

Robot Fine-Tuning Made Easy: Pre-Training Rewards and Policies for Autonomous Real-World Reinforcement Learning

Jingyun Yang*, Max Sobol Mark*, Brandon Vu, Archit Sharma, Jeannette Bohg, Chelsea Finn

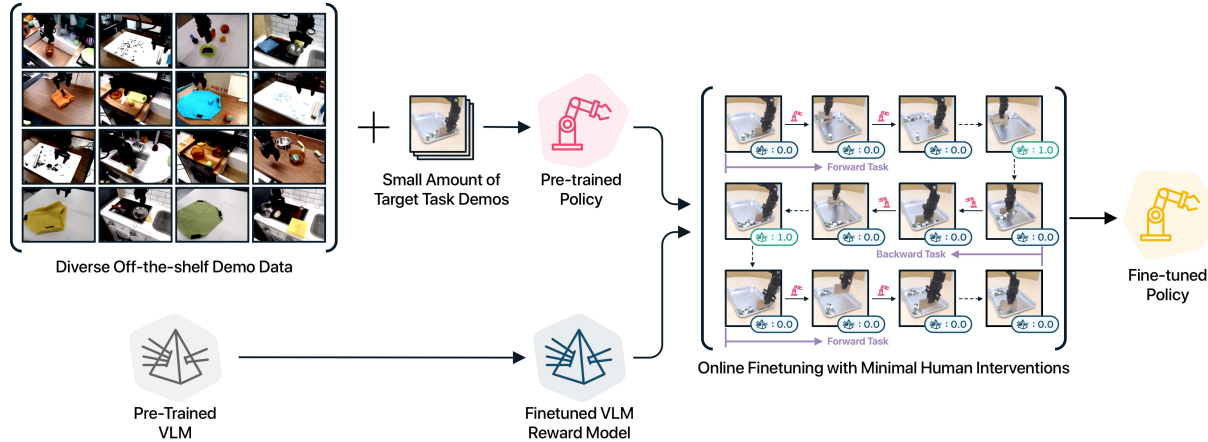


Fig. 1: We propose a system that enables autonomous and efficient real-world robot learning. First, we pre-train a **multi-task policy** and fine-tune a pre-trained **Vision-Language Model (VLM)** as a **reward model** using diverse off-the-shelf demonstration datasets and a small amount of target task demonstrations. Then, we fine-tune the **pre-trained policy** online reset-free with the **VLM reward model**.

Abstract—The pre-train and fine-tune paradigm in machine learning has had dramatic success in a wide range of domains because the use of existing data or pre-trained models on the internet enables quick and easy learning of new tasks. We aim to enable this paradigm in robotic reinforcement learning, allowing a robot to learn a new task with little human effort by leveraging data and models from the Internet. However, reinforcement learning often requires significant human effort in the form of manual reward specification or environment resets, even if the policy is pre-trained. We introduce ROBOFUME, a reset-free fine-tuning system that pre-trains a multi-task manipulation policy from diverse datasets of prior experiences and self-improves online to learn a target task with minimal human intervention. Our insights are to utilize calibrated offline reinforcement learning techniques to ensure efficient online fine-tuning of a pre-trained policy in the presence of distribution shifts and leverage pre-trained vision language models (VLMs) to build a robust reward classifier for autonomously providing reward signals during the online fine-tuning process. In a diverse set of five real robot manipulation tasks, we show that our method can incorporate data from an existing robot dataset collected at a different institution and improve on a target task within as little as 3 hours of autonomous real-world experience. We also demonstrate in simulation experiments that our method outperforms prior works that use different RL algorithms or different approaches for predicting rewards. Project website: <https://robofume.github.io>

I. INTRODUCTION

In many domains that involve machine learning, a widely successful paradigm for learning task-specific models is to first pre-train a general-purpose model from an existing diverse prior dataset, and then adapt the model with a

small addition of task-specific data [1]–[5]. This paradigm is attractive to real-world robot learning, since collecting data on a robot is expensive, and fine-tuning an existing model on a small task-specific dataset could substantially improve the data efficiency for learning a new task. Pre-training a policy with offline reinforcement learning and then fine-tuning it with online reinforcement learning is a natural way to implement this paradigm in robotics. However, numerous challenges arise when using this recipe in practice. First, off-the-shelf robot datasets often use different objects, fixture placements, camera viewpoints, and lighting conditions compared to the local robot platform. This causes non-trivial distribution shifts between pre-training and online fine-tuning data, which makes effectively fine-tuning a robot policy difficult. Indeed, most existing works only show the benefit of the pre-train and fine-tune paradigm where the robot uses the same hardware instance in both pre-training and fine-tuning phases [6], [7]. Second, training or fine-tuning a policy in the real world often requires extensive human supervision, which includes manually resetting the environment between trials [8]–[10] and engineering reward functions [7], [11], [12]. In this work, our goal is to address these two challenges and develop a practical framework that enables robot fine-tuning with minimal time and human effort.

