
RDD: Retrieval-Based Demonstration Decomposer for Planner Alignment in Long-Horizon Tasks

Mingxuan Yan¹ Yuping Wang^{1,2} Zechun Liu³ Jiachen Li^{1*}

¹University of California, Riverside ²University of Michigan ³Meta AI
{myan035, yuping.wang, jiachen.li}@ucr.edu zechunliu@meta.com

Abstract

To tackle long-horizon tasks, recent hierarchical vision-language-action (VLAs) frameworks employ vision-language model (VLM)-based planners to decompose complex manipulation tasks into simpler sub-tasks that low-level visuomotor policies can handle. Typically, the VLM planner needs finetuning to learn to decompose a new task, which requires target task demonstrations segmented into sub-tasks by either human annotation or heuristic rules. However, without prior knowledge, the heuristic sub-tasks can deviate significantly from the visuomotor policy’s training data, thereby degrading task performance. To address these issues, we propose a **Retrieval-based Demonstration Decomposer (RDD)** that automatically decomposes video demonstrations into sub-tasks with prior by aligning the visual features of the decomposed sub-task intervals with those from the training data of the low-level visuomotor policies. RDD outperforms the state-of-the-art sub-task decomposer on both simulation and real-world tasks, demonstrating robustness across diverse settings. Code and more results are available at [rdd-neurips.github.io](https://github.com/rdd-neurips).

1 Introduction

Developing generalist robots that are capable of executing complex, long-horizon tasks in unstructured environments has become one of the central goals of current robotics research. Traditional robotic programming and learning methods often struggle with the variability and complexity inherent in real-world scenarios. Building upon the success of Vision-Language Models (VLMs) and Large Language Models (LLMs), a new class of multi-modal foundation models known as Vision-Language-Action models (VLAs) [1, 2, 3, 4, 5] has emerged specifically for embodied AI applications. As recent studies [6, 7, 8, 9, 10, 11, 12, 13] have shown, integrating high-level planners above the low-level visuomotor policies vastly improves the performance for long-horizon robotic tasks. This has led to the hierarchical VLA paradigm [14, 15, 13, 16, 17, 18, 19, 20]. The planner, often a powerful VLM, performs task planning and reasoning to break down complex tasks into simpler sub-tasks with step-by-step language instructions. A learning-based visuomotor policy, trained on datasets with short-horizon sub-tasks and conditioned on the generated sub-task instructions, performs precise manipulation to complete the sub-tasks one by one, thereby completing long-horizon tasks.

Despite its versatility, a vanilla VLM planner typically needs to be finetuned with human demonstrations when deploying to a given task [18, 14, 16]. To build the dataset for planner finetuning, demonstrations are temporally decomposed to sub-tasks by human annotation [14, 16, 18, 19, 15] or heuristics [13, 15, 21, 22, 23, 24, 25]. However, these methods are neither scalable nor efficient, and, most importantly, they could generate sub-tasks that deviate significantly from the training data of the low-level visuomotor policy. Figure 1 illustrates this dilemma. The state-of-the-art sub-task decomposer UVD [25], which uses a heuristic decomposition rule based on visual feature change-point

*Corresponding author

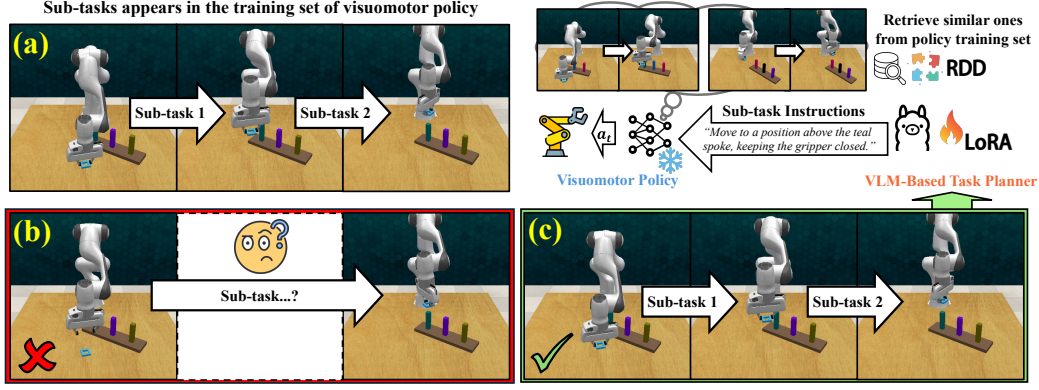


Figure 1: The core idea of RDD. (a) Two sub-tasks appear in the visuomotor policy’s training set, on which the policy has been optimized. (b) Existing sub-task decomposers, such as UVD [25], use heuristic decomposition rules and may generate “unfamiliar” sub-tasks that are difficult to handle for the low-level visuomotor policy. (c) In contrast, RDD decomposes the demonstration into sub-tasks that are visually similar to the ones in the training set of the visuomotor policy. The sub-tasks are then used to finetune the high-level planner, which gives sub-task instructions to the low-level visuomotor policy and guides it to finish the task step-by-step.

detection, generates sub-tasks that significantly deviate from the training data of the visuomotor policy. Finetuning the planner with these sub-tasks could make the planner generate sub-task instructions that the visuomotor policy is not optimized for, leading to compromised performance.

This gap motivates us to develop an automatic, training-free, and computationally efficient approach that *identifies sub-tasks from a video demonstration with prior*, i.e., the decomposed sub-tasks should be aligned with the training data of the low-level visuomotor policies. To achieve this, we propose a **Retrieval-based Demonstration Decomposer (RDD)** that decomposes the demonstration into sub-tasks visually similar to the ones in the training set of the visuomotor policy, as illustrated in Figure 1 (c). Inspired by previous work [25], we employ existing visual encoders [26, 27, 28, 29, 30] that encode images into a compact latent space where distance metrics (e.g., angular distance) are effective in describing the semantic relationship between images. To align the sub-tasks to the training data of the low-level visuomotor policy, we build a sub-task visual feature vector database with the visuomotor training set and design an effective sub-task similarity measure to ensure similar sub-task samples can be efficiently retrieved. We formulate sub-task identification as an optimal partitioning problem and employ a dynamic programming-based solver to optimize the sub-task partitioning strategy efficiently. The experiments show that RDD consistently outperforms state-of-the-art methods on both simulation and real-world benchmarks.

The main contributions of this paper are as follows:

- This work is the first to coordinate the high-level planner and low-level visuomotor policy in the hierarchical VLA framework by generating the planner’s finetuning dataset that is well aligned with the visuomotor policy to improve the long-horizon task performance.
- We propose RDD, a retrieval-based video sub-task identification algorithm with sub-task prior. Specifically, we model sub-task identification as an optimal partitioning problem, which can be solved efficiently with a dynamic programming solver.
- We evaluate RDD on both simulation and real-world benchmarks. Experimental results show that RDD outperforms the state-of-the-art heuristic decomposer and is robust across various settings.

2 Related Work

Hierarchical VLAs. While single-stage VLAs [1, 2, 3, 4, 5] achieve promising performance in short-horizon manipulation tasks, long-horizon tasks need an in-depth understanding of the task and general planning ability, which is hard to handle by a single-stage model. To this end, hierarchical structures have emerged as a compelling solution for long-horizon manipulation tasks [14, 15, 13, 16, 17, 18, 19, 20, 10]. As representative examples, Hi Robot [14] and $\pi_{0.5}$ [18] enhance their previous work on visuomotor policy [4, 3] with a VLM-based planner. According to image observation and the overall task goal, the planner provides sub-task instructions at each time step. The low-level

policy, conditioned on the instruction, outputs the final actions. Hierarchical structures also enable error correction and human intervention [13, 16, 14]. However, these methods rely on either human annotation or heuristic rules to identify sub-tasks when finetuning the planner, which is less efficient and could generate sub-tasks that are hard to handle by the visuomotor policy.

Sub-Task Identification. Finetuning the high-level planner in hierarchical VLAs requires demonstrations broken down into sub-tasks with associated labels. Manually performing this segmentation [14, 16, 18, 19, 15] is slow and expensive. Human subjectivity also leads to inconsistencies. Heuristic methods [13, 15, 21, 22, 23, 24], such as segmenting based on contact changes or end-effector velocity profiles, require task-specific knowledge for carefully designed rules. In contrast, UVD [25] leverages general visual representation and identifies sub-tasks by detecting frame-by-frame temporal change points of visual embedding distances. However, when applying to hierarchical VLAs, UVD can still sub-optimally decompose sub-tasks, which may deviate significantly from the training data of the visuomotor policy. In contrast, with the sub-task prior, RDD identifies sub-tasks by explicitly aligning the sub-tasks with the training set of the visuomotor policy, enabling seamless coordination between the planner and visuomotor policy.

Visual Representations. Considerable efforts have been made to develop visual encoders that embed RGB frames into compact latent vector spaces [26, 27, 28, 29, 30]. Some of these efforts are specially designed for robotics and manipulation scenarios. For instance, R3M [27] uses time-contrastive learning on large datasets of human videos; LIV [26] learns a value function conditioned on both language instructions and images. These visual representations are designed to capture meaningful information about the scene, objects, and potentially their relationships or temporal dynamics.

3 Retrieval-Based Demonstration Decomposer (RDD)

3.1 Problem Statement

Visuomotor-Planner Dataset Alignment. Hierarchical VLAs typically follow an imitation learning framework that trains a low-level visuomotor policy $\pi_\theta(a_t|s_t, o_t, l_t, L)$ and a high-level planner $p_\phi(l_t|s_t, o_t, l_{t-1}, L)$. The latter is usually a VLM. a_t denotes the waypoint action at timestep t , including 6-DoF pose and binary gripper state. Both policy π_θ and planner p_ϕ are conditioned on the RGB image observation o_t , proprioceptive states s_t , and the overall task objective description L in natural language, such as “put the cube in the drawer”. The policy π_θ is additionally conditioned on a sub-task instruction l_t like “first, pick up the cube”, which is determined by the planner p_ϕ at time t .

