

Retrieve, Don't Retrain: Extending Vision-Language-Action Models to New Tasks at Test Time

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Extending a vision-language-action (VLA) policy to a new task typi-
2 cally requires task-specific teleoperated demonstrations and per-task fine-tuning,
3 making adaptation costly in both data collection and compute. In this paper, we
4 show that this target-side per-task adaptation cost can be replaced by **retrieval**.
5 Our retrieval augmented policy is trained once on paired demonstrations from the
6 target embodiment (**query**) and a cheaper embodiment (**pool**, *e.g.*, human-hand
7 video), then frozen. New tasks are added at deployment by appending pool-side
8 demonstrations to a retrieval pool. The frozen policy conditions on retrieved tra-
9 jectories at every control step, so new tasks are absorbed by indexing data rather
10 than updating parameters. Fine-tuning is needed only to take on a new, unseen
11 embodiment, not for each new task. We show that retrieval improves policies
12 beyond a specific backbone, including standard VLA policies, but its effect is
13 especially pronounced in Cosmos Policy, a video-generation-based world-action
14 model (WAM). In this setting, retrieval supplies coarse task progression, while
15 the WAM's future-image objective provides an additional visual consistency sig-
16 nal that strengthens the retrieval-conditioned actions. On PushT, we study how
17 retrieval provides a reusable high-level motion prior for cross-embodiment gen-
18 eralization to unseen goal angles, while on RoboTwin 2.0 our method outperforms
19 cross-embodiment baselines on unseen tasks, and we additionally demonstrate the
20 method on a real robot.

21 **Keywords:** Robot foundation models, World-action models, Retrieval-augmented
22 policies, Vision-language-action models

23 1 Introduction

24 General-purpose robot policies [1, 2, 3, 4, 5, 6, 7, 8] aim to execute open-ended manipulation behav-
25 iors from natural-language instructions while generalizing across diverse environments, tasks, and
26 embodiments. Yet a new embodiment still requires its own teleoperated demonstrations and per-task
27 fine-tuning, so cost grows with each new task added. We argue that this per-task cost is avoidable:
28 behavioral knowledge from a cheap, data-rich source (*e.g.*, human-hand video demonstrations) can
29 transfer to the target embodiment through a **retrieval rather than retraining** paradigm.

30 The cost of the previous approach is twofold. On the data side, target-embodiment demonstrations
31 would be collected through teleoperation, which is slow, hardware-bound, and roughly $18\times$ slower
32 to acquire than equivalent human-hand demonstrations [9, 10]. On the compute side, modern vision-
33 language-action (VLA) models and robot foundation models operate over high-dimensional visual
34 and action sequences, so per-task fine-tuning of recent world-action models (WAM) [7, 11, 8] costs
35 roughly 24 GPU-hours per task and continues to scale with model size, context length, and action
36 horizon. Both costs compound with every new task introduced.

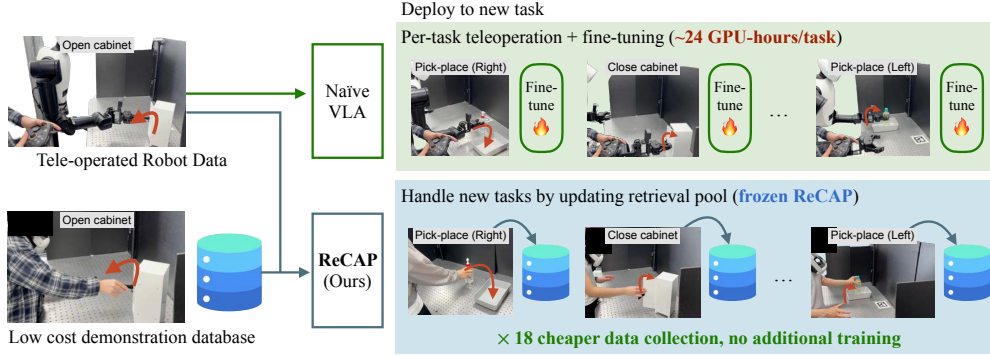


Figure 1: **RECAP overview.** Instead of teleoperating each new task and fine-tuning the policy (top, ~ 24 GPU-hours/task for Cosmos Policy [7]), ReCAP appends cheap human-hand demonstrations to a retrieval pool while keeping the policy frozen (bottom), $18\times$ cheaper [9, 10], no additional training.

37 We propose RECAP (**R**etrieval-**C**onditioned **A**ction **P**olicy), which shifts adaptation from repeated
 38 optimization to retrieval over a reusable pool of source-embodiment demonstrations. The policy is
 39 trained *once* to bridge the gap between source and target embodiments and is then frozen; behavioral
 40 coverage expands by *simply appending new demonstrations to the retrieval memory*.

41 RECAP builds on a world-action model (WAM) [7, 8, 12, 11, 13, 14, 15], specifically Cosmos Pol-
 42 icy [7]. We parameterize the action latents as a *residual* over retrieved trajectories: retrieval supplies
 43 the coarse high-level motion and task progression, while the policy learns only the embodiment-
 44 specific dynamics needed to execute the behavior on the target robot. Crucially, the WAM’s future-
 45 image prediction objective enforces consistency between the retrieved trajectory and the predicted
 46 evolution of the scene, a visual alignment signal that becomes informative *only* when paired with
 47 retrieval in unseen tasks, and that we find especially beneficial for long-horizon behaviors where
 48 high-level motion structure dominates.

