SAM2Act: Integrating Visual Foundation Model with A Memory Architecture for Robotic Manipulation

Anonymous CVPR submission

Paper ID *****

Abstract

Robotic manipulation systems operating in diverse, dy-001 002 namic environments must exhibit three critical abilities: multitask interaction, generalization to unseen scenarios, 003 and spatial memory. While significant progress has been 004 made in robotic manipulation, existing approaches often 005 fall short in generalization to complex environmental varia-006 tions and addressing memory-dependent tasks. To bridge 007 this gap, we introduce SAM2Act, a multi-view robotic 008 transformer-based policy that leverages multi-resolution 009 010 upsampling with visual representations from large-scale foundation model. SAM2Act achieves a state-of-the-art av-011 erage success rate of 86.8% across 18 tasks in the RLBench 012 benchmark, and demonstrates robust generalization on The 013 Colosseum benchmark, with only a 4.3% performance 014 gap under diverse environmental perturbations. Building 015 016 on this foundation, we propose SAM2Act+, a memorybased architecture inspired by SAM2, which incorporates 017 018 a memory bank, an encoder, and an attention mechanism to enhance spatial memory. To address the need for evaluating 019 020 memory-dependent tasks, we introduce MemoryBench, a 021 novel benchmark designed to assess spatial memory and 022 action recall in robotic manipulation. SAM2Act+ achieves competitive performance on MemoryBench, significantly 023 outperforming existing approaches and pushing the bound-024 aries of memory-based robotic systems. 025

026 1. Introduction

027 The world in which we live is diverse and constantly changing, encompassing a wide variety of objects, scenes, and 028 029 environmental conditions. Consider the seemingly simple 030 task of following a recipe when cooking: we can seamlessly 031 perform the action of picking it up and sprinkling it into the pan, recognize salt even if it comes in different types of 032 container, and remember whether we have already added 033 salt. Humans excel in such environments because they can 034 035 interact with their surroundings to achieve specific goals,

generalize to unseen scenarios, and retain knowledge from
past experiences [33]. These abilities—multitask interac-
tion, generalization, and memory—serve as guiding princi-
ples for developing robotic systems capable of operating in
similarly complex environments.036
037038
039

Significant progress has been made in robotic manipula-041 tion through prior work. Early methods, such as the Trans-042 porter Network [39] and CLIPort [31], demonstrated effec-043 tive 2D action-centric manipulation but were limited in their 044 ability to handle spatially complex tasks. More recent ap-045 proaches, such as PerAct [32] and RVT [9], have pushed to-046 ward 3D-based manipulation. PerAct employs a multitask 047 transformer that interprets language commands and predicts 048 keyframe poses, achieving strong results across a variety of 049 tasks. RVT builds on this foundation by adopting a 2.5D 050 representation, improving training efficiency and inference 051 speed. Its successor, RVT-2, further enhances performance 052 with a coarse-to-fine strategy, increasing precision for high-053 accuracy tasks. Despite these advances, important chal-054 lenges remain, including improving multitask performance, 055 enhancing generalization to novel environment configura-056 tions, and integrating memory mechanisms for tasks requir-057 ing episodic recall. 058

We introduce SAM2Act, a multi-view robotics 059 transformer-based policy that enhances feature repre-060 sentation by integrating multi-resolution upsampling with 061 visual embeddings from large-scale foundation models. 062 Built on the RVT-2 multi-view transformer, SAM2Act 063 achieves strong multitask success and generalization. 064 Building on this foundation, we introduce SAM2Act+, 065 which incorporates a memory-based architecture inspired 066 by SAM2's approach. Using a memory bank, an encoder, 067 and an attention mechanism, SAM2Act+ enables episodic 068 recall to solve spatial memory-dependent manipulation 069 tasks. We evaluate SAM2Act and SAM2Act+ using 070 MemoryBench, a new benchmark suite that tests policies' 071 spatial memory capabilities and the ability to retain and 072 recall past actions. SAM2Act+ achieves competitive 073 performance on MemoryBench, with an average accuracy 074 of 94.3%, outperforming next highest baseline by a huge 075

131



Figure 1. SAM2Act is a multi-view, language-conditioned behavior cloning policy trained with fewer demonstrations. Given a language instruction, it can execute high-precision tasks, such as turning the tiny knob on the lamp. It also generalizes to various environmental variations, such as changes in lighting conditions. Through further training with our proposed memory architecture, it now evolves into SAM2Act+, which is now capable of solving tasks that require implicit spatial memory—such as remembering where the robot previously stored the pliers, as depicted in the above figure.

076 margin of 39.3%. Furthermore, we assess the generalization capabilities of SAM2Act on The Colosseum 077 [26], a benchmark designed to test robotic manipulation 078 under various environmental perturbations. 079 SAM2Act 080 demonstrates robust performance on The Colosseum with an average decrease of 4.3% across all perturbations, 081 highlighting its ability to generalize effectively in diverse 082 and challenging scenarios. Lastly, our approach outper-083 forms the baseline methods in real-world evaluations while 084 085 exhibiting comparable generalization and spatial memory capabilities. 086

In summary, this work makes three key contributions. 087 First, we introduce a novel model formulation that lever-088 ages visual foundation models to solve high-precision, 089 memory-dependent manipulation tasks. Second, we pro-090 pose MemoryBench, a evaluation benchmark for assess-091 ing spatial memory in behavior cloning models. Finally, 092 we present empirical results and insights on the model's 093 performance across both simulation and real-world tasks. 094

