

RoboSSM: Scalable In-context Imitation Learning via State-Space Models

Youngju Yoo^{1,2} Jiaheng Hu¹ Yifeng Zhu¹ Bo Liu^{1,3}

Qiang Liu¹ Roberto Martín-Martín^{1,4} Peter Stone^{1,5}

¹The University of Texas at Austin ²KAIST ³FAIR at Meta ⁴Amazon ⁵Sony AI
youngju.yoo@kaist.ac.kr

{jiahengh,yifeng.zhu,bliu,lqiang,robertomm,pstone}@cs.utexas.edu

Abstract: In-context imitation learning (ICIL) enables robots to learn tasks from prompts consisting of just a handful of demonstrations. By eliminating the need for parameter updates at deployment time, this paradigm supports few-shot adaptation to novel tasks. However, recent ICIL methods rely on Transformers, which have computational limitations and tend to underperform when handling longer prompts than those seen during training. In this work, we introduce RoboSSM, a scalable recipe for in-context imitation learning based on state-space models (SSM). Specifically, RoboSSM replaces Transformers with Longhorn – a state-of-the-art SSM that provides linear-time inference and strong extrapolation capabilities, making it well-suited for long-context prompts. We evaluate our approach on the LIBERO benchmark and compare it against strong Transformer-based ICIL baselines. Experiments show that RoboSSM extrapolates effectively to varying numbers of in-context demonstrations, yields high performance on unseen tasks, and remains robust in long-horizon scenarios. These results highlight the potential of SSMs as an efficient and scalable backbone for ICIL.

Keywords: Imitation Learning, In-context Learning, State-Space Model

1 Introduction

Imitation Learning (IL) is a powerful framework that enables robots to learn behaviors from demonstrations without explicit programming or reward design [1, 2]. While IL has achieved notable success in manipulation and navigation tasks, a key limitation of conventional imitation learning lies in its restricted adaptation capability, particularly when faced with new tasks. Even with models trained on large multi-task datasets [3, 4, 5, 6], adapting to novel tasks still requires collecting a large amount of task-specific data and retraining, which can be computationally costly and often unstable [7, 8]. To address this challenge, In-Context Imitation Learning (ICIL) introduces a new paradigm, inspired by the success of large language models (LLMs) [9, 10, 11] in adapting to unseen language tasks through few-shot learning [12]. ICIL integrates the concept of prompting into imitation learning [13, 14, 15, 16, 12, 17, 18, 19], allowing the model to infer and perform tasks based on a prompt composed of demonstrations, with no post-demonstration training.

Given that ICIL formulates imitation learning as a sequence modeling problem, recent ICIL approaches have naturally adopted Transformer-based models as their primary architecture [17, 19, 18]. Although Transformers are the dominant architecture for sequence modeling [20], their time complexity scales quadratically with sequence length, and they struggle to extrapolate beyond training lengths [21, 22]. For ICIL to handle long prompts efficiently at test time, it is essential to adopt alternatives to Transformers that enhance scalability with input length. In this paper, we introduce **RoboSSM**, a scalable in-context learning framework that replaces Transformers with state-space models (SSMs). Specifically, RoboSSM utilizes Longhorn [23], a state-of-the-art SSM with linear inference time and strong extrapolation capability for long-context sequences. Leveraging these

properties, RoboSSM can process substantially longer prompts at test time compared to previous Transformer-based ICIL methods.

On the LIBERO [24] benchmark, RoboSSM uniquely benefits from using more in-context examples, maintaining high success rates on unseen tasks when trained with only a few demonstrations. For instance, on the task *pick up the plate and place it in the tray*, where the *plate* object was unseen during training, RoboSSM achieves its highest performance when prompted with 32 demonstrations, despite being trained on only two. Furthermore, our framework performs well on unseen long-horizon tasks, which we simulate by repeating frames in the demonstrations to create time-dilated scenarios. Consequently, RoboSSM handles test-time demonstration prompts up to 16 times longer than those seen in training while maintaining linear inference time, whereas Transformer-based ICIL methods sharply degrade once the test prompt exceeds the training length. These findings confirm that RoboSSM establishes a scalable in-context imitation learning framework by effectively leveraging long-range contextual information.

2 Related Work

In this section, we provide an overview of prior ICIL methods and state-space models (SSMs), along with their recent applications to robotics.

2.1 In-Context Imitation Learning

Imitation learning [1, 2, 25, 26] has long been a foundation for imparting skills to robots by learning from demonstration data. Standard behavior cloning approaches [27, 28] typically train a separate policy for each task or rely on large multi-task datasets to acquire broader skills. While multi-task imitation learning [3, 4] can handle diverse tasks, these methods still struggle to perform completely unseen tasks without additional data collection for fine-tuning.

