develop a discrete representation of object affordance by merging behavior cloning with
24 interaction mode identification. We focus on objects with multiple moving parts to explore a variety
25 of distinct and meaningful outcomes, which we define as *interaction modes*. These interaction modes
26 can be represented as a discrete set of options from which the agent can sample during inference
27 to determine the appropriate interaction mode for the object. Notably, we characterize interaction
28 mode as an affordance property of the object itself which can be learned from observation data. To
29 transfer such discrete interaction modes into robotic action, an action predictor is learned from play
30 data containing these diverse modes. Therefore, we argue that the manipulation policy is a joint
31 distribution over both interaction mode selector and action predictor.

32 The robotics field is replete with studies addressing such policy decomposed into modes (or skills)
33 and action distribution. However, most behavior cloning approaches [1, 2, 3, 4] require expert data
34 for task representation supervision, especially using language description as the task representation
35 such as RLBench [5], Calvin [6], SayCan [7], etc. This assumes that the distribution of the interaction
36 mode is known and can be represented as language prompting. Other works learn the interaction mode
37 (skill) distribution using unsupervised reinforcement learning [8, 9, 10, 11] with self-supervised
38 intrinsic reward as weak supervision for task representation. Moreover, there are other interaction
39 mode learning approaches [12, 13, 14] which learn a distinguishable interaction mode (skill) prior.
40 However, the distribution of interaction modes (skills) learned in these works does not specifically
41 map to correspondent outcomes, indicating that the agent cannot sample skills with reasonable and
42 limited options. Therefore, these works fail to find a structural and disentangled space of skill for an

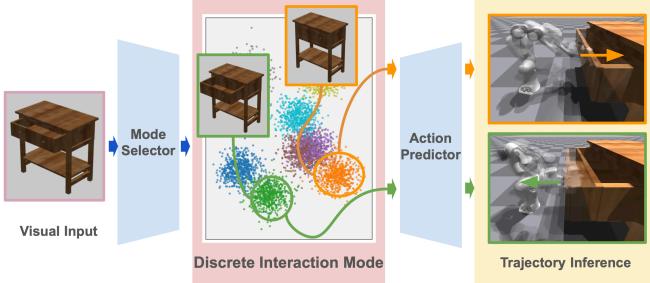


Figure 1: ActAIM2 identifies meaningful interaction modes such as open and close drawers from RGB-D images of articulated objects and robots. It represents these modes as discrete clusters of embeddings and trains a policy to generate control actions for each cluster-based interaction.

43 agent to sample the interaction mode discretely based on observation. To solve the issues above, we
 44 require that the policy 1) captures various interaction modes without necessitating expert labels or
 45 privileged information and 2) allows the agent to make discrete choices within this space, reflecting
 46 the finite interaction modes available in a given scene.

47 We introduce ActAIM2, which splits the policy into a discrete mode selector using a Gaussian
 48 Mixture Model for discrete sampling and an action predictor that processes sampled task embeddings
 49 to predict corresponding action sequences, trained via behavior cloning. Both components utilize self-
 50 supervised play data from a dataset collected through adaptive data collection and heuristic grasping,
 51 ensuring a balanced representation of interaction modes. Building upon ActAIM’s method [12],
 52 our data collection does not use privileged information like reward functions or part segmentation
 53 and enhances realism by employing a complete robot instead of a floating gripper. To summarize,
 54 our contributions are threefold: First, we introduce ActAIM2, a self-supervised learning approach
 55 that enables discrete sampling across different interaction modes. Second, We devise a novel data
 56 collection methodology and have constructed a dataset with diverse interaction modes for model
 57 training. Finally, we thoroughly evaluate our model against a spectrum of generative models and
 58 behavior cloning agents, demonstrating ActAIM2’s superior performance over existing baselines and
 59 the capability of performing the discrete sampling.

60 2 Related Work

61 **Articulated object manipulation** – The manipulation of articulated objects is challenging due to
 62 their complex geometries and kinematics, often understood through partially observed data like
 63 images. Common methods deduce the object’s kinematic structure via passive observation [15,
 64 16, 17] or interactive perception [18, 19]. These techniques are crucial for modeling articulated
 65 objects in robotics and incorporating these models into planning. While imitation learning relies
 66 on expert demonstrations [20, 21], it is limited by the need for extensive and costly data collection.
 67 Conversely, recent research explores obtaining actionable visual priors through direct interaction
 68 with objects [22, 23, 24], focusing on their geometric and semantic properties. Additionally, learning-
 69 based methods utilizing simulation supervised visual learning, and visual affordance learning are
 70 emerging [23, 24, 25]. ActAIM2 addresses these challenges by demonstrating the stability of our
 71 model in the experiments section.

72 **Transformer-based policies** – Many recent works have proposed Transformer-based policies for
 73 robotic control. For behavior cloning, Transformers model conditional action distributions [26, 3]
 74 and encode observations into latent vectors [2]. They also process language commands to specify and
 75 generalize tasks, demonstrating flexibility and robustness [27, 28, 29]. Additionally, Transformers are
 76 utilized to learn latent 3D representations from scenes, building on their success in computer vision
 77 [30, 31, 32], and to process voxelized scenes from RGB images and point clouds [2, 33]. However,
 78 challenges remain, such as high memory demands for high-resolution voxelizations that slow training,
 79 and susceptibility to noisy data when re-rendering views from point clouds [3]. Transformers also
 80 map language and visual inputs directly to robot actions, ranging from basic attention layers to
 81 complex pre-trained models [27, 34]. ActAIM2 builds on this foundation using the RVT model

82 but innovates by integrating a novel discrete interaction mode representation into the skill language
83 prompting.

84 ***Unsupervised Skill Learning*** – Unsupervised skill discovery enables agents to learn distinct behaviors
85 without a reward function but often struggles with inadequate state space exploration when using
86 variational Mutual Information maximization [35, 36]. This limitation hinders its effectiveness for
87 complex tasks. Some strategies counteract this by focusing the learning on smaller, Euclidean spaces
88 to diversify movements [37], although this often restricts the learning to navigation and simple
89 coordinate-based tasks. Alternative approaches propose auxiliary exploration mechanisms and novel
90 training methods to enhance state space exploration [36, 38]. ActAIM2 addresses these issues by
91 introducing the term “interaction mode,” defined by significant visual changes from changes in the
92 degrees of freedom (DoF) of an articulated object. By maximizing the mutual information loss and
93 using contrastive evaluation of visual changes, ActAIM2 encourages the model to identify discrete
94 interaction modes and develop a disentangled latent space for effective sampling.

95 3 Problem Formulation

96 The problem we are solving is how to use a parallel-jaw gripper robot to manipulate various articulated
97 objects and generalize such skills among all types of articulated objects. We adopt a two-phase
98 methodology to enhance a robot’s ability to manipulate articulated objects: data collection and model
99 training, inspired by [12]. We employ a structured predefined action primitive, executing a series
100 of actions across four heuristic phases (initiation, reaching, grasping, manipulation) to gather a
101 comprehensive dataset of observations O_i , including RGBD images and multi-view camera positions,
102 without relying on predefined inputs. Within each action phrase, we define the action a as the keypose
103 ($\mathbf{p}, \mathbf{R}, \mathbf{q}$) of the parallel-jaw gripper similar to [5]. The key pose is defined as the concatenation of
104 the 3 terms, which are $\mathbf{p} \in \mathbb{R}^3$ as the gripper position, $\mathbf{R} \in SO(3)$ as the gripper rotation quaternion,
105 and $\mathbf{q} \in \{0, 1\}$ as the binary parameter indicating whether the gripper is open or close.

106 Our model training aims to uncover the policy’s distribution $\mathbb{P}(a|o)$, with o representing the obser-
107 vation and $a = (\mathbf{p}, \mathbf{R}, \mathbf{q})$ the action, through a decomposition strategy that reconfigures the action
108 distribution as:

$$\mathbb{P}(a|o) = \int_{\text{action predictor}} \underbrace{\mathbb{P}(a|o, \epsilon)}_{\text{mode selector}} \underbrace{\mathbb{P}(\epsilon|o)}_{d\epsilon} \quad (1)$$

109 The mode selector $\mathbb{P}(\epsilon|o)$, contrasting with ActAIM [12]’s Gaussian space, utilizes a mixture of
110 Gaussian distributions to define distinct interaction modes, facilitated by a discrete, latent space $\epsilon \in \mathbb{Z}$
111 for enhanced action prediction and mode selection.

112 4 ActAIM2: Robotic Interaction Mode Discovery

113 We aim to derive the policy $\mathbb{P}(a|o)$ from Equation 1, employing the robot as the agent within an
114 environment populated by various articulated objects. Motivated by [12], our preliminary step was to
115 gather offline, self-supervised data via simulation, which served as the foundation for training the
116 policy $\mathbb{P}(a|o)$. Guided by Equation 1, we dissected the target policy into the action predictor $\mathbb{P}(a|o, \epsilon)$ and mode selector
117 $\mathbb{P}(\epsilon|o)$. For training efficacy, the action predictor $\mathbb{P}(a|o, \epsilon)$ and mode selector
118 $\mathbb{P}(\epsilon|o)$ were pre-trained individually before collectively fine-tuning the overarching pipeline $\mathbb{P}(a|o)$.

119 4.1 Iterative Data Collection

120 We collect trajectory data $T_j = \{(a_i, O_i) | i = 0, 1, 2, 3\}_j$ where O_i are RGBD observations from a
121 configuration of five cameras encircling the articulated object, and $a_i = (\mathbf{p}, \mathbf{R}, \mathbf{q})_i$ represents the key
122 pose and state of the gripper. Inspired by ActAIM [12], we use a similar GMM adaptive method
123 to collect diverse interaction modes. To collect the data with self-supervision, for each trajectory,
124 characterized by the initial observation $O_j^{init} = O_{0j}$ and final observation $O_j^{final} = O_{3j}$, we utilize

125 a pre-trained image encoder \mathcal{E}_O to transform the image observations into a latent vector v . The task
 126 embedding z_j for each trajectory T_j is defined as follows:

$$z_j = v_j^{init} - v_j^{final} = \mathcal{E}_O(O_j^{init}) - \mathcal{E}_O(O_j^{final}) \quad (2)$$

127 To determine the success or failure of a manipulation, we introduce a threshold \bar{z} , defining a trajectory
 128 T_j as successful if $z_j > \bar{z}$. Details of the dataset are presented in Appendix 7.2.

129 4.2 Learning Interaction Modes with Discrete Representations $\mathbb{P}(\epsilon|o)$

130 **Learning Generative Model using GMM Prior** –Based on Equation 1, we aim to identify an effective
 131 mode selector to generate a relevant task latent space $\epsilon \in \mathbb{Z}$. Inspired by the concept of visual
 132 affordances in [39], we learn such latent space encapsulates all conceivable future states using the
 133 conditional Gaussian Mixture Variational Autoencoder (GMVAE) generative model.

134 To learn the mode selector as prior $\mathbb{P}(\epsilon|o)$, we use the data from our collected trajectory and select data
 135 $(O^{init}, O^{final})_j$ for mode selector training. Knowing that z from Equation 2 contains the complete
 136 information of final states given initial observation, We pre-compute the ground-truth label z as our
 137 learning target and regard the initial observation O^{init} as the conditional variable to generate the task
 138 embeddings. We aim to learn a generative model which predicts all possible task embeddings z from
 139 the conditional initial observation O^{init} . Taking inspiration from [40], we construct the generative
 140 model as a Transformer-based Conditional GMVAE, as shown in Figure 2a. We defined our inference
 141 process as $O^{init} \rightarrow (c, y) \rightarrow x \rightarrow \epsilon$ where c is a categorical variable with $p(c) = \text{Multi}(\pi)$ and y is
 142 a Gaussian distribution with $p(y) = \mathcal{N}(0, \mathbf{I})$ which jointly forms a Gaussian Mixture distribution
 143 as the prior. We expect the output learned ϵ would have a similar distribution as the ground-truth
 144 task embedding z since z can be represented as a Mixture of Gaussian distribution fixing the initial
 145 observation (fixing the initial object and object state). Formally, we write the distribution of the
 146 variable under the following process; for simplification, we denote O^{init} as O^i ,

$$p_{\xi, \beta}(\epsilon, x, y, c|O^i) = p(y)p(c)p_{\xi}(x|y, c, O^i)p_{\beta}(\epsilon|x, O^i) \quad (3)$$

$$p_{\xi}(x|y, c, O^i) = \prod_{k=1}^K \mathcal{N}(\mu_{c_k}(y, O^i), \Sigma_{c_k}(y, O^i)) \quad (4)$$

$$p_{\beta}(\epsilon|x, O^i) = \mathcal{N}(\mu_{\beta}(x, O^i), \Sigma_{\beta}(x, O^i)) \quad (5)$$

147 where K represents the number of mixture components, a hyper-parameter within the training regime,
 148 and $\mu_{c_k}, \Sigma_{c_k}, \mu_{\beta}, \Sigma_{\beta}$ denote networks to be trained where ξ and β are the parameters of these
 149 networks. In our implementation, μ_{c_k}, Σ_{c_k} are instantiated as multi-layer ResNet [41] architectures,
 150 and $\mu_{\beta}, \Sigma_{\beta}$ as a transformer with 4 self-attention layers to enhance the stability of reconstruction.
 151 To improve the training stability, we model the categorical distribution $p(c)$ by the Gumbel-Softmax
 152 distribution [42].