095 2. Related Work

2.1. 3D-based Robotic Transformer for Manipula tion

098 2D-based methods [2, 5, 31, 39, 41] are effective for simple pick-and-place tasks due to fast training, low hardware 099 requirements, and minimal computational cost. However, 100 they depend on pretrained image encoders and fail in tasks 101 102 requiring high precision, robust spatial interaction, or resilience to environmental and camera variations [26]. Re-103 cent work addresses these limitations with 3D perception. 104 Methods like PolarNet [4], M2T2 [38], and Manipulate-105 Anything [7] reconstruct point clouds, while C2F-ARM 106 [15] and PerAct [32] use voxel-based 3D representations. 107 108 Act3D [8] and ChainedDiffuser [36] adopt multi-scale 3D

features. RVT [9] introduces 2.5D multi-view images for109faster training, refined by RVT-2 [10] with a coarse-to-fine110architecture for improved precision. Our work, SAM2Act,111combines RVT-2's spatial reasoning with enhanced virtual112images from the SAM2 visual encoder, achieving high pre-113cision and generalization across diverse tasks.114

2.2. Visual Representations for Robot Learning

Robotics research heavily relies on visual representations 116 from computer vision to process high-dimensional inputs 117 and improve policy learning. Visual representations are in-118 tegrated into robot learning through pre-training [23–25], 119 co-training [19, 20, 29, 37], or frozen encoders [28, 34, 40], 120 all of which effectively support policy training. These 121 representations also enhance invariance, equivariance, and 122 out-of-distribution generalization [6, 26, 35]. SAM-E [40] 123 demonstrates the use of a pre-trained SAM encoder for 124 robotic manipulation by leveraging image embeddings for 125 policy learning. Expanding on this, our approach employs 126 the SAM2 visual encoder to generate image embeddings 127 for robotic transformers and utilizes its multi-resolution fea-128 tures to improve convex upsampling for next-action predic-129 tion. 130

2.3. Memory in Robotics

Memory is a fundamental component of human cognition, 132 and equipping generalist robotic agents with episodic and 133 semantic memory is crucial for enabling them to perform 134 complex tasks effectively [17]. Early research on mem-135 ory in robotics primarily addressed navigation tasks, re-136 lying on semantic maps that were often constrained in 137 scope [1, 3, 11]. Recent advancements leverage represen-138 tations derived from vision-language models (VLMs) and 139 Large Vision Models (LVMs), utilizing voxel maps or neu-140 ral feature fields to encode, store, and retrieve information 141

142 [7, 13, 14, 22]. Alternative methods represent semantic memory for manipulation tasks using Gaussian splats to en-143 144 code spatial information [18, 30]. In contrast, our approach draws inspiration from the framework of Partially Observ-145 146 able Markov Decision Processes (POMDPs) [21], incorporating memory directly into the training process. By inte-147 grating spatial memory from past actions into the agent's 148 belief state, we enhance the robustness and adaptability of 149 150 learned policies.

151 3. MemoryBench: A Memory Benchmark for152 Robotic Manipulation

We introduce MemoryBench, a benchmark designed to
systematically evaluate the spatial memory capabilities of
robotic manipulation policies. In subsection 3.1, we begin
by outlining the logic and rules behind task design. We will
then describe the tasks we have developed in subsection 3.2.

158 3.1. Task Design

Unlike standard RLBench tasks [16], many of which in-159 volve long-horizon scenarios, our tasks are specifically de-160 signed to require spatial memory. Without such memory, 161 the agent would be forced to rely on random actions. To 162 create these tasks, we intentionally violate the Markov as-163 sumption, which states that in a Markov Decision Process 164 (MDP), the next observation depends solely on the current 165 observation and action: 166

167
$$P(o_{t+1} | o_1, a_1, \dots, o_t, a_t) = P(o_{t+1} | o_t, a_t)$$

168 This assumption implies that knowing only o_t and a_t is sufficient to predict o_{t+1} . However, in our tasks, we de-169 170 sign scenarios where two distinct action histories lead to the same observation o_t , but require different subsequent 171 172 actions. This forces the agent to recall which action history 173 led to o_t to perform the correct next action. Furthermore, 174 we standardized the language instructions to prevent unin-175 tentional leakage of spatial information that could aid the model in memory-based tasks. These principles guided the 176 development of our spatial memory-based tasks. 177

3.2. Spatial Memory-based Tasks

MemoryBench extends the RLBench simulator to pro-179 180 vide scripted demonstrations for three spatial memory tasks: reopen_drawer, put_block_back, and 181 rearrange_block. Each task is designed to evaluate a 182 specific aspect of spatial memory and adheres to the prin-183 184 ciples outlined in Section 3.1. To introduce complexity, these tasks include two to four variations and additional 185 steps-such as pressing a button mid-sequence-that dis-186 rupt the Markov property. This forces the agent to rely on 187 memory rather than solely on immediate observations. 188

The reopen_drawer task evaluates the agent's abilityto recall 3D spatial information along the z-axis. Initially,

one of three drawers (top, middle, or bottom) is open. The 191 agent must close the open drawer, press a button on the ta-192 ble, and then reopen the same drawer. After the button is 193 pressed, all drawers are closed, and the scene becomes vi-194 sually indistinguishable, requiring the agent to use memory 195 to identify the correct drawer. This task tests the agent's 196 ability to recall spatial states over a temporal sequence. 197 The put_block_back task tests the agent's ability to re-198 member 2D spatial information on the x-y plane. Four red 199 patches are placed on a table, with a block initially posi-200 tioned on one of them. The agent should move the block to 201 the center of the patches, press a button, and return the block 202 to its original position. The agent must rely on its memory 203 of the block's initial location to succeed, demonstrating its 204 capability to encode and retrieve 2D spatial information. 205