Inspired by the in-context learning paradigm in large language models (LLMs) [12], recent imitation learning methods aim to eliminate parameter updates at test time, instead prompting a multi-task policy with a few demonstrations of unseen tasks. Keypoint Action Tokens [19] introduce an ICIL framework that converts the visual observations and actions into tokens, which are fed into a pre-trained large language model. ICRT [17] performs in-context learning using a causal Transformer that predicts actions with next-token prediction, conditioned on a prompt consisting of a sequence of encoded teleoperated demonstrations. LipVQ-VAE [18] is an action tokenizer that uses vector quantization to address the lack of temporal smoothness in existing tokenizers and enable ICIL.

Although ICIL methods have been extensively studied and have achieved significant progress, they are typically trained and evaluated on test prompts that closely match the training prompt distribution in terms of length. For instance, ICRT learns from inputs containing five demonstrations and masks the first random k of them as the prompt, then at inference time is evaluated with three demonstration prompts. LipVQ-VAE is trained and evaluated using only a single full demonstration as the prompt. In contrast, RoboSSM explicitly aims to handle prompts that significantly deviate from the training distribution, such as those containing a larger number of demonstrations or long-horizon tasks.

2.2 State-Space Models

State-space models (SSMs) have emerged as a promising alternative to Transformers for sequence modeling in language tasks, addressing the quadratic time complexity of Transformers and their limitations in handling long contexts. SSMs originate from classical control theory and are particularly inspired by continuous-time linear dynamical systems. By discretizing the continuous-time formulation, SSMs can be expressed as discrete-time models that update the hidden state s_t via a linear recurrence:

$$s_t = As_{t-1} + Bx_t, \quad (1)$$

where x_t is the input and A and B are the state transition matrices. Recent SSMs aim to design the transition matrices A , B and the recurrence formulation. S4 [29], H3 [30], S5 [31], Mamba [32],

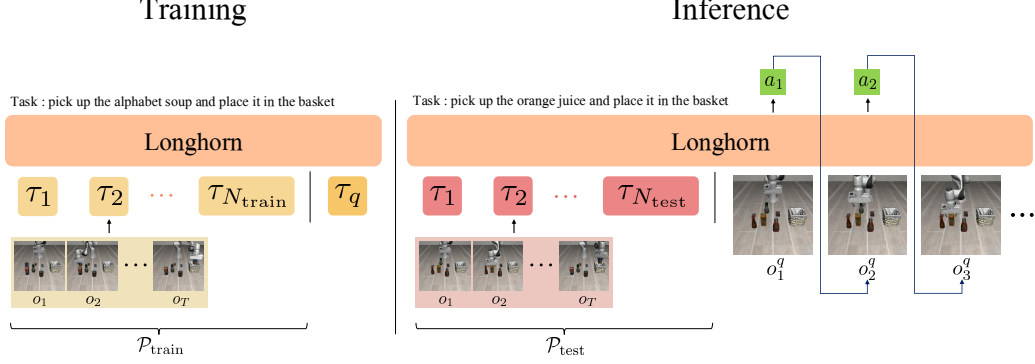


Figure 1: **Overview of RoboSSM.** **Training** (left): Longhorn receives N_{train} trajectories from $\mathcal{P}_{\text{train}}$ and a query trajectory for the same task. **Inference** (right): Given N_{test} trajectories from $\mathcal{P}_{\text{test}}$ containing unseen tasks, the model predicts actions and updates the environment iteratively from the initial observation embedding.

and Longhorn [23] have introduced structured state transition matrices and parallel computation schemes to enhance efficiency and capabilities. In particular, recent SSM architectures enable linear-time inference and demonstrate strong extrapolation capabilities over long-range contexts, while achieving comparable performance to Transformers in language modeling tasks.

As SSMs have evolved, their applications have expanded beyond language modeling to various other domains, including robotics. For instance, S5 is applied to reinforcement learning by allowing hidden state resets within a trajectory [33]. MAIL [34] proposes a novel imitation learning policy by leveraging Mamba. Building on these extensions, we explore using Longhorn, a recent state-of-the-art model, to perform in-context imitation learning.

3 Method

Our objective is to learn an ICIL policy π_{θ} that maximizes the success rate on unseen tasks when conditioned on few-shot demonstrations. In this section, we describe how RoboSSM implements π_{θ} with Longhorn [23] and how it processes trajectory prompts during training and inference.

3.1 Architecture

RoboSSM first processes the observations with multimodal encoders. The per-step encoded observation embeddings are then passed through the Longhorn state-space block to generate actions.

Input Encoding RoboSSM encodes multimodal observations at each time step of a demonstration. Visual data are processed by convolutional neural networks (CNNs), and proprioceptive data are embedded using multi-layer perceptrons (MLPs). The per-step features from each modality are concatenated and projected through an MLP to produce an observation embedding. To prevent the model from trivially copying the actions in the prompt, we exclude actions from the input representation. We further exclude any task language instructions to force the model to attend to the in-context demonstrations.