49 The main contribution of this paper is threefold: a paradigm that adapts a policy to new tasks entirely
 50 at test time, absorbing each new task by extending a retrieval pool with cheap pool-embodiment
 51 demonstrations while the policy stays frozen with no parameter updates; a retrieval-conditioned
 52 residual policy on a WAM (i.e., Cosmos Policy [7]) in which retrieval supplies the high-level mo-
 53 tion so the policy learns only the embodiment-specific correction, reinforced by the WAM’s future-
 54 image objective that is informative only when paired with retrieval; and consistent gains over cross-
 55 embodiment baselines on PushT [16] (34.9% vs. 6.0% on seven unseen angles) and RoboTwin
 56 2.0 [17] (31.5% vs. 26.0% on five unseen tasks), with a pool-progression study confirming mono-
 57 tonic coverage growth without parameter updates and a further real-robot validation.

58 2 Related Work

59 **World action models.** Recent VLA policies inherit either a pretrained language backbone with an
 60 added action head (OpenVLA [1], $\pi_{0.5}$ [3], GR00T N1.6 [4]) or a pretrained video model that folds
 61 actions into the same generative process (DreamZero [8], Cosmos Policy [7], mimic-video [12],
 62 Fast-WAM [11]). We call the latter family world-action models (WAMs); their video backbone
 63 is pretrained on internet-scale data and already encodes semantics and physical dynamics, so the
 64 policy learns only control on top, with action and future-observation prediction emerging from one
 65 shared video generation. These WAMs train on target-robot demonstrations alone, leaving their dy-
 66 namics priors unpaired with cheaper cross-embodiment supervision, which we address by building
 67 a WAM on Cosmos Policy that conditions on retrieved pool-embodiment trajectories at training and
 68 deployment (Section 4.1).

69 **Retrieval and Cross-embodiment policy transfer.** Retrieval-based imitation [18] adapts a policy
 70 with relevant demonstrations rather than training from scratch on each task. Within a single em-

71 bodiment, FlowRetrieval [19] co-trains the policy on optical-flow-retrieved data, and STRAP [20]
 72 trains a task-specific specialist at deployment from DTW-matched sub-trajectories. To cut data cost,
 73 attention has turned to cheaper human demonstrations [21, 22, 23]: Hong et al. [24] uses a single
 74 human-hand demonstration to retrieve matching robot sub-trajectories and efficiently fine-tunes a
 75 policy on them for fast adaptation. These methods still pay for each new task with training, whether
 76 by co-training or a test-time fine-tune. Other methods reduce this training cost by using human
 77 data more directly. R+X [25] retrieves everyday human videos and executes them directly without
 78 training, but it stays in the human domain and does not learn an embodiment-specific correction.
 79 In-context methods such as MimicDroid [9] condition on a human prompt, which a user must hand-
 80 pick for each task. We instead learn the cross-embodiment gap once on paired (query, pool) data
 81 and freeze it; the user then grows the retrieval pool with new pool-embodiment demonstrations at
 82 test time, which broadens the policy’s task coverage without retraining.

83 3 Problem Formulation

84 **Setting.** We study a cross-embodiment imitation setting with two sides: a *query* side that we want
 85 to control autonomously at deployment, carrying a *target embodiment* (e.g., a robot arm), and a *pool*
 86 side carrying a *pool embodiment* that is easier to collect demonstrations from but is never deployed
 87 (e.g., a human hand). The two embodiments differ in geometry, contacts, and dynamics, and differ
 88 sharply in data cost, since a query demonstration requires a teleoperation rig and operator while
 89 a pool demonstration only needs a lightweight tracker on a human. Transfer between them rests
 90 on two assumptions: a shared state/action representation (we use SE(3) end-effector pose plus a
 91 gripper signal where applicable), and motions that are semantically similar at the trajectory level, so
 92 a coarse plan derived from a pool trajectory is informative for the query.

93 **Train and test access.** At training time we assume that the model is given paired demonstrations
 94 $\mathcal{D}_{\text{train}}^{\text{query}}$ and $\mathcal{D}_{\text{train}}^{\text{pool}}$ on a fixed task distribution, each a set of state-action pairs $\{(s_t, a_t)\}$ collected on
 95 the corresponding embodiment. At test time the model faces *new* tasks outside this distribution; only
 96 pool demonstrations $\mathcal{D}_{\text{test}}^{\text{pool}}$ are provided, no additional query data is collected, and model parameters
 97 are not updated. The model is rolled out on the target embodiment, so the current query state s_t^{query}
 98 is observed at every control step.

99 **Retrieval-conditioned action prediction.** Let $\mathcal{D}^{\text{pool}}$ denote the active retrieval pool, $\mathcal{D}_{\text{train}}^{\text{pool}}$ at train-
 100 ing and $\mathcal{D}_{\text{test}}^{\text{pool}}$ at deployment. At each step t , we select $t' = \arg \min_{t'} d(s_t^{\text{query}}, s_{t'}^{\text{pool}})$ over $\mathcal{D}^{\text{pool}}$,
 101 where d is a feature-space distance specified in Section 4.2. The retrieved chunk $(s_{t':t'+H}^{\text{pool}}, a_{t':t'+H}^{\text{pool}})$
 102 of action chunk length H steps together with s_t^{query} is fed to the policy, which predicts the query
 103 action chunk $a_{t:t+H}^{\text{query}}$. The same rule applies at training and deployment with only $\mathcal{D}^{\text{pool}}$ changing;
 104 crucially, $\mathcal{D}_{\text{test}}^{\text{pool}}$ can be extended at any time without touching θ .