The rearrange_block task evaluates the agent's 206 ability to perform backward reasoning by recalling and re-207 versing prior actions. Initially, one block is placed on one 208 of two red patches, while the other patch remains empty. A 209 second block is positioned at the center of both patches. The 210 agent must move the second block to the empty patch, press 211 a button, and then relocate the first block off its patch. Suc-212 cessfully completing this task requires the agent to deter-213 mine which block to move without having interacted with 214 the correct one in previous actions, thereby testing its ca-215 pacity for backward spatial memory reasoning. These tasks 216 collectively evaluate both forward and backward spatial rea-217 soning across 3D (z-axis) and 2D (x-y plane) spaces. By 218 introducing non-Markovian elements, they emphasize the 219 need for memory representations to solve complex sequen-220 tial decision-making problems. 221

4. Method

Our method, SAM2Act, enables precise 3D manipulation 223 with strong generalization across environmental and object-224 level variations. Building upon the RVT-2 framework [10], 225 SAM2Act introduces key architectural innovations that en-226 hance visual feature representation and task-specific rea-227 soning. The architecture reconstructs a point cloud of the 228 scene, renders it from virtual cameras at orthogonal views, 229 and employs a two-stage multi-view transformer (coarse-230 to-fine) to predict action heatmaps. The coarse branch 231 generates zoom-in heatmaps to localize regions of inter-232 est, while the fine branch refines these into precise action 233 heatmaps. SAM2Act leverages the pre-trained SAM2 en-234 coder [27] to extract multi-resolution image embeddings, 235 which are further refined through the multi-resolution up-236 sampling technique to predict accurate translation heatmaps 237 with minimal information loss. To address tasks requiring 238 spatial memory, SAM2Act+ extends the SAM2Act archi-239 tecture by incorporating memory-based components. These 240 include Memory Bank, Memory Encoder, and Memory At-241 tention, enabling the model to encode historical actions and 242

CVPR 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 2. Simulation and Real Tasks. We demonstrate the effectiveness of SAM2Act+ in solving memory-based tasks by evaluating it against baselines on the three benchmark memory tasks (shown at the top). Additionally, we validate our approach using a Franka Panda robot on four real-world tasks (shown at the bottom), including tests under out-of-distribution perturbations.

condition current observations. This memory-based policy enhances the agent's ability to predict actions based on
past contextual information, significantly improving performance in tasks that require sequential decision-making.

In the following sections, we detail the SAM2Act architecture (subsection 4.1), including its multi-resolution
upsampling mechanism (Figure 4). We also present
the SAM2Act+ extension, which integrates memory-based
components for solving spatial memory tasks (subsection 4.2).

4.1. SAM2Act: Multi-Resolution Upsampling for Enhanced Visual Feature Representation

A distinctive feature of SAM2Act is the incorporation of 255 256 the SAM2Act Module into the manipulation backbone for training, as illustrated in Figure 4. The coarse and fine 257 SAM2Act Modules share the same architecture, with the 258 fine branch generating additional features to predict actions 259 beyond translation, while the coarse branch focuses exclu-260 261 sively on translation. Point-cloud representations are recon-262 structed from raw image inputs, and virtual images are generated from three viewpoints using virtual cameras. Instead 263 of directly inputting these images into the multi-view trans-264 former, their RGB channels are duplicated and processed 265 by the SAM2 [27] image encoder, which produces object-266 centric multi-resolution embeddings. These embeddings, 267 generated at three resolution levels, are combined with vir-268 tual images containing RGB, depth, 3D translation coordi-269 nates, and language instructions before being fed into the 270 multi-view transformer. 271

To adapt the SAM2 image encoder to our domain, 272 we fine-tune it using Low-Rank Adaptation (LoRA) [12] 273 with a default rank of 16, which enables domain adap-274 tation with minimal computational cost while maintain-275 ing model efficiency. Additionally, to fully leverage the 276 multi-resolution embeddings produced by the SAM2 im-277 age encoder, we introduce a multi-resolution upsampling 278 method. This method uses the embeddings as auxiliary 279 inputs to enhance the generation of translation heatmaps, 280 thereby improving spatial precision and overall system per-281 formance. The multi-resolution upsampling mechanism, 282 also detailed in Figure 4, leverages cascaded convex up-283



Figure 3. Overview of the SAM2Act (top) and SAM2Act+ (bottom) architectures. The SAM2Act architecture leverages the SAM2 image encoder to generate prompt-conditioned, multi-resolution embeddings, fine-tuned with LoRA for efficient adaptation to manipulation tasks. A multi-view transformer aligns spatial coordinates with language instructions, while a cascaded multi-resolution upsampling mechanism refines feature maps and generates accurate translation heatmaps. SAM2Act+ extends this architecture by incorporating memory-based components, including the Memory Encoder, Memory Attention, and Memory Bank, into the coarse branch. These components enable memory-driven reasoning by processing historical heatmaps and integrating prior observations, allowing the agent to predict actions based on stored contextual information. Observations are reconstructed into point clouds, rendered into three virtual images, and lifted into 3D translation points, enabling precise spatial reasoning across both architectures.



Figure 4. **SAM2Act Module and multi-resolution upsampling mechanism.** A cascade of three convex upsamplers processes feature maps at increasing resolutions, integrating multi-resolution embeddings from the SAM2 image encoder through elementwise addition and layer normalization. The upsamplers progressively refine features, doubling spatial dimensions at each stage, to generate accurate translation heatmaps while capturing fine-grained spatial details critical for manipulation tasks.

samplers to progressively refine feature maps across resolutions. Let $X^l \in \mathbb{R}^{B \times C^l \times H^l \times W^l}$ denote the feature maps at stage *l* and $E^l \in \mathbb{R}^{B \times C^l \times H^l \times W^l}$ the corresponding multiresolution embedding from SAM2. Also let $U(\cdot)$ denote the upsampling operator that doubles the spatial dimensions. 288 The feature maps are updated at each stage as follows: 289

$$X^{l+1} = \text{LayerNorm}(U(X^l) \oplus E^l),$$
 290

where \oplus represents element-wise addition. The upsampling 291 operator U is defined as: 292

$$U: \mathbb{R}^{B \times C^l \times H^l \times W^l} \to \mathbb{R}^{B \times (C^l/2) \times (2H^l) \times (2W^l)}.$$
 293