Longhorn state-space block The sequence of observation embeddings $\{x_t\}_{t=1}^T$ is fed into Longhorn, which recurrently updates a memory state matrix $s_t \in \mathbb{R}^{d \times m}$. At each time step, the input is interpreted as a key-value pair (k_t, x_t) , with $x_t \in \mathbb{R}^d$ and $k_t \in \mathbb{R}^m$ obtained via a linear projection of x_t , analogous to how Transformers use keys in the attention mechanism. Longhorn then performs a recurrent update:

$$s_t = A_t \odot s_{t-1} + B_t, \quad (2)$$

where \odot denotes the element-wise product, and $A_t, B_t : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times m}$ are functions of x_t , defined as

$$A_t = (\mathbf{1}_{d \times m} - \varepsilon_t \otimes k_t^{\odot 2}), \quad B_t = (\varepsilon_t \otimes \mathbf{1}_m) \otimes k_t, \quad \varepsilon_{t,i} = \frac{\beta_{t,i}}{1 + \beta_{t,i} k_t^\top k_t}, \quad (3)$$

where \otimes denotes the outer product and $\beta_t \in \mathbb{R}^d$ is a weighting vector. From the updated state, we compute a context vector:

$$r_t = s_t q_t \in \mathbb{R}^d, \quad (4)$$

where $q_t \in \mathbb{R}^m$ is a query vector derived from a linear projection of x_t . Finally, this context vector is passed through an output head to produce the corresponding action a_t .

Longhorn for In-Context Imitation Learning From an online-learning perspective, the recurrent form in (2) can be derived as the solution to the following online convex programming objective [35]:

$$s_t = \arg \min_{s \in \mathbb{R}^{d \times m}} \left\{ \|s - s_{t-1}\|_F^2 + \|s k_t - x_t\|_{\text{diag}(\beta_t)}^2 \right\}, \quad (5)$$

where $\|\cdot\|_F$ is the Frobenius norm and $\beta_t \in \mathbb{R}^d$ is a weighting vector. This objective balances two competing goals inherent in online learning: the first term encourages the updated state s_t to remain close to the previous state s_{t-1} , while the second term enforces that the current state s_t accurately reflects the new input, allowing the model to incorporate newly observed information. The weighting vector β_t modulates this trade-off by controlling the relative importance of the current observation embedding; it is obtained by applying a sigmoid activation to a linear projection of the input x_t . In this online regression view, the use of β_t in equation (5) naturally mitigates forgetting while integrating new information, thereby enabling efficient in-context learning with long contexts.

3.2 In-Context Imitation Learning

Following the standard ICIL formulation, RoboSSM conditions on a context of demonstration trajectories during both training and inference. At test time, the policy adapts to new tasks based solely on the provided demonstrations, without any parameter updates.

Each input to the policy consists of a prompt and a query trajectory. The prompt \mathcal{P} contains N trajectories that provide task context:

$$\mathcal{P} = [\tau_1, \tau_2, \dots, \tau_N], \quad (6)$$

where each demonstration trajectory τ_i is a sequence of T_i observation embeddings:

$$\tau_i = \left\{ o_1^{(i)}, o_2^{(i)}, \dots, o_{T_i}^{(i)} \right\}, \quad (7)$$

and $o_t^{(i)}$ denotes the t -th observation embedding from the i -th demonstration.

Figure 1 illustrates the training and inference procedure of the policy π_θ . At each time step t , given a prompt \mathcal{P} and the sequence of query observation embeddings $o_{1:t}^q$, the policy predicts the next action for the query trajectory as:

$$a_t = \pi_\theta(\mathcal{P}; o_1^q, \dots, o_t^q), \quad t = 1, \dots, T_q. \quad (8)$$

Both the prompt and the query trajectory are sampled from demonstrations of the same task. During training, the output actions of both the prompt and the query are supervised using the ground-truth actions. We adopt a multi-task learning approach following ICRT [17], enabling the policy to infer the task intent from the prompt and to generalize to unseen tasks.

At inference time, the query trajectory is initialized with the first observation embedding o_0^q . The policy then iteratively outputs the next action using the contextual information in \mathcal{P} and applies the action to the environment, gradually building on the resulting observations until the task is complete.



Figure 2: We evaluate RoboSSM on challenging manipulation tasks from the LIBERO benchmark. In LIBERO-90, the Study and Living Room suites each comprise 4 different scenes, while the Kitchen suite comprises 10 different scenes.

4 Empirical Results

In this section, we evaluate whether RoboSSM can execute unseen tasks based on demonstration prompts, comparing it to previous state-of-the-art ICIL methods. Section 4.1 describes the experimental setup, including dataset construction. We consider two regimes: out-of-distribution to assess length generalization with $|\mathcal{P}_{\text{test}}| > |\mathcal{P}_{\text{train}}|$ (Sec. 4.2), and in-distribution where $|\mathcal{P}_{\text{test}}| = |\mathcal{P}_{\text{train}}|$ (Sec. 4.3). We then compare RoboSSM against a multi-task learning (MTL) policy without in-context learning on unseen tasks (Sec. 4.4) and investigate the effect of β -scaling in Longhorn for ICIL (Sec. A.2). Additionally, we provide latent-space visualizations of trajectories (Sec. 4.5) and an inference runtime analysis (Sec. A.3).