105 4 Proposed Method

106 We propose RECAP to adapt a policy to new tasks at test time without retraining, by retrieving
 107 relevant demonstrations rather than updating weights. The intuition is that the target and pool em-
 108 bodiments largely agree on *what* a task requires and differ mainly in *how* to execute it. A retrieved
 109 pool trajectory therefore supplies the shared high-level plan cheaply, leaving the policy to learn only
 110 the embodiment-specific correction, so adapting to a new task becomes a matter of indexing data
 111 rather than updating parameters, much as retrieval-augmented generation externalizes knowledge
 112 into a searchable store (Fig. 2). Section 4.1 details the backbone and its retrieval conditioning, and
 113 Section 4.2 specifies how the retrieved chunk is selected and extended at test time.

114 4.1 Retrieval-Augmented World Action Model

115 Because the policy stays frozen and conditions on an external pool, new tasks can be absorbed at
 116 test time by extending that pool rather than by retraining; we expect coverage to grow with the pool,

117 and cheap pool-embodiment data, such as
 118 human-hand video, to partly substitute for
 119 target-robot teleoperation.

120 **Backbone and retrieval-conditioned in-**
 121 **put.** The backbone is the Cosmos Policy
 122 formulation [7], which emits query-side
 123 actions and future image observations as
 124 one denoised video sequence.

125 We extend its conditioning with the re-
 126 trieved pool-embodiment chunk:

$$\pi_{\theta} \left(s_t^{\text{query}}, (s_{t':t'+H}^{\text{pool}}, a_{t':t'+H}^{\text{pool}}) \right) \mapsto \hat{a}_{t:t+H}^{\text{query}}, \hat{s}_{t+H}^{\text{query}}. \quad (1)$$

127 The retrieved chunk and the observed
 128 query frame are encoded into clean latent
 129 frames and prepended along the tempo-
 130 ral axis as conditioning, while the query-
 131 side future actions and observations are
 132 denoised from noise (Fig. 2); the language
 133 instruction enters via cross-attention. The
 134 retrieved chunk thus extends the standard I2V
 conditioning, a single clean frame, to a
 clean state-action subsequence, with no
 architectural modification.

135 **Joint training objective (World action model).** Action and future-image latents are supervised
 136 jointly with a single flow-matching loss:

$$\mathcal{L}(\theta) = \lambda \mathcal{L}_{\text{act}}(\hat{a}_{t:t+H}^{\text{query}}, a_{t:t+H}^{\text{query}}) + \mathcal{L}_{\text{state}}(\hat{s}_{t+H}^{\text{query}}, s_{t+H}^{\text{query}}) \quad (2)$$

137 Joint training yields actions aligned with the predicted next state, producing more grounded action
 138 outputs. For a standard VLA, we lack $\mathcal{L}_{\text{state}}$ and $\hat{s}_{t+H}^{\text{query}}$.

139 **Residual action parameterization.** Because the retrieved pool action chunk $a_{t':t'+H}^{\text{pool}}$ already
 140 encodes a coarse motion the target should execute, we let the action latents represent only the
 141 embodiment-specific correction [26, 27] on top:

$$\hat{a}_{t:t+H}^{\text{query}} = a_{t':t'+H}^{\text{pool}} + \Delta a_{t:t+H}. \quad (3)$$

142 This narrows what the action latents must encode to “how the query action differs from the pool’s,”
 143 the variation actually caused by the embodiment gap. This variation is weakly reflected in action
 144 labels but clearly visible in pixels, e.g., how contact happens, how the gripper closes. Residual
 145 focuses the action latents on this correction, and state prediction provides the dense visual signal to
 146 learn it.

147 4.2 Retrieval

148 At each control step t , the policy retrieves a pool-embodiment index t' from $\mathcal{D}^{\text{pool}}$ whose surround-
 149 ing chunk best matches the current query context. We first form a candidate set $\mathcal{C}_t^{\text{traj}}$ by taking the
 150 top- K trajectories closest to the query under a composite initial-frame descriptor ψ_0 , a language
 151 embedding of the goal, initial task-relevant object positions (via SAM 3 [28]), and initial proprio-
 152 ception. Within $\mathcal{C}_t^{\text{traj}}$, the index distance d in Section 3 is a weighted sum of L_2 distances over object
 153 pose, proprioception, and the upcoming action chunk (training only; dropped at inference), and a
 154 cosine distance over a DINOv3 [29] image feature.

155 At inference, new pool-embodiment demonstrations $\mathcal{D}_{\text{test}}^{\text{pool}}$ replace the active pool and are reindexed
 156 under ψ_0 and the features above; retrieval re-runs every step, so $\mathcal{D}_{\text{test}}^{\text{pool}}$ can grow within a session.

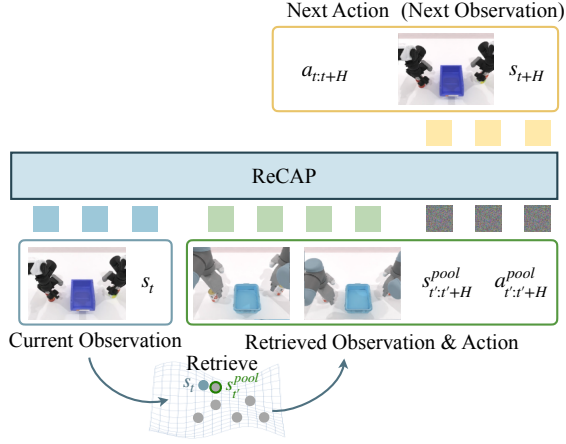


Figure 2: **RECAP framework.** The current observation retrieves a matching state-action chunk from the pool database; the retrieved chunk and the current observation then condition a world action model that denoises the next action and next observation in one video sequence.

