

---

# SAM2ACT: INTEGRATING VISUAL FOUNDATION MODEL WITH A MEMORY ARCHITECTURE FOR ROBOTIC MANIPULATION

Haoquan Fang<sup>1</sup> Markus Grotz<sup>1</sup> Wilbert Pumacay<sup>2</sup>

Yi Ru Wang<sup>1</sup> Dieter Fox<sup>\*1,3</sup> Ranjay Krishna<sup>\*1,4</sup> Jiafei Duan<sup>\*1,4</sup>

<sup>1</sup>University of Washington, <sup>2</sup>Universidad Católica San Pablo, <sup>3</sup>NVIDIA

<sup>4</sup>Allen Institute for Artificial Intelligence

[sam2act.github.io](https://sam2act.github.io)

## ABSTRACT

Robotic manipulation systems operating in diverse, dynamic environments must exhibit three critical abilities: generalization to unseen scenarios, multitask interaction, and spatial memory. While significant progress has been made in robotic manipulation, existing approaches often fall short in addressing memory-dependent tasks and generalization to complex environmental variations. To bridge this gap, we introduce **SAM2Act**, a multi-view robotic transformer that leverages multi-resolution upsampling and visual representations from large-scale foundation models. SAM2Act achieves a state-of-the-art average success rate of **86.8% across 18 tasks** in the RLBench benchmark, and demonstrates robust generalization on The Colosseum benchmark, with only a **4.3% performance gap** under diverse environmental perturbations. Building on this foundation, we propose **SAM2Act+**, a memory-augmented architecture inspired by SAM2, which incorporates a memory bank and attention mechanism for spatial memory. To address the need for evaluating memory-dependent tasks, we introduce MemoryBench, a novel benchmark designed to assess spatial memory and action recall in robotic manipulation. SAM2Act+ achieves **competitive performance on MemoryBench**, significantly outperforming existing approaches and pushing the boundaries of memory-enabled robotic systems.

## 1 INTRODUCTION

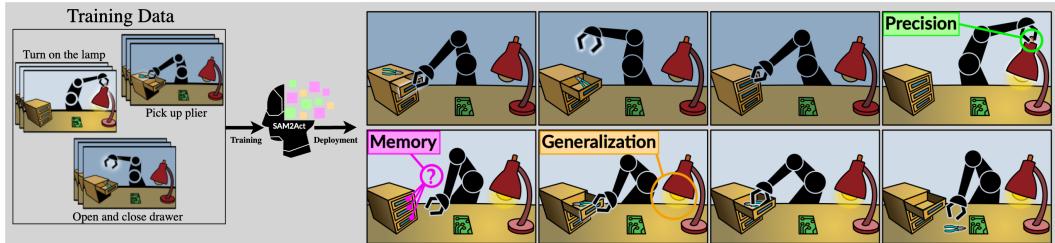


Figure 1: **Simulation and Real Tasks.** We demonstrate the effectiveness of SAM2Act+ in solving memory-based tasks by evaluating it against baselines on the three benchmark memory tasks (shown at the top). Additionally, we validate our approach using a Franka Panda robot on four real-world tasks (shown at the bottom), including tests under out-of-distribution perturbations.

The world in which we live is diverse and constantly changing, encompassing a wide variety of objects, scenes, and environmental conditions. Consider the seemingly simple task of following a recipe when cooking: we can recognize salt even if it comes in different types of container, seamlessly

---

\*Equal advising

---

perform the action of picking it up and sprinkling it into the pan, and remember whether we have already added salt. Humans excel in such environments due to their ability to generalize in unseen scenarios, interact with their surroundings to achieve specific goals, and retain knowledge from past experiences Smith & Gasser (2005); Duan et al. (2022). These abilities, generalization, multitask interaction, and memory, serve as guiding principles for the development of robotic systems capable of operating in similarly complex environments.

Significant progress has been made in robotic manipulation through prior work. Early methods, such as the Transporter Network Zeng et al. (2021) and CLIPort Shridhar et al. (2022), demonstrated effective 2D action-centric manipulation but were limited in their ability to handle spatially complex tasks. More recent approaches, such as PerAct Shridhar et al. (2023) and RVT Goyal et al. (2023), have pushed toward 3D-based manipulation. PerAct employs a multitask transformer that interprets language commands and predicts keyframe poses, achieving strong results across a variety of tasks. RVT builds on this foundation by adopting a 2.5D representation, improving training efficiency and inference speed. Its successor, RVT-2, further enhances performance with a coarse-to-fine strategy, increasing precision for high-accuracy tasks. Despite these advances, important challenges remain, including improving multitask performance, enhancing generalization to novel environment configurations, and integrating memory mechanisms for tasks requiring episodic recall.