Over the past few years, there has been a lot of progress in designing efficient and autonomous reinforcement learning algorithms. However, no existing system could both utilize diverse demonstration datasets and learn with minimal human supervision, without the need for human-engineered reward functions and manual environment resets. Some works propose to reduce the need for manual environment

*Equal contribution. All authors are affiliated with Department of Computer Science, Stanford University. This work was supported by Toyota Research Institute and ONR grant N00014-20-1-2675. Contact: jingyuny@stanford.edu, maxsobolmark@stanford.edu

resets using reset-free RL [7], [11], [13], where an agent alternates between running a task policy and a reset policy during training while updating both with online experience. However, these works do not leverage diverse off-the-shelf robot datasets. Recent advances in offline RL algorithms have enabled policies to leverage diverse offline data and improve further via online fine-tuning [14], [15], but these new methods have not been integrated into a system that aims at minimizing human supervision in the fine-tuning phase. There are also works that propose to eliminate the need for human-specified reward functions by learning reward prediction models [13], [16]–[18], but we found that many of these proposed models can be brittle when deployed in a real-world RL fine-tuning setup. In summary, although prior works have presented individual components that are vital to building a working system for efficient and human-free robot learning, it is not clear which components one should use to put together such a system and how.

We design ROBOFUME, a system that enables autonomous and efficient real-world robot learning by leveraging diverse offline datasets and online fine-tuning. Our system operates in two phases. In the pre-training phase, we assume access to a diverse prior dataset, a few task demonstrations and reset demonstrations of the target task, and a small collection of sample failure observations in the target task. From this data, we learn a language-conditioned, multi-task policy with offline RL. To cope with the distribution shift between the offline dataset and online interactions, we need an algorithm that could effectively digest diverse offline data, and display robust fine-tuning performance when placed into an environment different from those seen in the offline dataset. We find that calibrated offline RL techniques [15], by underestimating predicted values of the learned policy from offline data and correcting the scale of the learned Q-values, make sure that the pre-trained policy can effectively digest diverse offline data and continuously improve during online adaptation. To ensure the online fine-tuning phase requires minimal human feedback, we need to remove the need for reward engineering by learning a reward predictor. Our insight is to take a large vision-language model (VLM) to provide a robust pre-trained representation and fine-tune it with a small amount of in-domain data so that it is tailored for the reward classification setup. Pre-trained VLMs have already been trained on internet-scale visual and language data. This makes the model more robust to lighting and camera positioning variations than the models used in prior works. In the fine-tuning phase, a robot adapts the policy in the real world autonomously by alternating between attempting the task and attempting to reset the environment to the initial state distribution of the task. Meanwhile, the agent uses the pre-trained VLM model as a surrogate reward for updating the policy.

We evaluate our framework by pre-training it on the Bridge dataset [19] and testing it on a diverse set of real-world downstream tasks: cloth folding, cloth covering, sponge pick-and-place, placing lid on a pot, and putting a pot in a sink. We find that our system provides substantial improvements

over offline-only methods with as little as 3 hours of real-world training. We perform more quantitative experiments in a simulation setup, where we illustrate that our method outperforms imitation learning and offline RL methods that either do not perform fine-tuning online or do not incorporate diverse prior data.

Our main contributions include (1) a fully autonomous system for pre-training from a prior robot dataset and fine-tuning on an unseen downstream task with a minimal number of resets and learned reward labels; (2) a method for fine-tuning pre-trained vision-language models and using them to construct a surrogate reward for downstream RL training.

II. RELATED WORK

Offline RL. Offline RL algorithms [20]–[25] provide a framework for initializing robot manipulation policies from offline demonstrations or interaction datasets. Such algorithms can also be extended to include an online fine-tuning phase after training a policy offline [15], [26]–[32]. Our work utilizes a recent offline RL algorithm, calibrated Q-learning (CalQL) [15], a state-of-the-art method that effectively learns from offline data and continuously improves the policy’s performance online by explicitly correcting the scale of the learned Q-values. We show that integrating CalQL helps our framework effectively utilize diverse prior datasets that have large distribution shifts from real-world online interactions.