During the policy training phase, the raw training dataset $\mathcal{D}^{\text{train}} = \{(\mathcal{S}^i, L^i)\}_{i=1}^{N_{\text{train}}}$ is composed of N_{train} demonstrations where $\mathcal{S}^i = \{(a_t^i, s_t^i, o_t^i)\}_{t=1}^{T_i}$ and L^i represents the corresponding task objective description. To break the complex long-horizon tasks down to simple instructions required by the low-level policy π_θ , a demonstration \mathcal{S}^i is decomposed into a set of partitions $P^i = \{I_j^i\}_{j=1}^{B_i}$ based on task-specific rules or human annotations. The j -th interval $\mathcal{I}_j^i = \{\mathcal{S}^i[b_j^i], \dots, \mathcal{S}^i[e_j^i]\}$ ($b_j^i < e_j^i$) corresponds to a single coherent sub-task, where b_j^i, e_j^i are indexes of the starting and ending frames. All time steps t within the same interval share the same sub-task instruction $l_t^i = f_{\text{lang}}(\text{prompt}_j)$ labeled manually or generated by a powerful language model. As such, the demonstration is augmented with language descriptions l_t^i to $\mathcal{S}_{\text{aug}}^i = \{(a_t^i, s_t^i, o_t^i, l_t^i)\}_{t=1}^{T_i}$ and the augmented training set is denoted as $\mathcal{D}_{\text{aug}}^{\text{train}} = \{(\mathcal{S}_{\text{aug}}^i, L^i)\}_{i=1}^{N_{\text{train}}}$. The policy $\pi_\theta(a_t|s_t, o_t, l_t, L)$ is then optimized on $\mathcal{D}_{\text{aug}}^{\text{train}}$.

During the high-level planner finetuning phase, given M demonstrations ($M \ll N_{\text{train}}$) for each task, we construct a planner finetuning dataset $\mathcal{D}^{\text{demo}} = \{(\mathcal{S}^i, L^i)\}_{i=1}^M$ and predict the sub-task partitioning strategy $P \in \Pi(\mathcal{S}^i)$ for \mathcal{S}^i , where $\Pi(\mathcal{S})$ denotes all possible partitioning over a sequence \mathcal{S} :

$$\Pi(\mathcal{S}) = \left\{ P = \{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_K\} \mid \bigcup_{i=1}^K \mathcal{I}_i = \mathcal{S}, \mathcal{I}_i \cap \mathcal{I}_j = \emptyset \text{ for } i \neq j \right\}.$$

Sub-task Identification as Optimal Partitioning Problem. Finding the optimal sub-task partitioning strategy can be formulated as an optimal partitioning problem, as illustrated in Figure 2:

$$P^{i*} = \arg \max_{P \in \Pi(\mathcal{S}^i)} J(P), \quad (3.1)$$

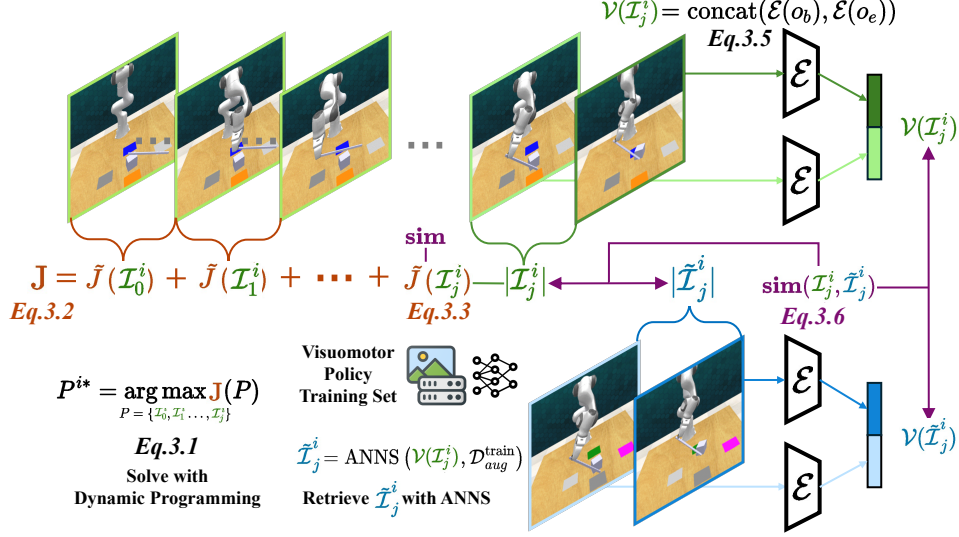


Figure 2: RDD formulates sub-task identification as an optimal partitioning problem. Intervals colored in green are proposed segments of the demonstration \mathcal{S}^i , and ones colored in blue are retrieved from the visuomotor policy’s training set $\mathcal{D}_{\text{aug}}^{\text{train}}$.

where $J(P)$ is the partitioning strategy scoring function defined on P that evaluates how close the strategy is to the low-level visuomotor policy’s training dataset $\mathcal{D}_{\text{aug}}^{\text{train}}$. Given the partitioning found, $\mathcal{D}^{\text{demo}}$ is augmented by f_{lang} and arranged to $\mathcal{D}_{\text{aug}}^{\text{demo}}$ following the same procedure as $\mathcal{D}_{\text{aug}}^{\text{train}}$. A pre-trained planner $p_\phi(l_t|s_t, o_t, l_{t-1}, L)$ is then finetuned on $\mathcal{D}_{\text{aug}}^{\text{demo}}$ with supervised learning to learn to decompose the new task.

3.2 Dynamic Programming Solver

Brute-force search of P^{i*} requires $O(2^{N-1})$ times of evaluation of J for a N frame’s demonstration, which is computationally intractable. Fortunately, [31] show that when J is interval-wise additive (as illustrated in Figure 2), i.e:

$$J(P) = \sum_{\mathcal{I} \in P} \tilde{J}(\mathcal{I}), \quad (3.2)$$

which implies $J(P) = J(P_1) + J(P_2)$, $(P_1, P_2) \in \{(P_1, P_2) | P = P_1 \cup P_2, P_1 \cap P_2 = \emptyset\}$, where \tilde{J} is the scoring function of a single interval. The following optimality holds:

Theorem 3.1 (Principle of Optimality [31]). *Given an additive scoring function J , any subset P' of an optimal partition P^* is the optimal partitioning strategy of the intervals it covers.*

This implies that if we find the partial optimal partitioning strategy for $\mathcal{S}^i[0 : j]$, it must be a subset of the global optimal P^{i*} . This optimality structure allows a dynamic programming algorithm [31] to find the optimal partition with $O(N^2)$ evaluations of the interval scoring function \tilde{J} .

In real-world robot learning scenarios, the duration of a sub-task is limited (typically tens of seconds) [32, 33, 34, 14], thus the complexity of the algorithm can be further improved by ignoring intervals excessively long. We show that: *if the length of the interval is bounded, the complexity can be further reduced to $O(N)$.* We provide the algorithm implementation in Appendix A.1, Algorithm 1, and draw the following conclusion:

Corollary 3.1.1. *If the length of every interval is in the range $[L_{\min}, L_{\max}]$, $0 < L_{\min} < L_{\max} \leq N$, Algorithm 1 finds the optimum with $O((L_{\max} - L_{\min}) \cdot \max(L_{\max} - L_{\min}, N - L_{\max}))$ evaluations of the interval scoring function \tilde{J} .*

We defer the proof to Appendix A.2. When the maximum sub-task interval length L_{\max} is bounded, which is common in robotics learning scenarios, a linear complexity $O(N)$ is achieved. Considering general cases, in this work, we make no assumption on L_{\max} and only mildly assume $L_{\min} = 2$ for

sanity (a valid interval must have both the starting and ending frame). We additionally remark that Algorithm 1 supports parallel evaluation of the scoring function, as the intervals to be evaluated are determined at the beginning.

Interval Scoring Function. Recall that \tilde{J} should reflect how well the proposed interval aligns with the intervals in the training set $\mathcal{D}_{\text{aug}}^{\text{train}}$, we define the interval scoring function \tilde{J} as:

Definition 3.1. *The scoring function \tilde{J} for an interval \mathcal{I} is defined as:*

$$\tilde{J}(\mathcal{I}_j^i) = |\mathcal{I}_j^i| \mathbf{sim}(\mathcal{I}_j^i, \text{ANNS}(\mathcal{V}(\mathcal{I}_j^i), \mathcal{D}_{\text{aug}}^{\text{train}})) = |\mathcal{I}_j^i| \mathbf{sim}(\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i), \quad (3.3)$$

where \mathcal{V} maps interval \mathcal{I} into a d -dimensional vector representation, $|\mathcal{I}|$ is the duration of the \mathcal{I} , and $\text{ANNS}(\mathcal{I}, \mathcal{D}_{\text{aug}}^{\text{train}})$ represents the approximate nearest neighbor of the interval proposal \mathcal{I} in the training set $\mathcal{D}_{\text{aug}}^{\text{train}}$ under some distance metric δ in \mathbb{R}^d . \mathbf{sim} is an interval similarity measure. For simplicity, we denote the result of approximate nearest neighbor search for \mathcal{I}_j^i as $\tilde{\mathcal{I}}_j^i$.

Eq. 3.3 essentially evaluates how close the proposed interval is to the training set of the visuomotor policy in the training set $\mathcal{D}^{\text{train}}$. Moreover, Def. 3.1 ensures the following notable property:

Proposition 3.1. *Suppose an interval \mathcal{I}_j^i can be split into K consecutive parts $\{\mathcal{I}_{j1}^i, \mathcal{I}_{j2}^i, \dots, \mathcal{I}_{jK}^i\}$, all of which have the same training set similarity score, i.e., $\mathbf{sim}(\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i) = \mathbf{sim}(\mathcal{I}_{j1}^i, \tilde{\mathcal{I}}_{j1}^i) = \dots = \mathbf{sim}(\mathcal{I}_{jK}^i, \tilde{\mathcal{I}}_{jK}^i)$. Given the interval scoring function \tilde{J} of Eq. 3.3, and an additive J , the following equality holds:*

$$J(\{\mathcal{I}_j^i\}) = J(\{\mathcal{I}_{j1}^i, \mathcal{I}_{j2}^i, \dots, \mathcal{I}_{jK}^i\}). \quad (3.4)$$

The proof is in Appendix B. This equality implies that J is ignorant of the number of intervals when evaluating nested partitionings with the same similarity score. An alternative way to interpret is that, in Eq. 3.3, \mathbf{sim} assigns scores to the sub-task assignment of each timestamp in an interval instead of assigning to the interval as a whole, thus the score summation is irrelevant to the number of intervals in the partitioning strategy.