We introduce SAM2Act, a multi-view robotics transformer that enhances feature representation by integrating multi-resolution upsampling with visual embeddings from large-scale foundation models. Built on the RVT-2 multiview transformer, SAM2Act achieves strong multitask success and generalization. Building on this foundation, we introduce SAM2Act+, which incorporates a memory-based architecture inspired by the SAM2 approach. Using a memory bank and an attention mechanism, SAM2Act+ enables episodic recall to solve more complex, memory-dependent manipulation tasks. We evaluate SAM2Act and SAM2Act+ using MemoryBench, a new benchmark suite that tests the spatial memory of policies and the ability to retain and recall actions. SAM2Act+ achieves competitive performance on MemoryBench, with an average accuracy of 91.3%, outperforming next highest baseline by a huge margin of 37.6%. Furthermore, we assess the generalization capabilities of SAM2Act on The Colosseum Pumacay et al. (2024), a benchmark designed to test robotic manipulation under various environmental perturbations. SAM2Act demonstrates robust performance on The Colosseum with an average decrease of 4.3% across all perturbations, highlighting its ability to generalize effectively in diverse and challenging scenarios. Lastly, our approach outperformed baseline methods in real-world evaluations while exhibiting comparable generalization and spatial memory capabilities. In summary, this work makes three key contributions. First, we introduce a **novel model formulation** that leverages visual foundation models to solve **high-precision, memory-dependent manipulation tasks**. Second, we propose MemoryBench, a dedicated evaluation benchmark for assessing **spatial memory in behavior cloning models**. Finally, we present **empirical results and insights** on the model’s performance across both simulation and real-world tasks.

## 2 RELATED WORK

### 2.1 3D-BASED ROBOTIC TRANSFORMER FOR MANIPULATION

2D-based methods Zhao et al. (2023); Chi et al. (2023); Zeng et al. (2021); Brohan et al. (2022); Shridhar et al. (2022) are effective for simple pick-and-place tasks due to fast training, low hardware requirements, and minimal computational cost. However, they depend on pretrained image encoders and fail in tasks requiring high precision, robust spatial interaction, or resilience to environmental and camera variations Pumacay et al. (2024). Recent work addresses these limitations with 3D perception. Methods like PolarNet Chen et al. (2023), M2T2 Yuan et al. (2023), and Manipulate-Anything Duan et al. (2024b) reconstruct point clouds, while C2F-ARM James & Abbeel (2022) and PerAct Shridhar et al. (2023) use voxel-based 3D representations. Act3D Gervet et al. (2023) and ChainedDiffuser Xian et al. (2023) adopt multi-scale 3D features. RVT Goyal et al. (2023) introduces 2.5D multi-view images for faster training, refined by RVT2 Goyal et al. (2024) with a coarse-to-fine architecture for improved precision. Our work, SAM2Act, combines RVT2’s spatial reasoning with enhanced virtual images from the SAM2 visual encoder, achieving high precision and generalization across diverse tasks.

---

## 2.2 VISUAL REPRESENTATIONS FOR ROBOT LEARNING

Robotics research heavily relies on visual representations from computer vision to process high-dimensional inputs and improve policy learning. Visual representations are integrated into robot learning through pretraining Majumdar et al. (2023); Ma et al. (2022); Nair et al. (2022), co-training Laskin et al. (2020b); Yarats et al. (2021); Laskin et al. (2020a); Shang et al. (2024); Duan et al. (2024a), or frozen encoders Shah & Kumar (2021); Wang et al. (2022a); Zhang et al. (2024), all of which effectively support policy training. These representations also enhance invariance, equivariance, and out-of-distribution generalization Wang et al. (2022b); Pumacay et al. (2024); Dasari et al. (2023). SAM-E Zhang et al. (2024) demonstrates the use of a pre-trained SAM encoder for robotic manipulation by leveraging image embeddings for policy learning. Expanding on this, our approach employs the SAM2 visual encoder to generate embeddings for robotic transformers and utilizes its multi-resolution features to improve convex upsampling for next-action prediction.

## 2.3 MEMORY IN ROBOTICS

Memory is a fundamental component of human cognition, and equipping generalist robotic agents with episodic and semantic memory is crucial for enabling them to perform complex tasks effectively Jockel et al. (2008). Early research on memory in robotics primarily addressed navigation tasks, relying on semantic maps that were often constrained in scope Henry et al. (2012); Bowman et al. (2017); Chaplot et al. (2020). Other works explicitly model the memory for a robot cognitive architecture Peller-Konrad et al. (2023). Recent advancements leverage representations derived from vision-language models (VLMs) and Large Vision Models (LVMs), utilizing voxel maps or neural feature fields to encode, store, and retrieve information Huang et al. (2024; 2023); Duan et al. (2024b); Liu et al. (2024). Alternative methods represent semantic memory for manipulation tasks using Gaussian splats to encode spatial information Kerbl et al. (2023); Shorinwa et al. (2024). In contrast, our approach draws inspiration from the framework of Partially Observable Markov Decision Processes (POMDPs) Lauri et al. (2022), incorporating memory directly into the training process. By integrating spatial memory from past actions into the agent’s belief state, we enhance the robustness and adaptability of learned policies.

