

# FLaRe: Achieving Masterful and Adaptive Robot Policies with Large-Scale Reinforcement Learning Fine-Tuning

Jiaheng Hu<sup>1:2</sup>, Rose Hendrix<sup>1</sup>, Ali Farhadi<sup>1:3</sup>, Aniruddha Kembhavi<sup>1:3</sup>, Roberto Martín-Martín<sup>2</sup>, Peter Stone<sup>2:4</sup>, Kuo-Hao Zeng<sup>1:y</sup>, and Kiana Ehsani<sup>1:y</sup>

<sup>1</sup> Allen Institute for Artificial Intelligence (AI2)

<sup>2</sup> University of Texas, Austin<sup>3</sup> University of Washington<sup>4</sup> Sony AI

<sup>y</sup> Equal Supervision.

**Abstract:** In recent years, the Robotics field has initiated several efforts toward building generalist robot policies through large-scale multi-task Behavior Cloning. However, direct deployments of these policies have led to unsatisfactory performance, where the policy struggles with unseen states and tasks. How can we break through the performance plateau of these models and elevate their capabilities to new heights? In this paper, we propose FLaRe, a large-scale Reinforcement Learning fine-tuning framework that integrates robust pre-trained representations, large-scale training, and gradient stabilization techniques. Our method aligns pre-trained policies towards task completion, achieving state-of-the-art (SoTA) performance both on previously demonstrated and on entirely novel tasks and embodiments. Specifically, on a set of long-horizon mobile manipulation tasks, FLaRe achieves an average success rate of 79.5% in unseen environments, with absolute improvements of +23.6% in simulation and +30.7% on real robots over prior SoTA methods. By utilizing only sparse rewards, our approach can efficiently master new capabilities beyond the pretraining data with minimal human effort. Moreover, we demonstrate rapid adaptation to new embodiments and behaviors with less than a day of fine-tuning, opening up possibilities for robots to continually adapt and improve when facing new tasks. Videos can be found on the project website [arobot-flare.github.io](https://arobot-flare.github.io)

## 1 INTRODUCTION

Foundation models in computer vision and natural language processing have recently achieved groundbreaking successes. Large transformer models, such as GPT-4 [SAM [29], have demonstrated the ability to perform an extensive range of tasks. Inspired by these advances, the robotics community has set its sights on training high-capacity, multi-task transformers for robotic applications.

One of the prominent methods in this pursuit is large-scale behavior cloning (BC) which leverages large datasets of real-world and simulated demonstrations (e.g. RT-2 [7], RT-X [42], and SPOC [4]) to train high-capacity policies that can perform many different tasks. While BC policies have shown promise, they remain fundamentally limited when directly deployed in the real world: models are constrained to the states observed during training, making it difficult to generalize beyond expert trajectories. Consequently, these policies often struggle when faced with unfamiliar states, and fail to recover from errors effectively.

On the other hand, reinforcement learning (RL) offers a complementary approach that directly optimizes the performance of the robot through trial-and-error learning, and RL algorithms have achieved many successes when a well-defined reward function is available [62, 65]. However, many RL algorithms are notoriously sample inefficient, requiring extensive training. As task horizon increases and action space grows, RL policies struggle to get off the ground due to the large search space. Moreover, RL's reliance on hand-crafted reward functions severely limits its scalability.

# FLaRe: Fine-tuning Large-Scale Robot Polices with RL

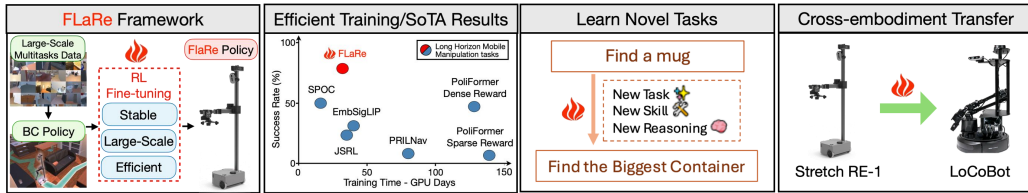


Figure 1: FLaRe is a simple but effective approach for large-scale fine-tuning of robotic policies. FLaRe achieves SoTA performance on simulation (+23.6%) and real-world (+30.7%) benchmarks, can generalize to unseen tasks, and adapts to new behaviors and embodiments.

Although insufficient for direct deployment, the policies trained through large-scale multi-task Behavior Cloning already possess extremely valuable features and behavior priors. How can we break through the performance plateau of these models and elevate their capabilities to new heights? Our key insight is that, through RL, we can align the behavior of these policies towards true objectives such as task completion (instead of the BC objective), thereby achieving masterful performance not only on tasks seen during BC training, but also on novel tasks and embodiments never seen by the pre-trained policy.

While attempts have been made to fine-tune BC policies with RL [58, 46, 45, 2], these works are only verified with small-scale networks and in single-task domains. Empirically, we find these methods ineffective as the capacity of the pre-trained policy increases, where the abrupt shift from BC to RL results in destructive gradient updates, leading to oscillations or even collapse during RL training.

In FLaRe, we introduce an effective, scalable, and robust solution for fine-tuning large-scale robot policies with RL. Illustrated in Fig. 1 top-left, FLaRe starts from a **multi-task** robotics policy, and fine-tunes it with **large-scale** RL through extensive use of simulation. To ensure the RL fine-tuning is **stable**, FLaRe introduces a set of simple yet highly effective techniques, detailed in Sec. 4.3, that drastically improve performance and reduce training time compared to previous methods.

FLaRe achieves SoTA performance on household mobile manipulation tasks. In established simulation benchmark [14], it achieves an average 79.5% success rate, +23.6% absolute improvements over the best baseline. In the real world, FLaRe achieves excellent results (80.7% SR on average), outperforming the best prior work by +30.7%. Furthermore, FLaRe offers several key advantages:

1. FLaRe enables efficient training with a 15x reduction in training time compared to the previous SoTA method, using a simple sparse reward without the need for handcrafted reward functions (Fig 1 top-right).
2. FLaRe allows for continual learning beyond the tasks seen during BC. Even for new tasks without expert trajectories or shaped rewards, FLaRe can be fine-tuned to achieve state-of-the-art performance (Fig 1 bottom-left).
3. FLaRe can rapidly adapt to new embodiments and behaviors, significantly enhancing the base model’s flexibility and applicability during lifelong deployment (Fig 1 bottom-right).

We find that FLaRe marks a promising achievement towards developing highly generalizable robotic systems that can handle a wide range of tasks in diverse environments in a lifelong fashion.