3.3 Interval Similarity and Overall Objective

Interval Similarity Measures. As introduced in Section 2, one can embed the RGB image observation o_t^i into a compact latent vector space for similarity measures. We define \mathcal{V} as:

$$\mathcal{V}(\mathcal{I}) = \text{concat}(\mathcal{E}(o_b), \mathcal{E}(o_e)). \quad (3.5)$$

As illustrated in Figure 2, o_b, o_e are image observations at the beginning and end of \mathcal{I} , and \mathcal{E} is the embedding function. This formulation is inspired by former studies [26, 35, 25] that the ending frame (i.e., the goal frame) contains rich information about the sub-task goal and thus can be a distinguishable representation. Eq. 3.5 also includes the starting frame, which is essentially the goal state of the previous sub-task, to aggregate context-related information into the vector representation.

Let the approximate nearest neighbor of \mathcal{I}_j^i be $\tilde{\mathcal{I}}_j^i = \text{ANNS}(\mathcal{I}_j^i, \mathcal{D}_{\text{aug}}^{\text{train}})$ we define the similarity measure \mathbf{sim} between $\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i$ as:

$$\mathbf{sim}(\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i) = - \left[\delta(\mathcal{V}(\mathcal{I}_j^i), \mathcal{V}(\tilde{\mathcal{I}}_j^i)) + \alpha \left| 1 - \frac{|\mathcal{I}_j^i|}{|\tilde{\mathcal{I}}_j^i|} \right| \right], \quad (3.6)$$

where the first term is the distance between the vector representations of \mathcal{I}_j^i and $\tilde{\mathcal{I}}_j^i$; the second evaluates the relative difference between the temporal durations of two intervals. α is a hyperparameter that controls the weights between temporal and visual similarity.

Considering OOD Sub-tasks. While the primary objective of RDD is to align the planner with the visuomotor policy’s existing capabilities, in real-world applications, out-of-distribution (OOD) sub-tasks not learned by the low-level visuomotor may exist. In such scenarios, the objective changes to: *aligning sub-task intervals to both existing visuomotor sub-tasks and general sub-tasks*, and the newly identified sub-tasks will be used to finetune both the visuomotor and the planner. Firstly, to detect the existence of new sub-tasks in demonstrations, one can quantify the novelty of a demonstration by $\Delta = \frac{1}{|P|} \sum_{\mathcal{I} \in P} \tilde{J}(\mathcal{I})$, the average similarity score of the optimal partition P found by the standard

RDD algorithm. A low value of Δ indicates a low averaged similarity, which signals novel sub-tasks. An alternate interval similarity measure **sim** for the OOD setting is defined:

$$\mathbf{sim}(\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i) = \underbrace{-\delta(\mathcal{V}_e(\mathcal{I}_j^i), \mathcal{V}_e(\tilde{\mathcal{I}}_j^i))}_{\text{retrieval}} + \underbrace{\beta G(\mathcal{I}_j^i)}_{\text{general}}, \quad (3.7)$$

where $\mathcal{V}_e(\mathcal{I}) = \mathcal{E}(o_e)$ and only the ending frame is used to calculate the semantic distance due to unpredictable OOD sub-task durations; G evaluates how well a proposed interval aligns with “general” sub-tasks. The hyperparameter β balances the trade-off between aligning with visuomotor sub-tasks and discovering novel, generalizable sub-tasks. G can be implemented using heuristic general sub-task identification functions like UVD [25] to measure how well an interval conforms to generic change-point detection heuristics:

$$G(\mathcal{I}) = -\frac{1}{|\mathcal{I}|} \mathbf{abs}(b - \mathbf{UVD}(e, \mathcal{I})), \quad (3.8)$$

where b, e represent the index of the beginning and ending frame of interval \mathcal{I} . $\mathbf{UVD}(e, \mathcal{I})$ gives the index of the UVD predicted beginning frame, given the goal frame on e .

Approximate Nearest Neighbor Search. Considering the vast number of intervals in $\mathcal{D}_{\text{aug}}^{\text{train}}$ and the high-dimensional vector space, we adopt approximate nearest neighbor search (ANNS) to implement the nearest neighbor searcher. We choose the popular random-projection-trees-based method Annoy [36] as the ANNS implementation, which is computationally efficient and shows good robustness on various data [37]. RDD can also work with GPU-accelerated ANNS libraries like FAISS [38] for further acceleration.

Overall Optimization Objective. By substituting Eq. 3.2 and Eq. 3.3 into Eq. 3.1, we have the complete definition of the optimization problem as:

$$P^{i*} = \arg \max_{P \in \Pi(\mathcal{S}^i)} \sum_{\mathcal{I}_j^i \in P} |\mathcal{I}_j^i| \mathbf{sim}(\mathcal{I}_j^i, \tilde{\mathcal{I}}_j^i), \quad (3.9)$$

where **sim** is defined by Eq. 3.6 or alternatively Eq. 3.7 for OOD settings. The optimal partitioning strategy P^{i*} of demonstration \mathcal{S}^i can be solved by Algorithm 1.

4 Experiments

Implementation and Parameter Settings. We adopt RACER [13] as the base hierarchical VLA framework, which uses RVT [39] as the low-level visuomotor policy π_θ and the recent LLaVa-based VLM llama3-llava-next-8B [40] as the pre-trained base model for planner p_ϕ . We use the pre-trained RVT policy π_θ provided by RACER [13] trained $\mathcal{D}_{\text{aug}}^{\text{train}}$ and the validation set of RL Bench (labeled with the same decomposition rule as in $\mathcal{D}_{\text{aug}}^{\text{train}}$). During the deployment phase, the planner is finetuned for two epochs on $\mathcal{D}_{\text{aug}}^{\text{demo}}$ using LoRA [41], with the rank of 128 and a scaling factor of 256 following RACER. The finetuning process takes about 5 minutes with 4 NVIDIA 6000 Ada GPUs. For base parameter settings, we set the weighting factor $\alpha = 1$ and interval similarity measure **sim** in Eq. 3.6 for non-OOD scenarios, and use LIV [26] as the visual encoder \mathcal{E} that is specifically designed for manipulation tasks. We use Gemini-1.5-flash [42] to generate sub-task language instructions for proposed intervals in $\mathcal{D}_{\text{aug}}^{\text{demo}}$.

Visuomotor Policy Training Dataset and Vector Database. We evaluate RDD on the RL Bench [32] robot manipulation benchmark. The visuomotor policy training set $\mathcal{D}_{\text{aug}}^{\text{train}}$ is adapted from [13]. $\mathcal{D}^{\text{train}}$ originally consists of 1908 teleoperated demonstrations from the RL Bench’s training set. When generating $\mathcal{D}_{\text{aug}}^{\text{train}}$, RACER additionally augmented it with heuristic failure-recovery samples, resulting in a training dataset with 10,159 demonstrations. In this work, we only use the original 1908 demonstrations to exclude interference. Demonstrations are decomposed into 12700 sub-task intervals using a task-specific heuristic decomposer based on motion and gripper states. Generally, the decomposer will mark a goal state of a sub-task whenever: 1) the gripper state closes or opens, 2) the arm stops for a pre-defined duration, and 3) the end of the demonstration. More details about this heuristic can be found in Section III.B of [13]; RACER uses GPT-4-turbo as the language labeling function f_{lang} to annotate the sub-task intervals, given the language descriptions of the robot movement and initial environment setup.

Table 1: Multi-task success rates (%) on RLBench.

Method	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Close Jar	Install Bulb	Meat off Grill	Open Drawer	Place Wine	Push Buttons
w/o Finetune	52.6 \pm 8.2	4.5 \pm 1.2	27.6 \pm 26.4	34.8 \pm 14.2	46.4 \pm 26.8	95.6 \pm 6.1	83.2 \pm 13.0	54.8 \pm 9.1
Uniform	71.3 \pm 5.4	3.1 \pm 1.2	46.4 \pm 29.9	51.2 \pm 19.2	76.4 \pm 22.4	100.0 \pm 0.0	80.8 \pm 14.5	82.0 \pm 7.8
UVD	71.4 \pm 5.1	3.0 \pm 1.3	44.0 \pm 28.7	54.8 \pm 20.0	85.2 \pm 20.6	100.0 \pm 0.0	80.8 \pm 15.3	67.2 \pm 13.6
RDD (Ours)	74.9 \pm 6.9	2.2 \pm 0.9	46.0 \pm 28.2	52.8 \pm 16.4	84.4 \pm 21.1	99.2 \pm 2.4	86.4 \pm 15.4	84.0 \pm 7.8
<i>Expert</i>	75.1 \pm 4.7	2.2 \pm 1.0	50.4 \pm 33.1	50.4 \pm 13.3	94.4 \pm 9.7	99.2 \pm 2.4	81.6 \pm 15.0	85.6 \pm 6.0
Method	Put in Cupboard	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Sweep to Dustpan	Turn Tap	
w/o Finetune	41.2 \pm 20.1	36.4 \pm 28.8	58.8 \pm 23.3	36.0 \pm 21.8	57.2 \pm 14.9	22.8 \pm 32.5	89.2 \pm 13.4	
Uniform	36.8 \pm 15.4	98.0 \pm 2.7	92.4 \pm 10.8	64.8 \pm 16.7	64.4 \pm 9.9	34.8 \pm 37.7	98.8 \pm 3.6	
UVD	35.2 \pm 12.1	90.4 \pm 8.6	96.8 \pm 6.6	74.4 \pm 29.2	66.8 \pm 21.2	43.6 \pm 24.6	89.6 \pm 11.1	
RDD (Ours)	41.2 \pm 17.1	97.2 \pm 3.1	98.4 \pm 3.2	68.0 \pm 25.0	65.2 \pm 14.3	57.2 \pm 29.7	94.0 \pm 5.1	
<i>Expert</i>	39.6 \pm 15.6	91.2 \pm 7.3	97.6 \pm 5.1	75.2 \pm 24.6	66.4 \pm 22.0	48.8 \pm 35.5	96.0 \pm 5.7	

Given $\mathcal{D}_{\text{aug}}^{\text{train}}$, we build a vector database following Eq. 3.5 and employ Annoy [36] as the ANNS algorithm to retrieve the approximate nearest neighbor. For each frame, to exclude the inference of occlusion, we concatenate the representation vectors of the front-view and gripper-view images into one. We apply the same configuration to UVD for fair comparison. For Annoy, we set the number of random-projection trees to 10 and let the searcher search through all trees at runtime. We empirically find that the choices of the ANNS algorithm or search parameters have a minor impact on the performance. We use angular distance as the distance measure δ , which is written as $\sqrt{2(1 - \cos(u, v))}$ for normalized vectors u, v . The finetuning dataset $\mathcal{D}_{\text{aug}}^{\text{demo}}$ is built on RLBench’s validation set following the same procedure except that the decomposition strategy is replaced by RDD. Each task has three demonstrations.