# 3 MEMORYBENCH: A MEMORY BENCHMARK FOR 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.

## 3.1 TASK DESIGN

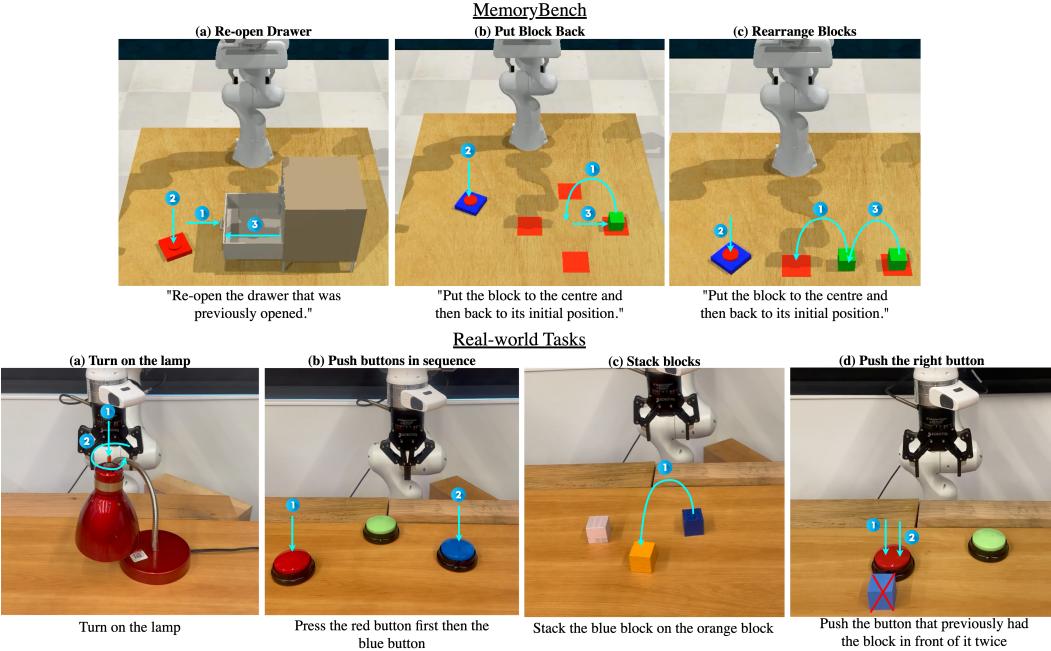
Unlike standard RLBench tasks James et al. (2020), many of which involve long-horizon scenarios, our tasks are specifically designed to require spatial memory. Without such memory, the agent would be forced to rely on random actions. To create these tasks, we intentionally violate the Markov assumption, which states that in a Markov Decision Process (MDP), the next observation depends solely on the current observation and action:

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

This assumption implies that knowing only  $o_t$  and  $a_t$  is sufficient to predict  $o_{t+1}$ . However, in our tasks, we design scenarios where two distinct action histories lead to the same observation  $o_t$ , but require different subsequent actions. This forces the agent to recall which action history led to  $o_t$  to perform the correct next action. These principles guided the development of our spatial memory-based tasks.

## 3.2 SPATIAL MEMORY-BASED TASKS

`MemoryBench` extends the `RLBench` simulator to provide scripted demonstrations for three spatial memory tasks: `reopen_drawer`, `put_block_back`, and `rearrange_block`. Each task is designed to evaluate a specific aspect of spatial memory and adheres to the principles outlined in Section 3.1. To introduce complexity, these tasks include two to four variations and additional



**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.

steps—such as pressing a button mid-sequence—that disrupt the Markov property. This forces the agent to rely on memory rather than solely on immediate observations.

The `reopen_drawer` task evaluates the agent’s ability to recall 3D spatial information along the z-axis. Initially, one of three drawers (top, middle, or bottom) is open. The agent must close the open drawer, press a button on the table, and then reopen the same drawer. After the button is pressed, all drawers are closed, and the scene becomes visually indistinguishable, requiring the agent to use memory to identify the correct drawer. This task tests the agent’s ability to recall spatial states over a temporal sequence. The `put_block_back` task tests the agent’s ability to remember 2D spatial information on the x-y plane. Four red patches are placed on a table, with a block initially positioned on one of them. The agent must move the block to the center of the patches, press a button, and return the block to its original position. The agent must rely on its memory of the block’s initial location to succeed, demonstrating its capability to encode and retrieve 2D spatial information.