At each stage, the output of the upsampler is combined with 294 the corresponding multi-resolution embedding E^{l} from the 295 SAM2 encoder, ensuring alignment between the multi-296 resolution features and the decoder's spatial refinement pro-297 cess. A layer normalization step follows each addition to 298 stabilize training and maintain feature coherence. This re-299 sults in direct integration of the embeddings into the trans-300 lation heatmap generation process. The cascading structure 301 refines features across multiple resolutions, capturing fine-302 grained spatial details critical for manipulation tasks. 303

379

Algorithm 1 Forward Pass of SAM2Act+ Module

1:	Initialize:	Number	of steps	<i>N</i> , 1	naximun	n number	of
	memories 1	M, number	er of viev	vs V,	empty m	nemory ba	nk
	Q with V s	eparate F	FIFO que	ues, i	nput X		
	· · ·						

2: for i = 1 to N do

- 3: for j = 1 to V do
- 4: Get embeddings \mathcal{E}_{raw} from MVT $T_{mv}(X_j)$
- Retrieve past memories \mathcal{M}_{old} from Q[j]5:
- Get memory-conditioned embeddings \mathcal{E}_{mem} 6: from Memory Attention $T_{mem}(\mathcal{E}_{raw}, \mathcal{M}_{old})$

7: Predict translation heatmap \mathcal{H} with upsampler $U(\mathcal{E}_{mem})$

Encode new memory \mathcal{M}_{new} using Memory En-8: coder $E_{mem}(\mathcal{H}, \mathcal{E}_{raw})$

9: Store new memory
$$Q[j] \leftarrow Q[j] \cup \{\mathcal{M}_{new}\}$$

if |Q[j]| = M then 10:

 $Q[j] \leftarrow Q[j]_{2:n}$ 11:

end if 12:

end for 13.

14: end for

304

305

306

4.2. SAM2Act+: Action Memory Architecture for Improved Spatial Awareness in Past Observations

To extend the SAM2Act architecture (subsection 4.1) with 307 308 memory-based capabilities inspired by SAM2, we intro-309 duce SAM2Act+, a task-specific variant designed for solving memory-based tasks. SAM2Act+ integrates the three 310 311 core memory components from SAM2-Memory Attention, Memory Encoder, and Memory Bank-into the coarse 312 313 branch of SAM2Act. Originally developed for object tracking in SAM2, these components are adapted to align with 314 the needs of SAM2Act+, enabling the agent to retain prior 315 actions and observations for sequential decision-making. 316 317 In SAM2, the Memory Encoder processes predicted ob-318 ject masks, while the Memory Attention module fuses im-319 age embeddings with positional information from previous 320 frames. SAM2Act+ adopts a similar structure: the predicted heatmaps, which serve as binary indicators of spa- § 5.3 Can SAM2Act generalize across object and environmen-321 tial positions in the image, function analogously to object 322 masks. This conceptual alignment ensures a seamless inte- § 5.4 Can SAM2Act+ solve spatial memory-based tasks that 323 324 gration of memory mechanisms, allowing the agent to leverage stored information to predict subsequent actions based 325 on historical context. 326

Architecture. The SAM2Act+ architecture is illustrated 327 328 in Figure 3. After pretraining SAM2Act in Stage 1, we 329 freeze the SAM2 image encoder and the multi-view trans-330 former in the coarse branch, as these components effectively generate robust embeddings for multi-view images in ma-331 nipulation tasks. We also freeze the entire fine branch, given 332 333 its proven ability to predict fine-grained actions accurately. 334 The reason why we only fine-tune the coarse branch is be-

cause it focuses on generating heatmaps that provide richer 335 contextual information for recalling past actions. The fine 336 branch, in contrast, primarily emphasizes small objects or 337 localized regions, which typically contain less information 338 relevant to memory-based tasks. 339

Training. To train SAM2Act+, we fine-tune the coarse 340 branch by integrating the three memory components (and 341 train them from scratch) with the multi-resolution up-342 sampling module. During fine-tuning, consecutive ac-343 tion keyframes are sampled as input, training the multi-344 resolution upsampler to predict new translations condi-345 tioned on memory. The memory components function sim-346 ilarly to their implementation in SAM2 for object track-347 ing, with one key distinction: the input to the Memory 348 Encoder. Instead of using image embeddings from the 349 SAM2 image encoder, we input feature embeddings gen-350 erated by the multi-view transformer (not conditioned by 351 memory). This adaptation ensures that memory encod-352 ing incorporates multi-view information while maintaining 353 independence in handling stored representations. Virtual 354 images are treated independently during memory encod-355 ing and attention, with each view's memory encoded sep-356 arately. Feature embeddings from each view are attended to 357 using their corresponding stored memories, preserving spa-358 tial and contextual alignment while leveraging fused multi-359 view information. This structured approach prevents cross-360 view interference and enhances the model's ability to reason 361 over sequential tasks. The memory-based forward pass for 362 SAM2Act+ is outlined in 1. By incorporating the memory 363 mechanism, SAM2Act+ enhances performance in scenarios 364 requiring long-term reasoning, enabling the agent to make 365 informed decisions based on historical context. 366

5. Experiments

We study SAM2Act and SAM2Act+ in both simulated and 368 real-world environments. Specifically, we are interested in 369 answering the following questions: 370 § 5.2 How does SAM2Act compare with state-of-the-art 3D 371 manipulation policies? 372 373 tal perturbations? 374 375 other baselines cannot? 376 § 5.5 How well does SAM2Act and SAM2Act+ perform on 377 real-world tasks? 378

5.1. Experimental Setup

We benchmark SAM2Act in both simulated and real-world 380 environments. The simulated environments serve as a con-381 trolled platform to ensure reproducible and fair compar-382 isons. The real-world experiments demonstrate the appli-383 cability of the method to real-world settings. Section 5.1 384 details our experimental setup and outlines the evaluation 385

426

427

428

429

449

Table 1. Multi-Task Performance on RLBench. We report the success rates for 18 RLBench tasks [16], along with the average success rate and ranking across all tasks. Our method, SAM2Act, outperforms all baselines, achieving a significant performance margin of 5.8% over RVT-2 [10], the current state-of-the-art 3D keyframe-based behavior cloning (BC) policy.

