

STRAP: ROBOT SUB-TRAJECTORY RETRIEVAL FOR AUGMENTED POLICY LEARNING

Marius Memmel^{*,1}, Jacob Berg^{*,1}, Bingqing Chen², Abhishek Gupta^{1†}, Jonathan Francis^{2,3,†}

¹Paul G. Allen School of Computer Science & Engineering, University of Washington

²Robot Learning Lab, Bosch Center for Artificial Intelligence

³Robotics Institute, Carnegie Mellon University

{mommelma, jacob33, abhgupta}@cs.washington.edu,

{bingqing.chen, jon.francis}@us.bosch.com

ABSTRACT

Robot learning is witnessing a significant increase in the size, diversity, and complexity of pre-collected datasets, mirroring trends in domains such as natural language processing and computer vision. Many robot learning methods treat such datasets as multi-task expert data and learn a multi-task, generalist policy by training broadly across them. Notably, while these generalist policies can improve the average performance across many tasks, the performance of generalist policies on any one task is often suboptimal due to negative transfer between partitions of the data, compared to task-specific specialist policies. In this work, we argue for the paradigm of training policies during deployment given the scenarios they encounter: rather than deploying pre-trained policies to unseen problems in a zero-shot manner, we non-parametrically retrieve and train models directly on relevant data at test time. Furthermore, we show that many robotics tasks share considerable amounts of low-level behaviors and that retrieval at the “*sub*”-trajectory granularity enables significantly improved data utilization, generalization, and robustness in adapting policies to novel problems. In contrast, existing full-trajectory retrieval methods tend to underutilize the data and miss out on shared cross-task content. This work proposes STRAP, a technique for leveraging pre-trained vision foundation models and dynamic time warping to retrieve sub-sequences of trajectories from large training corpora in a robust fashion. STRAP outperforms both prior retrieval algorithms and multi-task learning methods in simulated and real experiments, showing the ability to scale to much larger offline datasets in the real world as well as the ability to learn robust control policies with just a handful of real-world demonstrations. Project videos at <https://sites.google.com/view/strappaper/home>

1 INTRODUCTION

Robot learning techniques have shown the ability to shift the process of designing robot controllers from a large manual or model-based process to a data-driven one (Francis et al., 2022; Hu et al., 2023). A particularly promising paradigm is that of end-to-end imitation learning with large neural models (Chi et al., 2023; Haldar et al., 2024), which has shown considerable success with the proliferation of powerful neural architectures such as diffusion models (Chi et al., 2023; Wang et al., 2024) or transformers (Haldar et al., 2024; Zhao et al., 2023), as well as the availability of large, diverse robotics datasets (Khazatsky et al., 2024; Collaboration et al., 2023). While imitation learning can be effective for performing particular tasks,

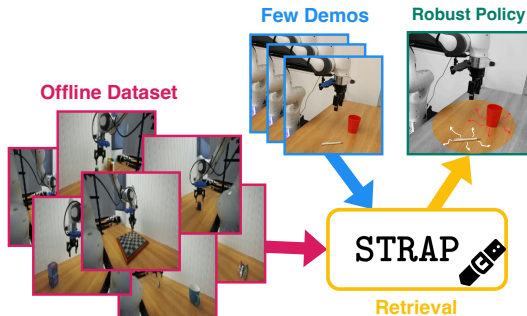


Figure 1: **STRAP**: Sub-trajectory retrieval for training robust policies during deployment.

with targeted in-domain data collection, this process can be expensive and time-consuming in terms of human effort. This becomes a challenge as we deploy robots into dynamic environments such as homes and offices, where new tasks and environments are commonplace and constant data collection is impractical.

Multi-task policy learning is often used in this situation, where data across multiple tasks is used to train a large task- or instruction-conditioned model that has the potential to generalize to new problems. While multi-task learning has seen successes in certain settings (Reed et al., 2022; Brohan et al., 2023), the performance of a multi-task, generalist policy is often lower than task-specific, specialist policies. This can be attributed to the model suffering from negative transfer and sacrificing per-task performance to improve the average performance across tasks. This challenge is exacerbated in unseen tasks or domains since zero-shot generalization is challenging and collecting large amounts of in-domain finetuning data can be expensive. In this work, we consider ways to better use pre-collected datasets and to enable few-shot finetuning of imitation learning models for new tasks.

In particular, we build on the paradigm of *non-parametric data retrieval*, where a small amount of in-domain data collected at test-time is used to retrieve a subset of particularly “relevant” data from the training corpus. This retrieved data can then be used for robust and performant model training and finetuning on new tasks. In this sense, the retrieved data can guide learned models towards desired behavior during test-time deployments; however, the question becomes: *How do we sub-select which data to retrieve from a large, pre-existing corpus?*

Several prior techniques have studied the problem of non-parametric retrieval, from the perspective of learning latent embeddings that encode states (Du et al., 2023), skills (Nasiriany et al., 2022), optical flow (Lin et al., 2024), and learned affordances (Kuang et al.). Most techniques are challenging to apply out of the box for two primary reasons. Firstly, they require training domain-specific encoders to embed states, skills, or affordances: this makes it challenging to apply to demonstrations collected in the open world, where visual appearance can show wide variations. Secondly, they often retrieve entire trajectories, limiting the policies’ ability to use data from other tasks that may share common components with the desired test-time behavior. These challenges limit both the broad applicability of these retrieval methods and the amount of cross-task data sharing. *How can we design easy-to-use off-the-shelf retrieval methods that maximally utilize the training data for test-time adaptation?*

The key insight in this work is that retrieval methods do not need to measure the similarity between entire trajectories (or individual states), but rather between *sub-trajectories* of the desired behavior at test-time and corresponding sub-trajectories of the training data. Notably, these sub-trajectories do not need to come from tasks that are similar in entirety to the desired test-time tasks. Instead, sub-components of many related tasks can be shared to enable robust, test-time policy training. For example, as shown in Fig. 1, for the multi-stage task of “*pick up the mug, put it in the drawer, and close it*”, both “*pick up the mug, put in on top of the drawer*” and “*close the bottom drawer, open the top drawer*” contain sub-tasks that when retrieved provide useful training data. Our proposed method, **Sub-sequence Trajectory Retrieval for Augmented Policy Learning (STRAP)**, uses a small amount of in-domain trajectories collected at test-time to retrieve and train on these relevant sub-trajectories across a large multi-task training corpus. The resulting policies show considerable improvements in robustness and generalization over previous retrieval methods, zero-shot multi-task policies, or policies that are trained purely on test-time in-domain data.

We show how STRAP can be used with minimal effort across training and evaluation domains with non-trivial visual differences. Our method first compares sub-trajectory similarity using features from off-the-shelf foundation models, *e.g.*, DINOv2 (Oquab et al.); these features capture strong notions of “object-ness”, discarding spurious visual differences such as lighting, texture, and local changes in object appearance. Secondly, our method leverages time-invariant alignment techniques, such as dynamic time warping (Giorgino, 2009), to compute the similarity between sub-trajectories of different lengths, removing requirements for retrieved trajectories to have a similar length and increasing the applicability of STRAP across tasks and domains. Lastly, we show how STRAP can be applied to arbitrary test corpora, with sub-trajectories being automatically extracted by our framework, thereby removing the requirement for manual segmentation of relevant sub-trajectories from the training corpus. We demonstrate how STRAP can be used out of the box to augment *any* few-shot imitation learning algorithm, providing significant gains in generalization at test-time, while avoiding expensive, test-time in-domain data collection. We instantiate STRAP with transformer-

based imitation learning policies and show the benefits of few-shot sub-trajectory retrieval on the LIBERO (Liu et al., 2024) benchmark in simulation and real-world imitation learning problems.

2 RELATED WORK

Retrieval for Behavior Replay: A considerable body of work has explored retrieval-based approaches for robotic manipulation, where the retrieval of relevant past demonstrations aids in replaying past experiences. The choices of embedding space hereby range from off-the-shelf models (Di Palo & Johns, 2024; Malato et al., 2024) like DINO (Caron et al., 2021), training encoders on the offline dataset (Pari et al., 2022) to abstract representation like object shapes (Sheikh et al.). Some works do not directly replay actions but add a layer of abstraction following sub-goals (Zhang et al., 2024), affordances (Kuang et al.) or keypoints (Papagiannis et al.). A key assumption of these methods is that the offline data either exactly resembles expert demonstrations collected in the test environment or that intermediate representations can bridge the gap. These drawbacks limit the usage of large multi-task datasets collected in various domains.