The `rearrange_block` task assesses the agent’s ability to perform backward reasoning by recalling and reversing prior actions. A block starts on one of two red patches, with the other patch empty. The agent must move the block to the empty patch, press a button, and then return the block to its original patch. This requires the agent to reconstruct past states and actions, testing its capacity for backward spatial memory reasoning. These tasks collectively evaluate both forward and backward spatial reasoning across 3D (z-axis) and 2D (x-y plane) spaces. By introducing non-Markovian elements, they emphasize the need for memory representations to solve complex sequential decision-making problems.

## 4 METHOD

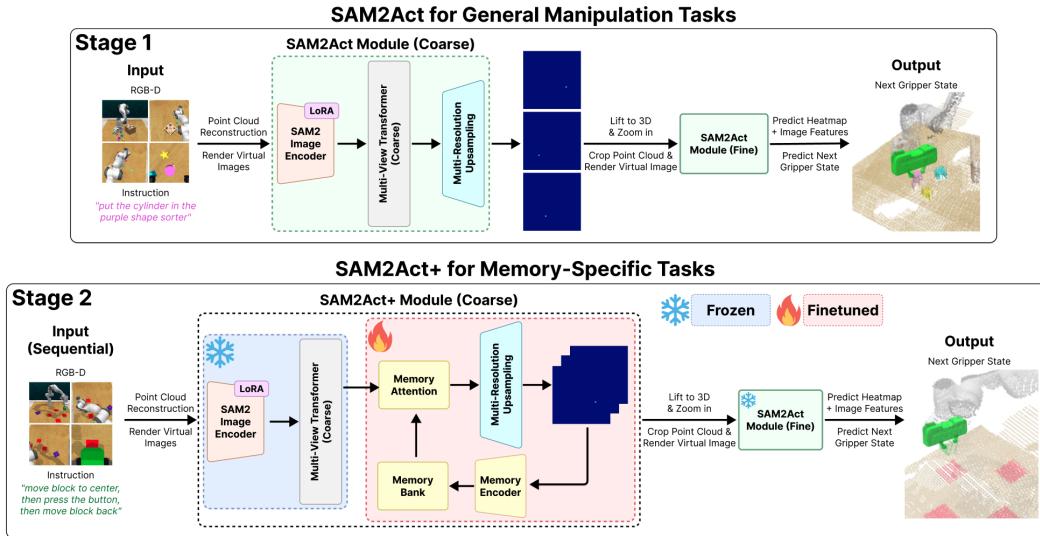
Our method, SAM2Act, enables precise 3D manipulation with strong generalization across environmental and object-level variations. Building upon the RVT-2 framework Goyal et al. (2024), SAM2Act introduces key architectural innovations that enhance visual feature representation and task-specific reasoning. The architecture reconstructs a point cloud of the scene, renders it from

virtual cameras at orthogonal views, and employs a two-stage multi-view transformer (coarse-to-fine) to predict action heatmaps.

The coarse branch generates zoom-in heatmaps to localize regions of interest, while the fine branch refines these into precise action heatmaps. SAM2Act leverages the pretrained SAM2 encoder Ravi et al. (2024) to extract multi-resolution image embeddings, which are further refined through advanced upsampling techniques to predict accurate translation heatmaps with minimal information loss. To address tasks requiring spatial memory, SAM2Act+ extends the SAM2Act architecture by incorporating memory-based reasoning components. These include a Memory Encoder, Memory Attention, and Memory Bank, enabling the model to process historical heatmaps and integrate prior observations. This memory-driven reasoning enhances the agent’s ability to predict actions based on past contextual information, significantly improving performance in tasks that require sequential decision-making.

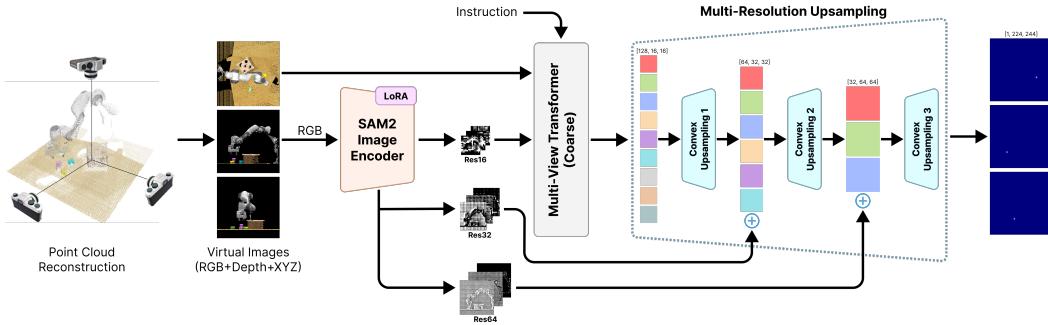
In the following sections, we detail the SAM2Act architecture (subsection 4.1), including its multi-resolution upsampling mechanism (Figure 5). We also present the SAM2Act+ extension, which integrates memory-based components for solving spatial memory tasks (subsection 4.2).