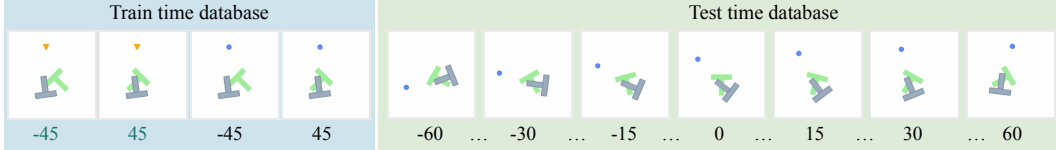


Figure 3: **PushT cross-embodiment pool database setting.** The training set pairs the *triangle* (target) and *disc* (pool) at $\pm 45^\circ$. The test set is a pool database of disc-pusher demonstrations spanning all goal angles, which the frozen triangle policy retrieves from on the seven unseen angles.

157 5 Experiments

158 We evaluate the proposed RECAP policy in three challenging cross-embodiment settings. PushT
 159 variant [16] (§5.2) provides a controlled environment for analyzing how retrieval improves gener-
 160 alization and why retrieval-conditioned WAMs are effective. RoboTwin 2.0 [17] (§5.3) evaluates
 161 whether the same paradigm scales to multi-task dual-arm manipulation, while real-robot experi-
 162 ments (§5.4) test whether unseen tasks can be absorbed through retrieval alone using human-hand
 163 demonstrations without additional robot training. Across all setups, we study the same hypothe-
 164 sis: retrieval supplies coarse task progression, while the policy learns only the embodiment-specific
 165 dynamics needed to execute the behavior on the target robot.

166 5.1 Experiment Setup

167 **PushT Environment.** In the 2D PushT benchmark [16, 30] an agent pushes a T-shaped block
 168 to a goal pose; we make it cross-embodiment with two pushers of different contact dynamics, a
 169 *triangle* (target) and a *disc* (pool), and take the goal rotation angle as the task axis. Training uses
 170 100 paired $\{\textit{triangle}, \textit{disc}\}$ demonstrations at $\pm 45^\circ$; the triangle is then evaluated on nine different
 171 angles ranging from -60° to $+60^\circ$ in 15° steps, seven of them unseen. The test-time pool holds
 172 disc-pusher demonstrations at 5° resolution over $[-60^\circ, +60^\circ]$, added without retraining.

173 **RoboTwin Simulation Environment.** On RoboTwin 2.0 [17], we take Aloha-Agilex [31] as the
 174 target and UR5 as the pool, training on five paired $\{\textit{target}, \textit{pool}\}$ tasks and evaluating on five unseen
 175 ones (Table 1). A test-time pool-progression eval grows the pool through five levels (i.e., 11, 17, 23,
 176 29, 35 tasks), each a strict superset of the previous, with the policy frozen throughout.

177 **Real Robot Setup.** On a physical robot, the pool embodiment is a human hand (video with wrist-
 178 pose tracked in VR) and the target is the teleoperated robot. We fine-tune on a single task, *open-*
 179 *cabinet*, with 25 paired demonstrations. Then we freeze the policy and evaluate on three tasks: the
 180 seen *open-cabinet* and two held out from fine-tuning, *place-bottle-in-plastic-box* and *close-cabinet*
 181 (Fig. 8a). The only test-time exposure to the held-out behaviors is 10 human-hand demonstrations
 182 per task added to the pool.

183 5.2 PushT Experiments

184 In this section, we study the following aspects of retrieval-conditioned generalization through PushT
 185 experiments [16]: (1) whether expanding the retrieval pool at test time improves unseen-task cov-
 186 erage without retraining, (2) whether retrieval benefits more from a WAM backbone than from an
 187 action-only policy, and (3) whether the retrieved trajectory acts as a reusable high-level motion prior
 188 that the policy adapts to the target embodiment. PushT is a controlled testbed whose generalization
 189 axis (i.e., the goal angle) is one-dimensional and densely measurable, which enables exploring these
 190 questions.

191 **Test-time pool progression.** Expanding the retrieval pool at test time without parameter updates
 192 steadily recovers the unseen angles. We grow the pool with disc-pusher demonstrations at inter-
 193 mediate goal angles and track per-angle success across five snapshots (Fig. 4). Specifically, the
 194 unseen-angle average rises monotonically from 6.0% without retrieval to 34.9% at the full pool.

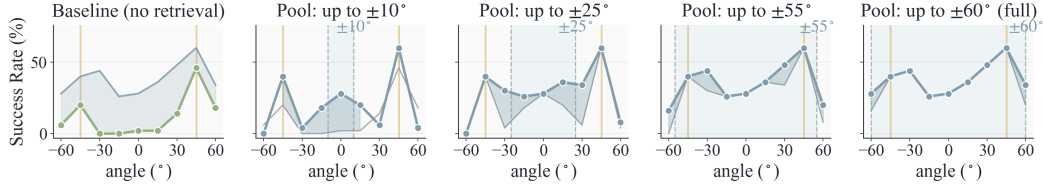


Figure 4: **Test-time pool progression on PushT.** The leftmost panel is the no-retrieval baseline with our full-pool curve overlaid (shaded gap). The other panels show per-angle success as the pool grows with no retraining, with the previous snapshot in gray and the incremental gain shaded.