Retrieval for Few-shot Imitation Learning: Retrieval for policy learning tries to mitigate these issues by learning policies from the retrieved data. While retrieval has shown to benefit policy learning from sub-optimal single-task data (Yin & Abbeel, 2024), most work focuses on retrieving from large multi-task datasets like DROID (Khazatsky et al., 2024) or OpenX (Collaboration et al., 2023) containing expert demonstrations. BehaviorRetrieval (BR) (Du et al., 2023) and FlowRetrieval (FR) (Lin et al., 2024) train an encoder-decoder model on state-action and optical flow respectively. Related to our work, SAILOR (Nasiriany et al., 2022) imposes skill constraints on the embedding space, clustering similar skills together to later retrieve those. A significant downside of training custom representations is that these methods do not scale well to the increasing size of available offline datasets and are unable to deal with significant visual and semantic differences. Moreover, techniques like BehaviorRetrieval and FlowRetrieval retrieve individual states, rather than sub-trajectories like our work, where sub-trajectory retrieval enables maximal data sharing between seemingly different tasks while capturing temporal information.

Learning from Sub-trajectories: Several works propose to decompose demonstrations into reusable sub-trajectories, e.g., based on end-effector-centric or full proprioceptive state-action transitions (Li et al., 2020; Belkhale et al., 2024; Shankar et al., 2022; Myers et al., 2024; Francis et al., 2022). Belkhale et al. (2024) propose to decompose demonstrations into end-effector-centric sub-tasks, e.g., "move forward" or "rotate left". The authors show that by decomposing and re-labeling the language instructions into a shared vocabulary, knowledge from multi-task datasets can be better shared when training multi-task policies. Myers et al. (2024) leverage VLMs to decompose demonstrations into sub-trajectories to better learn to imitate them. To our knowledge, we propose the first robot sub-trajectory retrieval mechanism, for partitioning large offline robotics datasets and for enabling cross-task positive transfer during policy learning.

3 PRELIMINARIES

3.1 DYNAMIC TIME WARPING

To match sequences of potentially variable length during retrieval, we build on an algorithm called dynamic time warping (DTW) (Müller, 2021). DTW methods compute the similarity between two time series that may vary in time or speed, e.g., different video or audio sequences. This algorithm aligns the varying length sequences by warping the time axis of the series using a set of step sizes to minimize the distance between corresponding points while obeying boundary conditions.

DTW algorithms are given two sequences, $X = \{x_1, x_2, \dots, x_n\}$ and $Y = \{y_1, y_2, \dots, y_m\}$, where $m \neq n$, and a corresponding cost matrix $C(x_i, y_j)$ that assigns the cost of assigning element x_i of sequence X to correspond with element y_j of sequence Y . The goal of DTW is to find a mapping between X and Y that minimizes the total cumulative distance between the assigned elements of both sequences while obeying boundary and continuity conditions. Dynamic time warping methods solve this problem efficiently using dynamic programming methods.

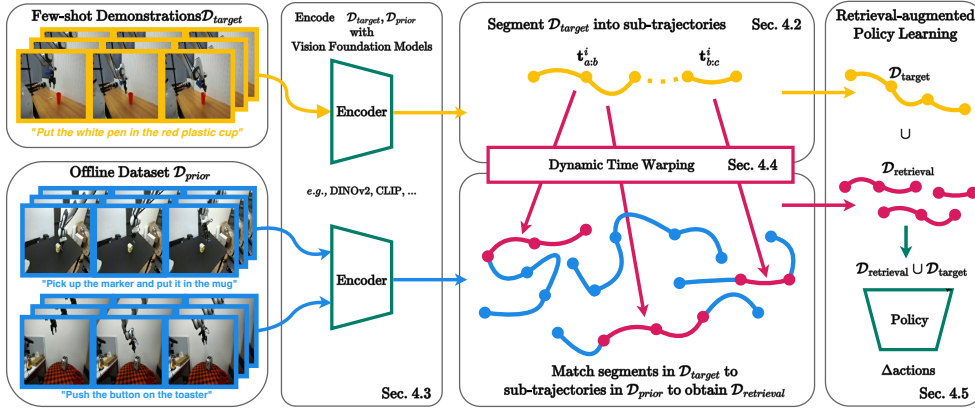


Figure 2: **Overview of STRAP:** 1) demonstrations $\mathcal{D}_{\text{target}}$ and offline datasets $\mathcal{D}_{\text{prior}}$ are encoded into a shared embedding space using a vision foundation model, 2) automatic slicing generates sub-trajectories which 3) S-DTW matches to corresponding sub-trajectories in $\mathcal{D}_{\text{prior}}$ creating $\mathcal{D}_{\text{retrieval}}$, 4) training a policy on the union of $\mathcal{D}_{\text{retrieval}}$ and $\mathcal{D}_{\text{target}}$ results in better performance and robustness.

A cumulative distance matrix D is computed via dynamic programming as follows: $D(0, 0) = C(0, 0)$, $D(n, 1) = \sum_{k=1}^n C(k, 1)$ for $n \in [1 : N]$ and $D(1, m) = \sum_{k=1}^m C(1, k)$ for $m \in [1 : M]$. Then the following dynamic programming calculation is performed:

$$D(i, j) = C(x_i, y_j) + \min\{D(i-1, j), D(i, j-1), D(i-1, j-1)\}, \quad (1)$$

where $C(x_i, y_j)$ is the distance between points x_i and y_j . We assume this cost matrix is pre-provided, and we describe how we compute this from raw camera images in Sec. 4.3. The optimal alignment between the sequences is found by backtracking from $D(n, m)$ to $D(0, 0)$, which gives the minimal alignment obeying the boundary condition. This guarantees that the start is matched to the start and the end is matched to the end or that the pairs (x_0, y_0) and (x_n, y_m) are the start and end of the path. This optimal pairing path consists of the best possible alignment between X and Y such that the cumulative cost between all matched pairs is minimized. DTW, as described, is widely used in time-series analysis, speech recognition, and other domains where temporal variations exist between sequences. In the context of our retrieval problem, DTW is used to go beyond retrieving exactly matched sequences to matching variable length subsequences, as we describe below.

Subsequence dynamic time warping (S-DTW) is an extension of the DTW algorithm for scenarios where a shorter query sequence must be matched to a portion of a longer reference sequence. Given a query sequence $X = \{x_1, x_2, \dots, x_n\}$ and a much longer reference sequence $Y = \{y_1, y_2, \dots, y_m\}$, the goal of S-DTW is to find a subsequence of Y (of a potentially different length from X), denoted $Y_{i:j}$ where $i \leq j$, that has the minimal DTW distance to X .

The cumulative cost matrix D for S-DTW is computed similarly to the traditional DTW described above, but with the distinction that it allows alignment to start and end at any point in R . D is initialized as $D(0, 0) = C(0, 0)$, $D(n, 1) = \sum_{k=1}^n C(k, 1)$ for $n \in [1 : N]$ and $D(1, m) = C(1, m)$ for $m \in [1 : M]$ and then completed using dynamic programming following Eq. (1).

This ensures that the query can match any sub-sequence of the reference. Once the cumulative cost matrix is computed, the optimal alignment is found by backtracking from the minimal value in the last row of the matrix, i.e., $\min(D(n, j))$ for $j \in \{1, \dots, m\}$. This gives the subsequence of Y that best aligns with X , obeying only temporality while relaxing the boundary condition. As we will show, using S-DTW for data retrieval enables the maximal retrieval of data across tasks in a retrieval-augmented policy training setting, as described in Sec. 4.3.

4 STRAP: SUB-SEQUENCE ROBOT TRAJECTORY RETRIEVAL FOR AUGMENTED POLICY TRAINING

In NLP, retrieval is a well-established paradigm for retrieving samples relevant to test-time scenarios using non-parametric similarity matching. These retrieval methods (Huang & Huang, 2024) can

retrieve high-quality, relevant data from the training corpus and use this to *target* the model for the particular test-time scenario with in-context learning. Informally, this transitions the model from being a jack-of-all-trades to a master-of-one. How can we adapt such a paradigm to policy learning for robotics? Prior retrieval methods (Lin et al., 2024; Du et al., 2023; Kuang et al.; Nasiriany et al., 2022) often require domain-specific training, or underutilize the training data.