We design our experiments to answer the following research questions:

- **Q1:** Can RoboSSM extrapolate to prompts composed of demonstrations that are longer than those used in training?
- **Q2:** How many training demonstrations are required for RoboSSM to effectively infer with long prompts?
- **Q3:** Can RoboSSM achieve comparable performance to Transformer-based baselines when the test-time prompt length is equal to that used in training?
- **Q4:** Can RoboSSM achieve superior performance on unseen tasks compared to multi-task learning?

4.1 Experimental Setup

Datasets We conduct experiments on the **LIBERO** benchmark [24], a challenging benchmark for visuomotor robot manipulation. LIBERO consists of five task suites: *LIBERO-Object*, *LIBERO-Goal*, *LIBERO-Spatial*, *LIBERO-Long*, and *LIBERO-90*. *LIBERO-90* consists of 90 tasks, while each of the other suites contains 10 tasks. Each task contains 50 demonstration trajectories.

In our experiments, we use the full *LIBERO-Object* suite and divide *LIBERO-90* into three task suites based on scene type, including *kitchen*, *living room*, and *study* scenes, as shown in Figure 2. For all experiments, the test set $\mathcal{D}_{\text{test}}$ contains tasks that are completely disjoint from the training set $\mathcal{D}_{\text{train}}$, ensuring that models are evaluated on entirely unseen tasks. We set $|\mathcal{D}_{\text{test}}| = 2$ for all task suites, and $|\mathcal{D}_{\text{train}}| = 8$ except for the *living room* suite, where $|\mathcal{D}_{\text{train}}| = 14$.

Baselines We compare **RoboSSM** with **ICRT** [17], a Transformer-based in-context imitation learning method that employs LLaMA2-Base [10] as its backbone. To ensure a fair comparison, we configure both backbones to have a similar number of parameters. Both baselines use 4 blocks from their respective backbones. We use LLaMA2-Base with 6 attention heads per layer and a hidden size of 512. For Longhorn, the input is projected to a value dimension $d = 512$ and the keys and

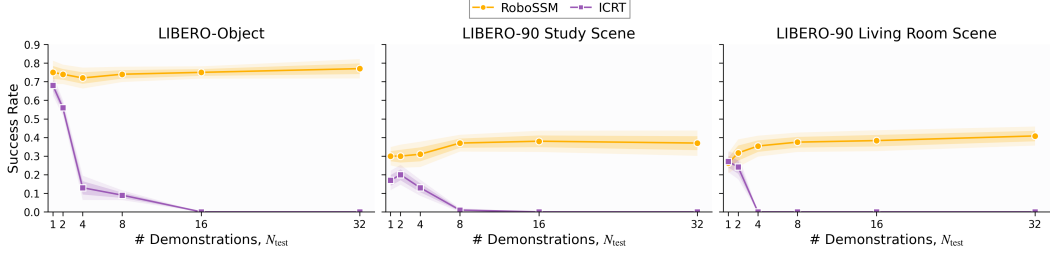


Figure 3: Comparison of RoboSSM and ICRT across test-time demonstrations (N_{test}), with both models trained with $N_{\text{train}}=2$. ICRT’s performance drops sharply once $N_{\text{test}} > N_{\text{train}}$.

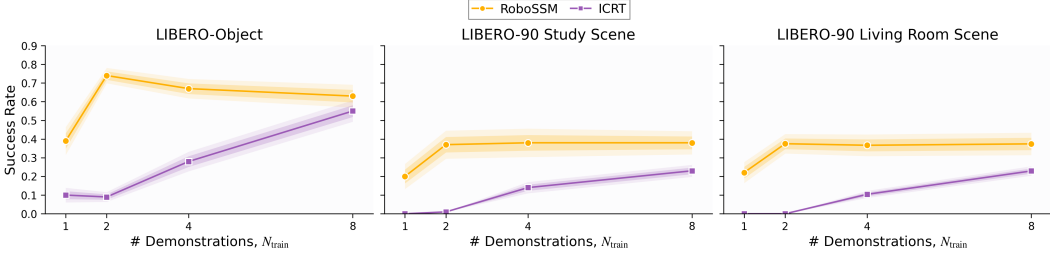


Figure 4: Results with a fixed test-time prompt $N_{\text{test}} = 8$ across models trained with different numbers of demonstrations (N_{train}).