#### 4.1 SAM2ACT: MULTI-RESOLUTION UPSAMPLING FOR ENHANCED VISUAL FEATURE REPRESENTATION



**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.

A distinctive feature of SAM2Act is the incorporation of the SAM2Act Module into the manipulation backbone for training, as illustrated in Figure 5. The coarse and fine SAM2Act Modules share the same architecture, with the fine branch generating additional features to predict actions beyond translation, while the coarse branch focuses exclusively on translation. Point-cloud representations are reconstructed from raw image inputs, and virtual images are generated from three viewpoints using virtual cameras. Instead of directly inputting these images into the multi-view transformer,



**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.

their RGB channels are extracted out and processed by the SAM2 Ravi et al. (2024) image encoder, which produces object-centric multi-resolution embeddings. These embeddings, generated at three resolution levels, are combined with virtual images containing RGB, depth, 3D translation coordinates, and language instructions before being fed into the multi-view transformer.

To adapt the SAM2 image encoder to our domain efficiently, we fine-tune it using Low-Rank Adaptation (LoRA) Hu et al. (2021), which enables domain adaptation with minimal computational cost while maintaining model efficiency. Additionally, to fully leverage the multi-resolution embeddings produced by the SAM2 image encoder, we introduce a multi-resolution upsampling method. This method uses the embeddings as auxiliary inputs to enhance the generation of translation heatmaps, thereby improving spatial precision and overall system performance. The multi-resolution upsampling mechanism, also detailed in Figure 5, leverages cascaded convex upsamplers 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 multi-resolution embedding from SAM2. Also let  $U(\cdot)$  denote the upsampling operator that doubles the spatial dimensions. The feature maps are updated at each stage as follows:

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

where  $\oplus$  represents elementwise addition. The upsampling operator  $U$  is defined as:

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

At each stage, the output of the upsample is combined with the corresponding multi-resolution embedding  $E^l$  from the SAM2 encoder, ensuring alignment between the multi-resolution features and the decoder’s spatial refinement process. A layer normalization step follows each addition to stabilize training and maintain feature coherence. This results in direct integration of the embeddings into the translation heatmap generation process. The cascading structure refines features across multiple resolutions, capturing fine-grained spatial details critical for manipulation tasks.

#### 4.2 SAM2Act+: ACTION MEMORY ARCHITECTURE FOR IMPROVED SPATIAL AWARENESS IN PAST OBSERVATIONS

To extend the SAM2Act architecture (subsection 4.1) with memory-based capabilities inspired by SAM2, we introduce SAM2Act+, a task-specific variant designed for solving memory-based tasks. SAM2Act+ integrates the three core memory components from SAM2—*Memory Attention, Memory Encoder, and Memory Bank*—into the coarse branch of SAM2Act. Originally developed for object tracking in SAM2, these components are adapted to align with the needs of SAM2Act+, enabling the agent to retain prior actions and observations for sequential decision-making. In SAM2, the Memory Encoder processes predicted object masks, while the Memory Attention module fuses image

---

**Algorithm 1** Forward Pass of SAM2Act+ Module

---

```
1: Initialize: Number of steps  $N$ , maximum number of memories  $M$ , number of views  $V$ , empty  
memory bank  $Q$  with  $V$  separate FIFO queues, input  $X$   
2: for  $i = 1$  to  $N$  do  
3:   for  $j = 1$  to  $V$  do  
4:     Get embeddings  $\mathcal{E}_{raw}$  from MVT  $T_{mv}(X_j)$   
5:     Retrieve past memories  $\mathcal{M}_{old}$  from  $Q[j]$   
6:     Get memory-conditioned embeddings  $\mathcal{E}_{mem}$  from Memory Attention  $T_{mem}(\mathcal{E}_{raw}, \mathcal{M}_{old})$   
7:     Predict translation heatmap  $\mathcal{H}$  with upsample U( $\mathcal{E}_{mem}$ )  
8:     Encode new memory  $\mathcal{M}_{new}$  using Memory Encoder  $E_{mem}(\mathcal{H}, \mathcal{E}_{raw})$   
9:     Store new memory  $Q[j] \leftarrow Q[j] \cup \{\mathcal{M}_{new}\}$   
10:    if  $|Q[j]| = M$  then  
11:       $Q[j] \leftarrow Q[j]_{2:n}$   
12:    end if  
13:  end for  
14: end for
```

---

embeddings with positional information from previous frames. SAM2Act+ adopts a similar structure: the predicted heatmaps, which serve as binary indicators of spatial positions in the image, function analogously to object masks. This conceptual alignment ensures a seamless integration of memory mechanisms, allowing the agent to leverage stored information to predict subsequent actions based on historical context. A detailed description of the Memory Attention and Memory Encoder modules can be found in the supplementary material.