Method	Avg. Success ↑	Avg. Rank \downarrow	Close Jar	Drag Stick	Insert Peg	Meat off Grill	Open Drawer	Place Cups	Place Wine	Push Buttons
PerAct [32]	49.4 ± 4.3	4.6	55.2 ± 4.7	89.6 ± 4.1	5.6 ± 4.1	70.4 ± 2.0	88.0 ± 5.7	2.4 ± 3.2	44.8 ± 7.8	92.8 ± 3.0
RVT [9]	62.9 ± 3.7	3.6	52.0 ± 2.5	99.2 ± 1.6	11.2 ± 3.0	88.0 ± 2.5	71.2 ± 6.9	4.0 ± 2.5	91.0 ± 5.2	$\textbf{100.0} \pm 0.0$
RVT-2 [10]	81.4 ± 3.1	1.9	$\textbf{100.0} \pm 0.0$	99.0 ± 1.7	40.0 ± 0.0	99.0 ± 1.7	74.0 ± 11.8	38.0 ± 4.5	95.0 ± 3.3	$\textbf{100.0} \pm 0.0$
SAM-E [40]	70.6 ± 0.7	2.6	82.4 ± 3.6	$\textbf{100.0}\pm0.0$	18.4 ± 4.6	95.2 ± 3.3	$\textbf{95.2} \pm 5.2$	0.0 ± 0.0	94.4 ± 4.6	$\textbf{100.0}\pm0.0$
SAM2Act (Ours)	$\textbf{86.8} \pm 0.5$	1.8	99.0 ± 2.0	99.0 ± 2.0	$\textbf{84.0} \pm 5.7$	98.0 ± 2.3	83.0 ± 6.0	$\textbf{47.0} \pm 6.0$	93.0 ± 3.8	$\textbf{100.0}\pm0.0$
Method	Put in Cupboard	Put in Drawer	Put in Safe	Screw Bulb	Slide Block	Sort Shape	Stack Blocks	Stack Cups	Sweep to Dustpan	Turn Tap
PerAct [32]	28.0 ± 4.4	51.2 ± 4.7	84.0 ± 3.6	17.6 ± 2.0	74.0 ± 13.0	16.8 ± 4.7	26.4 ± 3.2	2.4 ± 2.0	52.0 ± 0.0	88.0 ± 4.4
RVT [9]	49.6 ± 3.2	88.0 ± 5.7	91.2 ± 3.0	48.0 ± 5.7	81.6 ± 5.4	36.0 ± 2.5	28.8 ± 3.9	26.4 ± 8.2	72.0 ± 0.0	93.6 ± 4.1
RVT-2 [10]	66.0 ± 4.5	96.0 ± 0.0	96.0 ± 2.8	88.0 ± 4.9	92.0 ± 2.8	35.0 ± 7.1	$\textbf{80.0} \pm 2.8$	69.0 ± 5.9	$\textbf{100.0}\pm0.0$	99.0 ± 1.7
SAM-E [40]	64.0 ± 2.8	92.0 ± 5.7	95.2 ± 3.3	78.4 ± 3.6	95.2 ± 1.8	34.4 ± 6.1	26.4 ± 4.6	0.0 ± 0.0	$\textbf{100.0} \pm 0.0$	$\textbf{100.0} \pm 0.0$
SAM2Act (Ours)	$\textbf{75.0} \pm 3.8$	99.0 ± 2.0	$\textbf{98.0} \pm 2.3$	$\textbf{89.0} \pm 2.0$	86.0 ± 4.0	$\textbf{64.0} \pm 4.6$	76.0 ± 8.6	$\textbf{78.0} \pm 4.0$	99.0 ± 2.0	96.0 ± 5.7

386 methodology.

Simulation Setup. All simulated experiments were con-387 ducted in the CoppeliaSim environment via PyRep, using a 388 389 7-DoF Franka Emika Panda robot in a tabletop setting. Ob-390 servations were captured from five RGB-D cameras-front, left shoulder, right shoulder, overhead and wrist-each at 391 $128 \,\mathrm{px} \times 128 \,\mathrm{px}$. The robot receives a keyframe specify-392 ing translation and quaternion orientation and utilizes an 393 394 OMPL-based motion planner to move to the target pose.

Real-robot Setup. We validate SAM2Act in real-world 395 scenarios using a Franka Emika Panda robot with a Robotiq 396 2F-85 gripper and a exocentric Intel RealSense D455 depth 397 sensor. We study four manipulation tasks, aligning three 398 with RVT-2 for comparison and introducing a new memory-399 based task. For each task, we collect 10-15 demonstrations 400 via kinesthetic teaching and scripted execution with scene 401 and object variations. As in Figure 2, we evaluate SAM2Act 402 against RVT-2 for tasks (a)-(c) and SAM2Act+ for mem-403 ory task (d). Each task undergoes 10 in-distribution and 10 404 out-of-distribution trials, including environmental perturba-405 tions, measuring total success. 406

18 RLBench & MemoryBench Tasks. To evaluate the 407 general performance of SAM2Act and the memory capabil-408 ities of SAM2Act+, we conducted simulation experiments 409 on two benchmarks: a subset of 18 tasks from RLBench and 410 MemoryBench. RLBench is a standard multi-task manip-411 ulation benchmark, from which we selected 18 tasks well-412 studied in prior work. MemoryBench is a curated set of 413 three tabletop manipulation tasks in CoppeliaSim that re-414 quire the trained policy to have both semantic and spatial 415 416 memory of past scenes and actions. In both benchmarks, each task is defined by a language instruction with 2-60 417 variations (e.g., handling objects, locations, and colors). We 418 419 collected 100 demonstrations per task for training and held 420 out 25 unseen demonstrations per task for testing. All poli-421 cies are evaluated four times to obtain standard deviations.