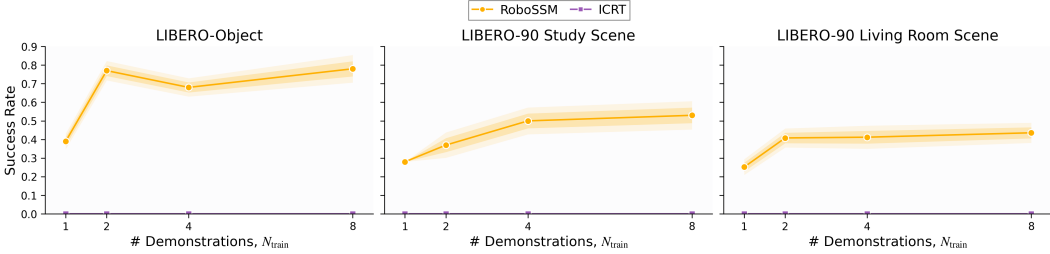


Figure 5: Results with a fixed test-time prompt $N_{\text{test}} = 32$ across models trained with different numbers of demonstrations (N_{train}). ICRT yields zero success across all N_{train} .

queries are projected to $m = 16$. We also implement multi-task learning (MTL) without in-context learning for each backbone, where **MTL-TF** uses LLaMA2-Base and **MTL-SSM** uses Longhorn. Unlike RoboSSM and ICRT, these MTL baselines take language instructions as input to specify the task.

4.2 Prompt Length Generalization

We investigate long-range in-context imitation learning with RoboSSM, focusing on prompt-length extrapolation to out-of-distribution contexts. In this section, we consider two approaches to making the test prompt substantially longer than the training prompt ($|\mathcal{P}_{\text{test}}| > |\mathcal{P}_{\text{train}}|$): (1) increasing the number of demonstrations and (2) applying temporal dilation to demonstrations. These experiments are conducted on LIBERO-Object, LIBERO-90 Study, and the Living Room Scene.

4.2.1 Number of demonstrations

To measure how performance changes when the number of test-time demonstrations differs from that of training, we train models with a small number of demonstrations $N_{\text{train}} = 2$. We then evaluate them on test prompts with $N_{\text{test}} \in \{1, 2, 4, 8, 16, 32\}$. Additionally, to assess how many training

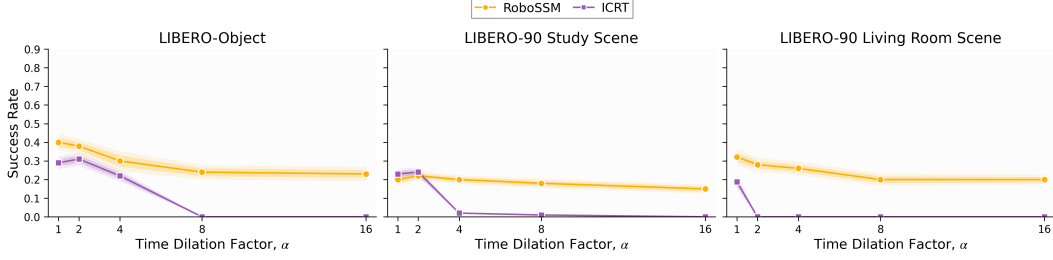


Figure 6: Comparison of RoboSSM and ICRT when evaluated on test-time prompts with temporally dilated demonstrations.

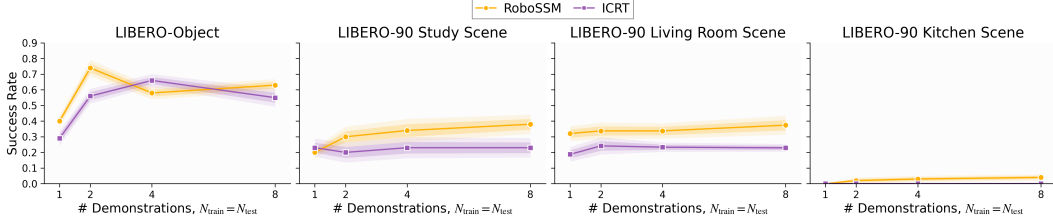


Figure 7: Results of RoboSSM and ICRT on in-context imitation learning, where the prompt consists of the same number of demonstrations during training and testing ($N_{\text{train}} = N_{\text{test}}$).

demonstrations are needed for long-range in-context learning, we train with $N_{\text{train}} \in \{1, 2, 4, 8\}$ and evaluate with fixed test-time prompts $N_{\text{test}} \in \{8, 32\}$.

As shown in Figure 3, RoboSSM maintains or even slightly improves its success rates as N_{test} increases beyond N_{train} . On LIBERO-Object, RoboSSM achieves its best performance at $N_{\text{test}} = 32$, which is 16 times longer than the training prompt length $|\mathcal{P}_{\text{train}}|$, despite never having observed such long prompts during training. This result answers **Q1** affirmatively, showing that RoboSSM extrapolates effectively to prompts far longer than the training horizon. In contrast, ICRT degrades sharply once $N_{\text{test}} > N_{\text{train}}$ and collapses on the long prompts. This result suggests that ICRT fails to generalize to longer prompts, performing reliably only when N_{test} is equal to or shorter than N_{train} .

According to Figure 4 and 5, RoboSSM achieves strong long-range in-context learning even when trained with few task examples, addressing **Q2**. With $N_{\text{test}} = 8$, ICRT is comparable to RoboSSM only when trained with demonstrations matching the test-time prompt length. Notably, with the longer test prompt ($N_{\text{test}} = 32$), RoboSSM maintains high success rates despite limited training demonstrations, whereas ICRT collapses to 0% success rate for all N_{train} .