Reset-free RL. Training an RL policy on a real robot typically requires manual environment resets. To eliminate such need to manually reset environments, prior works have studied approaches to learn robot policies in a ‘reset-free’ setup. Some work [11], [33]–[36] cast the ‘reset-free’ learning problem as a multi-task learning problem, observing that by learning a set of tasks where some of the tasks could reset others, an agent could then be trained to perform all of those tasks without needing manual resets. Other works [7], [12], [13], [17], [37]–[39] learn both a task policy and a reset policy for performing the task and resetting to the initial state distribution. Our work takes an approach between the two classes of approaches, learning a language-conditioned multi-task policy that can perform both the target task and the reset for the target task. Most of these prior works learn from scratch rather than incorporating prior data and assume that a reward function is available. ARIEL [7] combines incorporating prior data with reset-free learning but assumes a hand-crafted reward function for each environment. They also collect their own prior dataset on the same robot hardware set-up as their target task. MEDAL++ [13] learns a reward classifier with demonstration and online interaction data via adversarial training, but does not consider incorporating diverse prior data. Leveraging diverse, off-the-shelf prior demonstration datasets is desirable since these datasets are readily available to use and can help a system obtain a policy initialization for efficient fine-tuning on a target task. Our system offers an approach to both incorporate diverse prior data and improve the autonomy of the fine-tuning phase by learning a model for predicting rewards. In particular, we found out that by leveraging diverse demonstration data, our

system requires only about 3 hours of training in the real world compared to 10-30 hours in MEDAL++.

Reward learning. Early works have studied the problem of learning a reward or cost function in imitation learning. These works leverage inverse optimal control (IOC) or inverse reinforcement learning (IRL) to extract a reward function directly from expert demonstrations [40]–[42]. With the advent of deep neural networks, more recent works have explored learning a reward model for an imitation learning or RL policy [13], [17], [43]–[47]. When using classifier-based reward models in reinforcement learning, RL agents can exploit the learned model by exploring states unseen during classifier training, tricking the model to output incorrect rewards. To solve such an exploitation issue, many works that learn reward models leverage adversarial learning, where a system learns a discriminator that identifies states similar to those in demonstrations as positives and those visited by the policy as negatives [13], [17], [44], [47]. However, prior work has found this training objective to be sensitive to distribution shifts between offline and online setups, such as lighting and camera view changes [48]. In this work, we fine-tune vision language models (VLM), pre-trained on internet-scale data, to construct a reward model. Large scale pre-training can learn representations that are robust to natural variations such as lighting, camera shifts and distractors [16], [49].

Leveraging pre-trained representations as reward predictors. Several recent works have shown positive results in utilizing pre-trained vision models [16], [50], large language models (LLMs) [51] or vision language models (VLMs) [52] as reward predictors. We tried VIP [16], a method that pre-trains a visual representation for generating dense reward functions for novel robotic tasks, and found it insufficient for the real-world robot fine-tuning setup. In this work, we fine-tune a pre-trained VLM [53] and find that it performs most effectively as a reward model. Our proposed system is flexible and can easily be adapted to use other pre-trained visual representations and VLMs.

III. PRELIMINARIES

The goal of our method is to leverage diverse prior demonstration datasets and learn a novel target task autonomously in a robot hardware instance that is distinct from the one used to collect the datasets. Our method assumes access to a prior dataset $\mathcal{D}_{\text{prior}} = \cup_{j=1}^N \mathcal{D}_j = \cup_{j=1}^N \{(s_i^j, a_i^j, s_i'^j)\}_{i=1}^K$, which consists of demonstrations of N different tasks τ_1, \dots, τ_k . We assume that all demonstration data uses image observations. The method will be tested on a downstream task τ_f , which is different from any of the prior tasks.

To facilitate learning on the downstream task, we also assume the availability of a small set of target task demos \mathcal{D}_f , target task reset demos \mathcal{D}_b , and target task failure states \mathcal{D}_{\odot} . The reset demos \mathcal{D}_b come from the reset task τ_b which resets the environment from an end state of τ_f to the initial state distribution of τ_f . The failure states \mathcal{D}_{\odot} consist entirely of image observations that correspond to unsuccessful states and are collected to aid with the VLM reward learning. In

addition to all the given data ($\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\odot}$), each task τ is also accompanied with a language description l .

IV. ROBOFUME

Our work focuses on designing an efficient and scalable framework for pre-training on a diverse set of prior demonstrations and autonomously fine-tuning on target tasks. Our system consists of an offline pre-training phase and an online fine-tuning phase. In Section IV-A, we discuss how we pre-train a language-conditioned multi-task policy on diverse data that can be fine-tuned online efficiently. Online fine-tuning requires a reward function to label successes and failures. In Section IV-B, we introduce a VLM-based classifier for providing a reward signal to the policy in the fine-tuning phase. Finally, in Section IV-C, we describe how to autonomously adapt the pre-trained policy in the fine-tuning phase by utilizing the VLM-based reward classifier as a reward signal and chaining forward and backward behaviors to practice the task with minimal human interventions.