To address these challenges, we present STRAP, a scalable retrieval method that can target models to particular test-time distributions by finding semantically similar sub-trajectories through the use of subsequence dynamic time warping with metric spaces defined by features from off-the-shelf vision foundation models. Doing so makes STRAP both visually robust and able to maximally utilize relevant portions of the training data. Our key insights are as follows:

1. *Vision foundation models* offer powerful out-of-the-box representations for trajectory retrieval. They sufficiently encode scene semantics and offer visual robustness in contrast to brittle in-domain feature extractors from prior work.
2. *Sub-trajectory retrieval* can enable maximal re-use of prior data while capturing temporal information about tasks and dynamics.
3. Performing retrieval via *subsequence dynamic time warping* can find optimal sub-trajectory matches in offline datasets that are agnostic to segment length task horizon or fluctuations in demonstration frequency.

4.1 PROBLEM SETTING: RETRIEVAL-AUGMENTED POLICY LEARNING

We consider a few-shot learning setting where we’re given a target dataset $\mathcal{D}_{\text{target}} = \{(s_0^i, a_0^i, s_1^i, a_1^i, \dots, s_{H_i}^i, a_{H_i}^i, l^i)\}_{i=1}^N$ containing expert trajectories of states s (e.g., observations like camera views o and proprioception x), actions a (such as robot controls), and task-specifying language instructions l . This target dataset is collected in the test environment and task, but there is only a small set of N trajectories, which limits generalization for models trained purely on such a small dataset. Since $\mathcal{D}_{\text{target}}$ is often insufficient to solve the task alone, we posit that generalization can be accomplished by non-parametrically *retrieving* data from an offline dataset $\mathcal{D}_{\text{prior}}$. This offline dataset $\mathcal{D}_{\text{prior}} = \{(s_0^j, a_0^j, s_1^j, a_1^j, \dots, s_{H_j}^j, a_{H_j}^j, l^j)\}_{j=1}^M$ can contain data from different environments, scenes, levels of expertise, tasks, or embodiments. Notably, the set of tasks in the offline dataset do *not* need to overlap with the set of tasks in the target dataset. We assume that the offline dataset shares matching embodiment with the target dataset and consists of expert-level trajectories, but may consist of a diversity of scenes and tasks that vary widely from the target dataset $\mathcal{D}_{\text{target}}$.

Given $\mathcal{D}_{\text{prior}}$ and $\mathcal{D}_{\text{target}}$, the goal is to learn a language-conditioned policy $\pi_\theta(a|s, l)$ that can predict optimal actions a in the target environment when prompted with the current state s and language instruction l . Assuming we can obtain a measure of success (such as task completion), and a broad set of initial conditions $s_0 \sim \rho_{\text{test}}(s_0)$ in the test environment. The objective of policy learning is to determine the policy parameters θ to maximize the expected success metric when evaluated on test conditions, under the policy π_θ and test-time environment dynamics. Since we are only provided a limited corpus of data, $\mathcal{D}_{\text{target}}$, in the target domain, these policy parameters cannot be learned by simply performing maximum likelihood on $\mathcal{D}_{\text{target}}$. Instead, we will present an approach where a smaller, “relevant” subset of the offline dataset $\mathcal{D}_{\text{retrieval}} \subseteq \mathcal{D}_{\text{prior}}$ is retrieved non-parametrically and then mixed with the smaller in-domain dataset $\mathcal{D}_{\text{target}}$ to construct a larger, augmented training dataset, i.e., $\mathcal{D}_{\text{target}} \cup \mathcal{D}_{\text{retrieval}}$, which is most relevant to the desired test-time conditions $\rho_{\text{test}}(s_0)$. This can then be used for training policies via imitation learning, as we will describe in Sec. 4.5. Doing so avoids an expensive generalist training procedure and rather focuses the learned model to being a high-performing specialist in a particular setting of interest. The key questions becomes - *How can we define what subset of the offline dataset $\mathcal{D}_{\text{prior}}$ is relevant to construct $\mathcal{D}_{\text{retrieval}}$?*

To handle the unique nature of robotic data, e.g., multi-modal and temporally dependent, we design STRAP for retrieval-augmented policy learning. Firstly, we need to define the unit of retrieval. Rather than retrieving individual state-action pairs or entire trajectories, STRAP crucially retrieves sub-trajectories. We also propose a method to automatically segment trajectories in $\mathcal{D}_{\text{target}}$ into such sub-trajectories (Sec. 4.2). Secondly, we need to define a suitable distance metric for a pair of sub-trajectories (Sec. 4.3). Then, we need a computationally efficient algorithm to retrieve relevant sub-

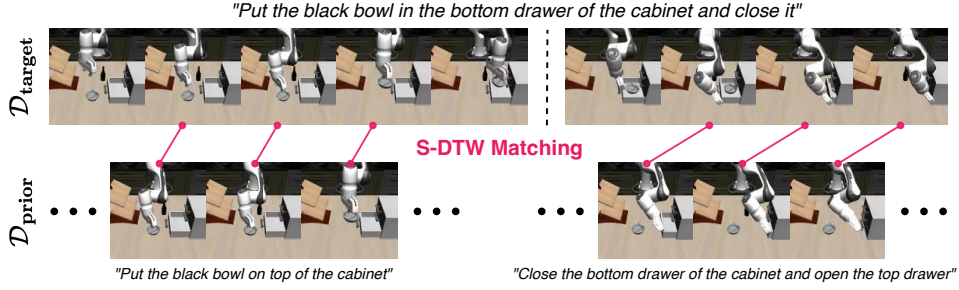


Figure 3: **Sub-trajectory matching:** S-DTW matches the sub-trajectories of $\mathcal{D}_{\text{target}}$ (top) to the relevant segments in $\mathcal{D}_{\text{prior}}$. A feature of S-DTW is that the start and end of the trajectories do not have to align, finding optimal matches for each pairing.

trajectories non-parametrically from the training set (Sec. 4.4). Finally, we put everything together and train policies based on retrieved data (Sec. 4.5).

4.2 SUB-TRAJECTORIES FOR DATA RETRIEVAL

To make the best use of the training dataset, while capturing temporal task-specific dynamics, we expand the notion of retrieval from being able to retrieve entire trajectories or single states to retrieving variable-length sub-trajectories. In doing so, retrieval can capture the temporal dynamics of the task, while still being able to share data between seemingly different tasks with potentially different task instruction labels. In particular, we define a sub-trajectory as a consecutive subset of a trajectory $t_{a:b}^i \subseteq T^i$ with the sub-trajectory $t_{a:b}^i = (s_a^i, s_{a+1}^i, \dots, s_b^i)$ including timestep a to b of the whole trajectory T^i of length H_i . Most long-horizon problems observed in robotics datasets (Liu et al., 2024; Khazatsky et al., 2024; Collaboration et al., 2023) naturally contain multiple such sub-trajectories. For instance, the task shown in Eq. 3 can be decomposed into “put the bowl in the drawer” and “close the drawer”. Note that we do not require these trajectories to explicitly have a specific semantic meaning, but semantically meaningful sub-trajectories often coincide with those most commonly encountered across tasks as we see in our experimental evaluation.