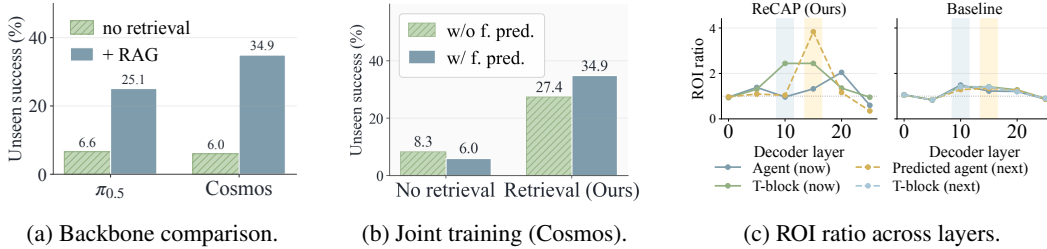


Figure 5: **Comparative analyses of ReCAP and baseline on PushT.** (a) Unseen-angle success with and without retrieval on a $\pi_{0.5}$ and a Cosmos (WAM) backbone; retrieval helps both, and the WAM benefits more. (b) The future-image objective improves unseen success only when paired with retrieval. (c) Action-slot attention across decoder layers, which peaks on the T-block and then on the predicted next position under retrieval but stays near uniform without it.

195 Notably, each angle reaches much of its final success *before* its matching pool angle is added, so the
 196 policy interpolates over neighboring demonstrations rather than memorizing the nearest one.

197 **Retrieval and the role of the WAM objective.** If retrieval already supplies the coarse motion
 198 plan, then the remaining learning problem is primarily adapting that plan to the target embodiment
 199 dynamics. We therefore hypothesize that a WAM should benefit more from retrieval than an action-
 200 only policy, since its future-image objective encourages the retrieved trajectory to remain consistent
 201 with the predicted future scene, providing a stronger learning signal for the embodiment-specific
 202 dynamics adaptation. Figure 5a reveals that retrieval improves both backbones, raising the unseen-
 203 angle success of the action-only $\pi_{0.5}$ [3] from 6.6% to 25.1%. However, the gain is larger with a
 204 WAM backbone. In Fig. 5b, removing the image-prediction objective reduces our model to 27.4%,
 205 comparable to retrieval-augmented $\pi_{0.5}$ (25.1%), while restoring it improves performance to 34.9%.

206 **How the retrieval-conditioned policy works.** The retrieved chunk acts as a coarse motion prior,
 207 while the policy adapts it to the target embodiment rather than planning from scratch. Decoder
 208 cross-attention analysis in Fig. 5c reveals a two-stage routing behavior: early layers attend to the
 209 manipulated object and retrieved trajectory, whereas later layers shift attention toward the policy’s
 210 own predicted next state, adapting the retrieved plan to the target embodiment dynamics. This
 211 structure does not emerge without retrieval, where the ROI ratio remains near 1.0 across layers,
 212 indicating near-uniform attention. Masking the action-to-retrieval attention further confirms that the
 213 retrieved trajectory causally influences the generated actions.

214 5.3 RoboTwin Simulation Experiments

215 In this section, we evaluate the proposed method in a multi-task, dual-arm manipulation simulation
 216 environment, RoboTwin 2.0 [17]. We compare against standard cross-embodiment baselines on seen
 217 and held-out unseen tasks, then test whether the test-time pool growth seen on PushT carries over to
 218 this multi-task regime.

219 **Baseline comparisons.** We compare against three cross-embodiment baselines that share our back-
 220 bone but differ in how they incorporate the UR5 (pool) data. *Baseline*, Cosmos Policy [7] is trained

Table 1: **Quantitative analysis on RoboTwin.** We report per-task success rate (%) on RoboTwin, with Aloha-Agilex as the target embodiment and UR5 as the retrieval pool. The left block shows seen tasks, and the right block shows unseen tasks.

Method	Seen tasks						Unseen tasks					
	PCB	OM	DB	MP	GR	Avg \uparrow	MPP	PBS	CB	HM	LP	Avg \uparrow
Baseline [7]	47.5	10.0	25.0	30.0	50.0	32.5	0.0	0.0	20.0	0.0	0.0	4.0
Retrieval Only	30.0	7.5	22.5	7.5	60.0	25.5	0.0	10.0	42.5	37.5	40.0	26.0
Co-training	0.0	5.0	7.5	50.0	72.5	27.0	0.0	0.0	40.0	0.0	10.0	10.0
RECAP (Ours)	60.0	12.5	40.0	35.0	70.0	43.5	5.0	12.5	47.5	47.5	45.0	31.5

PCB = Place Cans Plasticbox OM = Open Microwave DB = Pick Dual Bottles MP = Move Can Pot GR = Grab Roller
MPP = Move Pillbottle Pad PBS = Place Bread Skillet CB = Click Bell HM = Hand-over Mic LP = Lift Pot

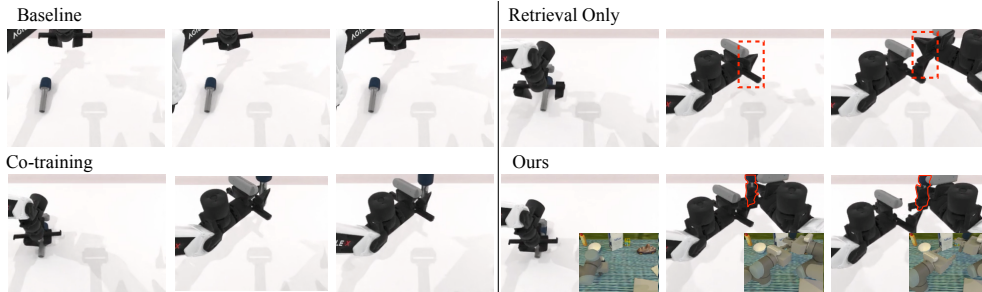


Figure 6: **Qualitative comparison on the held-out hand-over-mic task.** Baseline (top-left) and Co-training (bottom-left) fail to grasp the microphone; Retrieval Only (top-right) knocks it over (red box); Ours (bottom-right) grasps it successfully. Each inset shows the retrieved UR5 chunk that the policy conditions on.