Evaluation Metrics and Baselines. We evaluate the performance of RDD and baselines in terms of multi-task success rates and corresponding rankings across 13 RLBench tasks². We compare our approach with a variety of baselines that adopt different sub-task identification strategies:

- **Expert** [13]: The same expert heuristic decomposer used in $\mathcal{D}_{\text{aug}}^{\text{train}}$ as performance upper bound.
- **UVD** [25]: A task-agnostic decomposer that detects change points of learning-based visual features.
- **Uniform**: A decomposer that divides each demonstration into 10 partitions with equal duration.
- **w/o Finetune**: The planner p_ϕ is the pre-trained VLM model without finetuning on $\mathcal{D}^{\text{demo}}$.

4.1 Quantitative Results and Analysis

Multi-Task Performance on RLBench. Table 1 shows the overall performance of RDD and baseline methods on multiple manipulation tasks using the base setting in Section 4. Results are averaged over 10 random seeds. RDD achieves a *near-oracle performance* and only compromises the success rate of merely 0.2% compared with the expert decomposer, our performance upper bound. On the other hand, we observe that UVD performs similarly to the naive uniform sampling strategy. It implies that the change points of learning-based visual features are not always aligned with the samples in $\mathcal{D}_{\text{aug}}^{\text{train}}$. By aligning the high-level planner to the knowledge of low-level policy, RDD outperforms the baseline methods that blindly decompose the demonstrations without this knowledge. It also suggests that finetuning is necessary for VLM-based planners. All finetuning-based methods achieve over 35% improvement over the vanilla Llama model.

Choice of Visual Representation. As an important building block of RDD, the choice of visual representation is of great importance. Table 2 shows the performance of RDD when adopting different visual encoders \mathcal{E} , including robotics specialized encoders: LIV [26], R3M [27], VIP [35], VC-1 [28]; and encoders for general vision tasks: CLIP [43], DINOv2 [29] and ResNet [44] pre-trained for ImageNet-1k classification. Results are averaged over three random seeds.

² Tasks on which the low-level visuomotor policy has a decent performance (success rate $> 35\%$ with expert planner). It excludes the interference of poorly optimized visuomotor when evaluating planners. Performance on all 18 tasks can be found in Appendix C.

It can be seen that *RDD shows good robustness with various visual encoders* and consistently outperforms baselines with the majority of encoders except for VC-1 and VIP, which demonstrates the strong robustness of RDD. VC-1 and VIP, on the other hand, are the only models that do not involve any form of language integration during training and perform the worst among all encoders. *This implies the importance of language integration for visual encoders in VLA perception for semantic information retrieval.* For instance, subtle pixel differences, such as the change of gripper state, may have a significant difference in language description. Surprisingly, ResNet, whose training does not explicitly involve language supervision, demonstrates a strong performance. The reason may be that its training dataset, ImageNet-1k, implicitly correlates its latent space with the language image labels.

Weighting Parameters. Table 3 shows the impact of α on the performance of RDD. Results are averaged over three random seeds. When $\alpha = 0$, i.e., there is no temporal alignment, and the algorithm is confused about sub-tasks whose beginning and ending frames are similar (e.g., reciprocating motion). On the other hand, overly relying on the temporal similarity ignores the semantic relationship between intervals and leads to performance degradation. We also evaluate the impact of the β in Table 5 for OOD scenarios, and the result shows that RDD is less sensitive to β . The choice of β depends on specific applications and user needs.

Number of Demonstrations in $\mathcal{D}^{\text{demo}}$. To explore the data efficiency of RDD, Table 4 shows its averaged success rates under different numbers of demonstrations in $\mathcal{D}_{\text{aug}}^{\text{demo}}$. Results are averaged over three random seeds. Specifically, we break the three-demonstration base setting dataset into three non-overlapping datasets with one demonstration per task to avoid bias induced by varying demonstration qualities. This result shows a high data efficiency of RDD. We credit this efficiency to the less-noisy keyframes provided by RDD, which are more informative for VLM to learn the underlying decomposition rules.

Performance on Real-world and OOD sub-tasks. Here we demonstrate RDD’s performance on both real-world and settings where the OOD sub-task appears. We first evaluate RDD on the real-world manipulation benchmark AgiBotWorld-Alpha [33]. We test RDD and UVD on the “supermarket” task, using 152 demos to build the RDD database and 37 demos for testing. For OOD sub-tasks, we test RDD on the human-operated demonstration dataset from RoboCerebra [34], which features highly diverse demonstrations in terms of objects, task goals, and arrangements. We use 560 demos to build the RDD database and test on the remaining 140 demos. We use the similarity measure sim in Eq. 3.7 for the OOD setting.

We evaluated the quality of the decomposition against ground-truth segmentations using the mean intersection over union (mIoU). As shown in Table 5, RDD outperforms UVD on real-world data. Under OOD settings, RDD consistently outperforms UVD by leveraging potential similarity between sub-tasks.

Speed and Scalability. We test the running time of Algorithm 1 with different numbers of frames on AMD EPYC 9254 using **one** CPU core. Figure 3 plots the running time with/without the prior knowledge of the maximum length of interval L_{max} . The results show that the complexity with L_{max} grows linearly with the number of frames, which aligns with our conclusion in Corollary 3.1.1, which indicates that when L_{max} is determined, the complexity of Algorithm 1 will be $O(N)$.

Table 2: Results when using different visual encoders \mathcal{E} . Full results on all tasks can be found in Table 9 in the appendix.

Visu. Repr.	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)
LIV	81.1 \pm 0.9	3.7 \pm 1.6
R3M	80.0 \pm 3.5	3.9 \pm 1.7
VIP	75.3 \pm 3.4	4.1 \pm 2.0
VC-1	75.5 \pm 3.1	3.8 \pm 2.2
CLIP	78.2 \pm 2.1	4.7 \pm 2.0
DINOv2	78.4 \pm 2.4	4.5 \pm 1.8
ResNet	81.1 \pm 2.5	3.4 \pm 1.5

Table 3: Results when tuning the weighting parameter α . Full results on all tasks can be found in Table 10.

α	Avg. Succ.	Avg. Rank
0	75.0 \pm 2.5	3.0 \pm 1.0
0.5	75.7 \pm 2.4	2.5 \pm 0.7
1	81.1 \pm 0.9	2.3 \pm 1.4
2	76.2 \pm 3.0	2.2 \pm 0.8

Table 4: Results with different numbers of demonstrations per task in $\mathcal{D}_{\text{aug}}^{\text{demo}}$. Full results on all tasks are in Table 11.

Demo. Num.	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)
1 (RDD)	77.9 \pm 4.5	2.0 \pm 0.9
3 (RDD)	81.1 \pm 0.9	1.6 \pm 0.6
3 (UVD)	75.6 \pm 1.8	2.4 \pm 0.6

Table 5: Performance on real-world and OOD sub-tasks (IoU).

Method	AgiBot. (Real World)	LIBERO (OOD)
UVD	0.506	0.598
RDD	0.706	/
RDD ($\beta = 0.25$)	/	0.624
RDD ($\beta = 0.10$)	/	0.630
RDD ($\beta = 0.05$)	/	0.614

Note that *Algorithm 1* supports parallel evaluation of the scoring function \tilde{J} , and the latency can be significantly reduced with multiprocessing. Also, we demonstrate the scalability of RDD when working with GPU-accelerated ANNS algorithms like FAISS [38] in Appendix D.

Necessity of Finetuning on Target Tasks: One may ask if the planner can transfer to an unseen new task in zero-shot. We thus build a new planner finetuned before deployment on the training set of the following tasks: “Close Jar”, “Insert Peg”, and “Install Bulb” as the baseline, which learns the visual features but not the task decompositions. Then, we test its performance on the remaining tasks. Results in Table 6 are averaged across 10 random seeds, and we also exclude tasks where the visuomotor fails. The results prove the necessity of fine-tuning on target tasks.

Decompose with VLMs: VLMs pretrained on internet-scale data are promising to process a variety of video understanding tasks. In Table 7, we compared RDD with a Gemini-2.5-pro [42]-based decomposer with the following prompt:

There is a robot doing a task, which can be segmented into multiple steps. A keyframe is where the robot finishes the previous step and begins the next. Can you help me find ALL indexes of keyframes? Please return a list of indices, for example: [15, 65, 105, ...]. Note that the frame index starts from 0 instead of 1.

As shown, RDD outperforms Gemini-2.5-pro despite its powerful general video understanding abilities. This result highlights the necessity of the planner aligning and the effectiveness of RDD.

Extended Evaluations and Discussions. We provide extended evaluations results in C and further discussions in Appendix D. We also provide a conceptual speed evaluation of RDD when working with the GPU-accelerated ANNS method FAISS [38] in Appendix D.1.