153 **Mode Selector Training Loss** –The conditional-GMVAE generative model is optimized using the
 154 variational inference objective, combining reconstruction loss with the log-evidence lower bound
 155 (ELBO) loss. The ELBO loss is expressed as:

$$\mathcal{L}_{ELBO} = \mathbb{E}_q \left[\frac{p_{\xi, \beta}(\epsilon, x, y, c|O^i)}{q(x, y, c|\epsilon, O^i)} \right] \quad (6)$$

156 where the proxy posterior $q(x, y, c|\epsilon, O^i)$ is approximated as $q(x, y, c|\epsilon, O^i) =$
 157 $\prod_i q_{\psi_x}(x|\epsilon, O^i)q_{\psi_y}(y|\epsilon, O^i)q_{\psi_c}(c|x, y, O^i)$, with ψ_x, ψ_y representing the parameters of the
 158 q_{ψ_x} and q_{ψ_c} networks. Based on the decomposition, We refer to the terms in the lower bound as the
 159 reconstruction term, conditional prior term, y -prior term and x -prior term respectively. The prior
 160 term for variables c, y, x is computed as the Kullback-Leibler (KL) divergence, which penalizes the
 161 difference between the learned latent variable distribution and the prior distribution. We formalize
 162 the reconstruction term as the L2 loss between task embedding predictions ϵ and ground-truth z ,
 163 written as $\mathcal{L}_{reconstruct} = \|\epsilon - z\|^2$ in practice. Mathematical details of the conditional GMVAE are
 164 presented in Appendix 7.3.1.

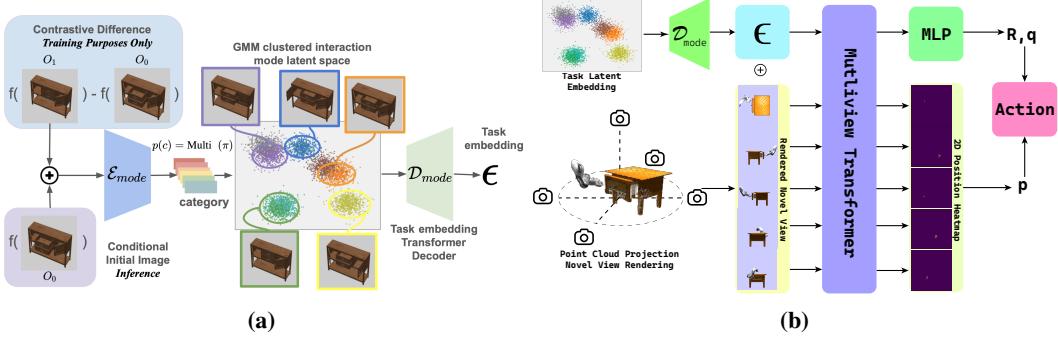


Figure 2: (a) GMM Model Selector The mode selector, a generative model, processes the differences between the initial and final image visual embeddings as generated data, using the initial image embeddings as the conditional variable. **(b) Behavior Cloning Action Predictor** Interaction mode ϵ is sampled from latent space embedding from model selector. 5 Multiview RGBD observations from circled cameras are back-projected and fused into a color point cloud to render novel views. Rendered image tokens and interaction mode token are contacted and fed through a multiview transformer to predict action $a = (\mathbf{p}, \mathbf{R}, \mathbf{q})$.

165 4.3 Supervised Action Predictor Learning $\mathbb{P}(a|o, \epsilon)$

166 Our final objective is to infer a sequence of low-level actions $a = (\mathbf{p}, \mathbf{R}, \mathbf{q})$ from the current
 167 observation O given the predicted task representation ϵ , ensuring the action sequence effectively
 168 accomplishes the articulated object manipulation task while aligning the specific given interaction
 169 mode ϵ . As shown in Figure 2b, The model inputs RGB-D images from encompass multi-view
 170 cameras, the present state of the robot gripper, and the task latent embedding ϵ . From the data
 171 collection phase, we collect successful manipulation trajectories T_j . We decomposed T_j to individual
 172 key-frame actions and observations $(a_i, O_i)_j \in T_j \in D$ as data points for training. To this end, we
 173 propose a multi-view transformer architecture as the basis for our behaviour cloning agent, aimed at
 174 learning the action distribution $\mathbb{P}(a|o, \epsilon)$.

175 **Novel View Rendering and Multiview Transformer** – Building upon the RVT [3] approach, we
 176 utilize novel view rendering from RGB-D multi-view cameras as our visual observation, strategically
 177 positioning five cameras around the robot and articulated objects to generate a merged RGB point
 178 cloud. This cloud is normalized to the scene center and projected onto orthogonal image planes,
 179 creating novel views from the top, front, behind, left, and right, each incorporating RGB color, XYZ
 180 position, and depth channels as model inputs. These rendered views are processed by a multiview
 181 transformer model, where images are patchified and encoded with MLPs and positional encoding,
 182 similar to the ViT process. Additionally, the task embedding ϵ and the current state of the gripper are
 183 encoded and their features are concatenated with image token representations. The transformer then
 184 processes the merged tokens, outputting per-view 2D heatmaps and global features which predict the
 185 current state action $a = (\mathbf{p}, \mathbf{R}, \mathbf{q})$.

186 **Action Predictor Training Loss** – We optimize our action predictor model utilizing a behaviour
 187 cloning loss framework. For the position heatmap, cross-entropy loss is employed, with the ground-
 188 truth image being synthesized from the 3D ground truth point $\hat{\mathbf{p}}$ through projection onto a 2D
 189 orthogonal view, following a Gaussian distribution for spatial smoothing. Similarly, the rotation
 190 heatmap is refined using cross-entropy loss for each Euler angle axis, translating the ground-truth
 191 rotation $\hat{\mathbf{R}}$ into an analogous one-hot vector representation. The binary classification loss, essentially
 192 a cross-entropy loss, updates the gripper jaw’s open-close state \mathbf{q} . Accordingly, the comprehensive
 193 training loss for the action predictor is articulated as:

$$\mathcal{L}_{action} = \mathcal{L}_{\mathbf{p}} + \mathcal{L}_{\mathbf{R}} + \mathcal{L}_{\mathbf{q}} = CE(\mathbf{p}, \hat{\mathbf{p}}) + CE(\mathbf{R}, \hat{\mathbf{R}}) + CE(\mathbf{q}, \hat{\mathbf{q}}) \quad (7)$$

194 4.4 Training Procedure

195 In our training process, we jointly train the policy $\mathbb{P}(a|o)$ as formulation:

$$\mathbb{P}(a|o) = \mathbb{P}(a|o, \text{sg}[\epsilon]) \mathbb{P}(\epsilon|o) \quad (8)$$

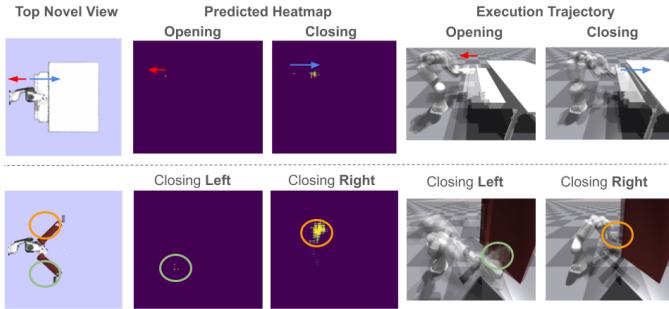


Figure 3: Given different task embedding, we see how action predictor produces actions representing distinct interaction modes. Here, we visualize the camera view and the prediction heatmap from the top for object instances. The first row shows heatmaps for pushing and pulling the handle, while the second row shows heatmaps for closing the left or right door. More qualitative results please see the appendix 7.4

196 In this context, $\text{sg}[\cdot]$ stands for the stop-gradient operator, which stops the flow of partial derivatives
 197 to the next network layers. During the training phase, we include the combined total loss from
 198 both the action predictor (behavior cloning loss) and the mode selector (ELBO loss), represented as
 199 $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{action}} + \mathcal{L}_{\text{ELBO}}$.

200 5 Experiments

201 Our experimental setup is designed to assess how well the proposed method, ActAIM2, performs
 202 in several important areas: 1) ActAIM2 proficiently handles various interaction modes, adapting to
 203 differences across object instances and categories. 2) Utilizing Gaussian Mixture Model (GMM)
 204 based priors facilitates a more structured latent space, enabling targeted searches for specific samples
 205 within the interaction modes. To evaluate ActAIM2, we report the average success rate, sample
 206 success rate, and the average reward achieved during interaction mode grounding iterations.

207 5.1 Experimental setup

208 Leveraging the success of previous studies [43, 2, 12], we use articulated objects from the SAPIEN
 209 dataset [44] for our experiments. Our training dataset includes six categories: faucets, tables, storage
 210 furniture, doors, refrigerators, and switches, while the testing dataset features three new categories:
 211 windows, boxes, and safes, with each category comprising 8 unique objects. We employ IsaacGym
 212 as our simulation platform [45], where we have designed a custom environment featuring a Franka
 213 Emika robot adjacent to an articulated object with 5 cameras circling the object. In our evaluations,
 214 we define a successful trajectory as follows: the model takes an initial observation and iteratively
 215 predicts a series of actions, executed across four steps. After these steps, we evaluate the object’s state.
 216 If the Degrees of Freedom (DoF) of any part of the object has changed by more than 30%, we deem
 217 the trajectory successful. This approach allows us to assess the efficacy of the robot’s manipulation
 218 capabilities in dynamically altering the object’s state based on the model’s iterative predictions.

219 **Baselines And Ablation Study** –We compare our results with the following baselines and ablation
 220 study. **Data Collection:** Adaptive data collection using a GMM-based heuristic grasping method to
 221 sample action sequences. **Where2Act** [43]: Calculates priors for discretized action primitives, using
 222 the complete object point cloud instead of segmented movable points. **ActAIM** [12]: Combines a
 223 Conditional Variational Autoencoder (CVAE) [46] with a transformer-based action predictor [47] to
 224 sample and forecast manipulation trajectories, employing an unrealistically simplified floating gripper.
 225 **Goal-RVT:** A supervised version of ActAIM2 that uses directly provided goal images, bypassing
 226 GMM prior sampling, to feed into the action predictor, demonstrating comparable performance to
 227 other supervised methods. **VQVAE-RVT:** Inspired by Genie [48], replacing the interaction mode
 228 selector with a VQVAE codebook, using a codebook size equal to the GMM cluster number in
 229 ActAIM2, and sampling task embeddings discretely during tests.

230 **Evaluation Metrics** –We conduct two types of evaluations to assess the ability of models for suc-
 231 cessful interaction and discrete sampling, testing across three object categories: Seen Objects (from
 232 the training set), Unseen Instances (same categories as Seen but different instances), and Unseen
 233 Categories (from outside the training set). **Interaction Mode Discovery** metric assesses the average

Table 1: Robotic Interaction Mode Discovery: We evaluate our design decisions through baseline comparisons and ablation studies using the sample-success rate (SSR) metric. We find that that ActAIM2 consistently surpasses competing models across various object types.

Test Set	Seen Objects							Unseen instances							Unseen Cats			
	SSR % ↑	fr.	tt	bb	ll	rr	AVG	fr.	tt	bb	ll	rr	AVG	bb	ll	rr	AVG	
Data Collection	12.8	8.9	9.4	16.9	10.4	8.9	11.2	11.4	9.5	14.0	13.2	11.9	12.9	12.2	20.4	17.3	15.0	17.6
Where2Act [43]	33.3	7.0	7.0	17.9	12.1	4.1	13.6	33.0	13.8	19.2	16.9	13.9	15.4	18.7	15.0	16.8	15.2	15.7
ActAIM [12]	49.3	41.4	36.2	28.6	24.5	19.7	33.3	22.0	38.1	35.5	21.0	18.2	16.2	25.2	38.4	24.1	31.8	31.5
Goal-RVT	58.4	48.9	51.2	72.1	23.2	33.3	47.8	40.2	43.4	39.2	65.1	18.3	25.3	38.6	32.3	23.1	33.9	29.8
ActAIM2	65.3	43.2	52.1	69.2	25.3	36.2	48.6	44.9	41.2	41.5	60.2	20.1	24.4	38.7	34.3	28.9	34.1	32.4

Table 2: Mode Sampling Evaluation ActAIM2 executes actions sampled uniformly across 8 clusters, while Goal-ActAIM and Goal-RVT use manually selected goal images as expert labels for mode-specific sampling. The best performance is highlighted in bold, underscoring ActAIM2’s consistent improvement across evaluations.