A. Pre-Training a Multi-Task Policy on Diverse Prior Data

Prior work has shown that training a policy using a conservative Q-value function is an effective way to obtain a good policy from an offline dataset [22], [24]. However, fine-tuning can be critical to learn competent policies as prior data may not provide sufficient coverage, especially for new tasks or scenes. We leverage CalQL [15] which modifies the conservative Q-learning algorithm CQL such that it enables efficient online fine-tuning by enforcing calibration on the Q-function (i.e. making the Q-value of the learned policy no lower than the Monte-Carlo returns in the prior dataset). CalQL allows us to improve the pre-trained policy efficiently with respect to online interactions.

CalQL requires the training of an actor and a critic. Since we use image observations, we additionally train an encoder $\phi(s_{\text{img}})$ that projects the images into a lower-dimensional space before giving them as inputs to the actor and critic. The encoder ϕ is a 4-layer CNN, and is optimized exclusively against the critic loss. To best utilize the multi-task data, we encode task descriptions l using pre-trained CLIP embeddings, resulting in an embedding $z = \text{CLIP}(l)$ which is used as the task representation. The policy then takes as inputs a concatenation of the encoded image observation $\phi(s_{\text{img}})$, task representation z , and proprioceptive information s_p , processes the concatenated vector through an MLP, and produces the output action a .

In addition to updating the policy using CalQL, we regularize policy learning with a behavior cloning (BC) loss, which encourages the behaviors to stay close to the seen demonstrations. Not only does this regularization improve performance of the offline pre-training, but we find that it also makes it less likely for the autonomous fine-tuning procedure to exploit false positive rewards from the VLM reward model. The weight of the BC regularization term is chosen such that the scales of the RL loss and the BC loss are similar throughout the pre-training phase. We train the policy π and the critic Q with datasets $\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b$.

After the offline learning phase, the policy and critic contain knowledge of all tasks in the prior data and the target task.

B. Fine-Tuning A Vision-Language Model for Rewards

To improve the autonomy of the policy fine-tuning phase, our agent needs to perform online fine-tuning without manually labeled or engineered reward functions. To achieve this, we propose to fine-tune off-the-shelf vision-language models as reward predictors. Leveraging existing vision-language models offers a number of benefits compared to utilizing a pre-trained visual representation or training a reward model from scratch using in-domain data: First, VLMs are trained on an Internet-scale dataset that contains diverse image and language contents. Such models possess better inductive biases and thus, can be more robust to natural shifts, such as perturbations to lighting conditions, or distractor objects that might be seen at test time. Second, since VLMs can take both visual and language information as input, they provide a natural interface for communicating the current observation and current task to the model when requesting a reward label.

We design a VLM-based reward model that takes the current observation and the task name as input and outputs a binary label of whether the current observation corresponds to a successful state or an unsuccessful state with respect to the task. Given a task name (eg. ‘put green cabbage into sink’), we first use GPT4 to convert the name to a short question that could serve as a prompt to know if the task has been completed or not (eg. ‘is green cabbage placed in the sink?’). Then, we pass the converted prompt to a VLM together with the current image of the environment. The VLM outputs a sparse binary reward, returning success if the ‘yes’ token has a higher probability than ‘no’ token.

We use MiniGPT4 [53] as the VLM for receiving (image, task prompt) pairs and answering whether the task has been successfully completed. We find the zero-shot performance of the pre-trained VLM to be unsatisfactory. To improve the VLM’s performance for reward modeling, we fine-tune it using the prior and target task data. In particular, for every demonstration, the last 3 states are used as success states and the ground truth answer is labeled as ‘Yes’; for all other states, we label the ground-truth answer as ‘No’. To provide the model with more information about failed states, we collect a small dataset \mathcal{D}_{\ominus} of images that correspond to unsuccessful states for the forward and backward target tasks. We fine-tune the VLM for 20 epochs using default fine-tuning hyperparameters provided by MiniGPT4. We find that fine-tuning leads to a more accurate reward model and is crucial for obtaining good policy learning performance.

C. Autonomous Online Fine-tuning

The offline pre-training phase produces a single language-conditioned policy $\pi(\cdot|s, l)$ that can perform the target and reset tasks when provided their respective language instructions l_f and l_b . The policy is then deployed in a hardware setup for further online fine-tuning. The outline of our pre-training and fine-tuning pipeline is presented in Algorithm 1.