**Architecture.** The SAM2Act+ architecture is illustrated in Figure 3. After pretraining SAM2Act in Stage 1, we freeze the SAM2 image encoder and the multi-view transformer in the coarse branch, as these components effectively generate robust embeddings for multi-view images in manipulation tasks. We also freeze the entire fine branch, given its proven ability to predict fine-grained actions accurately. Fine-tuning is applied only to the coarse branch, as it focuses on generating heatmaps that provide richer contextual information for recalling past actions. The fine branch, in contrast, primarily emphasizes small objects or localized regions, which typically contain less information relevant to memory-based tasks.

**Training.** To train SAM2Act+, we fine-tune the coarse branch by integrating the three memory components with the multi-resolution upsampling module. During fine-tuning, consecutive action keyframes are sampled as input, training the multi-resolution upsample to predict new translations conditioned on memory. The memory components function similarly to their implementation in SAM2 for object tracking, with one key distinction: the input to the Memory Encoder. Instead of using unconditioned image embeddings from the SAM2 image encoder, we input feature embeddings generated by the multi-view transformer. This adaptation ensures that memory encoding incorporates multi-view information while maintaining independence in handling stored representations. Virtual images are treated independently during memory encoding and attention, with each view’s memory encoded separately. Feature embeddings from each view are attended to using their corresponding stored memories, preserving spatial and contextual alignment while leveraging fused multi-view information. This structured approach prevents cross-view interference and enhances the model’s ability to reason over sequential tasks. The memory-augmented forward pass for SAM2Act+ is outlined in 1. By incorporating memory mechanisms, SAM2Act+ enhances performance in scenarios requiring long-term reasoning, enabling the agent to make informed decisions based on historical context.

## 5 EXPERIMENTS

We study SAM2Act and SAM2Act+ in both simulated and real-world environments. Specifically, we are interested in answering the following questions:

- § 5.2 How does SAM2Act compare with state-of-the-art 3D manipulation policies?
- § 5.3 Can SAM2Act+ solve spatial memory-based tasks that other baselines cannot?
- § 5.4 Can SAM2Act generalize across object and environmental perturbations?



also compare with RVT Goyal et al. (2023), PerAct Shridhar et al. (2023), and SAM-E Zhang et al. (2024).

## 5.2 PERFORMANCES ACROSS 18 RLBNCH TASKS

Table 1 compares SAM2Act with prior keyframe-based 3D BC methods on the RLBNch benchmark. Overall, SAM2Act achieves an average success rate of  **$86.8\% \pm 0.5$** , surpassing the previous best (RVT-2) by **5.4%**. A closer look at individual tasks reveals that SAM2Act ranks **first in 9 out of 18 tasks** and remains highly competitive in **7 others**, coming within **one successful attempt or 4%** of the best performance. These tasks include Close Jar, Drag Stick, Meat Off Grill, Place Wine, Screw Bulb, Sweep to Dustpan, and Turn Tap. The largest margin of improvement occurs in Insert Peg, where SAM2Act **exceeds RVT-2 by 44% (approximately 2.1x)**, and in Sort Shape, where it outperforms RVT-2 by 29%. Both tasks require precise manipulation, underscoring the effectiveness of SAM2Act’s multi-resolution upsampling strategy. These results establish SAM2Act as a **leading policy for complex 3D tasks**, highlighting its ability to handle high-precision manipulations - an area where prior methods have struggled.