3D Baselines. We benchmark SAM2Act and
SAM2Act+ against the current state-of-the-art 3D nextbest-pose prediction model, RVT-2. RVT-2 is a multi-

view robotics transformer that leverages a coarse-to-fine approach on the constructed point cloud to predict the next best action heatmap. We also compare with RVT [9], Per-Act [32], and SAM-E [40].

5.2. Performances Across 18 RLBench Tasks

Table 1 compares SAM2Act with prior keyframe-based 430 3D BC methods on the RLBench benchmark. Overall, 431 SAM2Act achieves an average success rate of $86.8\% \pm 0.5$, 432 surpassing the previous best (RVT-2) by 5.4%. A closer 433 look at individual tasks reveals that SAM2Act ranks first 434 in 9 out of 18 tasks and remains highly competitive in 435 7 others, coming within one successful attempt or 4% 436 of the best performance. These tasks include Close Jar, 437 Drag Stick, Meat Off Grill, Place Wine, Screw Bulb, Sweep 438 to Dustpan, and Turn Tap. The largest margin of im-439 provement occurs in Insert Peg, where SAM2Act ex-440 ceeds RVT-2 by 44% (approximately 2.1×), and in Sort 441 Shape, where it outperforms RVT-2 by 29%. Both tasks 442 require precise manipulation, underscoring the effective-443 ness of SAM2Act's multi-resolution upsampling strategy. 444 These results establish SAM2Act as a leading policy for 445 complex 3D tasks, highlighting its ability to handle high-446 precision manipulations - an area where prior methods have 447 struggled. 448

5.3. Semantic Generalization across Tasks

The results evaluated in subsection 5.2 were obtained 450 by training and testing models within the same environ-451 ment. However, to truly assess generalization perfor-452 mance, policies must remain robust against both environ-453 mental and object-level perturbations. We therefore trained 454 SAM2Act and the baseline methods on 20 tasks from The 455 Colosseum benchmark and tested them under 13 different 456 perturbation categories over three runs. SAM2Act exhibits 457 the smallest performance drop compared to the base-458 lines, with an average decrease of 4.3% (standard deviation 459 of 3.59%). Notably, it proves particularly robust to envi-460 ronmental perturbations - such as changes in lighting, table 461

Table 2. **The Colosseum results**. Task-average success rate percentage change for SAM2Act and other baselines across 13 perturbation factors from The Colosseum, relative to evaluations without perturbations. Our approach, SAM2Act, demonstrates the lowest average percentage change across all perturbations, with a minimal drop of $-4.3\pm3.6\%$, highlighting its robustness in handling various environmental and object-level perturbations.

Method	Average ↑	MO-Color ↑	RO-Color ↑	MO-Texture †	RO-Texture ↑	MO-Size ↑	RO-Size ↑
RVT-2 [10]	-19.5 ± 2.8	-20.7 ± 1.0	$-11.8 {\pm} 0.8$	-13.3±4.6	-11.4±3.7	-13.2±3.1	-17.7 ± 0.1
SAM2Act (SAM2 \rightarrow SAM)	-20.7 ± 1.2	-26.1 ± 0.7	-15.7 ± 2.9	-15.0 ± 3.3	-16.5 ± 6.2	-18.7 ± 1.9	-19.8 ± 1.3
SAM2Act (w/o Multi-res Input)	-19.1 ± 4.5	-15.5 ± 6.4	-13.5 ± 4.6	-20.4 ± 0.5	-16.6 ± 6.1	-21.3 ± 7.5	-12.6 ± 7.5
SAM2Act (Ours)	-4.3 ±3.6	-1.1±2.5	- 0.7 ±7.2	-3.3 ±2.4	24.72 ±6.1	-15.9 ± 5.0	0.9 ±6.8
Method	Light Color \uparrow	Table Color \uparrow	Table Texture \uparrow	Distractor \uparrow	Background Texture \uparrow	Camera Pose ↑	All Perturbations \uparrow
RVT-2 [10]	-15.6±1.3	-26.5 ± 4.4	-14.6 ± 4.4	-4.9 ± 5.3	$-4.4{\pm}4.0$	-19.5±2.8	-77.9±1.7
SAM2Act (SAM2 \rightarrow SAM)	-16.3 ± 1.2	-23.5 ± 5.3	-12.3 ± 3.1	$0.6{\pm}2.9$	-5.4 ± 3.2	-20.7 ± 1.2	-79.5±2.5
SAM2Act (w/o Multi-res Input)	-7.2 ± 3.6	-18.3 ± 6.1	-17.5 ± 3.3	-4.6 ± 3.5	-5.7 ± 3.5	-19.1 ± 4.5	-73.8 ± 2.2
SAM2Act (Ours)	4.5 ±4.4	1.1±2.5	-3.7 ±5.2	1.7 ±1.7	-1.5 ±2.7	- 4.3 ±3.6	-58.3 ±4.4

462 color/texture, the addition of distractors, and even camera
463 pose – while also maintaining competitive performance un464 der object-level perturbations.