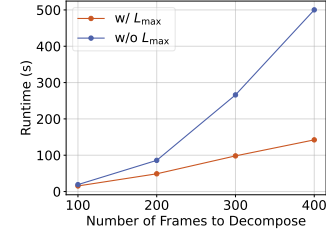


Figure 3: Linear scaling of running time of Algorithm 1 with L_{\max} . Tested with one CPU core.

Table 6: Vanilla Planner without finetuning on the target task. Full results on all tasks are in Table 12.

Method	Avg. Succ. (↑)	Avg. Rank (↓)
w/o finetuning on target task	77.9 ± 4.3	1.6 ± 0.5
RDD (Ours)	79.6 ± 7.2	1.4 ± 0.5

Table 7: Comparing RDD with Gemini-2.5-pro. Full results on all tasks are in Table 13.

Method	Avg. Succ. (↑)	Avg. Rank (↓)
Gemini-2.5-pro	72.6 ± 4.7	1.7 ± 0.4
RDD (Ours)	74.9 ± 6.9	1.3 ± 0.4

4.2 Qualitative Results and Analysis

Figure 4 visualizes the decomposition results of RDD and UVD on both real-world and simulation benchmarks. We can observe that RDD is robust to task-irrelevant interference and can robustly identify the sub-tasks that are close to the expert sub-task division. Also, RDD demonstrates strong robustness to nuanced arm movements, where the keyframe localization is challenging precisely. Conversely, UVD fails to locate keyframes precisely, and the generated sub-tasks largely deviate from expert sub-tasks.

5 Discussions and Future Works

Visuomotor Training Data Generation based on Source Dataset: While this work applies RDD to planner-visuomotor alignment, it can also be used to generate additional sub-task training data for visuomotor *aligned with a labeled source dataset*. By aligning the sub-task interval visual features with the existing source dataset, RDD may make the newly labeled data easier to learn, allowing the visuomotor reuse learned knowledge from the source dataset.

Specific Sub-task Interval Features: RDD measures sub-task interval similarity in the single-frame image feature space. Some applications, such as hierarchical vision-language navigation [19], which require the planner to use historical landmark images, may necessitate specialized designs of the similarity score function.

Data Quality of the Source Dataset and Data Curation: As a retrieval-based sub-task identification method, RDD’s primary objective is to let the high-level planner effectively utilize the skills that the low-level visuomotor policy *already possesses*. Therefore, RDD is agnostic to the “optimality” of the



Figure 4: Qualitative results of RDD and UVD functioning on both real-world (AgiBotWorld) and simulation (RLBench and LIBERO) benchmarks. Blocks outlined in black are sub-tasks decomposed by the same task-specific heuristic used in the visuomotor policy’s training set; blocks outlined in **green** are sub-tasks found by RDD; and blocks outlined in **red** are sub-tasks found by UVD.

skills themselves. This ensures the planner generates commands that the policy can reliably execute, rather than potentially “better” ones it cannot handle.

On the other hand, in scenarios where the visuomotor policy’s training data contains significant noisy samples that the policy fails to learn, RDD can be easily integrated with dataset curation techniques [45, 46]. These methods can serve as a pre-processing step to filter the visuomotor training set. For instance, CUPID [45] computes an “action influence” score for state-action pairs that can be used to evaluate each segment’s contribution to the policy’s final behavior. By applying a simple threshold, low-influence or flawed segments can be pruned from the dataset before RDD uses it as a reference. This would prevent catastrophic failures by ensuring RDD aligns demonstrations only with high-quality, influential sub-tasks.

6 Conclusion

In this work, we present the Retrieval-based Demonstration Decomposer (RDD), a training-free decomposition method that aligns the high-level task planner and low-level visuomotor policy in hierarchical VLAs. By retrieving and aligning sub-task segments with the low-level policy’s training data, RDD enables an effective planner that fully exploits the capability of the visuomotor policy. We formally formulate the sub-task identification task into an optimal partitioning problem, which can be efficiently solved by dynamic programming with our novel sub-task interval scoring function. Experiment results demonstrate that RDD outperforms state-of-the-art demonstration decomposers. RDD offers a scalable and promising solution for sub-task identification, opening new avenues for planner-policy coordination in hierarchical robot learning systems.

References

- [1] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.
- [2] 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.
- [3] 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.
- [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] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
- [6] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [7] Yingdong Hu, Fanqi Lin, Tong Zhang, Li Yi, and Yang Gao. Look before you leap: Unveiling the power of gpt-4v in robotic vision-language planning. *arXiv preprint arXiv:2311.17842*, 2023.
- [8] Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong, Paul Wohlhart, Sean Kirmani, Brianna Zitkovich, Fei Xia, et al. Open-world object manipulation using pre-trained vision-language models. *arXiv preprint arXiv:2303.00905*, 2023.
- [9] Hongyi Chen, Yunchao Yao, Ruixuan Liu, Changliu Liu, and Jeffrey Ichnowski. Automating robot failure recovery using vision-language models with optimized prompts. *arXiv preprint arXiv:2409.03966*, 2024.
- [10] 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.
- [11] Peiqi Liu, Yaswanth Orru, Jay Vakil, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024.
- [12] Fangchen Liu, Kuan Fang, Pieter Abbeel, and Sergey Levine. Moka: Open-vocabulary robotic manipulation through mark-based visual prompting. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024.
- [13] Yinpei Dai, Jayjun Lee, Nima Fazeli, and Joyce Chai. Racer: Rich language-guided failure recovery policies for imitation learning. *International Conference on Robotics and Automation (ICRA)*, 2025.
- [14] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, Karl Pertsch, Quan Vuong, James Tanner, Anna Walling, Haohuan Wang, Niccolo Fusai, et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025.
- [15] Yi Li, Yuquan Deng, Jesse Zhang, Joel Jang, Marius Memmel, Raymond Yu, Caelan Reed Garrett, Fabio Ramos, Dieter Fox, Anqi Li, et al. Hamster: Hierarchical action models for open-world robot manipulation. *arXiv preprint arXiv:2502.05485*, 2025.

- [16] Lucy Xiaoyang Shi, Zheyuan Hu, Tony Z Zhao, Archit Sharma, Karl Pertsch, Jianlan Luo, Sergey Levine, and Chelsea Finn. Yell at your robot: Improving on-the-fly from language corrections. *arXiv preprint arXiv:2403.12910*, 2024.
- [17] Weiyu Liu, Neil Nie, Ruohan Zhang, Jiayuan Mao, and Jiajun Wu. Learning compositional behaviors from demonstration and language. In *8th Annual Conference on Robot Learning*, 2025.
- [18] 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.
- [19] An-Chieh Cheng, Yandong Ji, Zhaojing Yang, Zaitian Gongye, Xueyan Zou, Jan Kautz, Erdem Biyik, Hongxu Yin, Sifei Liu, and Xiaolong Wang. Navila: Legged robot vision-language-action model for navigation. *arXiv preprint arXiv:2412.04453*, 2024.
- [20] Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang, Kiana Ehsani, Dieter Fox, and Ranjay Krishna. Manipulate-anything: Automating real-world robots using vision-language models. *arXiv preprint arXiv:2406.18915*, 2024.
- [21] Changyeon Kim, Minh Heo, Doohyun Lee, Jinwoo Shin, Honglak Lee, Joseph J Lim, and Kimin Lee. Subtask-aware visual reward learning from segmented demonstrations. *arXiv preprint arXiv:2502.20630*, 2025.
- [22] Wensheng Wang and Ning Tan. Hybridgen: Vlm-guided hybrid planning for scalable data generation of imitation learning. *arXiv preprint arXiv:2503.13171*, 2025.
- [23] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Ireteyo Akinola, Yashraj Narang, Linxi Fan, Yuke Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot learning using human demonstrations. *arXiv preprint arXiv:2310.17596*, 2023.
- [24] Tongzhou Mu, Minghua Liu, and Hao Su. Drs: Learning reusable dense rewards for multi-stage tasks. *arXiv preprint arXiv:2404.16779*, 2024.
- [25] Zichen Zhang, Yunshuang Li, Osbert Bastani, Abhishek Gupta, Dinesh Jayaraman, Yecheng Jason Ma, and Luca Weihs. Universal visual decomposer: Long-horizon manipulation made easy. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6973–6980. IEEE, 2024.
- [26] Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. *arXiv preprint arXiv:2306.00958*, 2023.
- [27] 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.
- [28] 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.
- [29] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Fernandez, et al. Dinov2: Learning robust visual features without supervision, 2023.
- [30] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2023.
- [31] Brad Jackson, Jeffrey D Scargle, David Barnes, Sundararajan Arabhi, Alina Alt, Peter Gioumoussis, Elyus Gwin, Paungkaew Sangtrakulcharoen, Linda Tan, and Tun Tao Tsai. An algorithm for optimal partitioning of data on an interval. *IEEE Signal Processing Letters*, 12(2):105–108, 2005.

- [32] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- [33] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Xindong He, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2025.
- [34] Songhao Han, Boxiang Qiu, Yue Liao, Siyuan Huang, Chen Gao, Shuicheng Yan, and Si Liu. Robocerebra: A large-scale benchmark for long-horizon robotic manipulation evaluation. *arXiv preprint arXiv:2506.06677*, 2025.
- [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] Erik Bernhardsson. ANNOY library. <https://github.com/spotify/annoy>. Accessed: 2025-05-05.
- [37] Martin Aumüller, Erik Bernhardsson, and Alexander Faithfull. Ann-benchmarks: A benchmarking tool for approximate nearest neighbor algorithms. *Information Systems*, 87:101374, 2020.
- [38] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024.
- [39] Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. In *Conference on Robot Learning*, pages 694–710. PMLR, 2023.
- [40] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild. URL <https://llava-vl.github.io/blog/2024-05-10-llava-next-stronger-llms>, 2024.
- [41] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- [42] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soriccut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [43] 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.
- [44] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [45] Christopher Agia, Rohan Sinha, Jingyun Yang, Rika Antonova, Marco Pavone, Haruki Nishimura, Masha Itkina, and Jeannette Bohg. Cupid: Curating data your robot loves with influence functions. *arXiv preprint arXiv:2506.19121*, 2025.
- [46] Joey Hejna, Suvir Mirchandani, Ashwin Balakrishna, Annie Xie, Ayzaan Wahid, Jonathan Tompson, Pannag Sanketi, Dhruv Shah, Coline Devin, and Dorsa Sadigh. Robot data curation with mutual information estimators. *arXiv preprint arXiv:2502.08623*, 2025.
- [47] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.