Given this definition of a sub-trajectory, our proposed retrieval technique only requires segmenting the target demonstrations into sub-trajectories $\mathcal{T}_{\text{target}} = \{t_{1:a}^i, t_{a:b}^i, \dots, t_{H_i-p_i:H_i}^i, \forall T^i \in \mathcal{D}_{\text{target}}\}$ but *not* the much larger offline training dataset $\mathcal{D}_{\text{prior}}$. Instead, appropriate sub-sequences will be retrieved from this dataset using a DTW based retrieval algorithm (Sec. 4.4). This makes the proposed methodology far more practical since $\mathcal{D}_{\text{prior}}$ is much larger than $\mathcal{D}_{\text{target}}$. While this separation into sub-trajectories can be done manually during data collection, we propose an automatic technique for sub-trajectory separation that yields promising empirical results. Building on techniques proposed by Belkhale et al. (2024), we split the demonstrations into atomic chunks, *i.e.*, lower-level motions, before retrieving similar sub-trajectories with our matching procedure (Sec. 4.4). In particular, we propose a simple proprioception-based segmentation technique that optimizes for changes in the robot’s end-effector motion indicating the transition between two chunks. For example, a Pick&Place task can be split into picking and placing separated by a short pause when grasping the object. Let x_t be a vector describing the end-effector position at timestep t . We define “transition states” where the absolute velocity drops below a threshold: $\|\dot{x}\| < \epsilon$ ¹. We empirically find that this proprioception-driven segmentation can perform reasonable temporal segmentation of target trajectories into sub-components. This procedure can certainly be improved further via techniques in action recognition using vision-foundation models (Team et al., 2023), or information-theoretic segmentation methods (Jiang et al., 2022).

4.3 FOUNDATION MODEL-DRIVEN RELEVANCE METRICS FOR RETRIEVAL

Given the definition and automatic segmentation of sub-trajectories, we must define a measure of similarity that allows for the retrieval of appropriate *relevant* sub-trajectory data from $\mathcal{D}_{\text{prior}}$, and at the same time is robust to variations in visual appearance, distractors, and irrelevant spurious

¹For trajectories involving “stop-motion”, this heuristic returns many short chunks as the end-effector idles, waiting for the gripper to close. To ensure a minimum length, we merge neighboring chunks until all are ≥ 20 .

Algorithm 1 STRAP ($\mathcal{D}_{\text{target}}, \mathcal{D}_{\text{prior}}, K, \epsilon, \mathcal{F}$)

Require: demos $\mathcal{D}_{\text{target}}$, offline dataset $\mathcal{D}_{\text{prior}}$, vision foundation model \mathcal{F} , # retrieved chunks K , chunking threshold ϵ ;

- 1: */* Pre-processing */*
- 2: $\mathcal{T}_{\text{target}} \leftarrow \text{SubTrajSegmentation}(\mathcal{D}_{\text{target}}, \epsilon)$; ▷ Heuristic demo chunking
- 3: $\mathcal{E}_{\text{prior}} \leftarrow \{\{\mathcal{F}(o_t)|o_t \in T\}|T \in \mathcal{D}_{\text{prior}}\}$; ▷ Embed $\mathcal{D}_{\text{prior}}$
- 4: $\mathcal{E}_{\text{target}} \leftarrow \{\{\mathcal{F}(o_t)|o_t \in T\}|T \in \mathcal{T}_{\text{target}}\}$; ▷ Embed chunked $\mathcal{D}_{\text{target}}$
- 5: */* Sub-trajectory Retrieval using S-DTW*/*
- 6: **for** $S_{\text{target}} \in \mathcal{D}_{\text{target}}$ **do**
- 7: $\mathcal{M} \leftarrow []$; ▷ Initialize empty match storage
- 8: **for** $T_{\text{prior}} \in \mathcal{D}_{\text{prior}}$ **do**
- 9: $D \leftarrow \text{computeCostMatrix}(\mathcal{E}_{\text{target}}, \mathcal{E}_{\text{prior}})$; ▷ Eq. (2)
- 10: $\mathcal{M}_{i,j} \leftarrow \text{extractSubTrajectory}(D, T_{\text{prior}})$; ▷ Dynamic Programming
- 11: **end for**
- 12: **end for**
- 13: $\mathcal{D}_{\text{retrieval}} \leftarrow \text{retrieveTopKMatches}(\mathcal{M}, K)$; ▷ Sec. 4.4
- 14: */* Policy Learning */*
- 15: **repeat**
- 16: sample $\mathcal{B} \sim \mathcal{D}_{\text{target}} \cup \mathcal{D}_{\text{retrieval}}$ to update policy π_{θ} with loss $\mathcal{L}(\mathcal{B}; \theta)$ ▷ Eq. (3)
- 17: **until** π_{θ} converged; **return** π_{θ}

features. While prior work has suggested objectives to train such similarity metrics through representation learning (Du et al., 2023; Lin et al., 2024; Kuang et al.), these methods are often trained purely in-domain, making them particularly sensitive to aforementioned variations. While using more lossy similarity metrics based on optical flow (*c.f.* (Lin et al., 2024)) or language (Zha et al., 2024) can help with this fragility, it often fails to capture the necessary task-specific or semantic details. This suggests the need for a robust, domain-agnostic similarity metric that can easily be applied out-of-the-box.

In this work, we will adopt the insight that vision(-language) foundation models (Oquab et al.; Radford et al., 2021) offer off-the-shelf solutions to this problem of measuring the semantic and visual similarities between sub-trajectories, capturing object- and task-centric affordances, while being robust to low-level variations in scene appearance. Trained on web-scale real-world image(-text) data, these models are typically robust to low-level perceptual variations, while providing semantically rich representations that naturally capture a notion of object-ness and semantic correspondence. Denoting a vision foundation model as $\mathcal{F}(\cdot)$, we can compute the pairwise distance of two camera views with an L2 norm² in embedding space, *i.e.*, $\|\mathcal{F}(o_i) - \mathcal{F}(o_j)\|_2$. While aggregation methods such as temporal averaging could be used to go from embedding of a single image to that of a sub-trajectory, they lose out on the actions and dynamics. We instead opt for a sub-trajectory matching procedure based on the idea of DTW (Giorgino, 2009) and use the embeddings for finding maximally relevant sub-trajectories. Given two sub-trajectories, t_i and t_j , we compute a pairwise cost matrix $C \in \mathbb{R}^{|t_i| \times |t_j|}$, where its value is as computed by:

$$C(i, j) = \|\mathcal{F}(o_i) - \mathcal{F}(o_j)\|_2 \quad (2)$$

4.4 EFFICIENT SUB-TRAJECTORY RETRIEVAL WITH SUBSEQUENCE DYNAMIC TIME WARPING

Given the above-mentioned definitions of sub-trajectories and foundation-model-driven similarity metrics, we instantiate an algorithm to find the K most relevant sub-trajectories $\mathcal{T}_{\text{match}}$ from the offline dataset $\mathcal{D}_{\text{prior}}$ for each sub-trajectory t segmented from $\mathcal{D}_{\text{target}}$. Sub-trajectories may have variable lengths and temporal positioning within a trajectory caused by varying tasks, platforms, or demonstrators. We employ S-DTW to match the target sub-trajectories $\mathcal{T}_{\text{target}}$ to appropriate segment $\mathcal{T}_{\text{match}}$ in $\mathcal{D}_{\text{prior}}$ (Sec. 3.1). S-DTW scales naturally with these challenges and allows for retrieval from diverse, multi-task datasets. On deployment, subsequence dynamic time warping accepts a query sub-sequences from the target dataset, *i.e.*, t_{target} , and uses dynamic programming to

²Other cost metrics such as (1-cosine similarity) could be used here as well.

compute matches that are maximally aligned with the query $\mathcal{T}_{\text{match}} = \{\text{SDTW}(t, \mathcal{D}_{\text{prior}}), \forall t \in \mathcal{T}_{\text{target}}\}$ along with matching costs, D . To construct $\mathcal{D}_{\text{retrieval}}$, we select the K matches with the lowest cost uniformly across the sub-trajectories in $\mathcal{T}_{\text{target}}$, *i.e.*, the same number of matches for each query until K matches are retrieved. We note that the resulting set of matches can contain duplicates if the demonstrations share similar chunks, but argue that if a chunk occurs multiple times in the demonstrations, it is important to the task and should be “*up-weighted*” in the training set – in this case through duplicated retrieval. For each match, we also retrieve its corresponding language instruction. The training dataset then contains a union of the target dataset $\mathcal{D}_{\text{target}}$ and the retrieved dataset $\mathcal{D}_{\text{retrieval}}$, $\mathcal{D}_{\text{target}} \cup \mathcal{D}_{\text{retrieval}}$. This significantly larger, retrieval-augmented dataset can then be used to learn policies via imitation learning, leading to robust, generalizable policies as we describe below.