465 5.4. Performance on MemoryBench

In Table 3, we evaluate SAM2Act+ against SoTA 3D 466 BC model, RVT-2 on MemoryBench, training all mod-467 els in a single-task setting to isolate memory-related 468 challenges (e.g., opening the wrong drawer rather than 469 unrelated mid-task failures). This setup ensures that 470 performance differences stem from memory capabili-471 ties. For a random agent, the expected success rates 472 are determined by the number of possible choices per 473 task: 33% for reopen_drawer (three drawers), 25% 474 for put_block_back (four patches), and 50% for 475 rearrange_block (two blocks). However, variations 476 in task complexity, fixed training data, and imbalanced 477 task distributions lead to slight deviations from these base-478 lines. Our proposed memory-based model, SAM2Act+, 479 demonstrates a strong understanding of spatial memory, 480 achieving an average success rate of 94.3% across all tasks. 481 It outperforms SAM2Act (without memory) by a huge 482 483 margin of 39.3% on MemoryBench, highlighting the sig-484 nificant impact of explicit memory modeling.

Table 3. **Performance on MemoryBench.** We report the success rates for the three spatial memory tasks in MemoryBench. Our method, SAM2Act+, significantly outperforms all baseline methods that lack an explicit memory mechanism, achieving an average improvement of 37.6% across all three tasks.

Methods / Tasks	Avg. Success \uparrow	(a) Reopen Drawer	(b) Put Block Back	(c) Rearrange Block
RVT-2	54.0 ± 5.3	60.0 ± 0.0	50.0 ± 2.3	52.0 ± 3.3
SAM2Act (Ours)	55.0 ± 24.3	48.0 ± 0.0	35.0 ± 3.8	82.0 ± 2.3
SAM2Act+ (Ours)	$\textbf{94.3} \pm \textbf{9.0}$	$\textbf{84.0}\pm0.0$	$\textbf{100.0}\pm0.0$	$\textbf{99.0} \pm 2.0$

485 5.5. Real-robot Evaluations

Table 4 presents our real-world experiment results, where
our method achieves a 75% task success rate, compared to
43% for RVT-2. SAM2Act significantly outperforms the

baseline in high-precision tasks (60% vs 0%). It excels489in memory-based tasks, such as (d) Push the same490button, which requires recalling the button's previous lo-
cation. Here, SAM2Act achieves 70% success, while RVT-
2, relying on random guessing, scores 40%. We also test
models' generalization against perturbations like lighting
changes, distractors, and position variations.491

Table 4. **Real-world results.** We compare RVT2 against SAM2Act for the first three tasks and SAM2Act+ on the last real-world tasks (indicated with *), evaluating performance both indistribution and out-of-distribution during test time.

	In-Di	stribution	Out-Distribution		
Task	RVT-2	SAM2Act	RVT-2	SAM2Act	
(a) turn on the lamp	0/10	6/10	0/10	6/10	
(b) push button sequence	4/10	9/10	1/10	9/10	
(c) stack cubes	8/10	8/10	3/10	3/10	
(d) push the same button *	4/10	7/10	2/10	6/10	

6. Conclusion & Limitation

We introduce SAM2Act, a multi-view, language-497 conditioned behavior cloning policy for 6-DoF 3D 498 enabling high-precision manipulations manipulation. 499 while generalizing effectively to unseen perturbations. 500 Building on this foundation, we propose SAM2Act+, a 501 memory-based multi-view language-conditioned robotic 502 transformer-based policy that equips the agent with spatial 503 memory awareness, allowing it to solve spatial memory-504 based tasks. While both SAM2Act and SAM2Act+ achieve 505 SOTA performance across multiple benchmarks, chal-506 lenges remain in extending them to dexterous continuous 507 control. Additionally, SAM2Act+ relies on a fixed memory 508 window length, which differs from task to task, limiting 509 its adaptability to tasks of varying length. Despite these 510 challenges, we believe SAM2Act+ is an important step 511 towards memory-based generalist manipulation policies. 512

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

513 References

- 514 [1] Sean L Bowman, Nikolay Atanasov, Kostas Daniilidis, and
 515 George J Pappas. Probabilistic data association for semantic
 516 slam. In 2017 IEEE international conference on robotics and
 517 automation (ICRA), pages 1722–1729. IEEE, 2017. 2
- 518 [2] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
 519 Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr520 ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al.
 521 Rt-1: Robotics transformer for real-world control at scale.
 522 arXiv preprint arXiv:2212.06817, 2022. 2
- 523 [3] Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Ab524 hinav Gupta, and Russ R Salakhutdinov. Object goal navi525 gation using goal-oriented semantic exploration. *Advances*526 *in Neural Information Processing Systems*, 33:4247–4258,
 527 2020. 2
- 528 [4] Shizhe Chen, Ricardo Garcia, Cordelia Schmid, and Ivan
 529 Laptev. Polarnet: 3d point clouds for language-guided
 530 robotic manipulation. *arXiv preprint arXiv:2309.15596*,
 531 2023. 2
- [5] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun
 Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song.
 Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page
 02783649241273668, 2023. 2
- 537 [6] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Ab538 hinav Gupta. An unbiased look at datasets for visuo-motor
 539 pre-training. In *Conference on Robot Learning*, pages 1183–
 540 1198. PMLR, 2023. 2
- [7] Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang,
 Kiana Ehsani, Dieter Fox, and Ranjay Krishna. Manipulateanything: Automating real-world robots using visionlanguage models. *arXiv preprint arXiv:2406.18915*, 2024.
 2, 3
- 546 [8] Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Ka547 terina Fragkiadaki. Act3d: Infinite resolution action detec548 tion transformer for robotic manipulation. *arXiv preprint*549 *arXiv:2306.17817*, 2023. 2
- Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. In *Conference on Robot Learning*, pages 694– 710. PMLR, 2023. 1, 2, 7
- [10] Ankit Goyal, Valts Blukis, Jie Xu, Yijie Guo, Yu-Wei Chao, and Dieter Fox. Rvt-2: Learning precise manipulation from few demonstrations. *arXiv preprint arXiv:2406.08545*, 2024.
 2, 3, 7, 8
- [11] Peter Henry, Michael Krainin, Evan Herbst, Xiaofeng Ren, and Dieter Fox. Rgb-d mapping: Using kinect-style depth cameras for dense 3d modeling of indoor environments. *The international journal of Robotics Research*, 31(5):647–663, 2012. 2
- 563 [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen564 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
 565 Lora: Low-rank adaptation of large language models. *arXiv*566 *preprint arXiv:2106.09685*, 2021. 4
- 567 [13] Haoxu Huang, Fanqi Lin, Yingdong Hu, Shengjie Wang, and
 568 Yang Gao. Copa: General robotic manipulation through

spatial constraints of parts with foundation models. *arXiv* preprint arXiv:2403.08248, 2024. 3