## 2 Related Work

### 2.1 Foundation model for robotics

Following the successes of foundation models in vision [29] and language [1], there has been a recent trend towards training robotics-specific foundation models [15, 22]. While these models focus on different robot applications, such as manipulation (e.g. RT-1 [6], RT-2 [7], RT-X[42], Octo [55], RoboCat [5], OpenVLA [28]), navigation (e.g. ViNT [50]), and mobile manipulation (e.g. SPOC[14]), they share a similar recipe of training a high-capacity transformer model through multi-task behavior cloning [48]. As a result, they generate the same end-product: a multi-task transformer policy, which FLaRe can use as a base model for fine-tuning.

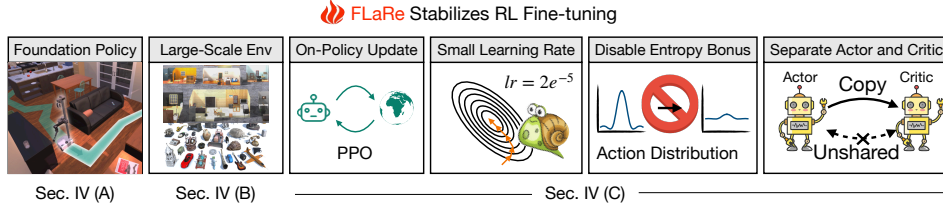


Figure 2: FLaRe introduces a series of design choices that help stabilize the RL training process, including 1) fine-tuning from a multi-task robotics policy, 2) large-scale fine-tuning in simulation, 3) using an on-policy algorithm as opposed to off-policy methods, 4) utilizing smaller learning rate than when performing RL from scratch, 5) disabling the entropy bonus objective that can potentially distort the policy at the start of the training, and 6) separating the actor and the critic network, so that the critic update will not influence the policy prediction.

## 2.2 RL training and fine-tuning of robotics models

While RL has achieved many successes in robotics[53], directly applying RL from scratch often requires extensive reward engineering and long training time [3, 21, 62, 65]. Hence, previous works have extensively explored leveraging pretrained models to facilitate RL [52, 27, 54, 38, 2, 19, 17, 24, 59, 30, 36, 45, 4, 46, 68, 60, 58, 20].

However, many of these approaches focus on fine-tuning models that have been pre-trained using either online RL [52, 27, 54] or offline RL [39, 32, 34, 33], which limits their applicability. This makes them unsuitable for fine-tuning most existing robotics foundation models, which are typically trained using Behavior Cloning. Many previous works also require access to the entire offline dataset during fine-tuning[38, 2, 19, 17, 24, 59, 30, 36, 45], which may be feasible for small-scale data and low-dimensional observations but is unlikely to be computationally feasible for large-scale data and image observations, as also noted by Ramrakhya *et al* [46].

In addition, the techniques proposed in many of these works are only evaluated on simple domains, with low-dimensional state spaces [45, 17, 38], small-scale network architecture (e.g. MLP)[45, 58, 2], single-task pretraining and fine-tuning [66, 30], and often no real robot experiments [2, 4, 17, 38]. PIRLNav [46] and JSRL [58] are two works that are closest to our setting, where only a pretrained policy is required for the fine-tuning phase. However, both of them focus on single-task setting with small-scale networks and no real robot experiments. In contrast, FLaRe explores fine-tuning from large robotics models, where both scalability and applicability to real robots are of critical concern.

## 3 Problem Formulation

We consider each robotics task  $T \in \mathcal{T}$  as a language-conditioned Partially Observable Markov Decision Process  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \mathcal{O}, \mathcal{L}, P(s_0), \gamma)$ , where  $\mathcal{S}$  is a state space,  $\mathcal{A}$  is an action space,  $\mathcal{O}$  is an observation space,  $\mathcal{P}$  is a Markovian transition model,  $\mathcal{L}$  is a set of natural language instruction,  $\gamma$  is a discount factor,  $P(s_0)$  is the initial state distribution, and  $R$  is a **sparse** reward function that takes in a natural language instruction  $l \in \mathcal{L}$  and a state  $s \in \mathcal{S}$  and outputs a binary value indicating whether a given instruction is completed. For the purpose of this paper, we assume that all tasks have the same action space (the actuators of the robot) and observation space (the robot’s sensors). Each task  $T \in \mathcal{T}$  defines a set of natural language instructions  $\mathcal{L}_T$  (e.g., for the task of Object Navigation, potential instructions can be “go to an apple”, “find a houseplant”, and more). At the start of every episode, an instruction  $l_T \in \mathcal{L}_T$  and an initial state  $s_0 \sim P(s_0)$  will be sampled. Every time a specific task  $T \in \mathcal{T}$  is given, our goal is to train a policy  $\pi_\theta^T$  that maximizes the expected return (i.e. success rate)  $\mathbb{E}_{\mathcal{L}_T, \pi} \sum_t R(s_t, l)$  for the given task over the possible language instructions  $\mathcal{L}_T$ .

## 4 Method

Considerable effort has been devoted to optimizing performance on robotics tasks via training high-capacity models  $\pi_\theta$  with large-scale, multi-task imitation learning[14, 6, 7, 42]. In practice, these efforts lead to unsatisfactory performance due to compounding errors [47], where small action prediction error leads to state distribution drift. Furthermore, for novel tasks and scenarios where no

Table 1: Success and Episode-length weighted Success (SEL) against baseline methods on the CHORES[14] benchmark (in-distribution tasks). Baselines with privileged information are *marked in blue*. FLARe significantly outperforms the previous SoTA methods.

Success (SEL) $\uparrow$	IL+RL: Sparse Reward			IL Only	RL Only		
	FLARe [Ours]	PIRLNav	JSRL	SPOC	Poliformer - Sparse	Poliformer - Dense	EmbSigLIP - Dense
ObjectNav	85.0 ( <b>67.6</b> )	20.0 (7.0)	21.0 (15.6)	55.0 (42.2)	14.5 (10.4)	<b>85.5</b> (61.2)	36.5 (24.5)
Fetch	<b>66.9</b> ( <b>54.7</b> )	0.0 (0.0)	2.9 (2.8)	14.0 (10.5)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
PickUp	<b>91.8</b> ( <b>90.4</b> )	0.0 (0.0)	50.9 (47.7)	90.1 (86.9)	0.0 (0.0)	90.1 (88.7)	71.9 (52.9)
RoomVisit	<b>70.4</b> ( <b>67.1</b> )	12.5 (11.0)	19.0 (18.6)	40.5 (35.7)	12.5 (12.5)	12.5 (10.9)	16.5 (11.9)

demonstration data is available, these models have shown limited generalization capabilities, likely due to the limited task coverage of the training data.

FLARe addresses both problems by fine-tuning the pre-trained model  $\pi_\theta$  with RL for each given task  $T \in \mathcal{T}$ . The key idea of FLARe is to achieve stable and effective RL fine-tuning through a series of design choices, including 1) utilize a large-scale multi-task model as the base model, 2) achieve large-scale fine-tuning through extensive use of simulations, and 3) a series of algorithmic design to stabilize the RL fine-tuning. Together, these design choices enable FLARe to effectively learn from **sparse** reward and achieve good performances. We elaborate in detail on each of these decisions in the following sections (Fig. 2).