Algorithm 1: RoboFuME

```

Initialize agent  $\mathcal{A} = \{\phi, \pi, Q\}$  and pre-trained VLM  $\hat{r}$ .
Initialize forward and backward tasks  $\tau_f, \tau_b$ .
// Prepare data and train VLM reward classifier.
 $\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\ominus} \leftarrow \text{load\_data}()$ .
 $\hat{r} \leftarrow \text{finetune\_vlm}(\hat{r}, \{\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\ominus}\})$ .
// Offline pre-training phase.
 $\mathcal{A}.\text{update\_buffer}(\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b)$ .
for  $t = 1$  to  $T_{\text{offline}}$  do
     $\mathcal{A}.\text{update\_params\_with\_calql}()$ .
// Online fine-tuning phase.
 $s \leftarrow \text{env}.\text{reset}(); l \leftarrow \tau_f.\text{get\_task\_lang}()$ .
for  $t = 1$  to  $T_{\text{online}}$  do
     $a \leftarrow \mathcal{A}.\text{act}(s, l); s' \leftarrow \text{env}.\text{step}(a)$ .
     $\mathcal{A}.\text{update\_buffer}(\{s, a, s', \hat{r}(s)\})$ .
    for  $i = 1$  to  $N_{\text{update\_ratio}}$  do
         $\mathcal{A}.\text{update\_params\_with\_calql}()$ .
    if switch then
        // Switch task after a fixed interval.
         $l \leftarrow \text{env}.\text{switch}(\tau_f, \tau_b).\text{get\_task\_lang}()$ .
    if interrupt then
        // Allow occasional human intervention.
         $s \leftarrow \text{env}.\text{reset}(); l \leftarrow \tau_f.\text{get\_task\_lang}()$ .
    else
         $s \leftarrow s'$ .

```

Since we aim for a fully autonomous setup, we roll out the policy in a reset-free manner, alternating between attempting the target task τ_f with $\pi(\cdot|s, l_f)$ and the reset task τ_b with $\pi(\cdot|s, l_b)$. We use the fine-tuned VLM from the previous subsection as the sparse reward function for the RL algorithm. When the VLM predicts the task has been completed successfully, we terminate the episode and switch the language instruction for the policy to complete the other task. In addition to switching tasks upon completion as predicted by the VLM, we switch after a fixed number of timesteps (150) to ensure the robot does not become stuck in bad states. As mentioned in Section IV-A, we fine-tune the agent using CalQL with an additional BC regularization term on the policy. We find that without a BC regularization term, behaviors degrade during training. By constraining the policy to stay close to the expert demos from the target and reset tasks, the agent becomes less likely to exploit false positives from the VLM reward model. We use the same fixed BC regularization weight throughout fine-tuning as we did in the offline pre-training phase. Using this regularization assumes that the agent could find an optimal policy close to the expert demo distribution. Exploring efficient fine-tuning without any BC regularization, which allows the agent to explore more widely without degrading performance, is a valuable direction for future work. Our fine-tuning pipeline is implemented on top of the implementation of MEDAL++ [13]. Please refer to this work for more details on our training procedure.

V. EXPERIMENTS

We design our experiments to answer the following questions: Is our method able to improve its performance through near autonomous online interactions? How does our proposed VLM reward function mechanism compare

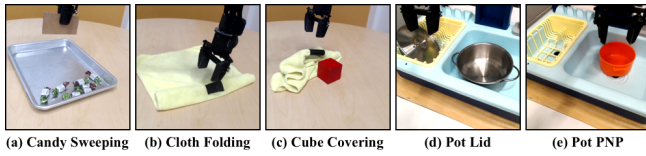


Fig. 2: **Illustrations of the five real-world evaluation tasks.** (a) Sweep candies to the top of the tray. (b) fold the yellow cloth. (c) cover a red wooden cube using the cloth. (d) place the lid on top of the metallic pot. (e) move the orange pot from the sink to the drying rack.

to existing alternatives? And, how does each component of ROBOFUME or data affect the performance of our method?

A. Real Robot Experiments

Setup. We evaluate ROBOFUME on five different real-robot manipulation tasks. We use a WidowX 250 robotic arm with a single third-person camera (Logitech C920, resizing images to 100x100 pixels). Figure 2 shows the five tasks we fine-tune and evaluate on. Our method runs autonomously executing back and forth the target task and the reset target task for a fixed number of steps or until the VLM predicts success. For tasks involving deformable objects (the two cloth tasks) we manually reset the object to the initial forward pose every 15-25 episodes, and for the rest of the tasks we reset every 30-35 episodes. Tasks that use the kitchen-sink environment (*pot lid* and *pot pnp*) frequently experience episode interruptions when the robot arm applies more than the maximum allowed torque, for example, when close to the sink borders. All tasks use 50 forward and 50 backward demos for the target task, and fewer than 20 combined trajectories of failures. We use demos from the BridgeDataV2 [19], [54] for pre-training our language-conditioned policy, selecting approximately 1,000 trajectories with relevant behaviors per task.