221 on Aloha-Agilex (target) demonstrations alone, with no access to the pool. *Retrieval Only* executes
222 the action sequence of the nearest pool demonstration without learning. *Co-training* is a common
223 cross-embodiment baseline, also used in EgoBridge [21] and STRAP [20], that jointly trains a single
224 policy on the union of target and pool trajectories. Table 1 shows that RECAP leads on both splits,
225 at 43.5% seen and 31.5% unseen versus 32.5% and 26.0% for the strongest baseline. Replaying
226 the nearest pool trajectory (i.e., Retrieval Only) is competitive only where that trajectory already
227 approximates the target action. Otherwise, the learned residual on top of retrieval is what closes the
228 gap. As Fig. 6 illustrates, the nearest UR5 trajectory collides and dislodges the object, while our
229 policy produces the finer grip orientation the task needs.

230 **Test-time pool progression.** Figure 7 shows that growing the re-
231 trieval pool at test time with the policy frozen raises unseen-task
232 success monotonically, from 9.0% to 31.5% at the full pool. Each
233 increase coincides with a held-out task becoming retrievable, and
234 once all five are in the pool, the frozen policy matches its supervised
235 unseen-task average. Cheap pool-embodiment data at deployment
236 can therefore stand in for new target-embodiment demonstrations on
237 tasks unseen during fine-tuning.

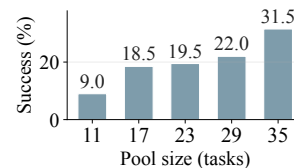


Figure 7: Test-time pool progression on RoboTwin.

238 5.4 Real Robot Experiments

239 We test whether the protocol transfers to real-world robots, despite the large embodiment gap be-
240 tween the human-hand demonstrations in the retrieval pool and the target robot (Fig. 8a). On two
241 held-out tasks, the no-retrieval baseline collapses to the trained open-cabinet motion regardless of
242 the target task, reaching only 10% and 0%, while retrieval enables the frozen policy to follow the
243 conditioned human trajectory and reach 80% and 30% on closing the cabinet and placing the bottle
244 (Fig. 8b). This indicates that human-hand demonstrations in the pool can partly substitute for addi-
245 tional target-robot teleoperation, even across a substantial embodiment gap. The qualitative results
246 are shown in Figure 9.

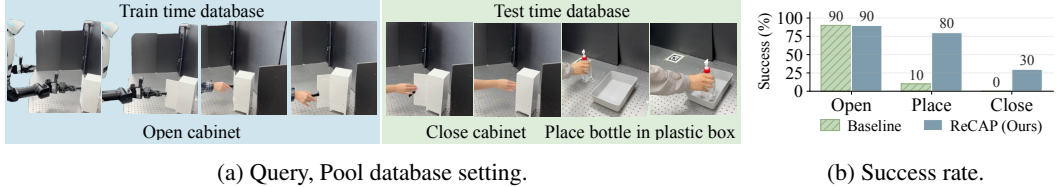


Figure 8: **Real-robot experiment.** (a) The training-time database pairs the robot (query) with a human-hand pool; the test-time database adds human-hand demonstrations for the held-out tasks. (b) Per-task success rate over 10 rollouts, Baseline vs RECAP (Ours).

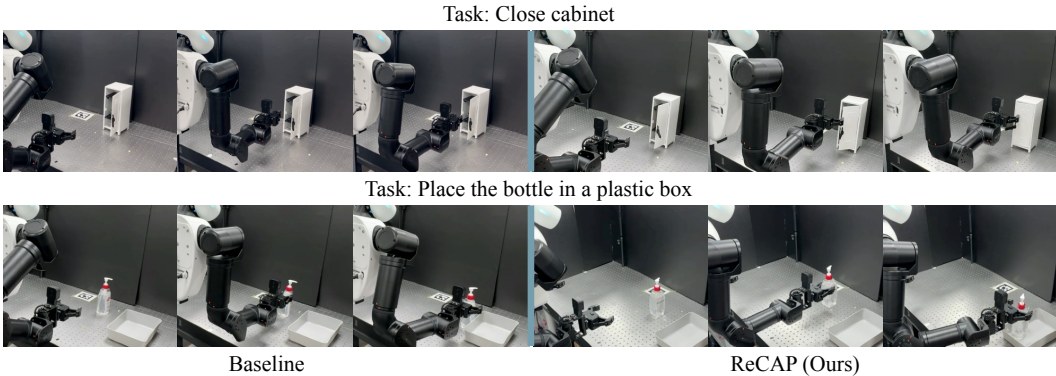


Figure 9: **Real-robot generalization to unseen tasks.** We show rollouts on the two held-out tasks (three frames each; baseline left, ours right). Trained only on *open-cabinet*, the baseline replays that trajectory and fails to close the cabinet (top) or grasp the bottle (bottom), whereas our policy follows the commanded behavior by conditioning on retrieved human-hand chunks.

247 6 Discussions

248 **Summary.** We extend a world-action-model policy to new tasks without retraining. Trained once
 249 to condition on a retrieval pool and predict an embodiment-specific residual on the retrieved tra-
 250 jectory, the frozen policy absorbs a new task by adding cheap pool-embodiment demonstrations at
 251 deployment. Across a cross-embodiment PushT variant, RoboTwin [17], and a physical robot, this
 252 improves generalization to unseen angles and tasks over cross-embodiment baselines, with success
 253 growing as the pool expands, and our analysis ties the gain to the WAM’s future-image objective
 254 acting together with retrieval. Indexing cheaper pool-embodiment data at deployment can thus stand
 255 in for collecting new target-robot demonstrations.