- [48] Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025.
- [49] Junming Zhang, Weijia Chen, Yuping Wang, Ram Vasudevan, and Matthew Johnson-Roberson. Point set voting for partial point cloud analysis. *IEEE Robotics and Automation Letters*, 6(2):596–603, 2021.
- [50] Mingxuan Yan, Ruijie Zhang, Xuedou Xiao, and Wei Wang. Detvpcc: Roi-based point cloud sequence compression for 3d object detection. *arXiv preprint arXiv:2502.04804*, 2025.
- [51] Zehao Wang, Yuping Wang, Zhuoyuan Wu, Hengbo Ma, Zhaowei Li, Hang Qiu, and Jiachen Li. Cmp: Cooperative motion prediction with multi-agent communication. *IEEE Robotics and Automation Letters*, 2025.
- [52] Yuping Wang and Jier Chen. Eqdrive: Efficient equivariant motion forecasting with multi-modality for autonomous driving. In *2023 8th International Conference on Robotics and Automation Engineering (ICRAE)*, pages 224–229. IEEE, 2023.
- [53] Yuping Wang and Jier Chen. Equivariant map and agent geometry for autonomous driving motion prediction. In *2023 International Conference on Electrical, Computer and Energy Technologies (ICECET)*, pages 1–6. IEEE, 2023.
- [54] Shuo Xing, Chengyuan Qian, Yuping Wang, Hongyuan Hua, Kexin Tian, Yang Zhou, and Zhengzhong Tu. Openemma: Open-source multimodal model for end-to-end autonomous driving. In *Proceedings of the Winter Conference on Applications of Computer Vision*, pages 1001–1009, 2025.
- [55] Yuping Wang, Xiangyu Huang, Xiaokang Sun, Mingxuan Yan, Shuo Xing, Zhengzhong Tu, and Jiachen Li. Uniocc: A unified benchmark for occupancy forecasting and prediction in autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, 2025.
- [56] Yuping Wang, Shuo Xing, Cui Can, Renjie Li, Hongyuan Hua, Kexin Tian, Zhaobin Mo, Xiangbo Gao, Keshu Wu, Sulong Zhou, et al. Generative ai for autonomous driving: Frontiers and opportunities. *arXiv preprint arXiv:2505.08854*, 2025.
- [57] Xu Liu, Tong Zhou, Chong Wang, Yuping Wang, Yuanxin Wang, Qijingwen Cao, Weizhi Du, Yonghuan Yang, Junjun He, Yu Qiao, et al. Toward the unification of generative and discriminative visual foundation model: A survey. *The Visual Computer*, pages 1–42, 2024.
- [58] Shuo Xing, Yuping Wang, Peiran Li, Ruizheng Bai, Yueqi Wang, Chengxuan Qian, Huaxiu Yao, and Zhengzhong Tu. Re-align: Aligning vision language models via retrieval-augmented direct preference optimization. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, 2025.
- [59] Shuo Xing, Zezhou Sun, Shuangyu Xie, Kaiyuan Chen, Yanjia Huang, Yuping Wang, Jiachen Li, Dezhen Song, and Zhengzhong Tu. Can large vision language models read maps like a human? *arXiv preprint arXiv:2503.14607*, 2025.
- [60] Congrui Hetang and Yuping Wang. Novel view synthesis from a single rgbd image for indoor scenes. In *2023 International Conference on Image Processing, Computer Vision and Machine Learning (ICICML)*, pages 447–450. IEEE, 2023.
- [61] Xiangbo Gao, Yuheng Wu, Xuwen Luo, Keshu Wu, Xinghao Chen, Yuping Wang, Chenxi Liu, Yang Zhou, and Zhengzhong Tu. Airv2x: Unified air-ground vehicle-to-everything collaboration. *arXiv preprint arXiv:2506.19283*, 2025.
- [62] Peiran Li, Xinkai Zou, Zhuohang Wu, Ruifeng Li, Shuo Xing, Hanwen Zheng, Zhikai Hu, Yuping Wang, Haoxi Li, Qin Yuan, et al. SafeFlow: A principled protocol for trustworthy and transactional autonomous agent systems. *arXiv preprint arXiv:2506.07564*, 2025.

NeurIPS Paper Checklist

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 claims made in the abstract and introduction accurately reflect the paper's contributions and scope.

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 the limitations of this work.

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: [\[Yes\]](#)

Justification: Assumptions are explicitly list in the conclusions, and proofs are listed in the appendix.

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 all the training and test details in the main text, appendix, and source code.

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: [Yes]

Justification: We will release the code and relevant data.

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, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide all the training and test details in the main text, appendix, and source code.

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: [Yes]

Justification: We repeated all experiments for multiple rounds and multiple random seeds and provided the statistics in the results.

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: We provide sufficient information on the computer resources in the main text and 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: The research conducted in the paper conformed, in every respect, with the NeurIPS 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: [Yes]

Justification: We discussed it in the appendix.

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., pre-trained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not release any new datasets or models. Our work only uses existing pre-trained models, which are widely available and do not introduce additional safety risks. Therefore, safeguards are not applicable in our case.

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 properly credited the creators or original owners of assets (e.g., code, data, models), used in the paper and conformed the license and terms.

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 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: [\[Yes\]](#)

Justification: We attach the details of the dataset/code/model along with the submission.

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: This work does not involve involve crowdsourcing nor research with human subjects.

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: This work does 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: [Yes]

Justification: Our method uses LLMs for dataset annotation, and we have described details in the main text.

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.

Appendix

A Algorithm Details

A.1 Dynamic Programming Solver to Problem 3.1

Algorithm 1 shows the dynamic programming solver. L_{\max} and L_{\min} are user-specified parameters that determine the minimum and maximum length of proposed sub-task intervals. \tilde{J} is the interval scoring function.

Algorithm 1 MaxSumPartition

Require: Sequence $u = [u_1, u_2, \dots, u_n]$, scoring function \tilde{J} , integer L_{\min} , integer L_{\max}
Ensure: Maximum score sum and partition of u

```

1: Initialize  $dp[0 \dots n] \leftarrow -\infty$ ,  $parts[0 \dots n] \leftarrow \emptyset$ 
2:  $dp[0] \leftarrow 0$ 
3: for  $i = L_{\min} + 1$  to  $n$  do
4:    $bestScore \leftarrow -\infty$ 
5:    $bestPartition \leftarrow \emptyset$ 
6:   for  $j = 0$  to  $i$  do
7:     if  $L_{\min} \leq i - j \leq L_{\max}$  then
8:        $segment \leftarrow u[j : i]$ 
9:        $s \leftarrow \tilde{J}(segment)$   $\triangleright$  can be evaluated in parallel before loops
10:      if  $dp[j] + s > bestScore$  then
11:         $bestScore \leftarrow dp[j] + s$ 
12:         $bestPartition \leftarrow parts[j] \cup \{segment\}$ 
13:      end if
14:    end if
15:  end for
16:  if  $bestPartition \neq \emptyset$  then
17:     $dp[i] \leftarrow bestScore$ 
18:     $parts[i] \leftarrow bestPartition$ 
19:  else
20:     $dp[i] \leftarrow dp[i - 1]$ 
21:     $parts[i] \leftarrow parts[i - 1]$ 
22:  end if
23: end for
24: return  $(dp[n], parts[n])$ 

```

A.2 Proof of Correctness and Complexity of Algorithm 1

Proof. The correctness when $L_{\min} = 1, L_{\max} = |\mathcal{S}^i|$ (we denote the algorithm under this case as Algorithm 0) has been proven in *Proof 2* of Jackson et al. [31]. It is sufficient to prove the cases when $1 < L_{\min} < L_{\max} < |\mathcal{S}^i|$. Notice that Algorithm 1 is equivalent to a special case of Algorithm 0 by constructing an adapted scoring function defined as Algorithm 2, where the score of invalid intervals

Algorithm 2 AdaptedScoreFunc

Require: Sub-sequence $u' = [u_1, u_2, \dots, u_m]$, scoring function \tilde{J} , integer L_{\min} , integer L_{\max}
Ensure: Adapted score of u'

```

1: if  $L_{\min} \leq |u'| \leq L_{\max}$  then
2:   return  $\tilde{J}(u')$ 
3: else
4:   return  $-\infty$ 
5: end if

```

is $-\infty$. Note that ADAPTEDSCOREFUNC preserves the additiveness of J , because if any interval in a strategy P violates the length assumption, P is also invalid, i.e., $J(P) = -\infty$. Given the facts: 1)

the correctness of Algorithm 0 has been proven by *Proof 2* [31]; 2) Algorithm 1 is equivalent to a special case of Algorithm 0, by implication, the correctness of Algorithm 1 is proved.