Results. Table I shows the results of our method after pretraining (labeled OFFLINE) and after autonomous fine-tuning (labeled FT 30K STEPS), comparing with a behavior cloning (BC) baseline. BC trains a language-conditioned policy on all the prior and target data. After 30k steps of autonomous online interaction, our method shows relative improvement of 51% upon the pre-trained performance, and outperforms BC by 58% on an average. For pick and place tasks (*pot lid* and *pot pnp*), the fine-tuned policy was more likely to retry the action if it initially failed to grasp the object. For candy sweeping, BC and the pre-trained policy were prone to overshooting and pushing on the border of the tray after the first sweep, whereas fine-tuning the policy enabled the policy to chain multiple sweeping attempts for higher success. Additionally, we find that policies learned by ROBOFUME (both offline and after fine-tuning) to be more robust to scene distractors on the *candy sweeping* task, as reported in Table III. The policies were trained without any distractors, but multiple objects not seen during training were placed in the background during evaluation. ROBOFUME policies retained 68% of its original performance, compared to BC which retained only 10% of its original performance. We hypothesize that BC might be more sensitive to spurious features, whereas ROBOFUME learns from more predictive features, leading to more robust policies.

Task	BC	ROBOFUME (OFFLINE)	ROBOFUME (FT 30K STEPS)
Cloth Covering	45%	60%	80%
Cloth Folding	60%	70%	85%
Candy Sweeping	31%	47%	66%
Pot Lid	60%	40%	95%
Pot PNP	45%	35%	55%

TABLE I: **Real-robot results on 5 manipulation tasks.** Our method significantly improves over both offline-only and BC performance after 30k steps of online interaction (2-4 hours). For the Candy sweeping we report the average percentage of candies out of a total of 7 that are placed in the top third of the tray by the end of the evaluation. For all other tasks, we report success rate over 20 trials.

B. Simulation Experiments and Ablations

We use a suite of simulated robotic manipulation environments to ablate contributions of different components of our algorithm. We test on three simulated environments used in [6]. We consider three bin-sorting tasks in which different objects (a vase, a tiny bench, and a dumbbell weight) have to be placed on the correct bin based on the language instruction, given only a sparse binary reward. We provide 10 forward and reset demonstrations for each task, 30 failure demos, and 10 demos each for 20 prior tasks that show picking and placing diverse objects on the same environment. For all methods that require online experience, we reset the environment every 1,000 environment steps, i.e. every 25 episodes of interactions. We compare our method against the following baselines: (1) *BC* behavior clones on all prior and target data; (2) *MEDAL++* learns separate forward and backward policies from target forward and backward task demonstrations and performs reset-free learning using an adversarially trained classifier as a reward signal; (3) *MEDAL++ with prior data* modifies *MEDAL++* to a single language-conditioned multi-task policy and adds all prior demonstration data into the replay buffer; (4) *ARIEL+VLM* modifies *ARIEL* [7] to use our VLM reward models as reward signal, instead of a handcrafted ground-truth reward. The results of our simulation experiments are presented in Figure 3. In all simulation tasks, our method ROBOFUME consistently outperforms prior methods, achieving success rates at least 20% higher than all baselines within 200k steps of online fine-tuning.

Ablations on RL Algorithm Design Choices. We evaluate our method trained with different critic and actor optimization procedures on the Vase simulated task, shown in Figure 4. Training with CalQL was the only method that yielded strong improvements in this task, with the other methods either failing completely or obtaining very poor performance. We find that training without the CalQL stabilizes training, while the losses for other methods would explode given the limited data.

Ablations on Reward Models. We compare our VLM reward function against other choices of automatic reward functions on the Vase simulated task in Figure 5. *VICE* [47] adversarially trains a binary classifier using positive samples from successful demonstrations, and labeling online experience as negative. We find that offline pre-training sufficiently limits the exploitation of the frozen VLM re-

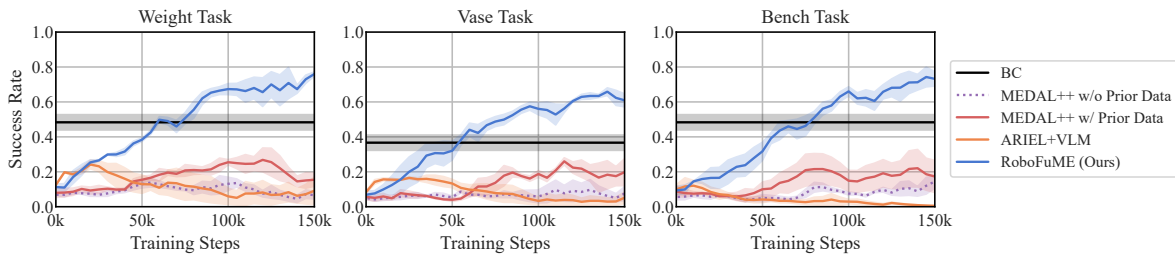


Fig. 3: **Performance of our method on three simulated environments.** We report the success rate over the course of training, averaged over three seeds. Our method ROBOFUME outperforms BC, ARIEL+VLM [7], and MEDAL++ [13] consistently on all three domains.