4.2.2 Time dilation

In real-world scenarios, robot demonstrations may vary in execution speed due to factors such as operator latency, hardware variability, or differing task conditions. To simulate such temporal variability, we evaluate whether the model can generalize to time-dilated demonstrations at test time. We create temporally stretched demonstrations by repeating each observation embedding in the original trajectory α times, resulting in a new trajectory of length $\alpha \cdot T$, where T is the original trajectory length and $\alpha \in \{1, 2, 4, 8, 16\}$ is the dilation factor. Although models are trained with the original prompt length ($\alpha = 1$), we evaluate their robustness under extended test-time prompts, where $|\mathcal{P}_{\text{test}}| = \alpha \cdot |\mathcal{P}_{\text{train}}|$, with α up to 16.

Figure 6 indicates that RoboSSM sustains competitive success rates across increasing dilation factors α , highlighting robustness to temporal stretching and long-context extrapolation (**Q1**). By contrast, ICRT exhibits a consistent performance decay with α , culminating in failure on the longest prompts.

Table 1: Comparison of RoboSSM and ICRT with multi-task learning policies using their respective backbones. *w/ lang* denotes that language instructions are included in the input, and *w/o lang* denotes that language instructions are excluded. Both ICIL frameworks consistently outperform the MTL baselines.

Method	LIBERO-Object	LIBERO-90		
		Study Scene	Living Room Scene	Kitchen Scene
MTL-TF	24.6 ± 1.7	14.6 ± 6.4	2.5 ± 3.5	0.0 ± 0.0
MTL-SSM	20.4 ± 3.3	15.0 ± 4.3	0.4 ± 0.9	0.0 ± 0.0
ICRT [17]				
w/o lang	66.3 ± 3.8	22.5 ± 6.1	23.3 ± 2.8	0.0 ± 0.0
w/ lang	48.3 ± 5.9	21.2 ± 3.8	27.1 ± 3.7	0.0 ± 0.0
RoboSSM				
w/o lang	57.9 ± 3.5	34.2 ± 3.7	33.8 ± 4.5	3.8 ± 1.3
w/ lang	45.4 ± 2.3	35.4 ± 4.7	29.2 ± 5.3	0.0 ± 0.0

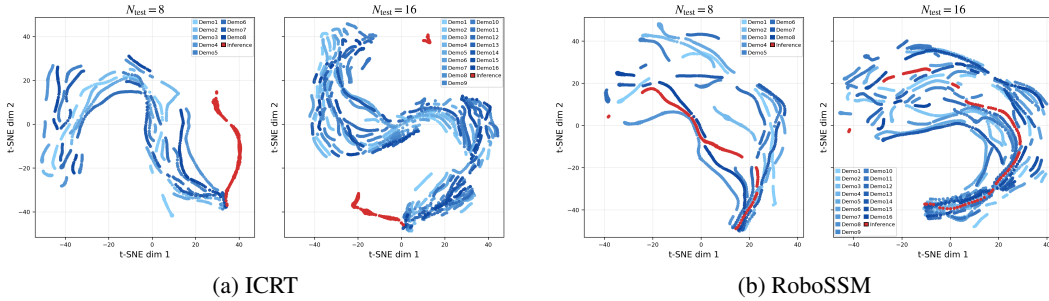


Figure 8: Visualization of latent representations for (a) ICRT and (b) RoboSSM with $N_{\text{test}} = 8$ (left) and $N_{\text{test}} = 16$ (right). Blue corresponds to prompt demonstrations and red to the predicted trajectory.

4.3 In-distribution ICIL

We evaluate RoboSSM and ICRT on in-distribution tasks, where $N_{\text{train}} = N_{\text{test}} \in \{1, 2, 4, 8\}$, covering varying numbers of demonstrations. This evaluation addresses **Q3** when both models operate under distributional conditions similar to those in training. Given that Longhorn exhibits performance parity with Transformers in language modeling tasks, we expect it to demonstrate comparable performance to Transformers in this setting.

Nevertheless, as shown in Figure 7, RoboSSM consistently achieves higher success rates than ICRT across most scenarios, particularly in the LIBERO-90 Study and Living Room scenes. This performance is likely due to the Longhorn architecture, whose formulation as an online regression problem enhances its in-context learning capability. However, both models struggle in the LIBERO-90 Kitchen scene, likely due to the inherent difficulty of the tasks in that suite.

4.4 Comparison to Multi-Task Learning

We compare RoboSSM against multi-task learning baselines, MTL-TF and MTL-SSM. During training, we set $N_{\text{train}} = 4$, and for evaluation, $N_{\text{test}} = 0$ for MTL baselines, while $N_{\text{test}} = 4$ for RoboSSM and ICRT. To enable a fair comparison with language-conditioned MTL, we additionally train and evaluate RoboSSM and ICRT with language instructions. In response to **Q4**, Table 1 shows that RoboSSM and ICRT reliably surpass the MTL baselines across their respective backbone architectures. However, the inclusion of language does not lead to improved performance on unseen tasks, as the language instructions serve to identify tasks [24]. These results demonstrate that RoboSSM effectively handles unseen tasks without any parameter updates and that RoboSSM achieves stronger performance compared to prior methods.