#### 4.1 Fine-tune from a multi-task robotics model

The first key design choice of FLARe is to start from a **multi-task** pre-trained large model (i.e. a foundational robotics model). Compared to fine-tuning from a single-task, small-scale network (as is often the case in previous works [58, 45, 46, 66]), starting from a robotics foundation model brings three key benefits. First, models pre-trained on diverse tasks can master more robust representations and more versatile behavior priors [63], which will benefit the fine-tuning process. Second, the highly capable network architecture (e.g. large transformer models) that comes with these foundational robotics models brings good inductive bias that can facilitate generalization [12], which is crucial to fine-tuning. Most importantly, the multi-task capability of these models allows us to reuse the same model for fine-tuning for many different tasks. In fact, as we will show in the experiments in Sec. 5.2, we can even fine-tune for tasks and embodiments that have never been seen by the pre-trained policy and still achieve good performance.

While our method can in principle work on any foundational robotics model, in this specific work, we focus on fine-tuning the SPOC model (Fig. 6) [14] — a multi-task transformer model for mobile manipulation tasks, trained on large-scale shortest path expert trajectories collected in Objaverse-Populated ProcTHOR houses[31, 11, 9]. We refer the reader to our supplementary material for more details regarding the SPOC model.

#### 4.2 Large-scale fine-tuning in simulation

The second key design choice of FLARe is to perform large-scale fine-tuning through extensive use of simulation. Recent advancements in robotics and embodied AI have given us a set of tools for simulating robotics tasks [31, 35, 23, 44, 67, 61]. In this work, we utilize AI2THOR [31] to perform large-scale simulated fine-tuning with diverse objects and scenes, which includes 150k procedurally generated PROCTHOR houses [9] and 800K+ annotated 3D objects [11].

When using simulation in robotics, addressing the sim-to-real gap [64] becomes a critical challenge. In FLARe, similar to Ehsani *et al.*[14], we employed two techniques to facilitate sim-to-real transfer. First, we perform extensive domain randomization, including color augmentation, applying random crops, and posterizing the images. Second, we extract visual features through DinoV2 [41], a pre-trained foundational vision model, which captures useful features that can generalize across simulation and the real world.

To ensure large-scale training of the transformer policy and value networks, we utilize the KV-cache technique[43] to reduce the computational costs during network inference, similar to Zeng *et al.*[62]. The KV-cache technique caches and reuses the keys and values of earlier observations within an episode. This reduces the inference complexity of the transformer network from quadratic to linear, which is crucial for affordable large-scale RL fine-tuning.

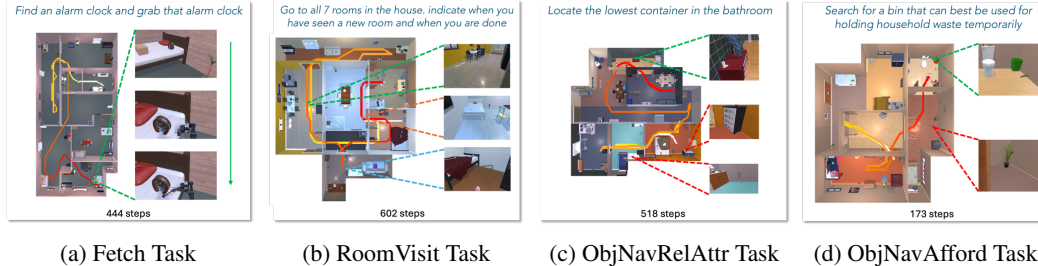


Figure 3: We evaluate FLaRe on mobile manipulation tasks. (a, b) In-distribution tasks, in unseen environments. (c, d) Novel tasks that require unseen capabilities from pretraining, in unseen environments. FLaRe excels in long-horizon tasks, showing strong object recognition, relational reasoning, and exploration abilities.

### 4.3 Stabilize RL fine-tuning

Finally, we introduce a set of simple but very critical algorithmic choices to ensure the stability of RL fine-tuning. While these techniques are relatively simple, as we will show in the ablation studies in Sec. 5.5, each choice is very important to ensure stable training and to obtain good performances.

**Using On-policy Algorithms.** Off-policy RL methods [37, 18] can utilize off-policy data during training, and thus bring the promise of sample-efficient RL. However, compared to on-policy methods, off-policy RL is often less stable and more sensitive to hyperparameters, both in theory and in practice, due to problems associated with the “deadly triad” [51]. In this work, since we perform fine-tuning entirely in simulation, we are less constrained by the sample efficiency of our RL algorithms, and therefore choose to use on-policy algorithms for stable fine-tuning. Specifically, we use PPO [49], a state-of-the-art on-policy policy gradient method.

**Taking Small Update Steps.** When setting the learning rate for RL, it is common practice to reuse a learning rate that has previously achieved success in the same/similar domains. However, what we found in FLaRe is that fine-tuning from an existing policy requires significantly lower learning rates than when starting from scratch. For example, in the object navigation task, the previous state-of-the-art result is achieved with PPO from scratch using a learning rate of  $2e - 4$ . In FLaRe, when fine-tuning on the exact same task, we have to reduce the learning rate by an order of magnitude to achieve stable learning. It is important to notice that we do not perform additional LR tuning in FLaRe - the same learning rate is used for all experiments and tasks.

**Disabling Entropy Bonus.** The PPO objective [49] contains an entropy bonus, which promotes the entropy of the action distribution predicted by the policy network to ensure sufficient exploration. However, we found that when fine-tuning from a pre-trained policy network, this additional entropy term can quickly distort the policy gradient update at the start of the training, leading to unlearning of the pre-trained policy. Hence, we remove this entropy bonus term from our PPO update in FLaRe.

**Disabling Feature Sharing.** When applying RL to high-dimensional observations such as images, a standard practice is to have a shared feature extractor between the actor and the critic network, which can facilitate the learning of useful features. However, we found that feature sharing during RL fine-tuning can actually hurt the performance since the gradient from the critic loss will change the pre-trained features and lead to the deterioration of the action prediction. Furthermore, during RL fine-tuning, since the pre-trained foundation model should already capture good representations, there is no need for the actor and the critic network to share the same feature extractor. Therefore, in FLaRe, we initialize the policy and the critic network as independent networks, both using the weight and architecture of the pre-trained transformer policy. The policy head of the critic network is replaced by a randomly initialized values.

We found that all four training components are important, and in Section 5.5, we show that removing any one of them results in training collapse.