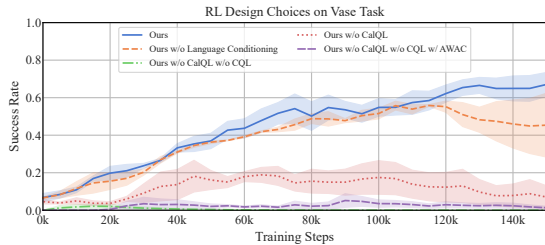


Fig. 4: **Performance of our method on the Vase simulated task with different actor-critic update objectives.** Fine-tuning with CalQL is critical to obtain stable improvements on this task, as training with CQL, AWAC, or SAC yields poor performance. We also find that language conditioned policies perform slightly better than one-hot task IDs in simulation.

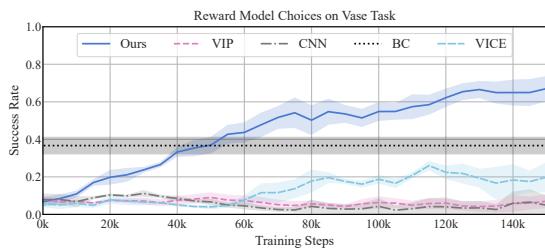


Fig. 5: **Performance of our method on the simulated Vase task using different reward functions.** Our method uses a fine-tuned VLM reward function and outperforms VICE rewards, whereas CNN and VIP rewards fail to improve online.

ward, outperforming VICE and thus, bypassing the need for adversarially trained reward functions. Such adversarial training can often learn to discriminate based on spurious shifts in the real world, such as lighting or scene changes, leading to instability in training outside simulation. VIP [16] trains a representation function such that the distance in representation space between the current observation and a goal image can be used to construct a dense reward function. We find that in the Vase simulated task, VIP fails to obtain good behaviors. Qualitatively, we observe VIP to be prone to false positives, which are exploited by the RL algorithm. To test the importance of the VLM large-scale pre-training compared to our fine-tuning procedure, we train a CNN classifier from scratch using the same data as we used to fine-tune the VLM, leading to unsatisfactory performance compared to fine-tuning a VLM.

How Accurate is the VLM Reward? We analyze the performance of the VLM reward over the course of fine-tuning for real-robot experiments. In Table II, we report the false positive rate, false negative rate, accuracy, and precision metrics for the VLM reward. The metrics are computed on the data collected during fine-tuning against a hand-engineered ground truth reward. We observe that while

Task	FP	FN	Accuracy	Precision
Cloth Covering	6.3%	80.9%	89.4%	15.3%
Cloth Folding	1.2%	59.8%	84.1%	92.0%
Pot PNP	6.1%	81.3%	86.9%	24.3%

TABLE II: **VLM reward model accuracy during real robot fine-tuning.** The low false positive (FP) rate indicates that online training has minimal reward exploitation.

Task	BC	ROBOFUME (offline)	ROBOFUME (fine-tuned @30k)
Candy Sweeping	31% → 3%	47% → 31%	66% → 45%

TABLE III: **Robustness of learned policy to distractors.** Entries in this table show the performance of the learned policy “before” → “after” adding distractors to the scene in the candy-sweeping task. Our system learns a policy that is much more robust to the distractors.

Task	ROBOFUME (offline)	ROBOFUME w/o Prior Data	ROBOFUME w/o Language Cond.
Candy Sweeping	47%	23%	13%

TABLE IV: **Evaluating effectiveness of prior data and language conditioned policies.** Results show that using prior data and using language conditioning positively affected the offline performance of our system.

false negative rates are high, false positive rates are low across all tasks. This asymmetry is crucial for successful RL fine-tuning, as RL policies can learn poor behaviors by exploiting false positives, but labeling some successful rollouts as negatives does not necessarily impede learning.

How Important is Diverse Prior Data and Language Conditioning? We ablate the contribution of diverse prior data and language-conditioned policies to ROBOFUME by evaluating the offline performance on the *candy sweeping* task, reported in Table IV. When pre-training without using prior data, that is, exclusively using target data, our method is able to sweep less than half the amount of candies on average. Similarly, we find that one-hot task encodings perform substantially worse than language-conditioned policies.

VI. CONCLUSION AND FUTURE WORK

We introduced an autonomous framework that leverages existing diverse prior robot demonstration datasets and improves performance in a new robot manipulation skill by fine-tuning online. By combining state-of-the-art offline-to-online RL algorithms, reset-free RL, and VLM-based reward models, our framework can fine-tune efficiently and nearly autonomously. Integrating this work with more recent VLM models and improving the reset efficiency of this framework are promising directions for future research.