As for complexity, let N be the length of the demonstration and M be the number of evaluations to \tilde{J} , we have:

$$\begin{aligned}
M &= \sum_{j=2}^{L_{\max}-L_{\min}+1} j + (N - L_{\max})(L_{\max} - L_{\min} + 1) \\
&= \frac{(L_{\max} - L_{\min} + 3)(L_{\max} - L_{\min})}{2} + (N - L_{\max})(L_{\max} - L_{\min} + 1) \\
&= O((L_{\max} - L_{\min}) \cdot \max(L_{\max} - L_{\min}, N - L_{\max}))
\end{aligned}$$

□

B Proof of Proposition 3.1

Proof. Let the identical similarity scores equal s , and let b_j^i, e_j^i be the starting and ending indexes of interval \mathcal{I}_j^i , respectively. By applying Eq. 3.2 and Eq. 3.3 we rewrite the left side of Eq. 3.1 to:

$$J(\{\mathcal{I}_j^i\}) = \tilde{J}(\mathcal{I}_j^i) = (e_j^i - b_j^i)s$$

And the right side:

$$\begin{aligned}
J(\{\mathcal{I}_{j1}^i, \mathcal{I}_{j2}^i, \dots, \mathcal{I}_{jK}^i\}) &= \sum_{k=1}^K \tilde{J}(\mathcal{I}_{jk}^i) \\
&= (e_{j1}^i - b_{j1}^i + e_{j2}^i - b_{j2}^i + \dots + e_{jk}^i - b_{jk}^i)s \\
&= \underbrace{(e_{j1}^i - b_j^i + e_{j2}^i - e_{j1}^i + \dots + e_j^i - e_{j(k-1)}^i)}_{\text{Since intervals are consecutive.}} s \\
&= (e_j^i - b_j^i)s \\
&= J(\{\mathcal{I}_j^i\})
\end{aligned}$$

□

C Additional Quantitative Results

Tables 8-13 provide the complete multi-task performances of the results in Section 4.1, including ones where the visuomotor policy fails.

D Discussions

D.1 Work with GPU-accelerated ANNS

The nearest neighbor (NN) search in RDD can be significantly accelerated using GPU-accelerated libraries like FAISS [38]. We conduct experiments on a typical database of 10 million entries (mainstream policy training dataset scale, as shown in Section D.2) of 2048 dimensions (same dimension as our main experiment in Table 1) As shown in Table 14, FAISS can achieve > 300 NN queries per second on one NVIDIA 4090 GPU. Under this setting, RDD only needs < 2 minutes to decompose a 500-frame video (5 fps), with a max interval length of 100 frames. (44549 NN queries in total). In other words, as part of the offline dataset building process, RDD can decompose demonstrations at a high speed of 4.3 fps, which shows the high scalability of RDD.

D.2 Scale of Mainstream Robotics Datasets

To support the aforementioned experiment settings, here we provide the scale of some of the most popular open-sourced robotics datasets. In summary, assuming each demonstration can be

Table 8: Main results with all RLBench Tasks.

Method	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Close Jar	Insert Peg	Install Bulb	Meat off Grill	Open Drawer	Place Cups	Sort Shape	Place Wine
w/o Finetune	39.7 \pm 6.5	4.3 \pm 1.3	27.6 \pm 26.4	5.6 \pm 6.7	34.8 \pm 14.2	46.4 \pm 26.8	95.6 \pm 6.1	3.2 \pm 4.3	16.0 \pm 11.7	83.2 \pm 13.0
Uniform	54.5 \pm 4.1	3.1 \pm 1.2	46.4 \pm 29.9	8.8 \pm 11.8	51.2 \pm 19.2	76.4 \pm 22.4	100.0 \pm 0.0	0.8 \pm 1.6	25.6 \pm 9.2	80.8 \pm 14.5
UVD [25]	54.3 \pm 3.9	3.2 \pm 1.2	44.0 \pm 28.7	10.4 \pm 14.1	54.8 \pm 20.0	85.2 \pm 20.6	100.0 \pm 0.0	1.2 \pm 1.8	25.2 \pm 11.0	80.8 \pm 15.3
Expert [13]	57.6 \pm 3.3	2.0 \pm 0.9	50.4 \pm 33.1	12.0 \pm 17.9	50.4 \pm 13.3	94.4 \pm 9.7	99.2 \pm 2.4	3.2 \pm 3.9	26.0 \pm 10.6	81.6 \pm 15.0
RDD (Ours)	57.3 \pm 5.3	2.4 \pm 1.1	46.0 \pm 28.2	16.8 \pm 18.6	52.8 \pm 16.4	84.4 \pm 21.1	99.2 \pm 2.4	2.0 \pm 2.0	32.4 \pm 10.2	86.4 \pm 15.4
Method	Push Buttons	Put in Cupboard	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
w/o Finetune	54.8 \pm 9.1	41.2 \pm 20.1	36.4 \pm 28.8	58.8 \pm 23.3	36.0 \pm 21.8	57.2 \pm 14.9	2.8 \pm 2.6	2.8 \pm 3.6	22.8 \pm 32.5	89.2 \pm 13.4
Uniform	82.0 \pm 7.8	36.8 \pm 15.4	98.0 \pm 2.7	92.4 \pm 10.8	64.8 \pm 16.7	64.4 \pm 9.9	13.6 \pm 7.8	5.2 \pm 4.7	34.8 \pm 37.7	98.8 \pm 3.6
UVD [25]	67.2 \pm 13.6	35.2 \pm 12.1	90.4 \pm 8.6	96.8 \pm 6.6	74.4 \pm 29.2	66.8 \pm 21.2	9.6 \pm 6.2	1.6 \pm 3.7	43.6 \pm 24.6	89.6 \pm 11.1
Expert [13]	85.6 \pm 6.0	39.6 \pm 15.6	91.2 \pm 7.3	97.6 \pm 5.1	75.2 \pm 24.6	66.4 \pm 22.0	14.8 \pm 11.2	5.2 \pm 4.4	48.8 \pm 35.5	96.0 \pm 5.7
RDD (Ours)	84.0 \pm 7.8	41.2 \pm 17.1	97.2 \pm 3.1	98.4 \pm 3.2	68.0 \pm 25.0	65.2 \pm 14.3	5.2 \pm 3.6	1.6 \pm 2.7	57.2 \pm 29.7	94.0 \pm 5.1

Table 9: Multi-task performances with different visual representations.

Visu. Repr.	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Close Jar	Insert Peg	Install Bulb	Meat off Grill	Open Drawer	Place Cups	Sort Shape	Place Wine
LIV [26]	61.0 \pm 0.4	3.6 \pm 1.7	68.0 \pm 17.3	4.0 \pm 3.3	41.3 \pm 21.7	96.0 \pm 5.7	100.0 \pm 0.0	1.3 \pm 1.9	32.0 \pm 5.7	96.0 \pm 3.3
R3M [27]	59.2 \pm 2.5	4.2 \pm 1.8	65.3 \pm 21.0	4.0 \pm 5.7	44.0 \pm 13.1	97.3 \pm 3.8	98.7 \pm 1.9	0.0 \pm 0.0	12.0 \pm 3.3	86.7 \pm 6.8
VIP [35]	56.5 \pm 2.0	4.0 \pm 2.0	72.0 \pm 14.2	2.7 \pm 1.9	38.7 \pm 15.4	93.3 \pm 9.4	100.0 \pm 0.0	5.3 \pm 5.0	22.7 \pm 10.0	89.3 \pm 8.2
VC-1 [28]	56.9 \pm 1.6	3.7 \pm 2.3	73.3 \pm 9.4	1.3 \pm 1.9	30.7 \pm 18.6	93.3 \pm 9.4	100.0 \pm 0.0	8.0 \pm 3.3	20.0 \pm 8.6	86.7 \pm 10.0
CLIP [43]	58.4 \pm 1.6	4.3 \pm 2.0	62.7 \pm 21.7	4.0 \pm 3.3	46.7 \pm 15.4	96.0 \pm 5.7	100.0 \pm 0.0	0.0 \pm 0.0	16.0 \pm 3.3	82.7 \pm 13.6
DINOv2 [29]	58.3 \pm 1.4	4.4 \pm 1.6	65.3 \pm 18.0	2.7 \pm 3.8	41.3 \pm 21.2	98.7 \pm 1.9	100.0 \pm 0.0	1.3 \pm 1.9	13.3 \pm 1.9	80.0 \pm 8.6
ResNet [44]	60.5 \pm 2.0	3.8 \pm 1.7	68.0 \pm 20.4	2.7 \pm 3.8	46.7 \pm 10.5	96.0 \pm 5.7	100.0 \pm 0.0	0.0 \pm 0.0	13.3 \pm 5.0	84.0 \pm 6.5
Visu. Repr.	Push Buttons	Put in Cupboard	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
LIV [26]	78.7 \pm 8.2	57.3 \pm 3.8	97.3 \pm 1.9	97.3 \pm 3.8	88.0 \pm 8.6	73.3 \pm 3.8	4.0 \pm 3.3	1.3 \pm 1.9	66.7 \pm 5.0	94.7 \pm 5.0
R3M [27]	89.3 \pm 5.0	50.7 \pm 10.0	85.3 \pm 5.0	94.7 \pm 5.0	94.7 \pm 5.0	82.7 \pm 12.4	8.0 \pm 3.3	1.3 \pm 1.9	53.3 \pm 8.2	97.3 \pm 3.8
VIP [35]	92.0 \pm 3.3	64.0 \pm 8.6	93.3 \pm 3.8	89.3 \pm 10.0	10.7 \pm 7.5	46.7 \pm 36.7	2.7 \pm 1.9	5.3 \pm 3.8	92.0 \pm 8.6	97.3 \pm 1.9
VC-1 [28]	93.3 \pm 5.0	65.3 \pm 8.2	93.3 \pm 6.8	92.0 \pm 8.6	9.3 \pm 6.8	52.0 \pm 31.5	4.0 \pm 3.3	9.3 \pm 7.5	92.0 \pm 8.6	100.0 \pm 0.0
CLIP [43]	89.3 \pm 5.0	46.7 \pm 13.2	81.3 \pm 3.8	94.7 \pm 5.0	94.7 \pm 5.0	81.3 \pm 10.5	10.7 \pm 3.8	5.3 \pm 5.0	52.0 \pm 8.6	88.0 \pm 14.2
DINOv2 [29]	88.0 \pm 3.3	50.7 \pm 15.4	85.3 \pm 5.0	94.7 \pm 5.0	94.7 \pm 5.0	78.7 \pm 6.8	9.3 \pm 1.9	4.0 \pm 3.3	46.7 \pm 5.0	94.7 \pm 7.5
ResNet [44]	93.3 \pm 1.9	61.3 \pm 9.4	98.7 \pm 1.9	90.7 \pm 8.2	86.7 \pm 10.5	73.3 \pm 6.8	17.3 \pm 11.5	1.3 \pm 1.9	56.0 \pm 5.7	100.0 \pm 0.0

Table 10: Multi-task performance with different weighting parameter α .