4.5 PUTTING IT ALL TOGETHER: STRAP

To start the retrieval process, we encode image observations in $\mathcal{D}_{\text{target}}$ and $\mathcal{D}_{\text{prior}}$ using a vision foundation model, *e.g.*, DINOv2 (Oquab et al.) or CLIP (Radford et al., 2021). To best leverage the multi-task trajectories in $\mathcal{D}_{\text{prior}}$, we split the demonstrations in $\mathcal{D}_{\text{target}}$ into atomic chunks based on a low-level motion heuristic. Then we generate matches between chunked $\mathcal{D}_{\text{target}}$ and $\mathcal{D}_{\text{prior}}$ and construct $\mathcal{D}_{\text{retrieval}}$ by selecting the top K matches uniformly across all chunks. Combining $\mathcal{D}_{\text{retrieval}}$ with $\mathcal{D}_{\text{target}}$ forms our dataset for learning a policy. In a standard policy learning setting, noisy retrieval data can lead to negative transfer, *e.g.*, when observations similar to the target data are labeled with actions that achieve a different task. Without conditioning, such contaminated samples hurt the policy’s downstream performance. We propose to use a language-conditioned policy to deal with this inconsistency. With conditioning, the policy can distinguish between samples from different tasks, separating misleading from expert actions while benefiting from positive transfer from the additional training data and context of the language conditioning.

We use behavior cloning (BC) to learn a visuomotor policy π similar to Haldar et al. (2024); Nasiriany et al. (2024). We choose a transformer-based (Vaswani, 2017) architecture feeding in a history of the last h observations $s_{t-h:t}$ and predicting a chunk of h future actions using a Gaussian mixture model action head. We sample batches from the union of $\mathcal{D}_{\text{target}}$ and $\mathcal{D}_{\text{retrieval}}$, as in $\mathcal{B} \sim \mathcal{D}_{\text{target}} \cup \mathcal{D}_{\text{retrieval}}$. As proposed in Haldar et al. (2024) we compute the multi-step action loss and add an L2 regularization term over the model weights θ , resulting in the following loss function:

$$\mathcal{L}(\mathcal{B}; \theta) = \frac{1}{|\mathcal{B}|} \sum_{(s_{i-h:i}, a_{i:i+h}, l) \in \mathcal{B}} -\log(\pi_{\theta}(a_{i:i+h} | s_{i-h:i}, l)) + \lambda \|\theta\|_2^2 \quad (3)$$

with policy π_{θ} and hyperparameter λ controlling the regularization.

5 EXPERIMENTS AND RESULTS

5.1 EXPERIMENTAL SETUP

Task Definition: We demonstrate the efficacy of STRAP in simulation on the LIBERO benchmark (Liu et al., 2024), and on a Pen-in-Cup manipulation task with a real world robot arm. Eq. 12 shows the target tasks and samples from the retrieval datasets.

- **LIBERO:** We evaluate STRAP on 10 long-horizon tasks of the LIBERO benchmark (Liu et al., 2024) which includes diverse objects, layouts, and backgrounds. The evaluation environments randomize the target object poses, providing an ideal test bed for robustness. We use the agent view (exocentric) observations for the retrieval and train policies on both agent view and in-hand observations. To showcase the benefits of sub-trajectory retrieval, we choose the 10 long-horizon tasks (LIBERO-10) as $\mathcal{D}_{\text{target}}$ and retrieve data from the other 90 short-horizon tasks (LIBERO-90), $\mathcal{D}_{\text{prior}}$. The following section features 5 tasks covering a variety of objectives and skills. While these tasks benefit the most from retrieval, we report results on the remaining ones in the appendix (Tab. 3). The tasks descriptions are as follows: *Stove-Moka* combines knob-turning and Pick&Place, *Bowl-Cabinet* combines Pick&Place with cabinet closing, *Soup-Cheese* and *Mug-Mug* both contain two consecutive Pick&Place tasks, and *Book-Caddy* involves Pick&Place and insertion. Each task comes with 50 demonstrations from which we select 5 random demonstrations in a few-shot imitation learning setting and retrieve data from all LIBERO-90 tasks, which amounts to 4500 total offline demonstrations.

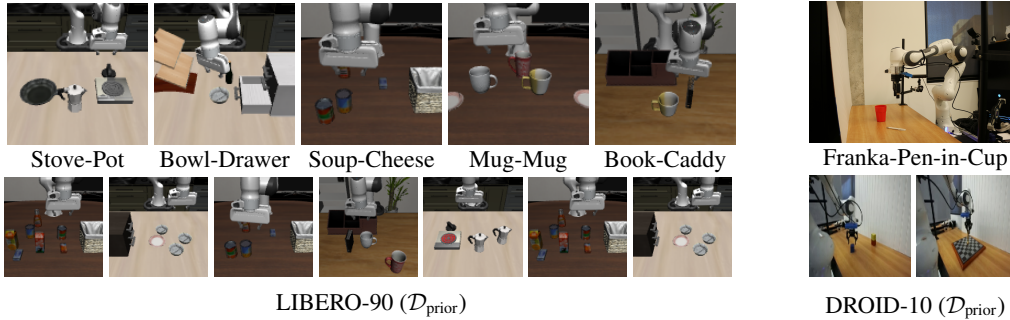


Figure 12: **Simulation and real-world tasks:** $\mathcal{D}_{\text{target}}$ tasks from LIBERO-10 and real-world Franka-Pen-in-Cup (top) and retrieval dataset $\mathcal{D}_{\text{prior}}$ (bottom).

Table 1: **Baselines:** Performance of baselines, ablations and variations of STRAP on the LIBERO 10 tasks (Eq. 12). DINOv2 and CLIP features perform similarly, making STRAP flexible in the encoder choice. **Bold** indicates best and underline runner-up results.

Task	Stove-Pot	Bowl-Cabinet	Soup-Cheese	Mug-Mug	Book-Caddy
BC	77.33% \pm 4.35	71.33% \pm 5.68	27.33% \pm 2.18	38.00% \pm 5.66	75.33% \pm 1.44
MT	0.00% \pm 0.00	0.00% \pm 0.00	0.00% \pm 0.00	0.00% \pm 0.00	88.00% \pm 1.89
BR (Du et al., 2023)	80.0% \pm 1.63	72.0% \pm 7.72	26.0% \pm 5.25	40.0% \pm 8.64	16.0% \pm 1.89
FR (Lin et al., 2024)	76.0% \pm 6.60	54.67% \pm 11.98	24.67% \pm 8.55	29.33% \pm 1.44	52.0% \pm 5.89
D-S	70.67% \pm 7.85	65.33% \pm 1.96	18.0% \pm 3.40	16.0% \pm 0.94	57.33% \pm 2.88
D-T	78.67% \pm 2.72	75.33% \pm 2.72	37.33% \pm 6.62	63.33% \pm 3.57	79.00% \pm 4.95
STRAP (CLIP)	86.00% \pm 4.10	90.67% \pm 2.18	<u>42.00% \pm 0.94</u>	54.67% \pm 3.31	83.33% \pm 3.03
STRAP (DINOv2)	85.33% \pm 2.18	91.33% \pm 2.18	42.67% \pm 7.20	57.33% \pm 7.68	<u>85.33% \pm 2.81</u>

- **Franka-Pen-in-Cup:** To demonstrate the efficacy of STRAP in a real-world setting, we solve a Pen-In-Cup task using the Franka Emika Panda robot. $\mathcal{D}_{\text{target}}$ contains 3 demonstrations of picking a pen and putting it in a cup next to it Eq. 12. $\mathcal{D}_{\text{prior}}$ consists of 100 demonstrations across 10 tasks in the same tabletop environment collected on the DROID (Khazatsky et al., 2024) hardware setup. For task details please refer to Appendix A.2. For retrieval, we average the embeddings per time-step across the left, right, and in-hand camera observations while training the policies on all three image observations.

Baselines and Ablation: We compare STRAP to the following baselines and ablations and refer the reader to Appendix A.1 for implementation details and Appendix A.3 for extensive ablations.