Test Set	Seen Objects							Unseen instances							Unseen Cats				
	Algorithm	Mode SSR % ↑	fr.	tt	bb	ll	rr	AVG	fr.	tt	bb	ll	rr	AVG	bb	ll	rr	AVG	
Goal-ActAIM	common mode	25.3	37.5	19.3	62.9	24.3	61.2	38.4	20.3	36.3	18.2	43.0	21.3	31.4	28.4	29.3	15.2	43.5	29.3
	rare mode	24.1	16.2	11.3	28.4	7.8	17.6	17.6	12.4	14.5	9.5	10.5	5.0	13.1	10.8	16.0	7.8	37.5	20.4
Goal-RVT	1st mode goal	42.4	57.4	76.5	75.4	44.9	65.4	60.3	40.2	48.8	72.4	74.3	30.2	50.3	52.7	35.5	35.5	49.2	40.1
	2nd mode goal	17.2	28.3	31.2	70.4	20.1	34.2	33.6	24.3	20.4	15.2	64.2	10.1	20.3	25.8	10.4	5.1	5.1	6.9
VQVAE-RVT	1st mode vector	64.4	44.9	56.3	64.3	27.4	35.2	48.75	45.6	39.3	40.2	51.8	29.1	30.2	39.4	34.2	26.9	36.4	32.5
	2nd mode vector	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
ActAIM2	1st mode cluster	95.4	45.3	83.4	80.4	58.8	65.4	71.5	70.4	45.4	72.4	78.5	40.3	58.9	61.4	43.8	34.9	65.8	48.2
	2nd mode cluster	15.4	24.5	20.4	60.4	10.2	31.2	27.0	10.1	15.8	15.2	58.3	5.0	20.3	20.9	10.4	5.1	5.1	6.9

sample success rate (SSR) defined as the ratio of successful trajectory samples to total samples. For Goal-RVT, expert-labeled goal images are used to ensure viable trajectories, whereas ActAIM2 samples from all 8 clusters to calculate SSR. Moreover, we measure SSR for dominant interaction modes within the selected 2 clusters as **Mode Sampling Evaluation**. This indicates that the evaluated method is required to have at least 2 interaction modes discovered and distributed in separate sample options. For ActAIM2 and VQVAE-RVT, we sample and execute actions from 8 distinct clusters, documenting the success rate of the primary interaction mode in the first and second clusters. Goal-RVT and Goal-ActAIM follow a similar approach but use ground truth goal images instead of sampled clusters, adjusting for changes in DoF.

5.2 Discussion of Results

Interaction Mode Discovery – Table 1 shows that ActAIM2 often outperforms baselines in interaction success rates for both trained and unseen instances. Where2Act and ActAIM use a simplified scenario with a heuristic floating gripper, whereas ActAIM2 utilizes a complete Franka Emika robot setup but still performs a higher success rate. This means ActAIM2 performs well even when it has to discard trajectories due to unreachable states or collisions. This highlights its ability to accurately predict gripper positions and efficiently sample from a learned latent space \mathbb{Z} , surpassing the performance of Goal-RVT’s broader goal image approach. We did not report the VQVAE-RVT results here since VQVAE-RVT outperforms around 5% in each test compared to ActAIM2 averagely. However, despite the average sample success rate, we show that VQVAE-RVT does not meet our requirement for discrete sampling in the following two experiments.

Mode Sampling Evaluation – ActAIM2 and Goal-RVT were evaluated for their ability to distinguish between interaction modes. ActAIM2 proved stable and successfully identified the primary interaction mode across various categories and conditions, matching Goal-RVT’s performance even for a secondary mode. Larger objects like doors and refrigerators consistently showed successful interaction trajectories, demonstrating ActAIM2’s adaptability to different object sizes and complexities. In contrast, VQVAE-RVT struggled to offer a viable second dominant interaction mode. Our tests showed that VQVAE-RVT tends to produce similar action heatmaps across different code vectors, complicating the identification of specific interaction modes. These issues are further explored in our reinforcement learning experiments on interaction mode grounding.

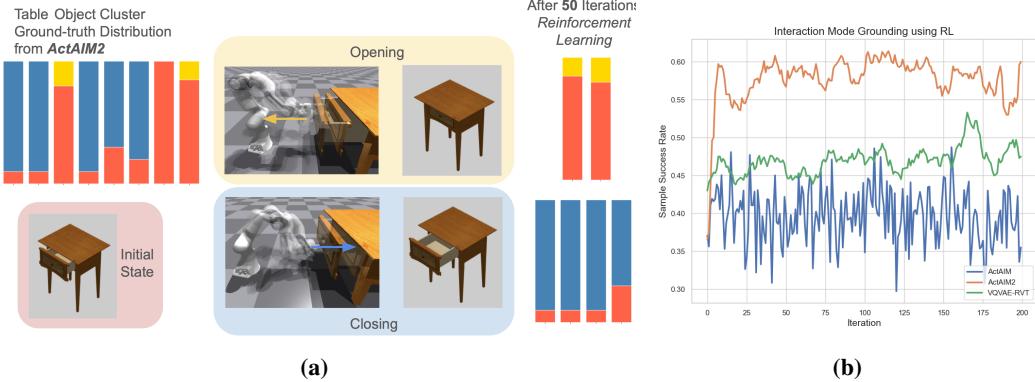


Figure 4: (a) Visualization of Reinforcement Learning on ActAIM2: ActAIM2’s grounding through reinforcement learning is demonstrated, highlighting successful openings (yellow), closings (blue), and failures (red) across eight interaction clusters. After 50 iterations, clusters are accurately identified, enabling consistent trajectory generation.

(b) Reward Optimizing Plot: This plot displays reward optimization over iterations (x-axis) and average reward (y-axis, indicative of success rate for the targeted mode). ActAIM2’s structured latent space shows clear advantages in sampling and convergence efficiency compared to ActAIM [12] and VQVAE-RVT.

263 5.3 Interaction Mode Grounding with ActAIM2 using Reinforcement Learning

264 We further explore ActAIM2’s capability as a prior in reinforcement learning to achieve specified
 265 interaction modes given a goal, iterating through all plausible interaction modes to effectively ground
 266 and distinguish each mode. In our setup, we used an unseen table (ID:20411) as the articulated object
 267 with a sparse reward function defined as $r = 1$ if $|d_i - d_g| < d_\epsilon$ else 0, where d_i is the selected i th
 268 degree of freedom (DoF), d_g is the target DoF, and d_ϵ is the threshold. Reinforcement learning was
 269 employed using ActAIM2 as the prior, framing the task as a Multi-Arm Bandit problem [49], which
 270 narrows the sampling space and enhances learning efficiency. For comparison, we used ActAIM
 271 and RVT-VQVAE as the prior in the baseline experiments, applying the same reward function but
 272 updating the sample task embedding via the cross-entropy method [50].

273 Qualitative and quantitative results are visualized in Figure 4a and Figure 4b. Figure 4a shows
 274 agents refining their understanding of interaction modes through repeated sampling, leading to self-
 275 supervised re-clustering. Figure 4b illustrates that using ActAIM2, the reward—equivalent to the
 276 average sample success rate—smoothly increases, whereas ActAIM struggles to ground the task
 277 embedding to a specific interaction mode representation, and VQVAE-RVT updates slowly due to its
 278 inability to provide distinct interaction mode representations.

279 6 Conclusion

280 ActAIM2 marks a major advancement in self-supervised learning for robotic control, enabling robust
 281 discrete sampling across diverse interaction modes. Our specialized data collection and dataset
 282 lay a solid foundation for future enhancements. Extensive comparisons show ActAIM2’s superior
 283 performance over existing models, enhancing discrete interaction mode learning and strengthening
 284 self-supervised discovery techniques for practical applications.

285 **Limitations.** ActAIM2 presents an exciting self-supervised approach to learning semantically
 286 meaningful priors as interaction modes. Yet, self-supervised discrete representation learning presents
 287 a technical challenge in balancing success rates across different modes, which results in partially
 288 successful rare interaction modes. This could be addressed with an iterative data collection focusing
 289 on less prevalent interaction modes. Further, ActAIM2 assumes interactions can be predefined
 290 as a simple primitive action sequence. Its reliance on self-supervised data collection limits naive
 291 application to complex multi-stage manipulation. We provide a more comprehensive discussion on
 292 ActAIM2’s assumption and future research in Appendix 7.1.

293 **References**

- 294 [1] Z. J. Cui, Y. Wang, N. M. M. Shafiullah, and L. Pinto. From play to policy: Conditional
295 behavior generation from uncurated robot data. In *11th International Conference on Learning
296 Representations (ICLR)*, 2023.
- 297 [2] M. Shridhar, L. Manuelli, and D. Fox. Perceiver-actor: A multi-task transformer for robotic
298 manipulation. In *6th Annual Conference on Robot Learning (CoRL)*, 2022.
- 299 [3] A. Goyal, J. Xu, Y. Guo, V. Blukis, Y.-W. Chao, and D. Fox. Rvt: Robotic view transformer for
300 3d object manipulation. *CoRL*, 2023.
- 301 [4] T. Gervet, Z. Xian, N. Gkanatsios, and K. Fragkiadaki. Act3d: 3d feature field transformers for
302 multi-task robotic manipulation, 2023.
- 303 [5] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark and
304 learning environment. In *IEEE Robotics and Automation Letters (RA-L)*, 2019.
- 305 [6] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-
306 conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and
307 Automation Letters (RA-L)*, 7(3):7327–7334, 2022.
- 308 [7] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan,
309 K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano,
310 K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine,
311 Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet,
312 N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, M. Yan, and A. Zeng.
313 Do as i can and not as i say: Grounding language in robotic affordances. In *arXiv preprint
314 arXiv:2204.01691*, 2022.
- 315 [8] M. Laskin, D. Yarats, H. Liu, K. Lee, A. Zhan, K. Lu, C. Cang, L. Pinto, and P. Abbeel.
316 URLB: Unsupervised reinforcement learning benchmark. In *Thirty-fifth Conference on Neural
317 Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL
318 https://openreview.net/forum?id=lwrPkQP_is.
- 319 [9] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills
320 without a reward function, 2018.
- 321 [10] A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for
322 reinforcement learning, 2020.
- 323 [11] M. Laskin, H. Liu, X. B. Peng, D. Yarats, A. Rajeswaran, and P. Abbeel. Cic: Contrastive
324 intrinsic control for unsupervised skill discovery, 2022.
- 325 [12] L. Wang, N. Dvornik, R. Dubeau, M. Mittal, and A. Garg. Self-supervised learning of action
326 affordances as interaction modes. *arXiv preprint arXiv:2305.17565*, 2023.
- 327 [13] N. M. M. Shafiullah, Z. J. Cui, A. Altanzaya, and L. Pinto. Behavior transformers: Cloning k
328 modes with one stone, 2022.
- 329 [14] S. Lee, Y. Wang, H. Etukuru, H. J. Kim, N. M. M. Shafiullah, and L. Pinto. Behavior generation
330 with latent actions, 2024.
- 331 [15] X. Li, H. Wang, L. Yi, L. J. Guibas, A. L. Abbott, and S. Song. Category-level articulated object
332 pose estimation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
333 (CVPR)*, pages 3703–3712, 2020. doi:[10.1109/CVPR42600.2020.00376](https://doi.org/10.1109/CVPR42600.2020.00376).