- [14] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv* preprint arXiv:2307.05973, 2023. 3
- [15] Stephen James and Pieter Abbeel. Coarse-to-fine q-attention with learned path ranking. *arXiv preprint arXiv:2204.01571*, 2022. 2
- [16] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. 3, 7
- [17] Sascha Jockel, Martin Weser, Daniel Westhoff, and Jianwei Zhang. Towards an episodic memory for cognitive robots. In Proc. of 6th Cognitive Robotics workshop at 18th European Conf. on Artificial Intelligence (ECAI), pages 68–74. Citeseer, 2008. 2
- [18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. ACM Trans. Graph., 42(4):139–1, 2023. 3
- [19] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. *Advances in neural information processing systems*, 33:19884–19895, 2020. 2
- [20] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. In *International conference on machine learning*, pages 5639–5650. PMLR, 2020. 2
- [21] Mikko Lauri, David Hsu, and Joni Pajarinen. Partially observable markov decision processes in robotics: A survey. *IEEE Transactions on Robotics*, 39(1):21–40, 2022. 3
- [22] Peiqi Liu, Zhanqiu Guo, Mohit Warke, Soumith Chintala, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Dynamem: Online dynamic spatio-semantic memory for open world mobile manipulation. arXiv preprint arXiv:2411.04999, 2024. 3
- [23] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. arXiv preprint arXiv:2210.00030, 2022. 2
- [24] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? Advances in Neural Information Processing Systems, 36:655–677, 2023.
- [25] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. arXiv preprint arXiv:2203.12601, 2022. 2
- [26] Wilbert Pumacay, Ishika Singh, Jiafei Duan, Ranjay Krishna, Jesse Thomason, and Dieter Fox. The colosseum: A benchmark for evaluating generalization for robotic manipulation. *arXiv preprint arXiv:2402.08191*, 2024. 2
- [27] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
 626

Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2:
Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 3, 4

- [28] Rutav Shah and Vikash Kumar. Rrl: Resnet as representation for reinforcement learning. *arXiv preprint arXiv:2107.03380*, 2021. 2
- [29] Jinghuan Shang, Karl Schmeckpeper, Brandon B May,
 Maria Vittoria Minniti, Tarik Kelestemur, David Watkins,
 and Laura Herlant. Theia: Distilling diverse vision
 foundation models for robot learning. *arXiv preprint arXiv:2407.20179*, 2024. 2
- [30] Olaolu Shorinwa, Johnathan Tucker, Aliyah Smith, Aiden
 Swann, Timothy Chen, Roya Firoozi, Monroe David
 Kennedy, and Mac Schwager. Splat-mover: Multi-stage,
 open-vocabulary robotic manipulation via editable gaussian
 splatting. In *8th Annual Conference on Robot Learning*,
 2024. 3
- 644 [31] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport:
 645 What and where pathways for robotic manipulation. In *Con-*646 *ference on robot learning*, pages 894–906. PMLR, 2022. 1,
 647 2
- 648 [32] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver649 actor: A multi-task transformer for robotic manipulation.
 650 In *Conference on Robot Learning*, pages 785–799. PMLR,
 651 2023. 1, 2, 7
- [33] Linda Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies. *Artificial life*, 11
 (1-2):13–29, 2005. 1
- [34] Che Wang, Xufang Luo, Keith Ross, and Dongsheng Li.
 Vrl3: A data-driven framework for visual deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:32974–32988, 2022. 2
- [35] Dian Wang, Robin Walters, Xupeng Zhu, and Robert Platt.
 Equivariant *q* learning in spatial action spaces. In *Conference* on *Robot Learning*, pages 1713–1723. PMLR, 2022. 2
- [36] Zhou Xian, Nikolaos Gkanatsios, Theophile Gervet, TsungWei Ke, and Katerina Fragkiadaki. Chaineddiffuser: Unifying trajectory diffusion and keypose prediction for robotic
 manipulation. In *7th Annual Conference on Robot Learning*,
 2023. 2
- [37] Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International conference on learning representations*, 2021. 2
- [38] Wentao Yuan, Adithyavairavan Murali, Arsalan Mousavian, and Dieter Fox. M2t2: Multi-task masked transformer for object-centric pick and place. *arXiv preprint arXiv:2311.00926*, 2023. 2
- [39] Andy Zeng, Pete Florence, Jonathan Tompson, Stefan
 Welker, Jonathan Chien, Maria Attarian, Travis Armstrong,
 Ivan Krasin, Dan Duong, Vikas Sindhwani, et al. Transporter
 networks: Rearranging the visual world for robotic manipulation. In *Conference on Robot Learning*, pages 726–747.
 PMLR, 2021. 1, 2
- [40] Junjie Zhang, Chenjia Bai, Haoran He, Wenke Xia, Zhigang
 Wang, Bin Zhao, Xiu Li, and Xuelong Li. Sam-e: Leveraging visual foundation model with sequence imitation for

 embodied manipulation.
 arXiv preprint arXiv:2405.19586,
 684

 2024.
 2,7
 685

[41] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
688 2