## 5 Results

We evaluate FLaRe on a set of navigation and manipulation tasks both in simulation and in the real world. Through our experiments, we seek to answer the following questions: **Q1**: Can FLaRe achieve state-of-the-art performance on tasks both within and outside the training data of the pre-trained

Table 2: FLaRe can fine-tune for tasks that are never seen by the base model, and achieve state-of-the-art performance. Baselines with privileged information are *marked in blue*. Sp: sparse reward. De: dense reward.

Success (SEL) $\uparrow$	FLaRe [ours]	Poliformer - Sp	SPOC++	Poliformer - De
ObjNavRelAttr	<b>71.0 (63.6)</b>	6.7 (6.7)	54.5 (44.6)	36.1 (32.4)
RoomNav	<b>91.6 (85.6)</b>	57.0 (51.8)	74.5 (59.9)	75.0 (62.4)
ObjNavAfford	<b>79.7 (70.6)</b>	35.5 (29.4)	62.4 (50.6)	53.8 (43.1)

policy? **Q2**: Can FLaRe learn new capabilities never seen during pre-training and generalize to unseen tasks? **Q3**: Can the policies learned by FLaRe transfer to the real world? **Q4**: Can FLaRe enable efficient adaptation to new robot embodiments and new behavior? **Q5**: Are the stabilization techniques in FLaRe necessary to ensure stable fine-tuning?

All of the experiments use the same hyperparameters, specified in the supplementary. Unless stated otherwise, results for FLaRe are obtained using sparse rewards that correspond to task completion. Visualizations of the robot trajectories are shown in Fig. 3 and on our project website.

## 5.1 FLaRe on seen capabilities

First, we evaluate the performance of FLaRe in comparison to prior behavior cloning (BC) and reinforcement learning (RL) baselines. Specifically, we test FLaRe on CHORES-S [14], a recently introduced simulation benchmark designed for household robot tasks. CHORES-S encompasses four task types that require various skills, including navigation, object recognition, object manipulation, and environment exploration. Similar to [14, 62], the policies use the agent’s RGB observations as input to predict discrete actions, which represent movements of the base, arm, gripper, and an *END* action to signify task completion. For further details on the action space, observation space, and task definitions, please refer to our project website.

CHORES tasks are very challenging due to their long-horizon nature, partial observability, RGB-only sensor, and diverse scenes and objects. Therefore, previous methods struggle to complete these tasks. Since CHORES tasks are in the training data of the SPOC model that FLaRe fine-tunes upon, our goal is to utilize FLaRe to improve performance on these in-distribution capabilities.

**Baselines.** Our baselines consist of prior works in imitation learning, reinforcement learning from scratch, and reinforcement learning fine-tuning from pre-trained policies. Aside from SPOC [14], the robot foundation model that we fine-tune upon, we additionally compare against **Poliformer** [62], a transformer-based RL-from-scratch method that achieved SOTA performance on object navigation; **EmbSigLIP** [26], a GRU-based RL-from-scratch method; **PIRLNav** [46], an RL fine-tuning method that employs learning rate scheduling to warm-start the value function; and **JSRL** [58], an off-policy RL fine-tuning method that gradually “roll in” experiences with the prior policy.

We compare against baselines in two settings: (1) a fair-comparison setting, where the baseline methods use the same sparse reward as FLaRe, and (2) an unfair-comparison setting, where the baseline methods use a privileged, task-specific dense reward that has been hand-coded by human experts. It is important to note that each new task demands significant researcher effort to design and curate a dense reward function that avoids collapsing during training and is not scalable to new tasks.

To demonstrate the superiority of FLaRe, all baseline methods are trained for more steps than FLaRe. Specifically, the fair-comparison baselines are trained for  $3x$  more steps on ObjectNav and RoomVisit, and  $2x$  more steps on Fetch and Pickup. The unfair-comparison baselines are trained until convergence to obtain the best possible result. Notice that this often means significantly longer training time. For example, for the Poliformer (Dense) on ObjectNav, the result is obtained after training for 300M steps - over  $15\times$  as many training steps that FLaRe uses on ObjectNav.

Results are shown in Table 1, Fig. 3 (a, b), and Fig. 4(a), where we evaluate on **unseen** simulated houses and report Success rate as well as Episode-length weighted Success (SEL [13]) which measures the efficiency of the policies. As shown by the table, FLaRe not only significantly outperforms the fair-comparison baselines, but also outperforms the unfair baseline on three out of the four tasks despite using significantly fewer training steps (**Q1**).

Table 3: Real-world results (total of 46 tasks) in an unseen apartment. For manipulation tasks, we report both full success (policy and heuristic grasping) and policy success (proximity) following [14].

Success Rate $\uparrow$	FLaRe [ours]	SPOC	Poliformer - Dense
ObjectNav	<b>94.4</b>	50.0	83.3
Fetch	<b>66.7</b> (55.6)	33.3 (11.1)	X
PickUp	<b>86.7</b> (66.7)	66.7 (46.7)	X
RoomVisit	<b>75.0</b>	50.0	X

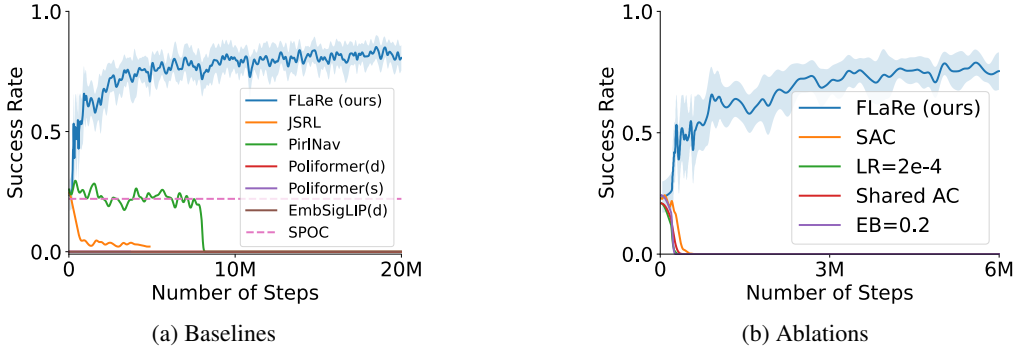


Figure 4: Baseline performances and ablation studies on the Fetch task. FLaRe is the only method that can achieve good performance on this challenging task.

## 5.2 FLaRe on novel capabilities

A well-trained robotics foundation model should learn features useful for all robotics tasks, not only applicable to in-distribution tasks appearing in its original training data. To investigate if FLaRe can take advantage of these pre-trained features, we examine the performance of FLaRe on a set of novel capabilities never seen by the foundation model. Specifically, we evaluate FLaRe on three navigation tasks that specify target objects/locations in different ways and require distinct types of explorations and skills, including 1) ObjNavRelAttr, which identifies target objects through relative object attributes comparison (e.g. “find the largest apple”); 2) RoomNav, which requires the robot to navigate to room types instead of objects (e.g. “go to the kitchen”); and 3) ObjNavAfford, which requires object affordance understanding (e.g. “find something I can sit on”). Note that new reasoning skills are required for these unseen tasks; for example, in ObjNavRelAttr, the agent must search the environment for all objects of the specified type, reason about their properties, and issue a completion action when it identifies the correct instance.