- 334 [16] B. AbbateMatteo, S. Tellex, and G. Konidaris. Learning to generalize kinematic models to
335 novel objects. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, *Proceedings of the*
336 *Conference on Robot Learning*, volume 100 of *Proceedings of Machine Learning Research*,
337 pages 1289–1299. PMLR, 30 Oct–01 Nov 2020. URL <https://proceedings.mlr.press/v100/abbatematteo20a.html>.
- 338
- 339 [17] A. Jain, R. Lioutikov, C. Chuck, and S. Niekum. Screwnet: Category-independent articulation
340 model estimation from depth images using screw theory. In *2021 IEEE International Conference*
341 *on Robotics and Automation (ICRA)*, pages 13670–13677, 2021. doi:[10.1109/ICRA48506.2021.9561132](https://doi.org/10.1109/ICRA48506.2021.9561132).
- 342
- 343 [18] R. Martín-Martín, S. Höfer, and O. Brock. An integrated approach to visual perception of
344 articulated objects. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*,
345 pages 5091–5097, 2016. doi:[10.1109/ICRA.2016.7487714](https://doi.org/10.1109/ICRA.2016.7487714).
- 346
- 347 [19] K. Hausman, S. Niekum, S. Osentoski, and G. S. Sukhatme. Active articulation model estimation
348 through interactive perception. In *2015 IEEE International Conference on Robotics and*
Automation (ICRA), pages 3305–3312, 2015. doi:[10.1109/ICRA.2015.7139655](https://doi.org/10.1109/ICRA.2015.7139655).
- 349
- 350 [20] T. Welschehold, C. Dornhege, and W. Burgard. Learning mobile manipulation actions from
351 human demonstrations. In *2017 IEEE/RSJ International Conference on Intelligent Robots and*
Systems (IROS), pages 3196–3201, 2017. doi:[10.1109/IROS.2017.8206152](https://doi.org/10.1109/IROS.2017.8206152).
- 352
- 353 [21] H. Xiong, Q. Li, Y.-C. Chen, H. Bharadhwaj, S. Sinha, and A. Garg. Learning by watching:
354 Physical imitation of manipulation skills from human videos. *2021 IEEE/RSJ International*
Conference on Intelligent Robots and Systems (IROS), pages 7827–7834, 2021. URL <https://api.semanticscholar.org/CorpusID:231632575>.
- 355
- 356 [22] R. Wu, Y. Zhao, K. Mo, Z. Guo, Y. Wang, T. Wu, Q. Fan, X. Chen, L. Guibas, and H. Dong.
357 VAT-mart: Learning visual action trajectory proposals for manipulating 3d ARTiculated objects.
358 In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=iEx3Pi0oLy>.
- 359
- 360 [23] K. Mo, L. Guibas, M. Mukadam, A. Gupta, and S. Tulsiani. Where2act: From pixels to actions
361 for articulated 3d objects. In *2021 IEEE/CVF International Conference on Computer Vision*
362 (*ICCV*), pages 6793–6803, 2021. doi:[10.1109/ICCV48922.2021.00674](https://doi.org/10.1109/ICCV48922.2021.00674).
- 363
- 364 [24] Z. Xu, Z. He, and S. Song. Universal manipulation policy network for articulated objects. *IEEE*
Robotics and Automation Letters, 7(2):2447–2454, 2022. doi:[10.1109/LRA.2022.3142397](https://doi.org/10.1109/LRA.2022.3142397).
- 365
- 366 [25] T. Mu, Z. Ling, F. Xiang, D. C. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su. Maniskill:
367 Generalizable manipulation skill benchmark with large-scale demonstrations. In *Thirty-fifth*
Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round
368 *2)*, 2021.
- 369
- 370 [26] N. M. M. Shafiullah, Z. J. Cui, A. Altanzaya, and L. Pinto. Behavior transformers: Cloning k
371 modes with one stone. In *Thirty-Sixth Conference on Neural Information Processing Systems*
(*NeurIPS*), 2022. URL <https://openreview.net/forum?id=agTr-vRQsa>.
- 372
- 373 [27] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan,
374 K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances.
arXiv preprint arXiv:2204.01691, 2022.
- 375
- 376 [28] S. Nair, E. Mitchell, K. Chen, b. ichter, S. Savarese, and C. Finn. Learning language-conditioned
377 robot behavior from offline data and crowd-sourced annotation. In A. Faust, D. Hsu, and
378 G. Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume 164 of
Proceedings of Machine Learning Research, pages 1303–1315. PMLR, 08–11 Nov 2022. URL
<https://proceedings.mlr.press/v164/nair22a.html>.
- 379

- 380 [29] Y. Zhang and J. Chai. Hierarchical task learning from language instructions with unified
 381 transformers and self-monitoring. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Findings of*
 382 *the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4202–4213, Online,
 383 Aug. 2021. Association for Computational Linguistics. doi:10.18653/v1/2021.findings-acl.368.
 384 URL <https://aclanthology.org/2021.findings-acl.368>.
- 385 [30] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. Bc-z:
 386 Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*,
 387 pages 991–1002. PMLR, 2022.
- 388 [31] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
 389 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. C.
 390 Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath,
 391 I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. S.
 392 Ryoo, G. Salazar, P. R. Sanketi, K. Sayed, J. Singh, S. A. Sontakke, A. Stone, C. Tan, H. Tran,
 393 V. Vanhoucke, S. Vega, Q. H. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich.
 394 Rt-1: Robotics transformer for real-world control at scale. *ArXiv*, abs/2212.06817, 2022. URL
 395 <https://api.semanticscholar.org/CorpusID:254591260>.
- 396 [32] Y. R. Wang, Y. Zhao, H. Xu, S. Eppel, A. Aspuru-Guzik, F. Shkurti, and A. Garg. Mytrans:
 397 Multi-view perception of transparent objects. In *International Conference on Robotics and*
 398 *Automation (ICRA) 2023*, 2023.
- 399 [33] Y. Ze, G. Yan, Y.-H. Wu, A. Macaluso, Y. Ge, J. Ye, N. Hansen, L. E. Li, and X. Wang. Gnfactor:
 400 Multi-task real robot learning with generalizable neural feature fields. In *CoRL2023 Oral*, 2023.
- 401 [34] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
 402 A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman,
 403 A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal,
 404 L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao,
 405 K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut,
 406 H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao,
 407 P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web
 408 knowledge to robotic control, 2023.
- 409 [35] M. Laskin, D. Yarats, H. Liu, K. Lee, A. Zhan, K. Lu, C. Cang, L. Pinto, and P. Abbeel.
 410 URLB: unsupervised reinforcement learning benchmark. *CoRR*, abs/2110.15191, 2021. URL
 411 <https://arxiv.org/abs/2110.15191>.
- 412 [36] D. Strouse, K. Baumli, D. Warde-Farley, V. Mnih, and S. Hansen. Learning more skills through
 413 optimistic exploration. *CoRR*, abs/2107.14226, 2021. URL <https://arxiv.org/abs/2107.14226>.
- 414 [37] S. Park, J. Choi, J. Kim, H. Lee, and G. Kim. Lipschitz-constrained unsupervised skill
 415 discovery. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=BGvt0ghNgA>.
- 416 [38] Z. Jiang, J. Gao, and J. Chen. Unsupervised skill discovery via recurrent skill training. In A. H.
 417 Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing
 418 Systems*, 2022. URL https://openreview.net/forum?id=sYDX_OxNNjh.
- 419 [39] M. Hassanin, S. Khan, and M. Tahtali. Visual affordance and function understanding: A
 420 survey. *ACM Comput. Surv.*, 54(3), apr 2021. ISSN 0360-0300. doi:10.1145/3446370. URL
 421 <https://doi.org/10.1145/3446370>.
- 422 [40] J. Bai, S. Kong, and C. P. Gomes. Gaussian mixture variational autoencoder with contrastive
 423 learning for multi-label classification. In *International Conference on Machine Learning*, pages
 424 1383–1398. PMLR, 2022.

- 427 [41] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016*
 428 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
 429 [doi:10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- 430 [42] I. A. M. Huijben, W. Kool, M. B. Paulus, and R. J. G. van Sloun. A review of the gumbel-max
 431 trick and its extensions for discrete stochasticity in machine learning, 2022.
- 432 [43] K. Mo, L. J. Guibas, M. Mukadam, A. Gupta, and S. Tulsiani. Where2act: From pixels to
 433 actions for articulated 3d objects. In *Proceedings of the IEEE/CVF International Conference on*
 434 *Computer Vision (ICCV)*, pages 6813–6823, October 2021.
- 435 [44] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, L. Yi,
 436 A. X. Chang, L. J. Guibas, and H. Su. SAPIEN: A simulated part-based interactive environment.
 437 In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 438 [45] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox. Gpu-accelerated
 439 robotic simulation for distributed reinforcement learning, 2018.
- 440 [46] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep condi-
 441 tional generative models. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett,
 442 editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates,
 443 Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf.
- 445 [47] A. Jaegle, S. Borgeaud, J.-B. Alayrac, C. Doersch, C. Ionescu, D. Ding, S. Koppula, D. Zoran,
 446 A. Brock, E. Shelhamer, et al. Perceiver io: A general architecture for structured inputs &
 447 outputs. 2021.
- 448 [48] J. Bruce, M. Dennis, A. Edwards, J. Parker-Holder, Y. Shi, E. Hughes, M. Lai, A. Mavalankar,
 449 R. Steigerwald, C. Apps, Y. Aytar, S. Bechtle, F. Behbahani, S. Chan, N. Heess, L. Gonzalez,
 450 S. Osindero, S. Ozair, S. Reed, J. Zhang, K. Zolna, J. Clune, N. de Freitas, S. Singh, and
 451 T. Rocktäschel. Genie: Generative interactive environments, 2024.
- 452 [49] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem.
 453 *Machine Learning*, 47(2):235–256, 2002.
- 454 [50] R. Y. Rubinstein. The cross-entropy method for combinatorial and continuous optimization.
 455 *Methodology and computing in applied probability*, 1(2):127–190, 1999.
- 456 [51] S. James and A. J. Davison. Q-attention: Enabling efficient learning for vision-based robotic
 457 manipulation. In *IEEE Robotics and Automation Letters (RA-L)*, 2022.
- 458 [52] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image
 459 recognition. pages 1–14. Computational and Biological Learning Society, 2015.
- 460 [53] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters
 461 in large spatial databases with noise. In *Proceedings of the Second International Conference on*
 462 *Knowledge Discovery and Data Mining*, KDD’96, page 226–231. AAAI Press, 1996.
- 463 [54] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet-1billion: A large-scale benchmark for
 464 general object grasping. In *2020 IEEE/CVF Conference on Computer Vision and Pattern*
 465 *Recognition (CVPR)*, pages 11441–11450, 2020. [doi:10.1109/CVPR42600.2020.01146](https://doi.org/10.1109/CVPR42600.2020.01146).
- 466 [55] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox. Contact-graspnet: Efficient 6-dof
 467 grasp generation in cluttered scenes. In *International Conference on Robotics and Automation*
 468 (*ICRA*) 2021, 2021.

469 **7 Supplementary**

470 **7.1 Assumptions and Future Research**

471 **7.1.1 Assumptions**

472 During the data collection phase, ActAIM2 operates under the assumption that: 1) the manipulations
473 are straightforward enough to be captured using a limited set of action primitives such as grasping,
474 pushing, or pulling; 2) an interaction mode is identified upon observing significant visual changes;
475 3) the interaction modes can be categorized into a few distinct types. A more detailed discussion of
476 these assumptions is provided below.

477 **Simple Action Space** – We employ a scripted, self-supervised method to collect actions that en-
478 compass diverse interaction modes. The action space is sufficiently simple, focusing primarily on
479 heuristic grasping and random actions. For more complex tasks, such as hammering, washing dishes,
480 or cooking, our current method fails to collect adequate data. Addressing these more intricate tasks
481 would require a more comprehensive and extensive dataset.

482 **Significant Visual Change** – Our data collection is entirely self-supervised, devoid of any expert data
483 or privileged information. We define an interaction as successful if it results in a significant visual
484 alteration to the targeted objects. This approach is effective for articulated objects in our studies, such
485 as doors, windows, or tables, which typically remain stationary except for their movable components.
486 However, challenges arise with objects like tools (e.g., hammers, cups, knives), where it is difficult to
487 discern visual changes either in the tools themselves or the targeted objects (e.g., nails, cup holders,
488 or deformable objects). Especially in tasks requiring repetitive actions, like continuously striking a
489 nail or repeatedly wiping dishes, a more nuanced and generalized method is necessary to determine if
490 meaningful interactions are occurring.

491 **Discrete Interaction Modes** – Articulated objects, by design, often have limited manipulation options.
492 However, when dealing with other objects such as tools, the number of potential interaction modes
493 significantly increases. The functionality of these objects can be diverse; for example, a hammer
494 might be used not only for hammering but also for hooking or reaching. Even the act of grasping
495 these objects presents countless variations, complicating the task of clustering them into discrete
496 modes.

497 **7.1.2 Future Research**

498 Based on the assumptions discussed earlier, we have identified two primary avenues for extending
499 our current research: long-horizon planning tasks and enhancing tool manipulation strategies.

500 **Long-horizon Planning Tasks** – Leveraging the discrete representation of interaction modes provided
501 by ActAIM2, we propose its application to long-horizon planning tasks. Examples of such tasks
502 include sequentially opening a table drawer, locating and opening a box within the drawer, and
503 finally pressing a button inside the box. These tasks illustrate the potential of ActAIM2 to serve as a
504 foundational prior, streamlining the process to discrete searches within complex sequences. To ensure
505 the robustness of our approach, it is crucial that the model accurately predicts all feasible interaction
506 modes based on the given scenario.