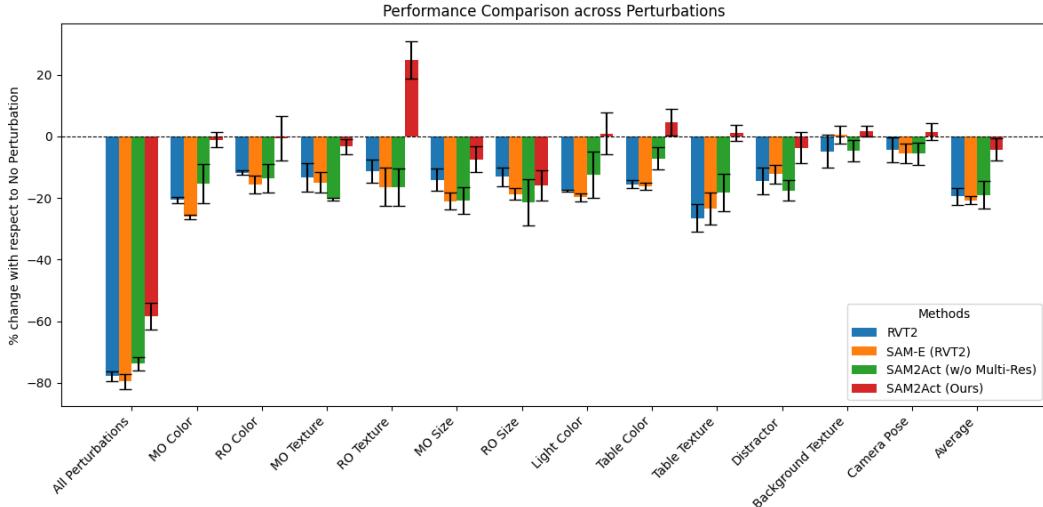


Figure 5: **The Colosseum results.** Task-average success rate percentage change for SAM2Act and other baselines across 13 perturbation factors from **The Colosseum**, relative to evaluations without perturbations. Our approach, SAM2Act, demonstrates the lowest average percentage change across all perturbations, with a minimal drop of  $-4.3 \pm 3.6\%$ , highlighting its robustness in handling various environmental and object-level perturbations.

## 5.3 PERFORMANCE ON MEMORYBENCH

In Table 2, we evaluate the state-of-the-art keyframe-based 3D BC model, RVT-2, alongside our method on MemoryBench, training all models in a single-task setting to isolate memory-related challenges (e.g., opening the wrong drawer rather than unrelated mid-task failures). This setup ensures that performance differences stem from memory capabilities. For a random agent, the expected success rates are determined by the number of possible choices per task: 33% for Reopen Drawer (three drawers), 25% for Put Block Back (four patches), and 25% for Rearrange Block (two blocks). However, variations in task complexity, fixed training data, and imbalanced task distributions lead to slight deviations from these baselines. Our proposed memory-based model, SAM2Act+, demonstrates a **strong understanding of spatial memory**, achieving an average success rate of 91.3% across all tasks. It **outperforms SAM2Act (without memory) by a huge margin of 37.6% on MemoryBench**, highlighting the impact of explicit memory modeling.

---

## 5.4 SEMANTIC GENERALIZATION ACROSS TASKS

The results evaluated in Section 5.2 were obtained by training and testing models within the same environment. However, to truly assess **generalization performance**, policies must remain robust against both environmental and object-level perturbations. We therefore trained SAM2Act and the baseline methods on 20 tasks from The Colosseum benchmark and tested them under 13 different perturbation categories over three runs. **SAM2Act exhibits the smallest relative performance drop compared to the baselines**, with an average decrease of just 4.3% (standard deviation of 3.59%). Notably, it proves particularly robust to environmental perturbations – such as changes in lighting, table color/textured, the addition of distractors, and even camera pose – while also maintaining competitive performance under object-level perturbations.

## 5.5 REAL-ROBOT EVALUATIONS

Table 3 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%) and consistently matches or surpasses RVT-2 in tabletop manipulation. It excels in memory-based tasks, such as (d) Push same button, which requires recalling the button’s previous location. Here, SAM2Act achieves 70% success, while RVT-2, relying on random guessing, scores 50%. We also test models’ generalization against perturbations like lighting changes, distractors, and position variations. Additional details are in the Appendix, with real-world rollout videos available on our project website.

**Table 3: Real-world results.** Comparison of RVT2 against SAM2Act for the first three tasks and SAM2Act+ for the last task (\*), evaluating in-distribution and out-of-distribution performance during test time.

Task	In-Distribution		Out-Distribution	
	RVT2	SAM2Act	RVT2	SAM2Act
(a) Turn on lamp	0/10	<b>6/10</b>	0/10	<b>6/10</b>
(b) Push button sequence	4/10	<b>9/10</b>	1/10	<b>9/10</b>
(c) Stack cubes	8/10	8/10	3/10	3/10
(d) Push same button *	5/10	<b>7/10</b>	2/10	<b>6/10</b>

## 6 CONCLUSION & LIMITATION

We introduce SAM2Act, a multi-view, language-conditioned behavior cloning policy for 6-DoF 3D manipulation, enabling high-precision manipulations while generalizing effectively to unseen perturbations. Building on this foundation, we propose SAM2Act+, a memory-based multi-view language-conditioned robotic transformer that equips the agent with spatial memory awareness, allowing it to solve spatial memory-based tasks with reduced uncertainty. While both SAM2Act and SAM2Act+ achieve state-of-the-art performance across multiple benchmarks, challenges remain in extending them to dexterous continuous control. Additionally, SAM2Act+ relies on a fixed memory window length based on task keyframes, limiting its adaptability to tasks of varying length. Despite these challenges, we believe SAM2Act+ marks an important step towards a memory-based generalist manipulation policy.