We compare against the Poliformer [62] baseline described in Sec. 5.1, as well as SPOC++, a BC baseline that has the same network architecture as SPOC but uses additional expert demonstrations (1M frames per aforementioned task). Note that these demonstrations are not available to FLaRe, nor to the SPOC model that FLaRe fine-tunes.

We show the results in Table 2 and Fig. 3 (c, d). On these out-of-distribution tasks that require novel capabilities, FLaRe achieves state-of-the-art performance without any additional hyperparameter tuning (**Q2**), even where the baselines have unfair advantages. It is worth noting that, since specifying each of these new tasks  $T_n$  is as simple as specifying a success criteria  $R_n$  and the associated language instructions  $L_n$ , these results imply that we can apply FLaRe to on-the-fly tasks without much engineering effort. This suggests a path towards continual adaptation.

## 5.3 FLaRe on real robots

To examine the performance of FLaRe on real robots, we evaluate the policies learned by FLaRe in a real-world apartment on a Stretch RE-1 [25]. This layout (Fig. 5) is never seen by the robot during training. We directly deploy policies without any adaptation or real-world fine-tuning, and leverage a heuristic object grasping model following SPOC [14]. We compare against SPOC and

Poliformer<sup>1</sup> with dense reward, and report the results in Table 3. Sim-to-real approaches introduced in Sec. 4.2 enable the successes of FLaRe in simulation to directly transfer to the real world, achieving state-of-the-art performances on a set of real world navigation and mobile manipulation tasks (Q3).

#### 5.4 FLaRe for adaptation

FLaRe opens up the possibilities for learning behaviors not captured by the demonstration data (and thus the foundation robotics model). We examine this in two setups, cross-objective and cross-embodiment capabilities of FLaRe (Q4).

##### 5.4.1 Adaptation to New Embodiment

We use FLaRe to fine-tune SPOC (which is trained only on Stretch-RE1) to adapt to Locobot [16]. Locobot has different action space and camera parameters: it lacks the manipulation degrees-of-freedom that Stretch possesses, but has a rotatable, narrow field-of-view camera mounted lower. To facilitate cross-embodiment transition, we simply mask out the invalid actions output by the policy, and repurpose two of the invalid actions to control the camera. FLaRe effectively utilizes the pre-trained policy to adapt to the new embodiment on ObjectNav, as shown by the table below:

New Embodiment	Success Rate $\uparrow$	SEL $\uparrow$
FLaRe	<b>72.0</b>	<b>47.2</b>
Poliformer zero-shot	57.5	30.1
Poliformer (Sparse Reward)	44.0	29.7

##### 5.4.2 Adaptation to New Behavior

We investigate whether FLaRe can be used to shape a robot’s behavior after the policy is trained, using only a few training steps. We test two new behaviors: 1) encouraging the agent to be more efficient (+step penalty  $-0.01/\text{step}$ ), and 2) reducing the number of unwanted collisions with the environment (+collision penalty  $-0.5/\text{collision}$ ). By adding a reward term tailored to each behavior, the policy adapts to these new behaviors after just 6 hours of training, with minimal impact on the success rate. The following table presents the results for the Fetch task:

New Behaviors	Success Rate $\uparrow$	Episode Length $\downarrow$	# of Collisions $\downarrow$
FLaRe	<b>66.9</b>	258.2	10.0
+ Step Pen.	65.7	<b>222.8</b>	10.0
+ Coll. Pen.	66.7	251.2	<b>3.1</b>

#### 5.5 Ablation studies

To evaluate whether the techniques proposed in Sec. 4.3 are necessary for the performance of FLaRe, we evaluate four ablation variants of our method. To evaluate whether using on-policy methods is important, we tested switching the PPO algorithm by Soft Actor-Critic [18] (SAC). To evaluate whether a small learning rate is necessary, we tested setting the learning rate to  $2e - 4$ , 10 times our original learning rate. To evaluate the importance of having separated actor and critic, we tested **Shared AC**, where the actor and critic share the transformer encoder and decoder trunk. Finally, we tested **EB=0.2**, which set the coefficient of the entropy bonus in PPO to 0.2. We show the training curves in Fig. 4(b).

Perhaps surprisingly, removing any single one of the stabilizing techniques in FLaRe results in the success rate quickly collapsing to 0 on the fetch task, while FLaRe learns very robustly with the same set of hyperparameters across variety of tasks and experiment setups (Q5). This showcases the importance of all of the techniques introduced in FLaRe.

<sup>1</sup>Poliformer reported real-world results only for the ObjectNav Task.

## 6 Conclusion

FLaRe is an efficient and scalable approach for fine-tuning large-scale robot policies using RL. It enables effective adaptation to unseen tasks and achieves state-of-the-art performance in both simulation and real-world settings. FLaRe’s adaptability to new embodiments and behaviors unlocks the potential for flexible deployment across a wide range of robotic platforms in a lifelong fashion. FLaRe’s main limitation lies in its reliance on simulation environments for fine-tuning. While leveraging recent work in simulation generation [10, 56] offers a promising direction, tackling tasks where robust simulations are unavailable—such as those involving liquids or soft objects—remains challenging and may require fine-tuning directly in the real world.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Belle-mare. Reincarnating reinforcement learning: Reusing prior computation to accelerate progress. *Advances in neural information processing systems*, 35:28955–28971, 2022.
- [3] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [4] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022.
- [5] Konstantinos Bousmalis, Giulia Vezzani, Dushyant Rao, Coline Manon Devin, Alex X Lee, Maria Bauza Villalonga, Todor Davchev, Yuxiang Zhou, Agrim Gupta, Akhil Raju, et al. Robocat: A self-improving generalist agent for robotic manipulation. *Transactions on Machine Learning Research*, 2023.
- [6] 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.
- [7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [8] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, et al. Robothor: An open simulation-to-real embodied ai platform. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3164–3174, 2020.
- [9] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Kiana Ehsani, Jordi Salvador, Winson Han, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Proctor: Large-scale embodied ai using procedural generation. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022.
- [10] Matt Deitke, Rose Hendrix, Ali Farhadi, Kiana Ehsani, and Aniruddha Kembhavi. Phone2proc: Bringing robust robots into our chaotic world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9665–9675, 2023.
- [11] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023.