507 **Extension to Tool Manipulation Tasks** – Another direction for expansion involves applying our work
508 to tool manipulation. Here, defining the interaction modes for various tools will be pivotal. A robust
509 dataset specifically tailored for tool manipulation is essential to support this endeavour. Additionally,
510 a more sophisticated scene descriptor is required to effectively determine which objects to manipulate
511 and which to designate as targets. This development would facilitate more nuanced and effective tool
512 interactions in automated systems.

513 **7.2 Dataset Generation**

514 **7.2.1 Iterative Data Collection Method**

515 When collecting data, we employ a strategy of random sampling, subsequently filtering successful
 516 actions as determined by our vision model without resorting to any privileged information. Drawing
 517 inspiration from [43, 51], we delineate the task of manipulating articulated objects into four fundamen-
 518 tal poses: initiation, reaching, grasping, and manipulating. Throughout these stages, we capture the
 519 robot's key action poses $a_i = (\mathbf{p}, \mathbf{R}, \mathbf{q})_i$ and RGBD observations O_i from a configuration of five cam-
 520 eras encircling the articulated object. Upon collecting the trajectory $T_j = \{(a_i, O_i) | i = 0, 1, 2, 3\}_j$,
 521 we also archive the initial and final observations, O_j^{init} and O_j^{final} , respectively, captured from the
 522 multi-view cameras with the robot occluded, to facilitate manipulation success evaluation.

523 We introduced our method of identifying successful interacted trajectories, which can be purely from
 524 vision data, specifically the initial and final observation. For each trajectory, characterized by the
 525 initial observation O_j^{init} and final observation O_j^{final} , we utilize a pre-trained image encoder \mathcal{E}_O to
 526 transform the image observations into a latent vector v . The task embedding z_j for each trajectory T_j
 527 is defined as follows:

$$z_j = v_j^{init} - v_j^{final} = \mathcal{E}_O(O_j^{init}) - \mathcal{E}_O(O_j^{final}) \quad (9)$$

528 In our implementation, we employ a pre-trained VGG-19 network [52], without the final fully
 529 connected layers, to serve as our image encoder \mathcal{E}_O . To determine the success or failure of a
 530 manipulation, we introduce a threshold \bar{z} , defining a trajectory T_j as successful if $z_j > \bar{z}$. It is
 531 important to note that this process does not rely on any privileged information. To illustrate the
 532 validity of our method, we define the trajectory's success as a 30% change in the ground-truth DoF
 533 value. The efficacy of this criterion is validated against the ground-truth DoF values, demonstrating a
 534 97.4% accuracy rate across our training and testing dataset. The collected trajectories must exhibit
 535 the diversity of interaction modes of the articulated objects. Thus, we employ three distinct methods
 536 of action sampling, as outlined below. The final dataset is a composite of these three methods.

537 **1. Random Sampling** – We generate play data for manipulation without prior interaction. First,
 538 we select an interaction point $p_1 \in \mathbb{R}^3$ on the articulated object, ensuring it lies within the robot's
 539 workspace. Subsequently, we sample a uniformly random manipulation rotation $\mathbf{R}_0 \in SO(3)$ and
 540 a manipulation position p_2 within the valid area, applying filters to exclude any configurations that
 541 would result in a collision. The robot's initial position p_0 is also determined through random sampling,
 542 which is a specified distance from the interaction point p_1 , ensuring a feasible starting position for the
 543 manipulation task. Based on the previous sampling, we define the randomly sampled action sequence
 544 as $\{(p_0, \mathbf{R}_0, 0), (p_1, \mathbf{R}_0, 0), (p_1, \mathbf{R}_0, 1), (p_2, \mathbf{R}_0, 1)\}$.

545 **2. Heuristic Grasping Sampling** – Heuristic grasping sampling is employed to select interaction
 546 points on the articulated object to enhance the precision of grasping actions. Utilizing the RGB-D
 547 observations, we crop the articulated object and transform it into an RGB point cloud, which under-
 548 goes preprocessing with DBSCAN clustering [53], aimed at identifying segments with significant
 549 geometric features, such as handles or buttons. After clustering, each segment is analyzed by a
 550 pre-trained GraspNet model [54] to generate a set of potential grasps. From this set, grasps with
 551 the highest scores are selected, with the grasp point designated as the interaction point and the grasp
 552 orientation as the gripper rotation for the trajectory. The initial and manipulation poses are determined
 553 using the previously described random sampling approach. This heuristic approach to grasping not
 554 only bolsters the stability of grasp actions but also enriches the dataset with a higher proportion of
 555 complex interaction modes, such as "grasp to open", enhancing the dataset's diversity and utility for
 556 training models to manipulate articulated objects in 'hard' interaction scenarios.

557 **3. GMM-based Adaptive Sampling** – To foster a wide array of interaction modes within our
 558 dataset, we implement GMM-based adaptive sampling inspired by the methodology outlined in [12].
 559 Following the acquisition of M trajectory datasets $\{T_j | j = 1, 2, \dots, M\}$ through random and heuristic
 560 grasping sampling from previous interactions, we compute the task embeddings $\{z_j | j = 1, 2, \dots, M\}$
 561 based on Equation 9. A Gaussian Mixture Model (GMM) prior is constructed from these task

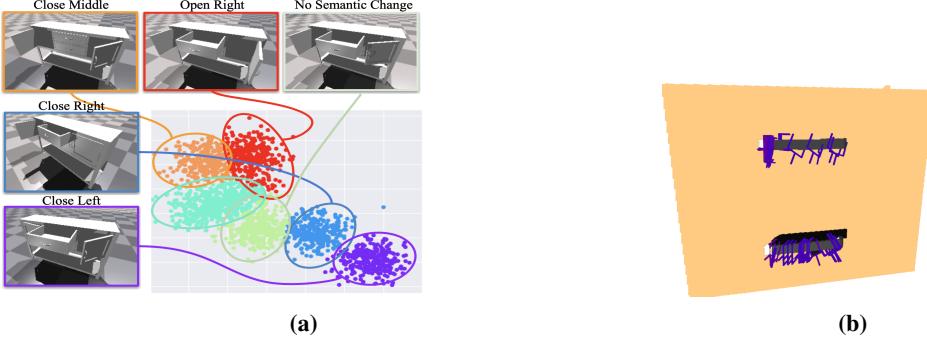


Figure 5: (a) GMM clustering adaptive sampling: his figure illustrates the visualization of using GMM to represent different interaction modes.

(b) Visualization of Heuristic Grasping: We illustrate the proposed grasping using our predefined heuristic with ContactGraspNet [55].

embeddings, denoted as $\mathbb{P}(z|\theta) = \sum_{k=1}^K \kappa_k p(a|\theta_k)$, where θ_k represents the parameters of each Gaussian component within the mixture. The choice of K , the number of clusters, is a hyper-parameter that reflects the presumed number of interaction modes inherent to the object.

Subsequently, we cluster the task embeddings z_j , assigning a unique cluster label to each corresponding trajectory. We found that the task embeddings from different trajectories grouped within the same cluster indicate a similar interaction mode, as they share proximate visual characteristics from initial and final observation. Upon clustering, a new GMM is formulated for each cluster, based on the action sequences, represented as $\mathbb{P}_k(a|\phi) = \sum_{l=1}^L \beta_l p_k(a|\phi_l)$. We then aim to sample an equal number of actions from each cluster, ensuring that the representation of actions—and, by extension, interaction modes—within the dataset are as diverse as possible, thus facilitating a comprehensive exploration of the articulated object’s potential interactions.

Utilizing these sampling methodologies, we concurrently collect data across all articulated objects within our dataset, culminating in a dataset denoted as:

$$D = \{T_j\}_{\text{random}} \cup \{T_j\}_{\text{grasp}} \cup \{T_j\}_{\text{GMM}} \quad (10)$$

$$= \{(a_i, O_i)_j\}_{\text{random}} \cup \{(a_i, O_i)_j\}_{\text{grasp}} \cup \{(a_i, O_i)_j\}_{\text{GMM}} \quad (11)$$

$$= \{(O_i, a_i)\}_{\text{random} \cup \text{grasp} \cup \text{GMM}} \quad (12)$$

After data collection, we enrich each trajectory within our dataset by associating the respective task embedding with the data tuple (O, a) , thereby forming atomic training data instances represented as $(O, a, \epsilon)_j$.

7.2.2 Data Collection Algorithm

The dataset we developed for training purposes is available on our official website. Our dataset was constructed through a combination of random sampling, heuristic grasp sampling, and Gaussian Mixture Model (GMM)-based adaptive sampling, featuring the Franka Emika robot engaging with various articulated objects across multiple interaction modes. It encompasses categories such as faucets, tables, storage furniture, doors, refrigerators, and switches, with 8 unique instances per category. For each instance, we collected 150 trajectories, ensuring comprehensive coverage of the objects’ interaction modes. Objects were scaled to realistic size and initialized in a ‘half-open’ state, denoting a median value for each degree of freedom (DoF). The data collection methodology is detailed in Algorithm 1.

7.3 Model Architecture and Implementation Details

This section outlines the detailed implementation of the model architecture, encompassing both the mode selector and the action predictor components.

Algorithm 1 Data Collection Algorithm

Require: Initial observation O^i , Number of GMM component K , hyper-paramter M for GMM in each cluster

Ensure: All sampled trajectories are filtered successful by evaluating $\epsilon > \bar{\epsilon}$

```

 $D \leftarrow \emptyset$                                 ▷ Set the initial dataset to be empty
while  $D$  not have enough data do
     $D_r = \{(a, o)_i\} \sim \text{RandomSampling}$           ▷ Random Sampling
     $G = \{g_i\} \sim \text{GraspNet}(O^i)$                   ▷ Sample Grasp using GraspNet
     $D_g = \{(a, o)_i\} \sim \text{GenerateTraj}(G)$         ▷ Generate trajectory based on grasp
     $D \leftarrow D \cup D_r \cup D_g$ 
     $\epsilon_i \sim D$                                          ▷ Compute task embedding in current  $D$ 
    Cluster  $\epsilon_i$  with GMM, assign cluster label on each trajectory
     $\{D_j | j = 1, \dots, K\} \leftarrow D$ 
     $D_{GMM} \leftarrow \emptyset$ 
    for  $j$  in range  $K$  do
        Extract  $D_j$  in  $D$  based on cluster label
         $p(D_j | \pi, \mu, \Sigma) = \prod_{n=1}^N \left( \sum_{m=1}^M \pi_m \mathcal{N}(\mathbf{x}_n | \mu_m, \Sigma_m) \right)$       ▷ fit GMM
         $\hat{D}_j \leftarrow \{(a, o)_i\} \sim p(D_j | \pi, \mu, \Sigma)$                                      ▷ Sample action from GMM
         $D_{GMM} \leftarrow D_{GMM} \cup \hat{D}_j$ 
     $D \leftarrow D \cup D_{GMM}$ 

```

591 **7.3.1 Mode Selector Architecture and Implementation Detail**

592 This section revisits the stochastic variables' definitions and distributions, as previously emphasized.
593 The distributions of the model parameters are formalized as follows:

$$p(c) = \text{Multi}(\pi) \quad (13)$$

$$p(y) = \mathcal{N}(0, \mathbf{I}) \quad (14)$$

$$p_{\xi, \beta}(\epsilon, x, y, c | O^i) = p(y)p(c)p_{\xi}(x|y, c, O^i)p_{\beta}(\epsilon|x, O^i) \quad (15)$$

$$p_{\xi}(x|y, c, O^i) = \prod_{k=1}^K \mathcal{N}(\mu_{c_k}(y, O^i), \Sigma_{c_k}(y, O^i)) \quad (16)$$

$$p_{\beta}(\epsilon|x, O^i) = \mathcal{N}(\mu_{\beta}(x, O^i), \Sigma_{\beta}(x, O^i)) \quad (17)$$

594 Here, $\mu_{c_k}, \Sigma_{c_k}, \mu_{\beta}, \Sigma_{\beta}$ are the model parameters to be optimized. Furthermore, we delineate the
595 generative model and compute the inference at test time by defining the posterior as follows:

$$q(x, y, c | \epsilon, O^i) = \prod_i q_{\psi_x}(x | \epsilon, O^i)q_{\psi_y}(y | \epsilon, O^i)q_{\psi_c}(c | x, y, O^i) \quad (18)$$

596 This necessitates the computation of three additional network parameters: $q_{\psi_x}, q_{\psi_y}, q_{\psi_c}$. We then
597 elaborate on deriving the posterior $q_{\psi_c}(c | x, y, O^i)$ for categorical variables c , employing the Gumbel
598 Softmax for the representation of categorical distributions.