4.5 Latent Space Analysis

Figure 8 displays visualizations of per-timestep latent representations using two-dimensional t-SNE [36], covering both the prompt demonstrations and the subsequent predicted trajectory. For ICRT, we plot the MLP output from the final LLaMA2-Base block. For RoboSSM, we plot the query-projected state from the final Longhorn block, as defined in equation (4). RoboSSM generates trajectories that stay close to the prompt demonstration clusters, indicating stable long-context extrapolation. In contrast, ICRT produces trajectories that drift far outside the prompt manifold, with the deviation becoming more pronounced as the number of prompt demonstrations increases.

5 Conclusion

In this work, we introduce RoboSSM, a scalable in-context imitation-learning framework built on state-space models (SSMs). RoboSSM executes unseen tasks and exhibits strong prompt-length extrapolation, handling few-shot prompts with linear-time inference. Across the LIBERO benchmarks, it processes prompts up to 16 times longer than those seen in training and outperforms Transformer-based ICIL methods, highlighting SSMs as a promising backbone for long-context robotics. Building on these results, RoboSSM can support continual adaptation for lifelong learning by simply being fed demonstration prompts for new tasks, without any task-specific parameter updates. Despite this potential, our study is limited in its coverage of more complex, compositional tasks. Moreover, comprehensively addressing novel tasks will require broader and more diverse training corpora. Future work may explore scaling datasets in both size and diversity to enable more effective generalization to novel tasks and complex task suites.

References

- [1] B. D. Argall, S. Chernova, M. Veloso, and B. Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- [2] A. Billard, S. Calinon, R. Dillmann, and S. Schaal. Robot programming by demonstration. In *Springer handbook of robotics*, pages 1371–1394. Springer, 2008.
- [3] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- [4] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- [5] A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.
- [6] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024.
- [7] M. J. Kim, C. Finn, and P. Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025.
- [8] J. Hu, R. Hendrix, A. Farhadi, A. Kembhavi, R. Martín-Martín, P. Stone, K.-H. Zeng, and K. Ehsani. Flare: Achieving masterful and adaptive robot policies with large-scale reinforcement learning fine-tuning. *arXiv preprint arXiv:2409.16578*, 2024.
- [9] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [10] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [11] G. Comanici, E. Bieber, M. Schaekermann, I. Pasupat, N. Sachdeva, I. Dhillon, M. Blistein, O. Ram, D. Zhang, E. Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- [12] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [13] Y. Duan, M. Andrychowicz, B. Stadie, O. Jonathan Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba. One-shot imitation learning. *Advances in neural information processing systems*, 30, 2017.
- [14] Z. Mandi, F. Liu, K. Lee, and P. Abbeel. Towards more generalizable one-shot visual imitation learning. In *2022 International conference on robotics and automation (ICRA)*, pages 2434–2444. IEEE, 2022.
- [15] E. Valassakis, G. Papagiannis, N. Di Palo, and E. Johns. Demonstrate once, imitate immediately (dome): Learning visual servoing for one-shot imitation learning. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8614–8621. IEEE, 2022.

- [16] M. Xu, Y. Shen, S. Zhang, Y. Lu, D. Zhao, J. Tenenbaum, and C. Gan. Prompting decision transformer for few-shot policy generalization. In *international conference on machine learning*, pages 24631–24645. PMLR, 2022.
- [17] L. Fu, H. Huang, G. Datta, L. Y. Chen, W. C.-H. Panitch, F. Liu, H. Li, and K. Goldberg. In-context imitation learning via next-token prediction. *arXiv preprint arXiv:2408.15980*, 2024.
- [18] A. D. Vuong, M. N. Vu, D. An, and I. Reid. Action tokenizer matters in in-context imitation learning. *arXiv preprint arXiv:2503.01206*, 2025.
- [19] N. Di Palo and E. Johns. Keypoint action tokens enable in-context imitation learning in robotics. *arXiv preprint arXiv:2403.19578*, 2024.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [21] Y. Zhou, U. Alon, X. Chen, X. Wang, R. Agarwal, and D. Zhou. Transformers can achieve length generalization but not robustly. *arXiv preprint arXiv:2402.09371*, 2024.
- [22] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35: 16344–16359, 2022.
- [23] B. Liu, R. Wang, L. Wu, Y. Feng, P. Stone, and Q. Liu. Longhorn: State space models are amortized online learners. *arXiv preprint arXiv:2407.14207*, 2024.
- [24] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023.
- [25] W. Sun, A. Venkatraman, G. J. Gordon, B. Boots, and J. A. Bagnell. Deeply aggravated: Differentiable imitation learning for sequential prediction. In *International conference on machine learning*, pages 3309–3318. PMLR, 2017.
- [26] B. Wang, E. Adeli, H.-k. Chiu, D.-A. Huang, and J. C. Nibbles. Imitation learning for human pose prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7124–7133, 2019.
- [27] M. Bain and C. Sammut. A framework for behavioural cloning. In *Machine intelligence 15*, pages 103–129, 1995.
- [28] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.
- [29] A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [30] D. Y. Fu, T. Dao, K. K. Saab, A. W. Thomas, A. Rudra, and C. Ré. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- [31] J. T. Smith, A. Warrington, and S. W. Linderman. Simplified state space layers for sequence modeling. *arXiv preprint arXiv:2208.04933*, 2022.
- [32] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [33] C. Lu, Y. Schroecker, A. Gu, E. Parisotto, J. Foerster, S. Singh, and F. Behbahani. Structured state space models for in-context reinforcement learning. *Advances in Neural Information Processing Systems*, 36:47016–47031, 2023.