α	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Close Jar	Insert Peg	Install Bulb	Meat off Grill	Open Drawer	Place Cups	Sort Shape	Place Wine
0	57.3 \pm 2.1	2.8 \pm 1.0	74.7 \pm 10.0	0.0 \pm 0.0	32.0 \pm 18.2	52.0 \pm 14.2	98.7 \pm 1.9	6.7 \pm 5.0	29.3 \pm 5.0	81.3 \pm 5.0
0.5	57.6 \pm 2.2	2.7 \pm 0.8	73.3 \pm 10.5	0.0 \pm 0.0	33.3 \pm 11.5	49.3 \pm 21.0	100.0 \pm 0.0	5.3 \pm 3.8	29.3 \pm 6.8	92.0 \pm 5.7
1	61.0 \pm 0.4	2.3 \pm 1.4	68.0 \pm 17.3	4.0 \pm 3.3	41.3 \pm 21.7	96.0 \pm 5.7	100.0 \pm 0.0	1.3 \pm 1.9	32.0 \pm 5.7	96.0 \pm 3.3
2	58.0 \pm 2.3	2.2 \pm 0.8	76.0 \pm 9.8	0.0 \pm 0.0	33.3 \pm 10.5	48.0 \pm 11.3	100.0 \pm 0.0	8.0 \pm 3.3	29.3 \pm 3.8	88.0 \pm 6.5
α	Push Buttons	Put in Cupboard	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
0	90.7 \pm 5.0	62.7 \pm 12.4	96.0 \pm 5.7	77.3 \pm 3.8	76.0 \pm 11.8	58.7 \pm 9.4	18.7 \pm 12.4	1.3 \pm 1.9	78.7 \pm 16.4	96.0 \pm 3.3
0.5	85.3 \pm 6.8	62.7 \pm 13.2	96.0 \pm 3.3	80.0 \pm 0.0	76.0 \pm 14.2	58.7 \pm 6.8	17.3 \pm 13.6	0.0 \pm 0.0	80.0 \pm 15.0	97.3 \pm 1.9
1	78.7 \pm 8.2	57.3 \pm 3.8	97.3 \pm 1.9	97.3 \pm 3.8	88.0 \pm 8.6	73.3 \pm 3.8	4.0 \pm 3.3	1.3 \pm 1.9	66.7 \pm 5.0	94.7 \pm 5.0
2	88.0 \pm 8.6	60.0 \pm 9.8	100.0 \pm 0.0	81.3 \pm 6.8	77.3 \pm 12.4	58.7 \pm 6.8	14.7 \pm 10.0	1.3 \pm 1.9	84.0 \pm 17.3	96.0 \pm 5.7

decomposed into 10 sub-tasks, the mainstream policy training datasets typically have 10 million sub-tasks. (\approx 10 million entries in the database). **The Open X-Embodiment (OXE) Dataset [47]:** A landmark collaboration among 21 institutions, OXE provides over 1 million robot trajectories from 22 different robot embodiments. Its explicit goal is to foster the development of generalist models, demonstrating that the community is actively removing the proprietary data barriers of the past. The explicit purpose of OXE is to provide a standardized, large-scale resource to train generalist models that have demonstrated significant performance gains by training on this diverse data. **Hugging Face**

Table 11: Multi-task performance with different numbers of demonstrations.

Demo. Num.	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Close Jar	Insert Peg	Install Bulb	Meat off Grill	Open Drawer	Place Cups	Sort Shape	Place Wine
1 (RDD)	59.1 \pm 3.4	1.8 \pm 0.8	75.6 \pm 11.1	6.2 \pm 5.7	35.6 \pm 9.5	65.8 \pm 20.9	100.0 \pm 0.0	5.8 \pm 4.7	25.8 \pm 3.8	91.1 \pm 7.7
3 (RDD)	61.0 \pm 0.4	1.8 \pm 0.7	68.0 \pm 17.3	4.0 \pm 3.3	41.3 \pm 21.7	96.0 \pm 5.7	100.0 \pm 0.0	1.3 \pm 1.9	32.0 \pm 5.7	96.0 \pm 3.3
3 (UVD [25])	57.1 \pm 0.3	2.3 \pm 0.6	66.7 \pm 13.2	4.0 \pm 5.7	37.3 \pm 19.1	93.3 \pm 9.4	100.0 \pm 0.0	2.7 \pm 1.9	21.3 \pm 10.5	77.3 \pm 11.5
Demo. Num.	Push Buttons	Put in Cupboard	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
1 (RDD)	86.7 \pm 5.7	60.4 \pm 11.8	97.8 \pm 2.0	78.2 \pm 13.3	61.8 \pm 28.7	79.6 \pm 16.2	11.6 \pm 8.1	2.2 \pm 2.7	87.1 \pm 16.2	92.9 \pm 14.7
3 (RDD)	78.7 \pm 8.2	57.3 \pm 3.8	97.3 \pm 1.9	97.3 \pm 3.8	88.0 \pm 8.6	73.3 \pm 3.8	4.0 \pm 3.3	1.3 \pm 1.9	66.7 \pm 5.0	94.7 \pm 5.0
3 (UVD [25])	62.7 \pm 12.4	44.0 \pm 6.5	84.0 \pm 6.5	96.0 \pm 5.7	85.3 \pm 13.2	82.7 \pm 12.4	16.0 \pm 3.3	1.3 \pm 1.9	60.0 \pm 16.3	93.3 \pm 5.0

Table 12: Multi-task performance of Vanilla Planner without finetuning on the target task.

Method	Avg. Succ. (\uparrow)	Avg. Rank (\downarrow)	Meat off Grill	Open Drawer	Place Wine	Push Buttons	Put in Cupboard
w/o finetuning on target task	77.9 \pm 4.3	1.6 \pm 0.5	99.2 \pm 2.4	99.6 \pm 1.2	86.4 \pm 8.8	70.4 \pm 8.0	61.2 \pm 16.8
RDD (Ours)	79.6 \pm 7.2	1.4 \pm 0.5	84.4 \pm 21.1	99.2 \pm 2.4	86.4 \pm 15.4	84.0 \pm 7.8	41.2 \pm 17.1

Method	Put in Drawer	Put in Safe	Drag Stick	Slide Block	Sweep to Dustpan	Turn Tap
w/o finetuning on target task	86.0 \pm 14.3	94.8 \pm 9.0	74.0 \pm 23.3	62.4 \pm 16.8	30.0 \pm 15.3	92.4 \pm 14.4
RDD (Ours)	97.2 \pm 3.1	98.4 \pm 3.2	68.0 \pm 25.0	65.2 \pm 14.3	57.2 \pm 29.7	94.0 \pm 5.1

Table 13: Comparing RDD with Gemini-2.5-pro.

Method	Avg.	Avg.	Close	Install	Meat off	Open	Place	Push
	Succ. (\uparrow)	Rank (\downarrow)	Jar	Bulb	Grill	Drawer	Wine	Buttons
Gemini-2.5-pro	72.6 \pm 4.7	1.7 \pm 0.4	41.2 \pm 30.1	40.8 \pm 16.5	83.2 \pm 15.2	99.6 \pm 1.2	86.4 \pm 11.1	82.4 \pm 8.6
RDD (Ours)	74.9 \pm 6.9	1.3 \pm 0.4	46.0 \pm 28.2	52.8 \pm 16.4	84.4 \pm 21.1	99.2 \pm 2.4	86.4 \pm 15.4	84.0 \pm 7.8
Method	Put in	Put in	Put in	Drag	Slide	Sweep to	Turn	
	Cupboard	Drawer	Safe	Stick	Block	Dustpan	Tap	
Gemini-2.5-pro	38.4 \pm 10.6	94.0 \pm 6.8	93.6 \pm 9.2	73.6 \pm 22.3	63.6 \pm 14.4	48.4 \pm 14.9	99.2 \pm 2.4	
RDD (Ours)	41.2 \pm 17.1	97.2 \pm 3.1	98.4 \pm 3.2	68.0 \pm 25.0	65.2 \pm 14.3	57.2 \pm 29.7	94.0 \pm 5.1	

Table 14: Performance of FAISS nearest neighbor search and RDD time on NVIDIA 4090.

Hardware	Dim	Vec Num	QPS	L_{max}	L_I	RDD Time (s)
NVIDIA 4090	2048	10M	386	100	500	115 (4.3 fps)

SmolVLA Dataset [48]: The emergence of models like SmolVLA, a capable vision-language-action model trained entirely on 23k episodes from 487 open-sourced community datasets through the LeRobot framework, outperforms the closed-source-dataset policy π_0 [4]. **AgiBot World [33]:** AgiBot World provides not just datasets but complete open-source toolchains and standardized data collection pipelines, further enriching the public ecosystem. It has collected over 1 million trajectories on over 100 homogeneous robots. Their proposed model GO-1, entirely trained on this open-sourced dataset, outperforms the closed-source dataset policy π_0 [4].

E Broader Impacts

The potential negative societal impacts of our work are consistent with those commonly observed in robotics research, including risks related to privacy, labor displacement, and unintended misuse in sensitive contexts. While our method is primarily designed to enhance the scalability and efficiency of robotic systems, such advancements may accelerate deployment in real-world settings, amplifying both positive and negative consequences. In parallel, advances in point cloud analysis [49, 50],

cooperative motion prediction [51], autonomous driving frameworks [52, 53, 54, 55], and generative AI for driving [56] highlight both the promise and the risks of deploying increasingly capable vision-action models. Broader surveys of visual foundation models [57] and new work on multimodal alignment [58, 59] further strengthen the importance of trustworthy design and governance, especially for safety-critical applications such as transportation and human-robot interaction [60, 61]. To mitigate these risks, we emphasize alignment with ethical guidelines, including fairness, accountability, transparency, and safety, and encourage interdisciplinary collaboration to monitor societal impacts as these technologies evolve [62].