256 **Limitations and Future Work.** Several constraints point to future work. The target and pool em-
 257 bodiments must share an end-effector action space, since the residual refines a retrieved chunk in a
 258 common low-level representation; structurally different action spaces (e.g., a dexterous hand versus
 259 a parallel gripper) would require an embodiment-agnostic interface or learned action translator. The
 260 pool must also contain trajectories rather than video alone, so video-only sources such as raw hu-
 261 man or web video would first need to be lifted into a state-action representation. What constitutes an
 262 effective representation for cross-embodiment retrieval also remains an open question; our current
 263 descriptor combines language, object pose, proprioception, and visual features, but more scalable or
 264 embodiment-invariant representations may further improve transfer.

265 In addition, the residual formulation becomes less reliable when retrieved motions differ substan-
 266 tially in execution speed or temporal scale, particularly for larger chunks where errors can accumu-
 267 late over time. Developing retrieval and adaptation mechanisms that remain robust under significant
 268 temporal or dynamical mismatch is an important direction for future work. Finally, scaling retrieval
 269 beyond curated trajectory pools to in-the-wild video sources such as raw YouTube video remains a
 270 promising step toward broadly reusable robot experience.

References

- 271
- 272 [1] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. P.
273 Foster, P. R. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine,
274 P. Liang, and C. Finn. OpenVLA: An open-source vision-language-action model. In *Proc. of*
275 *the 8th Annual Conference on Robot Learning (CoRL)*, 2024. URL [https://openreview.](https://openreview.net/forum?id=ZMnD6QZAE6)
276 [net/forum?id=ZMnD6QZAE6](https://openreview.net/forum?id=ZMnD6QZAE6).
- 277 [2] J. Lee, J. Duan, H. Fang, Y. Deng, S. Liu, B. Li, B. Fang, J. Zhang, Y. R. Wang, S. Lee,
278 et al. Molmoact: Action reasoning models that can reason in space. *arXiv preprint*
279 *arXiv:2508.07917*, 2025.
- 280 [3] Physical Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail,
281 M. Equi, C. Finn, N. Fusai, et al. $\pi_{0.5}$: a vision-language-action model with open-world
282 generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- 283 [4] J. Bjorck et al. Gr00t n: An open foundation model for generalist humanoid robots. *arXiv*
284 *preprint arXiv:2503.14734*, 2025.
- 285 [5] J. Zheng, J. Li, Z. Wang, D. Liu, X. Kang, Y. Feng, Y. Zheng, J. Zou, Y. Chen, J. Zeng,
286 T. Wang, Y.-Q. Zhang, J. Liu, and X. Zhan. X-VLA: Soft-prompted transformer as scal-
287 able cross-embodiment vision-language-action model. In *Proc. of the Fourteenth International*
288 *Conference on Learning Representations (ICLR)*, 2026. URL [https://openreview.net/](https://openreview.net/forum?id=kt51kZH4aG)
289 [forum?id=kt51kZH4aG](https://openreview.net/forum?id=kt51kZH4aG).
- 290 [6] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu. RDT-1b: a
291 diffusion foundation model for bimanual manipulation. In *Proc. of the Thirteenth International*
292 *Conference on Learning Representations (ICLR)*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=yAzN4tz7oI)
293 [forum?id=yAzN4tz7oI](https://openreview.net/forum?id=yAzN4tz7oI).
- 294 [7] M. J. Kim, Y. Gao, T.-Y. Lin, Y.-C. Lin, Y. Ge, G. Lam, P. Liang, S. Song, M.-Y. Liu, C. Finn,
295 and J. Gu. Cosmos policy: Fine-tuning video models for visuomotor control and planning. In
296 *Proc. of the Fourteenth International Conference on Learning Representations (ICLR)*, 2026.
297 URL <https://openreview.net/forum?id=wPEIStHxYH>.
- 298 [8] S. Ye, Y. Ge, K. Zheng, S. Gao, S. Yu, G. Kurian, S. Indupuru, Y. L. Tan, C. Zhu, J. Xiang,
299 et al. World action models are zero-shot policies. *arXiv preprint arXiv:2602.15922*, 2026.
- 300 [9] R. Shah, S. Liu, Q. Wang, Z. Jiang, S. Kumar, M. Seo, R. Martín-Martín, and Y. Zhu. Mim-
301 icdroid: In-context learning for humanoid robot manipulation from human play videos. *arXiv*
302 *preprint arXiv:2509.09769*, 2025.
- 303 [10] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mim-
304 icplay: Long-horizon imitation learning by watching human play. In *Proc. of the 7th Annual*
305 *Conference on Robot Learning (CoRL)*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=hRZ1YjDZmTo)
306 [hRZ1YjDZmTo](https://openreview.net/forum?id=hRZ1YjDZmTo).
- 307 [11] T. Yuan, Z. Dong, Y. Liu, and H. Zhao. Fast-wam: Do world action models need test-time
308 future imagination? *arXiv preprint arXiv:2603.16666*, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2603.16666)
309 [abs/2603.16666](https://arxiv.org/abs/2603.16666).
- 310 [12] J. Pai, L. Achenbach, V. Montesinos, B. Forrai, O. Mees, and E. Nava. mimic-video: Video-
311 action models for generalizable robot control beyond vlas. *arXiv preprint arXiv:2512.15692*,
312 2025.
- 313 [13] L. Li, Q. Zhang, Y. Luo, S. Yang, R. Wang, F. Han, M. Yu, Z. Gao, N. Xue, X. Zhu, et al.
314 Causal world modeling for robot control. *arXiv preprint arXiv:2601.21998*, 2026.
- 315 [14] S. Wang, J. Shi, Z. Fu, X. He, F. Liu, C. Yang, Y. Zhou, Z. Fei, J. Gong, J. Fu, et al. World
316 action models: The next frontier in embodied ai. *arXiv preprint arXiv:2605.12090*, 2026.