- [34] X. Jia, Q. Wang, A. Donat, B. Xing, G. Li, H. Zhou, O. Celik, D. Blessing, R. Lioutikov, and G. Neumann. Mail: Improving imitation learning with mamba. *arXiv preprint arXiv:2406.08234*, 2024.
- [35] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936, 2003.
- [36] L. Metz, C. D. Freeman, S. S. Schoenholz, and T. Kachman. Gradients are not all you need. *arXiv preprint arXiv:2111.05803*, 2021.
- [37] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

A Appendix

A.1 Implementation Details

We use a single NVIDIA A100 GPU for all models and task suites. During training, we use the AdamW [37] optimizer with weight decay $1e-4$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. The learning rate follows a cosine decay schedule from $1e-4$ to $1e-5$. We use front-view and hand-view RGB images for visual observations, and the robot’s joint angles and gripper state for proprioception. For visual observations, data augmentation includes color jitter with brightness, contrast, and saturation factors of 0.3, and a random masking scheme that applies up to 8 square masks of size 16×16 . Models are trained for 200 epochs with a batch size of 4. During evaluation, we execute 20 rollouts per task, generate trajectories of up to 200 time steps, and report the average success rate over all tasks in the suite across 6 random seeds.

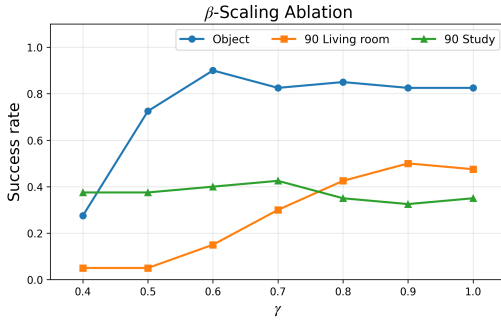


Figure 9: Effectiveness of test-time β -scaling across task suites.

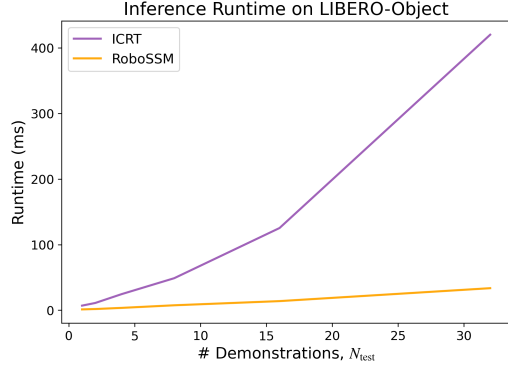


Figure 10: Inference runtime on LIBERO-Object for ICRT and RoboSSM.

A.2 β -Scaling of Longhorn for ICIL

We rescale the weight vector β in equation (5) at test time as $\beta'_t = \gamma \beta_t$ to investigate how regulating the trade-off between reliance on the previous state and incorporation of new inputs affects RoboSSM. Figure 9 reports results for $N_{\text{test}} = 8$ and $N_{\text{train}} = 2$ with $\gamma \in [0.4, 1.0]$. Decreasing β improves performance on several task suites by biasing the model toward attending to prior state information.

A.3 Runtime Analysis

Figure 10 reports the inference runtime on LIBERO-Object. Each run uses a prompt constructed from N_{test} demonstrations and executes $T_q = 200$ action-update steps. Empirically, RoboSSM achieves lower runtime, whereas ICRT’s runtime increases rapidly as the number of demonstrations grows. RoboSSM scales nearly linearly with the prompt length L , yielding total work $\mathcal{O}(L + T_q)$. In contrast, ICRT, with a LLaMA-2 backbone and a KV cache, incurs $\mathcal{O}(L^2)$ prefill over the prompt and $\mathcal{O}(T_q)$ decoding for next-action prediction. Its runtime therefore increases markedly more rapidly as L grows, widening the runtime gap in favor of RoboSSM at longer prompt lengths. Overall, RoboSSM maintains efficient inference with long-context demonstrations.