- **Behavior Cloning (BC)** behavior cloning using a transformer-based policy trained on $\mathcal{D}_{\text{target}}$;
- **Multi-task Policy (MT)** transformer-based policy trained on $\mathcal{D}_{\text{prior}}$;
- **BR** (BehaviorRetrieval) (Du et al., 2023) prior work that trains a VAE on state-action pairs for retrieval and uses cosine similarity to retrieve single state-action pairs;
- **FR** (FlowRetrieval) (Lin et al., 2024) same setup as BR but VAE is trained on pre-computed optical flow from GMFlow (Xu et al., 2022);
- **D-S** (DINO state) same as BR and FR but uses off-the-shelf DINOv2 (Oquab et al.) features instead of training a VAE;
- **D-T** (DINO trajectory) retrieves *full* trajectories (rather than sub-trajectories) with S-DTW and DINOv2 features;

5.2 EXPERIMENTAL EVALUATION

Our evaluation aims to address the following questions: (1) Does *sub-trajectory retrieval* improve performance in few-shot imitation learning? (2) How effective are the representations from *vision-foundation models* for retrieval? (3) What types of matches are identified by *S-DTW*?

Does sub-trajectory retrieval improve performance in few-shot imitation learning? STRAP outperforms the retrieval baselines BR and FR on average by +12.20% and +12.47% across all 10 tasks (Tab. 1). These results demonstrate the policy’s robustness to varying object poses. BC represents a strong baseline on the LIBERO task as the benchmark’s difficulty comes from pose variations during evaluation. By memorizing the demonstrations, BC achieves high success rates, outperforming BR and FR by +4.53% and +4.80% across all 10 tasks. The multi-task baseline trained on LIBERO-90

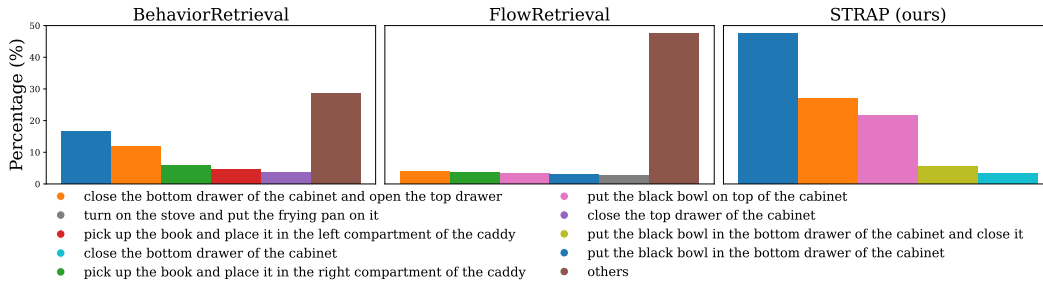


Figure 13: **Tasks distribution** in $\mathcal{D}_{\text{retrieval}}$ for different retrieval methods with target task “put the black bowl in the bottom drawer of the cabinet and close it”.

struggles to generalize to unseen language instructions, failing on 9/10 tasks, only succeeding on the one with an almost exact match in LIBERO-90 (*c.f.* Tab. 1). To prove that the robustness benefits are not unique to the LIBERO benchmark we perform a real-world evaluation in Tab. sec. 5.2. While BC and STRAP solve the Franka-Pen-in-Cup demonstrated in $\mathcal{D}_{\text{target}}$ (*base*), BC lacks robustness to out-of-distribution (*OOD*) scenarios. The policy replays the trajectories observed in $\mathcal{D}_{\text{target}}$. STRAP retrieves relevant sub-trajectories from $\mathcal{D}_{\text{prior}}$, *e.g.*, the robot putting the screwdriver in the cup or picking up pens in various poses. Augmented policy learning then distills this knowledge into a policy, resulting in generalization to an *OOD* scenario.

To further investigate the efficacy of sub-trajectories, we compare sub-trajectory retrieval with S-DTW (STRAP) to retrieving full trajectories with S-DTW (D-T) in Tab. 1. We find sub-trajectory retrieval to improve performance by +4.17% across all 10 tasks. We hypothesize that full trajectories can contain segments irrelevant to the task, effectively hurting performance and reducing the accuracy of the cumulative cost.

Pen-in-Cup	<i>base</i>		<i>OOD</i>	
	Pick	Place	Pick	Place
BC	100%	100%	0%	0%
STRAP	100%	90%	100%	100%

Table 2: **Real-world results:** Franka-Pen-in-Cup task

How effective are the representations from

vision-foundation models for retrieval? Next, we ablate the choice of foundation model representation in STRAP. We compare CLIP, a model trained through supervised learning on image-text pairs, with DINOv2, a self-supervised model trained on unlabeled images. We don’t find any representation to significantly outperform the other with DINOv2 separated from CLIP by only +0.73% across all 10 tasks. To show the efficacy of vision-foundation models for retrieval, we replace the in-domain feature extractors from prior work (BR, FR) trained on $\mathcal{D}_{\text{prior}}$ with an off-the-shelf DINOv2 encoder model (D-S). Comparing them in their natural configuration, *i.e.*, state-based retrieval using cosine similarity, allows for a side-by-side comparison of the representations. Tab. 1 shows the choice of representation to depend on the task with no method outperforming the others on all tasks. Since D-S has no notion of dynamics and task semantics due to single-state retrieval, BR and FR outperform it by +5.00% and +4.73%, respectively. We want to highlight that vision foundation models don’t have to be trained on $\mathcal{D}_{\text{prior}}$ and, therefore, scale much better with increasing amounts of trajectory data and on unseen tasks.

What types of matches are identified by S-DTW? To understand what data STRAP retrieves, we visualize the distribution over tasks as a function of $\mathcal{D}_{\text{retrieval}}$ proportion in Figure 13. The figure visualizes the top five tasks retrieved and accumulates the rest into the “others” category. It becomes clear that STRAP retrieves semantically relevant data – each task shares at least one sub-task with the target task. For example, “put the black bowl in the bottom drawer of the cabinet”, “close the bottom drawer of the cabinet ...” (Eq. 3). Furthermore, STRAP’s retrieval is sparse, only selecting data from 5/90 semantically relevant tasks and ignoring irrelevant ones. We observe that DINOv2 features are surprisingly agnostic to different environment textures, retrieving data from the same task but in a different environment (*c.f.* Eq. 13, “put the black bowl in the bottom drawer of the cabinet and close it”). Furthermore, DINOv2 is robust to object poses retrieving sub-trajectories that “close the drawer” with the bowl either on the table or in the drawer (*c.f.* Eq. 25, “close the

bottom drawer of the cabinet and open the top drawer”). Trained on optical flow, FR has no notion of visual appearance, failing to retrieve most of the semantically relevant data.

6 CONCLUSION

We introduce STRAP as an innovative approach for leveraging visual foundation models in few-shot robotics manipulation, eliminating the need to train on the entire retrieval dataset and allowing it to scale with minimal compute overhead. By focusing on sub-trajectory retrieval using S-DTW, STRAP improves data utilization and captures dynamics more effectively. Overall, it outperforms state-of-the-art methods BehaviorRetrieval and FlowRetrieval by 12.20% and 12.47%, respectively, across all 10 LIBERO tasks.

7 REPRODUCIBILITY STATEMENT

To ensure reproducibility, we report runs over multiple seeds (1234, 42, 4325), seeding the retrieval procedure Tab. 6 as well as the training Tab. 1 and Tab. 3. This comprehensive approach allowed us to verify the consistency of our results across various runs ensuring reproducibility. We conduct all baseline and ablation experiments on the LIBERO-10 simulated benchmark and report hyperparameters in Appendix A.1 and Sec. 5.1. We will include a code release with our final paper, providing detailed instructions for reproducing our experiments exactly. This release will encompass all necessary components, including data preprocessing scripts, vision foundation model inference, hyperparameters, and evaluation scripts.