## REFERENCES

Sean L Bowman, Nikolay Atanasov, Kostas Daniilidis, and George J Pappas. Probabilistic data association for semantic slam. In *2017 IEEE international conference on robotics and automation (ICRA)*, pp. 1722–1729. IEEE, 2017.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

- 
- Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. *Advances in Neural Information Processing Systems*, 33:4247–4258, 2020.
- Shizhe Chen, Ricardo Garcia, Cordelia Schmid, and Ivan Laptev. Polarnet: 3d point clouds for language-guided robotic manipulation. *arXiv preprint arXiv:2309.15596*, 2023.
- 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*, pp. 02783649241273668, 2023.
- Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at datasets for visuo-motor pre-training. In *Conference on Robot Learning*, pp. 1183–1198. PMLR, 2023.
- Jiafei Duan, Samson Yu, Nicholas Tan, Li Yi, and Cheston Tan. Boss: A benchmark for human belief prediction in object-context scenarios. *arXiv preprint arXiv:2206.10665*, 2022.
- Jiafei Duan, Wilbert Pumacay, Nishanth Kumar, Yi Ru Wang, Shulin Tian, Wentao Yuan, Ranjay Krishna, Dieter Fox, Ajay Mandlekar, and Yijie Guo. Aha: A vision-language-model for detecting and reasoning over failures in robotic manipulation. *arXiv preprint arXiv:2410.00371*, 2024a.
- Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang, Kiana Ehsani, Dieter Fox, and Ranjay Krishna. Manipulate-anything: Automating real-world robots using vision-language models. *arXiv preprint arXiv:2406.18915*, 2024b.
- Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: Infinite resolution action detection transformer for robotic manipulation. *arXiv preprint arXiv:2306.17817*, 2023.
- 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*, pp. 694–710. PMLR, 2023.
- 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.
- Markus Grotz, Mohit Shridhar, Yu-Wei Chao, Tamim Asfour, and Dieter Fox. Peract2: Benchmarking and learning for robotic bimanual manipulation tasks. In *CoRL 2024 Workshop on Whole-body Control and Bimanual Manipulation: Applications in Humanoids and Beyond*, 2024. URL <https://openreview.net/forum?id=nIU0ZFmptX>.
- 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.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Haoxu Huang, Fanqi Lin, Yingdong Hu, Shengjie Wang, and Yang Gao. Copa: General robotic manipulation through spatial constraints of parts with foundation models. *arXiv preprint arXiv:2403.08248*, 2024.
- 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.
- Stephen James and Pieter Abbeel. Coarse-to-fine q-attention with learned path ranking. *arXiv preprint arXiv:2204.01571*, 2022.
- 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.

- 
- 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)*, pp. 68–74. Citeseer, 2008.
- 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.
- Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. In *International conference on machine learning*, pp. 5639–5650. PMLR, 2020a.
- 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, 2020b.
- 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.
- 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.
- 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.
- 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.
- 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.
- Fabian Peller-Konrad, Rainer Kartmann, Christian RG Dreher, Andre Meixner, Fabian Reister, Markus Grotz, and Tamim Asfour. A memory system of a robot cognitive architecture and its implementation in armarmx. *Robotics and Autonomous Systems*, 164:104415, 2023.
- 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.
- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- Rutav Shah and Vikash Kumar. Rrl: Resnet as representation for reinforcement learning. *arXiv preprint arXiv:2107.03380*, 2021.
- 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.
- Olaolu Shorinwa, Johnathan Tucker, Aliyah Smith, Aiden Swann, Timothy Chen, Roya Firooz, 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.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on robot learning*, pp. 894–906. PMLR, 2022.
- Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, pp. 785–799. PMLR, 2023.

- 
- Linda Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies. *Artificial life*, 11(1-2):13–29, 2005.
- 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, 2022a.
- Dian Wang, Robin Walters, Xupeng Zhu, and Robert Platt. Equivariant  $q$  learning in spatial action spaces. In *Conference on Robot Learning*, pp. 1713–1723. PMLR, 2022b.
- Zhou Xian, Nikolaos Gkanatsios, Theophile Gervet, Tsung-Wei Ke, and Katerina Fragkiadaki. Chaineddiffuser: Unifying trajectory diffusion and keypose prediction for robotic manipulation. In *7th Annual Conference on Robot Learning*, 2023.
- 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.
- 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.
- 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*, pp. 726–747. PMLR, 2021.
- 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*, 2024.
- 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.