- [12] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [13] Ainaz Eftekhari, Kuo-Hao Zeng, Jiafei Duan, Ali Farhadi, Ani Kumbhavi, and Ranjay Krishna. Selective visual representations improve convergence and generalization for embodied ai. *arXiv preprint arXiv:2311.04193*, 2023.
- [14] Kiana Ehsani, Tanmay Gupta, Rose Hendrix, Jordi Salvador, Luca Weihs, Kuo-Hao Zeng, Kunal Pratap Singh, Yejin Kim, Winson Han, Alvaro Herrasti, et al. Spoc: Imitating shortest paths in simulation enables effective navigation and manipulation in the real world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16238–16250, 2024.
- [15] Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. Foundation models in robotics: Applications, challenges, and the future. *arXiv preprint arXiv:2312.07843*, 2023.
- [16] Abhinav Gupta, Adithyavairavan Murali, Dhiraj Prakashchand Gandhi, and Lerrel Pinto. Robot learning in homes: Improving generalization and reducing dataset bias. *Advances in neural information processing systems*, 31, 2018.
- [17] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv preprint arXiv:1910.11956*, 2019.
- [18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR, 2018.
- [19] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan, John Quan, Andrew Sendonaris, Ian Osband, et al. Deep q-learning from demonstrations. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [20] Hengyuan Hu, Suvir Mirchandani, and Dorsa Sadigh. Imitation bootstrapped reinforcement learning. *arXiv preprint arXiv:2311.02198*, 2023.
- [21] Jiaheng Hu, Peter Stone, and Roberto Martín-Martín. Causal policy gradient for whole-body mobile manipulation. *arXiv preprint arXiv:2305.04866*, 2023.
- [22] Yafei Hu, Quanting Xie, Vidhi Jain, Jonathan Francis, Jay Patrikar, Nikhil Keetha, Seungchan Kim, Yaqi Xie, Tianyi Zhang, Zhibo Zhao, et al. Toward general-purpose robots via foundation models: A survey and meta-analysis. *arXiv preprint arXiv:2312.08782*, 2023.
- [23] Emily Jin, Jiaheng Hu, Zhuoyi Huang, Ruohan Zhang, Jiajun Wu, Li Fei-Fei, and Roberto Martín-Martín. Mini-behavior: A procedurally generated benchmark for long-horizon decision-making in embodied ai. *arXiv preprint arXiv:2310.01824*, 2023.
- [24] Ryan Julian, Benjamin Swanson, Gaurav S Sukhatme, Sergey Levine, Chelsea Finn, and Karol Hausman. Never stop learning: The effectiveness of fine-tuning in robotic reinforcement learning. *arXiv preprint arXiv:2004.10190*, 2020.
- [25] Charles C Kemp, Aaron Edsinger, Henry M Clever, and Blaine Matulevich. The design of stretch: A compact, lightweight mobile manipulator for indoor human environments. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 3150–3157. IEEE, 2022.
- [26] Apoorv Khandelwal, Luca Weihs, Roozbeh Mottaghi, and Aniruddha Kumbhavi. Simple but effective: Clip embeddings for embodied ai. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14829–14838, 2022.
- [27] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual reinforcement learning: A review and perspectives. *Journal of Artificial Intelligence Research*, 75: 1401–1476, 2022.

- [28] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- [29] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [30] Jens Kober, Betty Mohler, and Jan Peters. Imitation and reinforcement learning for motor primitives with perceptual coupling. In *From motor learning to interaction learning in robots*, pages 209–225. Springer, 2010.
- [31] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- [32] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- [33] Aviral Kumar, Anikait Singh, Frederik Ebert, Mitsuhiko Nakamoto, Yanlai Yang, Chelsea Finn, and Sergey Levine. Pre-training for robots: Offline rl enables learning new tasks from a handful of trials. *arXiv preprint arXiv:2210.05178*, 2022.
- [34] Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, pages 1702–1712. PMLR, 2022.
- [35] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabrael Levine, Michael Lingelbach, Jiankai Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR, 2023.
- [36] Yao Lu, Karol Hausman, Yevgen Chebotar, Mengyuan Yan, Eric Jang, Alexander Herzog, Ted Xiao, Alex Irpan, Mohi Khansari, Dmitry Kalashnikov, et al. Aw-opt: Learning robotic skills with imitation and reinforcement at scale. *arXiv preprint arXiv:2111.05424*, 2021.
- [37] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [38] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- [39] Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [40] Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith B Hall, Daniel Cer, and Yinfei Yang. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. *arXiv preprint arXiv:2108.08877*, 2021.
- [41] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [42] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya 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.

- [43] Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling transformer inference. *Proceedings of Machine Learning and Systems*, 5:606–624, 2023.
- [44] Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey, Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. Habitat 3.0: A co-habitat for humans, avatars and robots. *arXiv preprint arXiv:2310.13724*, 2023.
- [45] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- [46] Ram Ramrakhya, Dhruv Batra, Erik Wijmans, and Abhishek Das. Pirlnav: Pretraining with imitation and rl finetuning for objectnav. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17896–17906, 2023.
- [47] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- [48] Stefan Schaal, Auke Ijspeert, and Aude Billard. Computational approaches to motor learning by imitation. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 358(1431):537–547, 2003.
- [49] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- [50] Dhruv Shah, Ajay Sridhar, Nitish Dashora, Kyle Stachowicz, Kevin Black, Noriaki Hirose, and Sergey Levine. Vint: A foundation model for visual navigation. *arXiv preprint arXiv:2306.14846*, 2023.
- [51] Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- [52] Adrien Ali Taiga, Rishabh Agarwal, Jesse Farebrother, Aaron Courville, and Marc G Bellemare. Investigating multi-task pretraining and generalization in reinforcement learning. In *The Eleventh International Conference on Learning Representations*, 2023.
- [53] Chen Tang, Ben Abbatematteo, Jiaheng Hu, Rohan Chandra, Roberto Martín-Martín, and Peter Stone. Deep reinforcement learning for robotics: A survey of real-world successes. *arXiv preprint arXiv:2408.03539*, 2024.
- [54] Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(7), 2009.
- [55] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- [56] Marcel Torne, Anthony Simeonov, Zechu Li, April Chan, Tao Chen, Abhishek Gupta, and Pulkit Agrawal. Reconciling reality through simulation: A real-to-sim-to-real approach for robust manipulation. *arXiv preprint arXiv:2403.03949*, 2024.
- [57] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [58] Ikechukwu Uchendu, Ted Xiao, Yao Lu, Banghua Zhu, Mengyuan Yan, Joséphine Simon, Matthew Bennis, Chuyuan Fu, Cong Ma, Jiantao Jiao, et al. Jump-start reinforcement learning. In *International Conference on Machine Learning*, pages 34556–34583. PMLR, 2023.
- [59] Mel Vecerik, Todd Hester, Jonathan Scholz, Fumin Wang, Olivier Pietquin, Bilal Piot, Nicolas Heess, Thomas Rothörl, Thomas Lampe, and Martin Riedmiller. Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards. *arXiv preprint arXiv:1707.08817*, 2017.