REFERENCES

- Suneel Belkhal, 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.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S. Ryoo, Grecia Salazar, Pannag R. Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong T. Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: robotics transformer for real-world control at scale. In Kostas E. Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin Yu (eds.), *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023. doi: 10.15607/RSS.2023.XIX.025. URL <https://doi.org/10.15607/RSS.2023.XIX.025>.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.
- Open X-Embodiment Collaboration, Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram Madhukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin

-
- Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minh Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaesan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Bajjal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.
- Norman Di Palo and Edward Johns. Dinobot: Robot manipulation via retrieval and alignment with vision foundation models. *arXiv preprint arXiv:2402.13181*, 2024.
- Maximilian Du, Suraj Nair, Dorsa Sadigh, and Chelsea Finn. Behavior retrieval: Few-shot imitation learning by querying unlabeled datasets. *arXiv preprint arXiv:2304.08742*, 2023.
- Jonathan Francis, Nariaki Kitamura, Felix Labelle, Xiaopeng Lu, Ingrid Navarro, and Jean Oh. Core challenges in embodied vision-language planning. *Journal of Artificial Intelligence Research*, 74: 459–515, 2022.
- Toni Giorgino. Computing and visualizing dynamic time warping alignments in R: The dtw package. *Journal of Statistical Software*, 31(7):1–24, 2009. doi: 10.18637/jss.v031.i07.
- Siddhant Haldar, Zhuoran Peng, and Lerrel Pinto. Baku: An efficient transformer for multi-task policy learning. *arXiv preprint arXiv:2406.07539*, 2024.
- Yafei Hu, Quanting Xie, Vidhi Jain, Jonathan Francis, Jay Patrikar, Nikhil Keetha, Seungchan Kim, Yaqi Xie, Tianyi Zhang, Zhibo Zhao, et al. Toward general-purpose robots via foundation models: A survey and meta-analysis. *arXiv preprint arXiv:2312.08782*, 2023.
- Yizheng Huang and Jimmy Huang. A survey on retrieval-augmented text generation for large language models. *CoRR*, abs/2404.10981, 2024. doi: 10.48550/ARXIV.2404.10981. URL <https://doi.org/10.48550/arXiv.2404.10981>.
- Yiding Jiang, Evan Zheran Liu, Benjamin Eysenbach, J. Zico Kolter, and Chelsea Finn. Learning options via compression. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/8567a53e58a9fa4823af356c76ed943c-Abstract-Conference.html.
- Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis,

-
- Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pannag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Bajjal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Foster, Jensen Gao, David Antonio Herrera, Minh Heo, Kyle Hsu, Jiaheng Hu, Donovan Jackson, Charlotte Le, Yunshuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, Abigail O’Neill, Rosario Scalise, Derick Seale, Victor Son, Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick Yin, Yunchu Zhang, Osbert Bastani, Glen Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman, Joseph J Lim, Jitendra Malik, Roberto Martín-Martín, Subramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, Yuke Zhu, Thomas Kollar, Sergey Levine, and Chelsea Finn. Droid: A large-scale in-the-wild robot manipulation dataset. 2024.
- Yuxuan Kuang, Junjie Ye, Haoran Geng, Jiageng Mao, Congyue Deng, Leonidas Guibas, He Wang, and Yue Wang. Ram: Retrieval-based affordance transfer for generalizable zero-shot robotic manipulation. In *8th Annual Conference on Robot Learning*.
- Chengshu Li, Fei Xia, Roberto Martin-Martin, and Silvio Savarese. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators. In *Conference on Robot Learning*, pp. 603–616. PMLR, 2020.
- Li-Heng Lin, Yuchen Cui, Amber Xie, Tianyu Hua, and Dorsa Sadigh. Flowretrieval: Flow-guided data retrieval for few-shot imitation learning. In *8th Annual Conference on Robot Learning, 2024*.
- Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Federico Malato, Florian Leopold, Andrew Melnik, and Ville Hautamäki. Zero-shot imitation policy via search in demonstration dataset. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7590–7594. IEEE, 2024.
- Vivek Myers, Chunyuan Zheng, Oier Mees, Kuan Fang, and Sergey Levine. Policy adaptation via language optimization: Decomposing tasks for few-shot imitation. In *8th Annual Conference on Robot Learning, 2024*.
- Meinard Müller. *Fundamentals of Music Processing: Using Python and Jupyter Notebooks*. Springer Cham, 2 edition, 2021. ISBN 978-3-030-69807-2. URL <https://doi.org/10.1007/978-3-030-69808-9>.
- Soroush Nasiriany, Tian Gao, Ajay Mandlekar, and Yuke Zhu. Learning and retrieval from prior data for skill-based imitation learning. In *Conference on Robot Learning, 2022*.
- Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of everyday tasks for generalist robots. *arXiv preprint arXiv:2406.02523*, 2024.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*.
- Georgios Papagiannis, Norman Di Palo, Pietro Vitiello, and Edward Johns. R+ x: Retrieval and execution from everyday human videos. In *RSS 2024 Workshop: Data Generation for Robotics*.
- Jyothish Pari, Nur Muhammad Mahi Shafiullah, Sridhar Pandian Arunachalam, and Lerrel Pinto. The surprising effectiveness of representation learning for visual imitation. In *18th Robotics: Science and Systems, RSS 2022*. MIT Press Journals, 2022.

-
- 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*, pp. 8748–8763. PMLR, 2021.
- Scott E. Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent. *Trans. Mach. Learn. Res.*, 2022, 2022. URL <https://openreview.net/forum?id=likK0kHjvj>.
- Tanmay Shankar, Yixin Lin, Aravind Rajeswaran, Vikash Kumar, Stuart Anderson, and Jean Oh. Translating robot skills: Learning unsupervised skill correspondences across robots. In *International Conference on Machine Learning*, pp. 19626–19644. PMLR, 2022.
- Jannik Sheikh, Andrew Melnik, Gora Chand Nandi, and Robert Haschke. Language-conditioned semantic search-based policy for robotic manipulation tasks. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Lirui Wang, Jialiang Zhao, Yilun Du, Edward H. Adelson, and Russ Tedrake. Poco: Policy composition from and for heterogeneous robot learning. *CoRR*, abs/2402.02511, 2024. doi: 10.48550/ARXIV.2402.02511. URL <https://doi.org/10.48550/arXiv.2402.02511>.
- Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezafofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8121–8130, 2022.
- Zhao-Heng Yin and Pieter Abbeel. Offline imitation learning through graph search and retrieval. *arXiv preprint arXiv:2407.15403*, 2024.
- Lihan Zha, Yuchen Cui, Li-Heng Lin, Minae Kwon, Montserrat Gonzalez Arenas, Andy Zeng, Fei Xia, and Dorsa Sadigh. Distilling and retrieving generalizable knowledge for robot manipulation via language corrections. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 15172–15179. IEEE, 2024.
- Yuying Zhang, Wenyan Yang, and Joni Pajarinen. Demobot: Deformable mobile manipulation with vision-based sub-goal retrieval. *arXiv preprint arXiv:2408.15919*, 2024.
- Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In Kostas E. Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin Yu (eds.), *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023. doi: 10.15607/RSS.2023.XIX.016. URL <https://doi.org/10.15607/RSS.2023.XIX.016>.

A APPENDIX

A.1 SIM EVALUATION

Table 3: **Baselines (sim)**: Performance of different methods on LIBERO-10 tasks in simulation

Method	Mug-Microwave	Moka-Moka	Soup-Sauce	Cream-Cheese-Butter	Mug-Pudding
BC	28.00% \pm 0.94	0.00% \pm 0.00	17.33% \pm 4.46	26.67% \pm 4.25	18.00% \pm 2.49
MT	0.00% \pm 0.00	0.00% \pm 0.00	0.00% \pm 0.00	0.00% \pm 0.00	0.00% \pm 0.00
BR (Du et al., 2023)	28.67% \pm 3.93	0.0% \pm 0.0	13.33% \pm 3.81	<u>32.0% \pm 4.32</u>	26.0% \pm 1.89
FR (Lin et al., 2024)	27.33% \pm 1.44	0.0% \pm 0.0	11.33% \pm 3.03	41.33% \pm 5.52	14.67% \pm 1.09
D-S	30.0% \pm 3.4	0.0% \pm 0.0	4.67% \pm 0.54	16.0% \pm 5.66	6.0% \pm 0.94
D-T	34.67% \pm 1.96	0.0% \pm 0.0	4.67% \pm 1.09	27.33% \pm 4.46	14.0% \pm 3.4
STRAP (CLIP)	30.00% \pm 2.49	0.00% \pm 0.00	8.67% \pm 6.28	29.33% \pm 10.51	<u>24.00% \pm 4.32</u>
STRAP (DINO)	29.33% \pm 2.72	0.00% \pm 0.00	<u>16.67% \pm 1.97</u>	29.33% \pm 11.34	18.67% \pm 1.44

Remaining results on LIBERO-10 Tab. 3 shows the results for the remaining LIBERO-10 task not reported in the main sections. Both FR and BR outperform STRAP on the Cream-Cheese-Butter task. We hypothesize that our chunking heuristic generates sub-optimal sub-trajectories (too long) causing them to contain multiple different semantic tasks, leading to worse matches in our retrieval datasets and eventually in decreasing downstream performance.