599 Notice that c is a categorical parameter that $c \sim \text{Multi}(\pi)$. We defined that $c \in \mathcal{C} = \{c_1, c_2, \dots, c_k\}$
600 and the each class probability is described as $\{\pi_1, \pi_2, \dots, \pi_k\}$. We use the Gumbel Softax trick which
601 provides a simple and efficient way to draw samples c from a categorical distribution with class
602 probabilities $\{\pi_1, \pi_2, \dots, \pi_k\}$. The following form represents the categorical c as,

$$c = \text{one-hot}(\text{argmax}_i[g_i + \log \pi_i]) \quad (19)$$

603 where $\{g_1, g_2, \dots, g_k\}$ are i.i.d samples drawn from $\text{Gumbel}(0,1)$. Assuming that categorical samples
604 c are encoded as k -dimensional one-hot vectors ω lying on the corners of the $(k - 1)$ -dimensional
605 simplex Δ^{k-1} We use the softmax function as a continuous, differentiable approximation to arg max ,
606 and generate k -dimensional sample vectors $\omega \in \Delta^{k-1}$. We defined ω as

$$\omega_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{i=1}^k \exp((\log(\pi_i) + g_i)/\tau)} \quad (20)$$

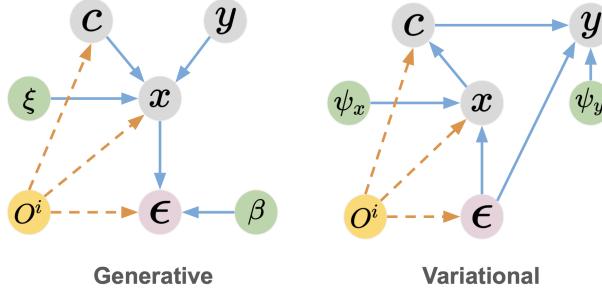


Figure 6: The graphical representations elucidate the Conditional Gaussian Mixture Variational Autoencoder (CGMVAE) framework, showcasing two distinct models: the generative model on the left and the variational family on the right. These graphical models serve to visually communicate the structural and functional relationships between variables within the CGMVAE, illustrating the data generation process and the approximation strategy employed by the variational family to infer latent variable distributions.

607 Where τ is the temperature as the hyperparameter. Therefore, we define the density of the Gumbel-
 608 Softmax distribution as,

$$p(c) = p_{\pi, \tau}(\omega_1, \dots, \omega_k) = \Gamma(k)\tau^{k-1} \left(\sum_{i=1}^k \frac{\pi_i}{\omega_i^\tau} \right)^{-k} \prod_{i=1}^k \frac{\pi_i}{\omega_i^\tau} \quad (21)$$

609 Now, given the representation of the categorical distribution of c from [Equation 21](#), we derive how
 610 we compute the posterior q_{ψ_c} for c . We consider the posterior $q_{\psi_c}(c = c_j | x, y, O^i)$ given $c = c_j$,

$$q_{\psi_c}(c = c_j | x, y, O^i) = \frac{p(c = c_j)p(x|c = c_j, y, O^i)}{\sum_{l=1}^k p(c = c_l)p(x|c = c_l, y, O^i)} \quad (22)$$

$$= \frac{\pi_j p(x|c = c_j, y, O^i)}{\sum_{l=1}^k \pi_l p(x|c = c_l, y, O^i)} \quad (23)$$

611 Therefore, we derive the posterior q_{ψ_c} directly and leave 2 posterior network q_{ψ_x}, q_{ψ_y} to be trained.
 612 Based on the following discussion, we draw the generative model and variational model view as
 613 graphical models in the [Figure 6](#).

614 In the implementation detail, we write parameters $p_\beta = (\mu_\beta, \Sigma_\beta)$ and $p_\xi = (\mu_{c_k}, \Sigma_{c_k})$ to generate
 615 a Gaussian distribution with each representing the mean and variance. We implement the network
 616 $\mu_{c_k}, \Sigma_{c_k}, \psi_x, \psi_y$ with a multi-layer ResNet and implement the network μ_β, Σ_β as a multi-view
 617 transformer since both O^i and ϵ represent multi-view information with the same number on the
 618 channel as the correspondent view number. We show our model μ_β, Σ_β architecture in [Figure 7](#).

619 7.3.2 Mode Selector Training and Inference

620 We illustrate the functionality and application of our mode selector through two distinct plots,
 621 highlighting both the training process and the inference mechanism for task embedding generation.
 622 Figure [Figure 9a](#) depicts the model's operation during training, where it processes the conditional
 623 variable O^i along with the ground truth data ϵ , to accurately reconstruct the task embedding.
 624 Conversely, Figure [Figure 9b](#) demonstrates the inference stage, where the model, requiring only the
 625 initial observation O^i and a discretely sampled cluster (employing an 8-cluster configuration for
 626 implementation), successfully generates the corresponding task embedding ϵ .

627 7.3.3 Action Predictor

628 We provide the architecture of the action predictor which is a joint transformer that takes in task
 629 embedding ϵ and novel view as input. The detailed implementation is shown at [Figure 8](#).

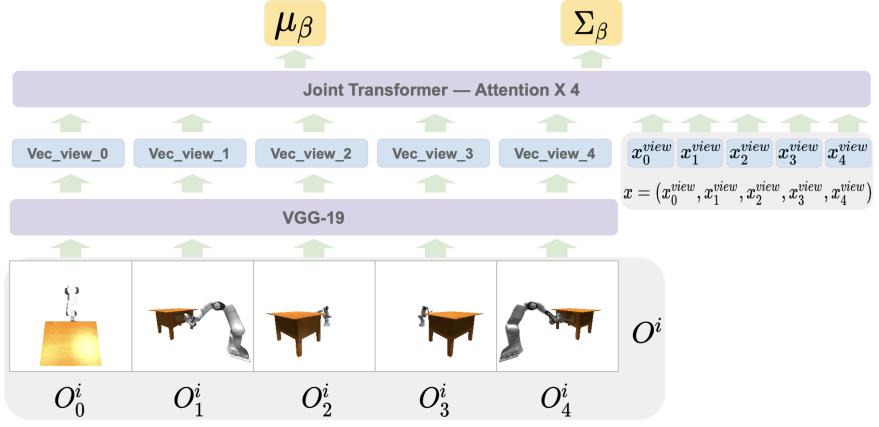


Figure 7: Mode Selector Decoder Architecture: The depicted architecture highlights the functionality of the mode selector decoder, which is designed to process two primary inputs: multi-view RGBD images $O^i = (O_0^i, O_1^i, O_2^i, O_3^i, O_4^i)$, and the Mixture of Gaussian (GMM) variable x . It is important to note that x can be represented as a multi-view feature vector, with our encoding approach preserving the separation of multi-view channels. Initially, the multi-view RGBD images are passed through a pre-trained VGG-19 image encoder to extract feature vectors for each view. Subsequently, these feature vectors, along with the GMM variable x , are inputted into a joint transformer. This transformer, featuring four attention layers, is tasked with producing the means and variances associated with the reconstructed task embedding ϵ .

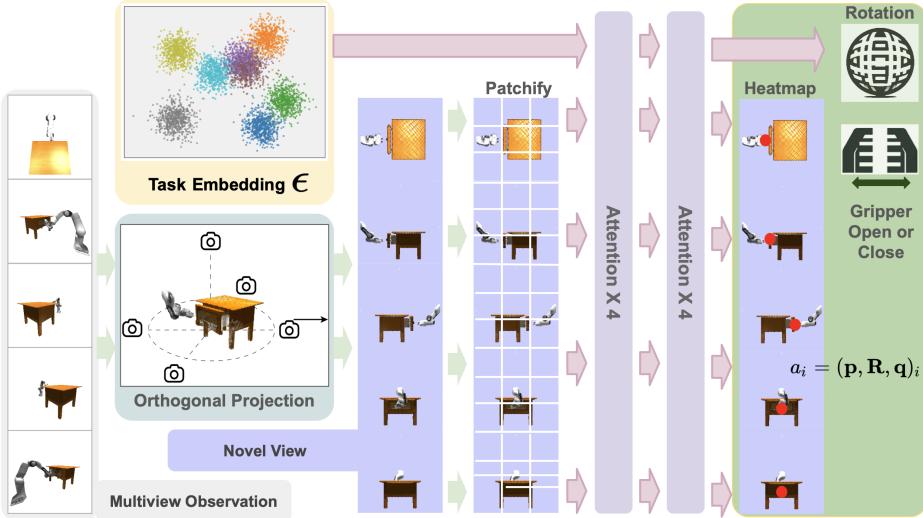
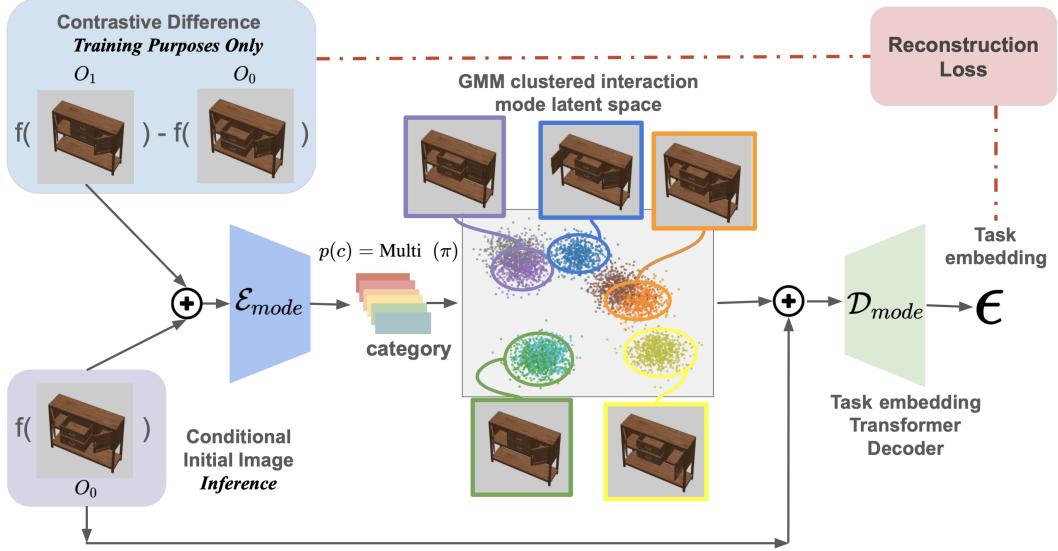


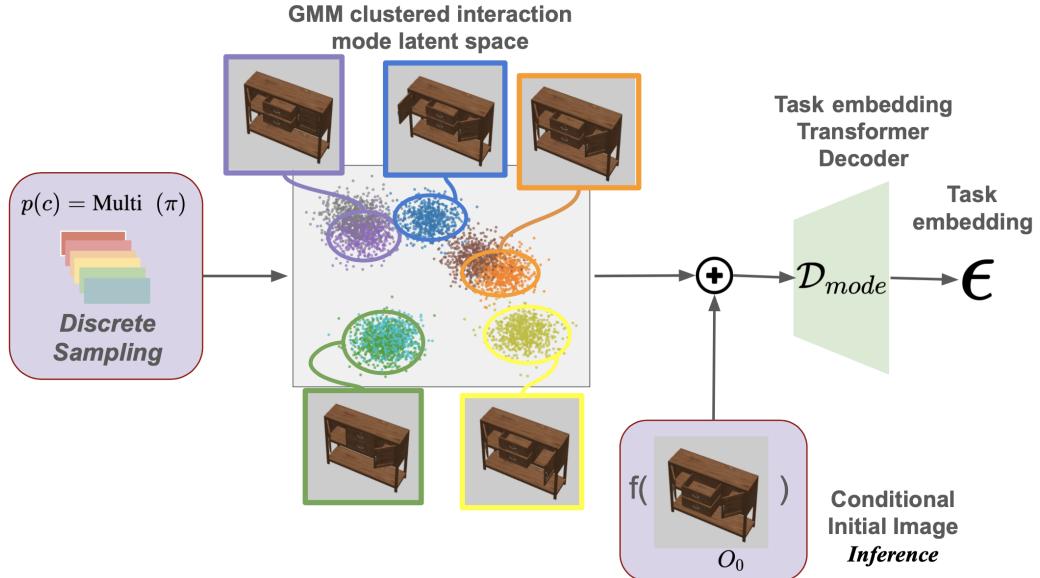
Figure 8: Action Predictor Architecture: This model integrates multi-view observations directly as input, sourced from predefined cameras within the scene. The process begins with the extraction of five RGBD images, which are subsequently transformed into RGB point clouds. These are then subject to orthogonal projection to generate five novel view images. Subsequently, these novel views are partitioned into smaller patches and fed into a joint transformer. This transformer, characterized by four attention layers, integrates the sampled task embedding derived from a Mixture of Gaussian distribution. The architecture of the joint transformer encompasses eight attention layers, culminating in the production of a heatmap. This heatmap delineates the action’s translation, the discretized rotation, and a binary variable indicating the gripper’s state—open or closed.