- 317 [15] S. Li, Y. Gao, D. Sadigh, and S. Song. Unified video action model. In *Proc. of the Robotics:
318 Science and Systems (RSS), 2025, 2025.*
- 319 [16] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion
320 policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics
321 Research*, 44(10-11):1684–1704, 2025.
- 322 [17] T. Chen, Z. Chen, B. Chen, Z. Cai, Y. Liu, Z. Li, Q. Liang, X. Lin, Y. Ge, Z. Gu, et al.
323 Robotwin 2.0: A scalable data generator and benchmark with strong domain randomization
324 for robust bimanual robotic manipulation. *arXiv preprint arXiv:2506.18088*, 2025.
- 325 [18] M. Du, S. Nair, D. Sadigh, and C. Finn. Behavior retrieval: Few-shot imitation learning by
326 querying unlabeled datasets. In *Proc. of the Robotics: Science and Systems (RSS), 2024, 2024.*
- 327 [19] L.-H. Lin, Y. Cui, A. Xie, T. Hua, and D. Sadigh. Flowretrieval: Flow-guided data retrieval for
328 few-shot imitation learning. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL),
329 2024*. URL <https://openreview.net/forum?id=FHnVRmeqxf>.
- 330 [20] M. Memmel, J. Berg, B. Chen, A. Gupta, and J. Francis. STRAP: Robot sub-trajectory retrieval
331 for augmented policy learning. In *Proc. of the Thirteenth International Conference on Learning
332 Representations (ICLR), 2025*. URL <https://openreview.net/forum?id=4VHiptx7xe>.
- 333 [21] R. Punamiya, D. Patel, P. Aphiwetsa, P. Kuppili, L. Y. Zhu, S. Kareer, J. Hoffman, and
334 D. Xu. Egobridge: Domain adaptation for generalizable imitation from egocentric human
335 data. In *Proc. of the Thirty-ninth Annual Conference on Neural Information Processing Sys-
336 tems (NeurIPS), 2026*. URL <https://openreview.net/forum?id=FGMBxzipgis>.
- 337 [22] V. Jain, M. Attarian, N. J. Joshi, A. Wahid, D. Driess, Q. Vuong, P. R. Sanketi, P. Sermanet,
338 S. Welker, C. Chan, I. Gilitschenski, Y. Bisk, and D. Dwibedi. Vid2robot: End-to-end video-
339 conditioned policy learning with cross-attention transformers. In *Proc. of the Robotics: Science
340 and Systems (RSS), 2024, 2024.*
- 341 [23] A. Bahety, P. Mandikal, B. Abbatematteo, and R. Martín-Martín. Screwmimic: Bimanual
342 imitation from human videos with screw space projection. In *Proc. of the Robotics: Science
343 and Systems (RSS), 2024, 2024.*
- 344 [24] M. Hong, A. Liang, K. Kim, H. Rajaprakash, J. Thomason, E. Bıyık, and J. Zhang. Hand me
345 the data: Fast robot adaptation via hand path retrieval. *arXiv preprint arXiv:2505.20455*, 2025.
- 346 [25] G. Papagiannis, N. Di Palo, P. Vitiello, and E. Johns. R+ x: Retrieval and execution from
347 everyday human videos. In *Proc. of the 2025 IEEE International Conference on Robotics and
348 Automation (ICRA)*, pages 8284–8290. IEEE, 2025.
- 349 [26] S. Sha, Y. Wang, B. Huang, A. Loquercio, and Y. Li. Efficient and reliable teleoperation
350 through real-to-sim-to-real shared autonomy. *arXiv preprint arXiv:2603.17016*, 2026.
- 351 [27] C. Schaff and M. R. Walter. Residual policy learning for shared autonomy. In *Proc. of the
352 Robotics: Science and Systems (RSS), 2020, 2020.*
- 353 [28] N. Carion, L. Gustafson, Y.-T. Hu, S. Debnath, R. Hu, D. Suris, C. Ryali, K. V. Alwala,
354 H. Khedr, A. Huang, J. Lei, T. Ma, B. Guo, A. Kalla, M. Marks, J. Greer, M. Wang, P. Sun,
355 R. Rädle, T. Afouras, E. Mavroudi, K. Xu, T.-H. Wu, Y. Zhou, L. Momeni, R. Hazra, S. Ding,
356 S. Vaze, F. Porcher, F. Li, S. Li, A. Kamath, H. K. Cheng, P. Dollár, N. Ravi, K. Saenko,
357 P. Zhang, and C. Feichtenhofer. Sam 3: Segment anything with concepts, 2025. URL <https://arxiv.org/abs/2511.16719>.
- 359 [29] O. Siméoni, H. V. Vo, M. Seitzer, F. Baldassarre, M. Oquab, C. Jose, V. Khalidov,
360 M. Szafraniec, S. Yi, M. Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*,
361 2025.

- 362 [30] R. Cadene, S. Alibert, A. Soare, Q. Gallouedec, A. Zouitine, S. Palma, P. Kooijmans, M. Ar-
363 actingi, M. Shukor, D. Aubakirova, M. Russi, F. Capuano, C. Pascal, J. Choghari, J. Moss,
364 and T. Wolf. Lerobot: State-of-the-art machine learning for real-world robotics in pytorch.
365 <https://github.com/huggingface/lerobot>, 2024.
- 366 [31] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation
367 with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.