Hyperparameters for sim results: All results are reported over 3 training and evaluation seeds (1234, 42, 4325). We fixed both the number of segments retrieved to 100, the camera viewpoint to the agent view image for retrieval, and the number of expert demonstrations to 5. Our transformer policy was trained over all input images for 300 epochs with batch size 32 and an epoch every 200 gradient steps.

Baseline implementation details: Following Lin et al. (2024), we retrieve single-state action pairs for the state-based retrieval baselines (BR, FR, D-S) and pad them by also retrieving the states from $t - h$ to $t + h - 1$ to make the samples compatible with our transformer-based policy. We refer the reader to Appendix A.3 for extensive ablation.

A.2 REAL EXPERIMENTS

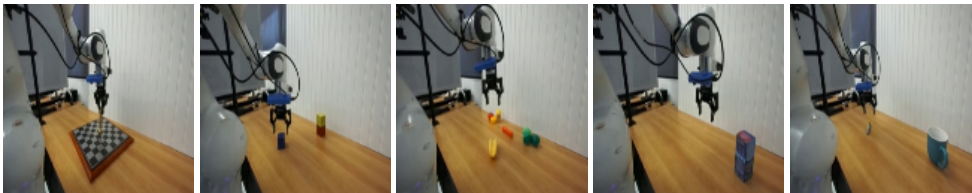


Figure 14: chess

Figure 15: cube_stacking

Figure 16: hotdog

Figure 17: knock_over_box

Figure 18: marker_in_mug



Figure 19: medicine_pnp

Figure 20: dispense_soap

Figure 21: pull_cable_right

Figure 22: pen_next_to_pens

Figure 23: screwdriver

Figure 24: Environment setup for the real-world tasks

Environment Name	Language Instruction
chess	Move the king to the top right of the chess board
cube_stacking	Stack the blue cube on top of the tower
hotdog	Put the hotdog in the bun
knock_over_box	Knock over the box
marker_in_mug	Put the marker in the mug
medicine_pnp	Pick up the medicine box on the right and put it next to the other medicine boxes
dispense_soap	Press the soap dispenser
pull_cable_right	Pull the cable to the right
pen_next_to_pens	Put the pen next to the markers
screwdriver	Pick up the screwdriver and put it in the cup

Table 4: Task/language instructions for the real-world dataset $\mathcal{D}_{\text{prior}}$

A.3 ABLATIONS

Table 5: **Ablations - Retrieval Method:** We explore different approaches for trajectory-based retrieval. Besides the heuristic reported in the main paper, we experiment with a sliding window approach that segments a trajectory into sub-trajectories of equal length (here: 30). We use S-DTW for both sliding window sub-trajectories and full trajectories.

Method	Stove-Moka	Bowl-Cabinet	Mug-Microwave	Moka-Moka	Soup-Cream-Cheese
Sub-traj (sliding window)	76.0% \pm 4.71	75.33% \pm 2.72	26.0% \pm 1.89	0.0% \pm 0.0	37.33% \pm 6.62
Full traj	78.67% \pm 2.72	68.67% \pm 1.44	34.67% \pm 1.96	0.0% \pm 0.0	28.67% \pm 3.81
Method	Soup-Sauce	Cream-Cheese-Butter	Mug-Mug	Mug-Pudding	Book-Caddy
Sub-traj (sliding window)	40.00% \pm 0.94	27.33% \pm 2.18	63.33% \pm 3.57	30.00% \pm 3.40	79.0% \pm 4.95
Full traj	4.67% \pm 1.09	27.33% \pm 4.46	43.33% \pm 1.09	14.0% \pm 3.4	68.0% \pm 5.66

Table 6: **Ablations - Retrieval Seeds:** We run STRAP on different retrieval seeds on a subset of LIBERO-10 tasks. We report results over all possible combinations of 3 training and 3 retrieval seeds

Method	Stove-Moka	Mug-Cabinet	Book-Caddy
BC Baseline	93.11% \pm 1.57	83.11% \pm 2.69	93.11% \pm 1.57
STRAP	98.0% \pm 1.04	88.67% \pm 2.11	98.0% \pm 1.04

Table 7: **Ablations - amount data retrieved:** We explore the effect of increasing the size of $\mathcal{D}_{\text{retrieval}}$. We evaluate performance on LIBERO-10 tasks in simulation on 2 different retrieval and 3 training seeds. We randomly sample 10 demos from $\mathcal{D}_{\text{target}}$ and retrieve 1500 segments. This demonstrates STRAP’s robustness over multiple seeds, as well as scalability to more data even leading to performance gains

Task	Stove-Pot	Bowl-Cabinet	Soup-Cheese	Mug-Mug	Book-Caddy
BC	86.33% \pm 2.18	76.0% \pm 3.97	41.67% \pm 3.72	59.0% \pm 2.25	92.67% \pm 1.81
STRAP (DINO)	88.67% \pm 3.42	95.67% \pm 1.19	45.67% \pm 7.41	67.67% \pm 1.59	93.71% \pm 1.87
Method	Mug-Microwave	Pots-On-Stove	Soup-Sauce	Cream cheese-Butter	Mug-Pudding
BC	47.67% \pm 4.75	0.00% \pm 0.00	23.0% \pm 3.42	57.33% \pm 0.77	32.0% \pm 1.33
STRAP (DINO)	31.33% \pm 3.73	0.00% \pm 0.00	45.0% \pm 5.09	58.67% \pm 9.58	38.33% \pm 3.38

Table 8: **Ablations - Diffusion Policies:** Performance on LIBERO-10 tasks using diffusion policies without language conditioning for BR and FR. These experiments replicate the training setup for BR and FR. Both methods fall short of the baselines reported in the rest of the paper.

Task	Stove-Pot	Bowl-Cabinet	Soup-Cheese	Mug-Mug	Book-Caddy
Diffusion Behavior Retrieval	36.67% \pm 1.44	68.0% \pm 2.49	34.0% \pm 2.49	55.33% \pm 1.44	42.0% \pm 1.63
Diffusion Flow Retrieval	68.67% \pm 2.37	56.0% \pm 4.32	18.0% \pm 3.4	56.0% \pm 3.4	35.33% \pm 6.28
Method	Mug-Microwave	Pots-On-Stove	Soup-Sauce	Cream cheese-Butter	Mug-Pudding
Diffusion Behavior Retrieval	30.67% \pm 0.54	0.00% \pm 0.00	10.67% \pm 1.96	24.0% \pm 0.94	9.33% \pm 1.44
Diffusion Flow Retrieval	32.67% \pm 3.31	68.0% \pm 2.49	6.0% \pm 0.0	35.33% \pm 0.54	8.0% \pm 1.89

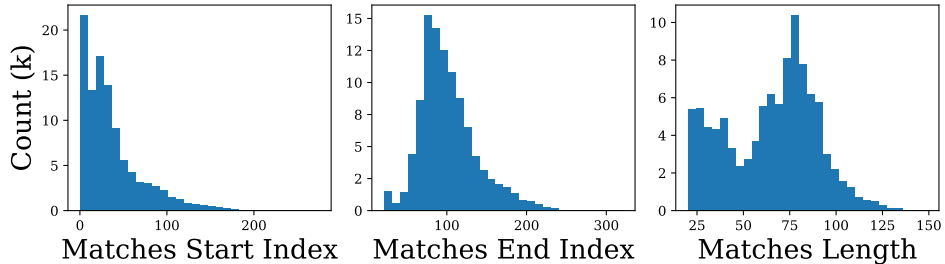


Figure 25: Match distribution $\mathcal{D}_{\text{prior}}$ for STRAP with target task: *”put the black bowl in the bottom drawer of the cabinet and close it”*. S-DTW finds the best matches regardless of start and end points or trajectory length. This results in a distribution over start and end points as well as a variety of trajectory lengths retrieved.