630 7.4 More Qualitative Results

631 We supplement our presentation with additional qualitative results, further elucidating the model’s
 632 proficiency in learning the disentanglement of interaction modes. Initially, we demonstrate the
 633 efficacy of the mode selector through a t-SNE plot. This choice of visualization is motivated by our
 634 methodology of training the mode selector and action predictor independently, allowing for a focused
 635 examination of the mode selector’s performance.



(a) Training Process of the Mode Selector: This figure illustrates the training procedure of the mode selector, mirroring the approach of a conditional generative model. It highlights the contrastive analysis between the initial and final observations—the latter serving as the ground truth for task embedding—to delineate generated data against the backdrop of encoded initial images as the conditional variable. The process involves inputting both the generated task embedding data and the conditional variable into a 4-layer Residual network-based mode encoder, which then predicts the categorical variable c . Following the Gaussian Mixture Variational Autoencoder (GMVAE) methodology, the Gaussian Mixture Model (GMM) variable x is computed and introduced alongside the conditional variable to the task embedding transformer decoder. This model is tasked with predicting the reconstructed task embedding, sampled from the Gaussian distribution as outlined in the architecture of the mode selector decoder, and calculating the reconstruction loss against the input ground truth data.



(b) Inference Process: In the inference phase, the agent discretely samples a cluster from the trained Gaussian Mixture Variational Autoencoder (GMVAE) model to calculate the Mixture of Gaussian variable x . This variable x , in conjunction with the conditional variable (initial image observation), is then inputted into the mode selector transformer decoder. The objective is to reconstruct the task embedding for inference, effectively translating the conditional information and sampled cluster into actionable embeddings.

636 Subsequently, we extend our qualitative analysis with figures akin to those presented in the main
637 paper, offering a comprehensive view of the model's capabilities. These additional figures serve to
638 reinforce the insights gained from the initial results, showcasing the model's nuanced understanding
639 of interaction modes through the distinct visual representations of the data.

640 **7.4.1 Mode selector TSNE plot Figure 14**

Utilizing our pre-trained Conditional Gaussian Mixture Variational Autoencoder (CGMVAE) mode selector, we conduct disentanglement learning visualization on our comprehensive dataset. Specifically, we focus on the "single drawer" object (object ID: 20411), employing the mode selector to delineate the generated clusters and compare them with the ground truth task embeddings. The data for this visualization is derived from our dataset, and we calculate the task embedding ϵ_j for each data point as the difference between the initial and final object states, represented by

$$\epsilon_j = v_j^{init} - v_j^{final} = \mathcal{E}_O(O_j^{init}) - \mathcal{E}_O(O_j^{final})$$

641 .

642 Subsequently, we employ a t-SNE plot to simultaneously visualize the ground truth and generated
643 task embeddings. In this visualization, distinct colors within the ground truth plot indicate data points
644 originating from different interaction modes. Similarly, varied colors in the generated plot correspond
645 to data points arising from disparate clusters within the Mixture of Gaussians model. Through this
646 approach, we demonstrate that:

- 647 1. The ground truth task embeddings ϵ are distinctly clustered based on the interaction modes.
- 648 2. The CGMVAE model effectively generates clusters that categorize data points by their respective
649 categories c .
- 650 3. The reconstructed data closely aligns with the ground truth data points, with the majority of the
651 clustered data encompassed within the respective ground truth clusters.

652 This visualization underscores the efficacy of our generative model mode selector in extracting task
653 embeddings for further application in the action predictor, highlighting the model's capability to
654 discern and categorize interaction modes accurately.

655 **7.4.2 Action Predictor Qualitative Results**

656 We present extensive qualitative results in [Figure 15a](#), [Figure 15b](#), [Figure 16a](#), and [Figure 16b](#),
657 demonstrating the model's ability to predict distinct interaction modes through discrete sampling. For
658 each object, we explore three different clusters, each representing a unique interaction mode. The
659 initial state of the robot and the articulated object is depicted from three perspectives: top-down, front,
660 and side views. The heatmaps, derived from the top view during manipulation steps, highlight the
661 variance in action space corresponding to different sampled interaction modes. Subsequent imagery
662 illustrates the robot's movement within the simulator and the outcome following interaction with the
663 articulated objects. It is important to note that comprehensive **video demonstrations** accompany this
664 document and are accessible on our website, <https://actaim2.github.io/>.

665 **7.4.3 Comparison of ActAIM2 and VQVAE-RVT**

666 Inspired by the Genie [48] approach, we have compared our ActAIM2 with VQVAE-RVT to assess
667 the efficacy of these models in discerning discrete interaction modes in robotic manipulation tasks.
668 Our primary objective was to evaluate the distinction between interaction modes using a simplified
669 scenario, a single-drawer table, which naturally exhibits two distinct interaction modes: opening and
670 closing.

671 In our experiments, we visualized the latent spaces generated by both ActAIM2 and VQVAE-RVT.
672 Particularly for VQVAE-RVT, the latent space visualization involved examining the distribution of
673 eight code vectors. As depicted in Figure 10, these vectors clustered into two categories, which ideally
674 should correspond to the two expected interaction modes of the drawer. This clustering pattern was
675 anticipated and desired as it suggests a clear demarcation between the distinct modes of interaction.

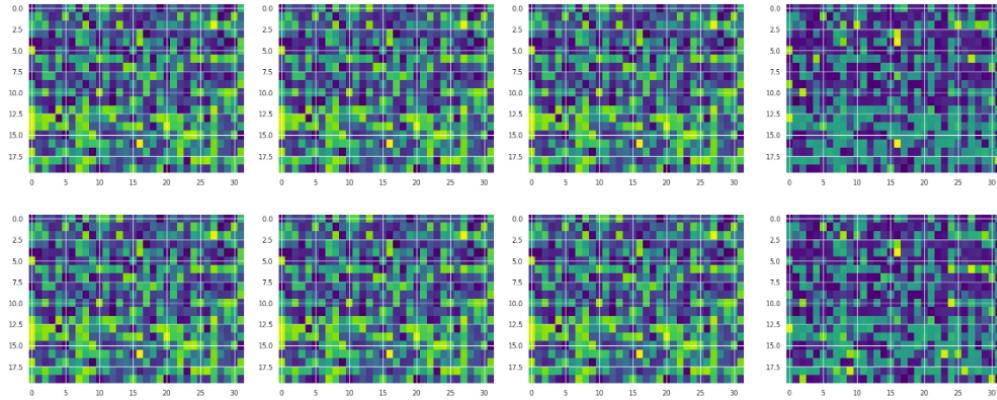


Figure 10: Visualization of Latent Space Clustering in VQVAE-RVT: This figure illustrates the distribution of eight code vectors within the latent space, categorized into two distinct clusters. These clusters are intended to represent the discrete interaction modes of opening and closing a drawer. The spatial arrangement highlights the expected separation of code vectors, symbolizing the potential for mode-specific action mapping in robotic manipulation tasks. Despite this apparent clustering, subsequent heatmaps (see Figure 11) reveal a lack of diversity in the action predictions, undermining the practical utility of this model configuration.

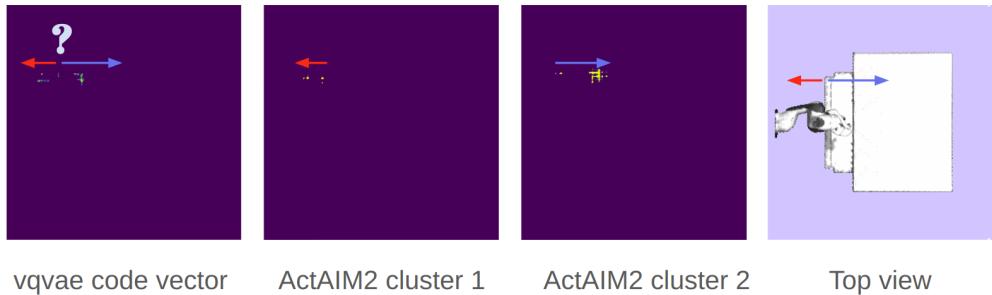


Figure 11: Comparative Visualization of Action Heatmaps and Observational Data From left to right: (1) VQVAE-RVT action heatmap synthesized using all eight code vectors, showing identical outcomes across the board, indicating a failure to differentiate interaction modes. (2) Action heatmap generated by \algoName when sampling from one cluster, demonstrating a specific interaction mode. (3) Action heatmap from \algoName when sampling from a different cluster, showcasing another distinct mode of interaction. (4) Top-view observation of the drawer, correlating with the spatial contexts of the heatmaps, providing a visual reference for the interaction zones mapped by the heatmaps. This series highlights \algoName's capability to discern and represent distinct action strategies through targeted cluster sampling.

676 However, subsequent visualizations raised concerns about the practical efficacy of the VQVAE-RVT
 677 model in our application context. When we explored the heatmaps generated by the VQVAE-RVT
 678 model, we observed a critical limitation: all 8 code vectors produced essentially the same heatmap,
 679 despite their differing positions in the latent space. This heatmap, illustrated in Figure [Y], consistently
 680 depicted all plausible interaction modes for the drawer, regardless of the specific code vector used.
 681 This outcome was in stark contrast to the results from ActAIM2, where distinct heatmaps clearly
 682 indicated specific interaction actions like pushing or pulling, depending on the sampled cluster within
 683 the latent space.

684 These findings led us to conclude that merely replacing the GMVAE component with a VQVAE in
 685 the setup did not achieve the desired disentanglement of interaction modes. The VQVAE-RVT model
 686 failed to map the code vectors to unique, mode-specific interaction strategies, instead converging on a
 687 generalized representation that was not useful for distinguishing between the actionable options of
 688 opening and closing the drawer. Consequently, ActAIM2's ability to discriminate between distinct
 689 interaction modes via cluster-specific sampling proves superior in contexts demanding discrete and
 690 distinguishable action representations.

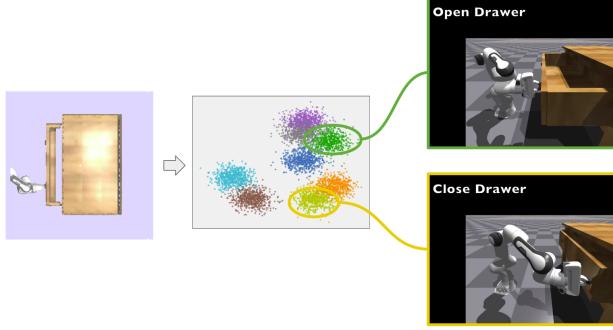


Figure 12: Opening and Closing a Drawer: This figure demonstrates the effective action sequence generated by ActAIM2 for a drawer. The left part of the image shows the drawer being opened, showcasing the robot's approach and grip adjustment. The right part of the image captures the drawer in a fully closed position, illustrating the final state after the action sequence execution.

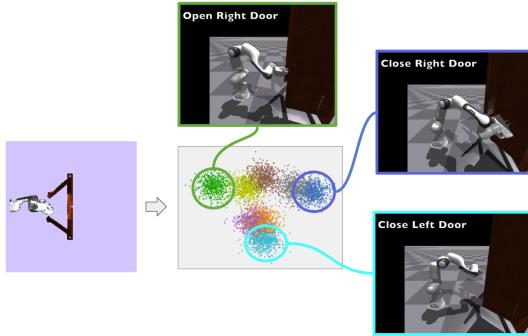


Figure 13: Opening and Closing a Door: This figure illustrates the ActAIM2's manipulation capability with a door. The left image displays the door being opened, highlighting the robot's positioning and the initial interaction phase. The right image shows the door completely closed, detailing the end of the manipulation sequence and the effectiveness of the action predictor.

691 8 Generation of Demonstration Videos

692 To illustrate the practical applications and effectiveness of ActAIM2, we generated demonstration
 693 videos by employing its inference mechanism. The process involves several key steps:

- 694 1. **Generative Mode Selection:** Initially, observations are inputted into the generative mode selector
 695 of ActAIM2. This component is responsible for reconstructing the task's latent space, which is
 696 modeled as a Mixture of Gaussians. This structure enables discrete sampling of clusters, which
 697 represent distinct interaction modes that the robotic system can execute.
- 698 2. **Sampling and Action Prediction:** From the reconstructed latent space, we sample the task em-
 699 beddings by selecting a cluster within the Gaussian Mixture Model (GMM) and its corresponding
 700 Gaussian distribution. This sampled task embedding is then forwarded to the action predictor. The
 701 action predictor generates the specific actions needed to interact with the environment effectively.
- 702 3. **Simulation and Recording:** As depicted in Figure 12 and Figure 13, ActAIM2 reconstructs
 703 an object-based GMM and samples different task embeddings. Depending on the sampled task
 704 embedding, different interactions are reconstructed and executed within a simulator. We recorded
 705 the manipulation processes, which are detailed in the video provided in the supplementary files.
 706 Each video showcases how ActAIM2 navigates through different interaction scenarios, reflecting
 707 the diverse capabilities of the model in real-time applications.