- [60] Maciej Wołczyk, Bartłomiej Cupiał, Mateusz Ostaszewski, Michał Bortkiewicz, Michał Zając, Razvan Pascanu, Łukasz Kuciński, and Piotr Miłoś. Fine-tuning reinforcement learning models is secretly a forgetting mitigation problem, 2024. URL <https://arxiv.org/abs/2402.02868>.
- [61] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11097–11107, 2020.
- [62] Kuo-Hao Zeng, Zichen Zhang, Kiana Ehsani, Rose Hendrix, Jordi Salvador, Alvaro Herrasti, Ross Girshick, Aniruddha Kembhavi, and Luca Weihs. Poliformer: Scaling on-policy rl with transformers results in masterful navigators. *arXiv preprint arXiv:2406.20083*, 2024.
- [63] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE transactions on knowledge and data engineering*, 34(12):5586–5609, 2021.
- [64] Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE symposium series on computational intelligence (SSCI)*, pages 737–744. IEEE, 2020.
- [65] Henry Zhu, Justin Yu, Abhishek Gupta, Dhruv Shah, Kristian Hartikainen, Avi Singh, Vikash Kumar, and Sergey Levine. The ingredients of real-world robotic reinforcement learning. *arXiv preprint arXiv:2004.12570*, 2020.
- [66] Yuke Zhu, Ziyu Wang, Josh Merel, Andrei Rusu, Tom Erez, Serkan Cabi, Saran Tunyasuvunakool, János Kramár, Raia Hadsell, Nando de Freitas, et al. Reinforcement and imitation learning for diverse visuomotor skills. *arXiv preprint arXiv:1802.09564*, 2018.
- [67] Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, and Yifeng Zhu. robosuite: A modular simulation framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020.
- [68] Zhuangdi Zhu, Kaixiang Lin, Anil K Jain, and Jiayu Zhou. Transfer learning in deep reinforcement learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

## A APPENDIX

### A.1 Results Visualization

We encourage the reader to visit our website ([robot-flare.github.io](http://robot-flare.github.io)) for visualizations of trajectories generated by FLARe both in simulation and in the real world, including performances visualization, behavior analysis, and failure mode analysis.

### A.2 Hyperparameter

Training and Model Details	
Parameter	Value
Total Rollouts	32
Learning Rate	0.0002
Mini Batch per Update	1
Update Repeats	4
Max Gradient Norm	0.5
Discount Value Factor $\gamma$	0.99
GAE $\lambda$	0.95
PPO Surrogate Objective Clipping	0.1
Value Loss Weight	0.5
Entropy Loss Weight	0.0
Steps for PPO Update	128
Transformer State Encoder Layers	3
Transformer State Encoder Hidden Dims	512
Transformer State Encoder Heads	8
Causal Transformer Deocder Layers	3
Causal Transformer Deocder Hidden Dims	512
Causal Transformer Deocder Heads	8

Table 4: Hyperparameters for training and model architecture.

### A.3 Number of Training Steps

The base SPOC model that we evaluated and fine-tuned upon is trained for 50k gradient update steps on a total of 100k episodes of demonstrations across the CHORES tasks, where the training hyperparameter and training data is exactly the same as in the original SPOC paper.

For navigation tasks that do not involve manipulating objects (i.e. ObjectNav and RoomVisit), FLARe performs RL fine-tuning for 20M steps, while all other fair-comparison baseline methods perform RL training for 60M steps. For mobile manipulation tasks (i.e. Fetch and Pickup), FLARe performs RL fine-tuning for 50M steps, while all other fair-comparison baseline methods perform RL training for 100M steps. For adaptation tasks, we run FLARe fine-tuning for 50M steps on ObjNavRelAttr and ObjNavAfford, and 20M steps on RoomNav. For cross-embodiment, we run FLARe for 30M steps.

All of the aforementioned experiments use the same base SPOC mode, with exactly the same set of hyperparameters.

### A.4 CHORES Benchmark

A big portion of FLARe’s evaluation is carried out on the CHORES benchmark. We provided detailed information about this benchmark, including the observation space, action space, and task specifications.

#### A.4.1 Observation Space

The observation space of CHORES consists of two ego-centric  $384 \times 224$  RGB camera pointing towards orthogonal directions, where one points towards the direction of navigation and the other

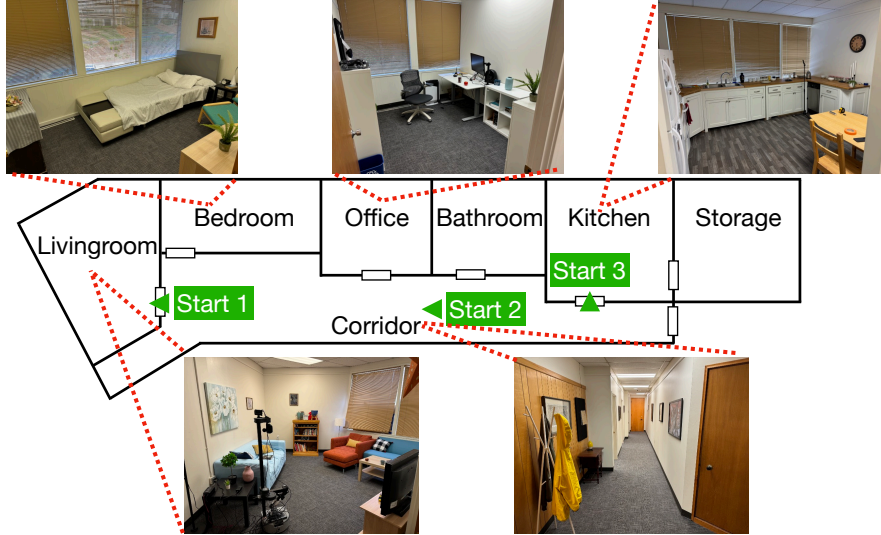


Figure 5: The real-world layout that we tested upon

Task	Description & Example	Max Steps
ObjectNav	Locate an object category: "find a mug"	600
PickUp	Pick up a specified object in agent line of sight: "pick up a mug"	600
Fetch	Find and pick up an object: "locate a mug and pick up that mug"	600
RoomVisit	Traverse the house. "Visit every room in this 5-room house. Indicate when you have seen a new room and when you are done."	1000

Table 5: CHORES tasks.

points at the arm. Additionally, a natural language text instruction is re-sampled at the start of each episode and appended to the observation to specify what the robot should be doing.