708 This comprehensive demonstration not only validates the functionality of ActAIM2 but also provides
 709 a visual understanding of its potential in diverse robotic manipulation tasks. The videos highlight the
 710 nuanced interactions achievable through targeted sampling within the model's structured latent space.

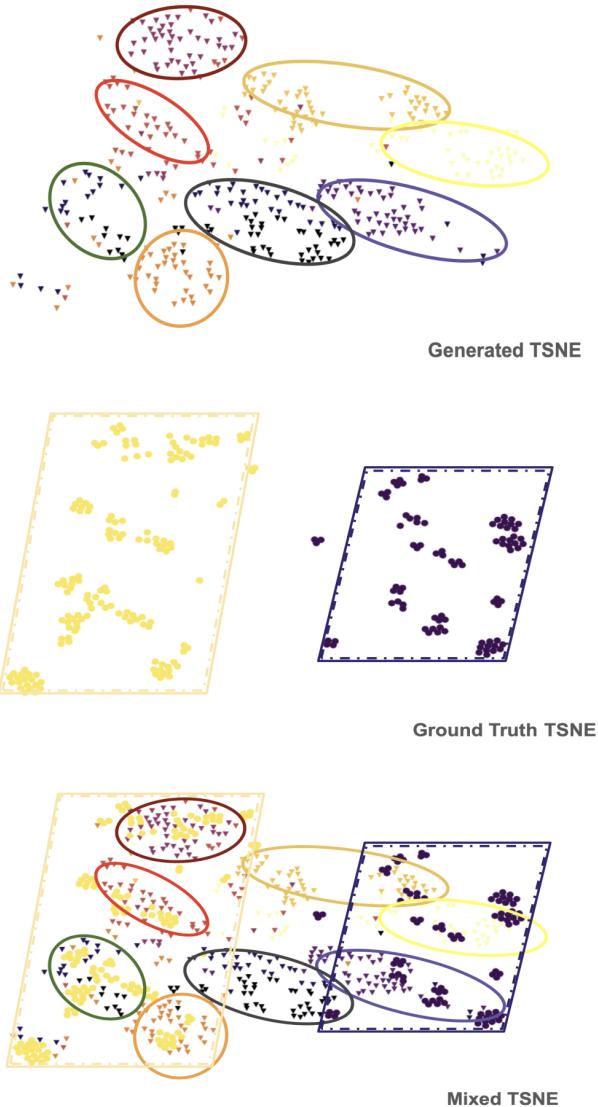
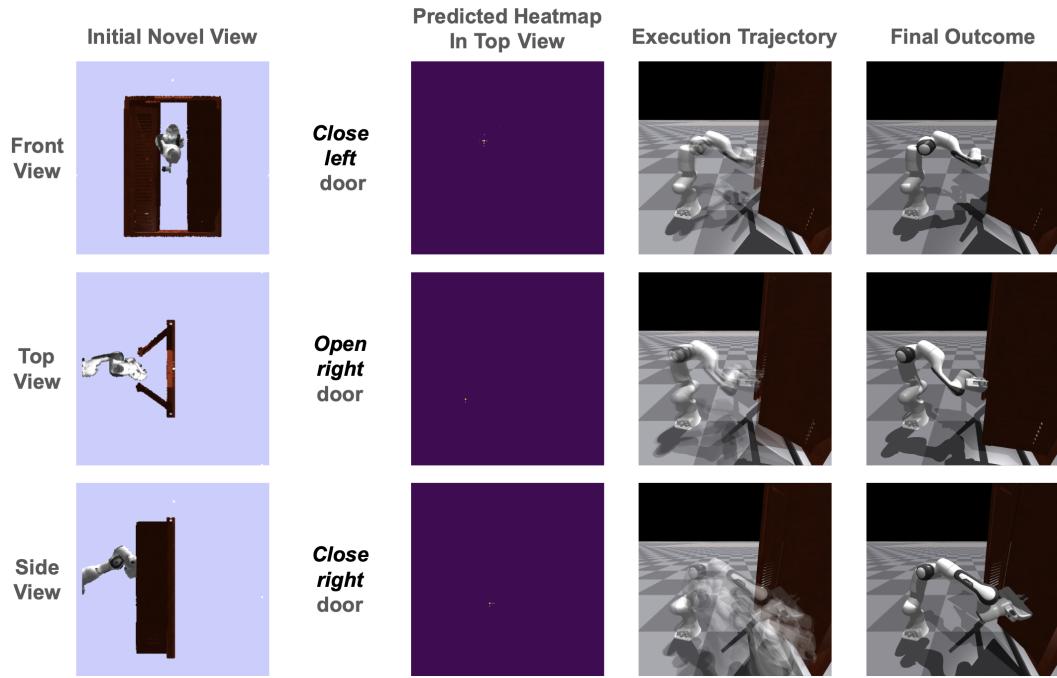
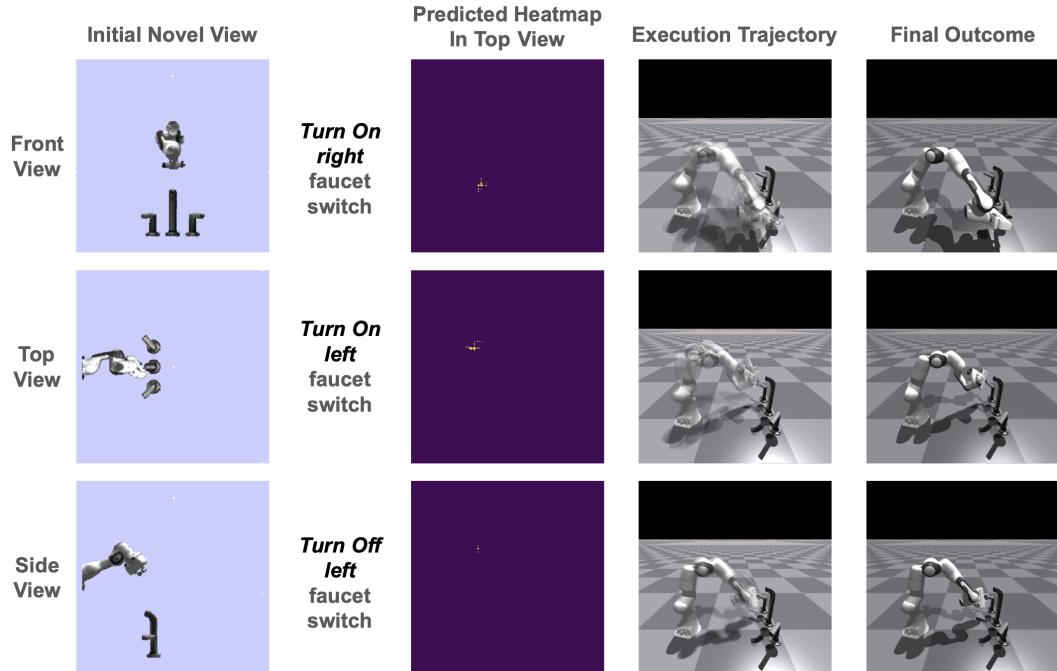


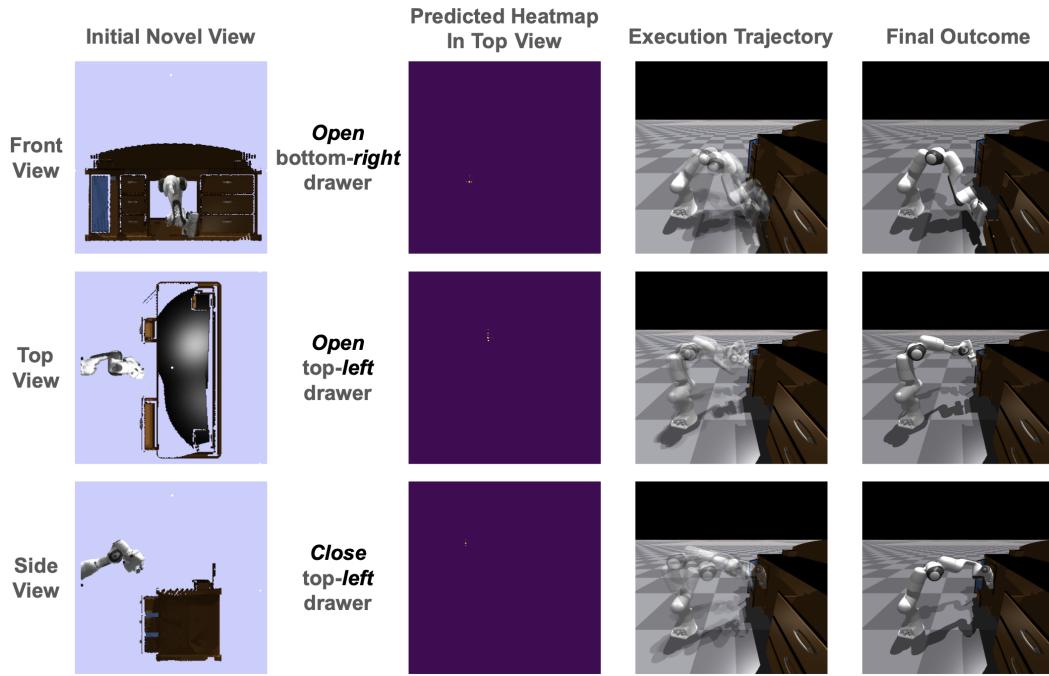
Figure 14: Disentanglement Visualization with CGMVAE: This figure illustrates the efficacy of the Conditional Gaussian Mixture Variational Autoencoder (CGMVAE) in disentangling interaction modes for the "single drawer" object (ID: 20411), using a t-SNE plot for visualization. Task embeddings ϵ_j , defined by the variance between initial and final object states, are visualized in distinct colors to denote various interaction modes and clusters. The sequence of figures demonstrates the CGMVAE's precision in clustering and aligning data points with their respective interaction modes: (1) Generated clusters from the CGMVAE mode selector reveal distinct groupings. (2) Ground truth task embeddings confirm the model's capacity for accurate interaction mode classification. (3) A combined visualization underscores the alignment between generated clusters and ground truth, showcasing the model's ability to consistently categorize tasks within identical interaction modes.



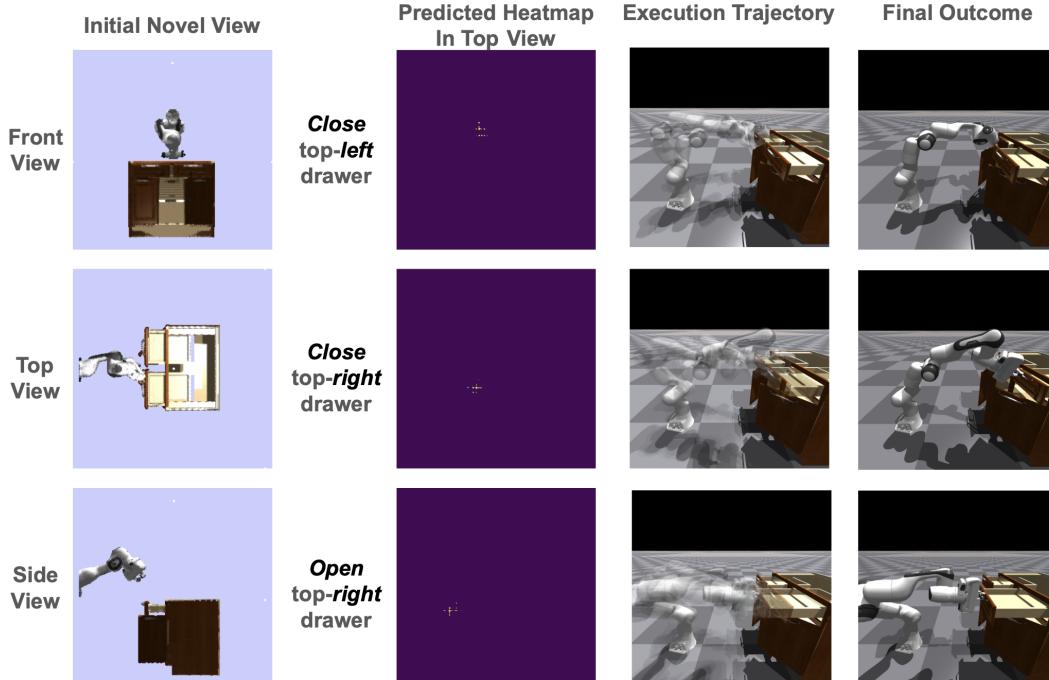
(a) Door, Object ID: 8961



(b) Faucet, Object ID: 154



(a) Table, Object ID: 19898



(b) Table, Object ID: 41083