#### A.4.2 Action Space

The action space of CHORES consists of 20 discrete actions: Move Base ( $\pm 20$  cm), Rotate Base ( $\pm 6^\circ$ ,  $\pm 30^\circ$ ), Move Arm ( $x, z$ ) ( $\pm 2$  cm,  $\pm 10$  cm), Rotate Grasper ( $\pm 10^\circ$ ), pickup, dropoff, done with subtask, and terminate.

#### A.4.3 Tasks Specifications

We describe the CHORES tasks in Table 5. For each task, if the robot exceeds the maximum steps, the episode is terminated and marked as failed.

For each task, we split a total of 191,568 houses from ProcThor[9] into training and testing sets with a ratio of 10:1, to ensure that the test evaluation is conducted in unseen houses.

### A.5 The SPOC Model

In this work, we use a slightly modified version of the SPOC model [14] inspired by Poliformer [62], where the transformer decoder block in SPOC is replaced by the decoder from Llama 2 LLM [57] to speed up training and inference. At each step, the SPOC model takes in the new observations consisting of two RGB images and a text instruction. Each of these images are separately passed through a frozen **vision transformer model** (DinoV2[41]) to extract a set of visual tokens. These tokens, along with an embedding of the natural language instructions using a pre-train text encoder T5[40], are summarized by a **transformer state encoder** to produce the observation representation. A **causal transformer decoder** then decodes the observations feature across all steps within the current episode into a belief vector that is passed through an actor head to generate the action prediction. We provide a visualization of our model in Fig. 6, and explain each of these components in detail below.

### A.5.1 Vision Transformer Model

We use DINOv2 as the visual foundation backbone because of its remarkable ability to make dense predictions that generalize across sim and real. Our input to the visual backbone are two RGB observations  $i_a$  and  $i_b$ .  $i_a \in \mathbb{R}^{H \times W \times 3}$  is captured by the navigation camera and  $i_b \in \mathbb{R}^{H \times W \times 3}$  is captured by manipulation camera, where  $H$  and  $W$  are the height and width of the image. The visual backbone then produces a patch-wise representation  $r \in \mathbb{R}^{\frac{H}{14} \times \frac{W}{14} \times h}$ , where  $h$  is the hidden dimensions of the visual representations.  $r$  is then reshaped and projected to generate visual tokens  $v_{\text{raw}} \in \mathbb{R}^{n_{\text{patch}} \times d_{\text{encoder}}}$ . A learnable camera-type embedding is then added to this visual tokens to ensure the model can differentiate between the navigation and the manipulation cameras, resulting in the final visual features  $v$ . To ensure sim-to-real transfer, we freeze the DinoV2 weight throughout training.

### A.5.2 Transformer State Encoder

This module summarizes the observations at each timestep as a vector  $s \in \mathbb{R}^d$ . The input to this encoder includes the visual representation  $v$ , the text feature  $g$ , and a learnable STATE token  $f$ . We concatenate these features together and feed them to a non-causal transformer encoder. This encoder then returns the output corresponding to the STATE token as the state feature vector. The transformer state encoder digests both visual and text features, and can thus be seen as generating a text-conditioned visual state representation.

### A.5.3 Causal Transformer Decoder

To deal with partial observability and handle long-horizon tasks, SPOC uses a causal transformer decoder to perform explicit memory modeling over time. The causal transformer decoder consumes the visual representations generated by the transformer state encoder, additively combines them with sinusoidal temporal position encodings and learned previous time step action embeddings, and generates the belief vector used for action generation.

## A.6 Real Robot Setup

Following SPOC [14], we equipped our Stretch RE-1 robot with two identical Intel RealSense 455 fixed cameras, namely the navigation and the manipulation camera. These cameras have a vertical field of view of  $59^\circ$  and are capable of capturing  $1280 \times 720$  RGB-D images. Both of these cameras point slightly down, with the horizon at a nominal  $30^\circ$ , to optimize the agent’s perspective of its functional workspace. The images returned by these cameras are first resized to  $396 \times 224$ , and the cropped to  $384 \times 224$ , to match the image observations during training.

Same as SPOC, we assess the performance of our models on ObjectNav and Fetch in a 6-room apartment also used in Phone2Proc [10], Pickup in RoboThor [8], and RoomVisit in both environments. The 6-room apartment contains environment variations wholly unseen at train time, including a new configuration (multiple rooms off a long corridor), two new room types (office and corridor), rooms with non-orthogonal wall alignment, and many unseen object instances. For each object in ObjectNav and Fetch, we tested three starting positions: once from the living room, once from the middle of the corridor, and once from the kitchen. We visualize these starting locations in Fig. 5. Below, we provide objects that we tested upon in the real world for each tasks.

### A.6.1 ObjectNav

Target objects are Sofa, Bed, Chair, Apple, Vase, and Houseplant, each from three starting positions.

### A.6.2 Fetch

Target objects are Apple, Vase, and Houseplant from the same three starting positions. In one small change from ObjectNav episodes, object instances are replaced with instances which better fit into Stretch’s grasping envelope and in some cases at a better height for interaction, but availability and placement are nearly identical.

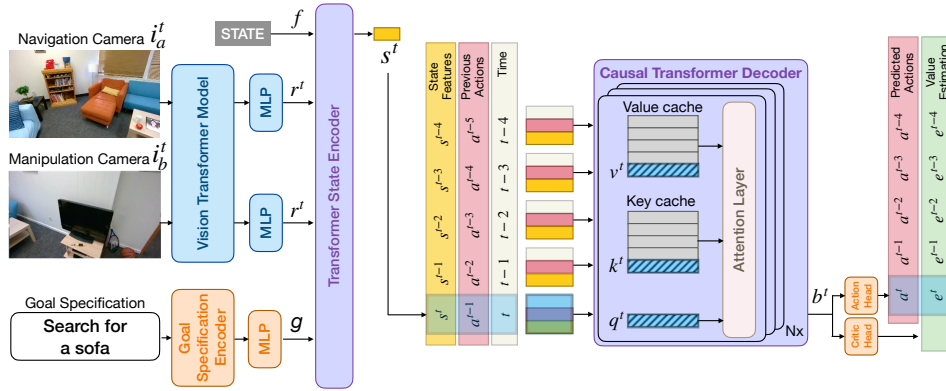


Figure 6: A visualization of the network architecture of the transformer-based SPOC model that FLaRe fine-tunes upon.

### A.6.3 PickUp

Objects are placed on three different surfaces (coffee table, desk, and nightstand) at three different heights. Objects are Apple, Houseplant, Spray Bottle, Mug, and Vase.

### A.6.4 RoomVisit

The full 6-room apartment is explored, and then partitioned into two 3-room apartments to evaluate the ability of SPOC to explore large and small spaces. We additionally explore a section of RoboTHOR and attached workroom as a novel 3-room apartment.