REFERENCES

- [1] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, *et al.*, “Improving language understanding by generative pre-training,” 2018.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [3] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [4] 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*, vol. 33, pp. 1877–1901, 2020.
- [5] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 000–16 009.
- [6] A. Kumar, A. Singh, F. Ebert, Y. Yang, C. Finn, and S. Levine, “Pre-training for robots: Offline rl enables learning new tasks from a handful of trials,” *arXiv preprint arXiv:2210.05178*, 2022.
- [7] H. R. Walke, J. H. Yang, A. Yu, A. Kumar, J. Orbik, A. Singh, and S. Levine, “Don’t start from scratch: Leveraging prior data to automate robotic reinforcement learning,” in *Conference on Robot Learning*. PMLR, 2023, pp. 1652–1662.
- [8] V. Kumar, E. Todorov, and S. Levine, “Optimal control with learned local models: Application to dexterous manipulation,” in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 378–383.
- [9] A. Ghadirzadeh, A. Maki, D. Kragic, and M. Björkman, “Deep predictive policy training using reinforcement learning,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 2351–2358.
- [10] K. Ploeger, M. Lutter, and J. Peters, “High acceleration reinforcement learning for real-world juggling with binary rewards,” in *Conference on Robot Learning*. PMLR, 2021, pp. 642–653.
- [11] A. Gupta, J. Yu, T. Z. Zhao, V. Kumar, A. Rovinsky, K. Xu, T. Devlin, and S. Levine, “Reset-free reinforcement learning via multi-task learning: Learning dexterous manipulation behaviors without human intervention,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 6664–6671.
- [12] C. Sun, J. Orbik, C. M. Devin, B. H. Yang, A. Gupta, G. Berseth, and S. Levine, “Fully autonomous real-world reinforcement learning with applications to mobile manipulation,” in *Conference on Robot Learning*. PMLR, 2022, pp. 308–319.
- [13] A. Sharma, A. M. Ahmed, R. Ahmad, and C. Finn, “Self-improving robots: End-to-end autonomous visuomotor reinforcement learning,” *arXiv preprint arXiv:2303.01488*, 2023.
- [14] I. Kostrikov, A. Nair, and S. Levine, “Offline reinforcement learning with implicit q-learning,” *arXiv preprint arXiv:2110.06169*, 2021.
- [15] M. Nakamoto, Y. Zhai, A. Singh, M. S. Mark, Y. Ma, C. Finn, A. Kumar, and S. Levine, “Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning,” *arXiv preprint arXiv:2303.05479*, 2023.
- [16] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang, “Vip: Towards universal visual reward and representation via value-implicit pre-training,” *arXiv preprint arXiv:2210.00030*, 2022.
- [17] A. Sharma, R. Ahmad, and C. Finn, “A state-distribution matching approach to non-episodic reinforcement learning,” *arXiv preprint arXiv:2205.05212*, 2022.
- [18] Y. J. Ma, W. Liang, V. Som, V. Kumar, A. Zhang, O. Bastani, and D. Jayaraman, “Liv: Language-image representations and rewards for robotic control,” *arXiv preprint arXiv:2306.00958*, 2023.
- [19] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and S. Levine, “Bridge data: Boosting generalization of robotic skills with cross-domain datasets,” *arXiv preprint arXiv:2109.13396*, 2021.
- [20] Y. Wu, G. Tucker, and O. Nachum, “Behavior regularized offline reinforcement learning,” *arXiv preprint arXiv:1911.11361*, 2019.
- [21] X. B. Peng, A. Kumar, G. Zhang, and S. Levine, “Advantage-weighted regression: Simple and scalable off-policy reinforcement learning,” *arXiv preprint arXiv:1910.00177*, 2019.
- [22] S. Levine, A. Kumar, G. Tucker, and J. Fu, “Offline reinforcement learning: Tutorial, review, and perspectives on open problems,” *arXiv preprint arXiv:2005.01643*, 2020.
- [23] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma, “Mopo: Model-based offline policy optimization,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 14 129–14 142, 2020.
- [24] A. Kumar, A. Zhou, G. Tucker, and S. Levine, “Conservative q-learning for offline reinforcement learning,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 1179–1191, 2020.
- [25] W. Zhou, S. Bajracharya, and D. Held, “Plas: Latent action space for offline reinforcement learning,” in *Conference on Robot Learning*. PMLR, 2021, pp. 1719–1735.
- [26] N. Ashvin, D. Murtaza, G. Abhishek, and L. Sergey, “Accelerating online reinforcement learning with offline datasets,” *CoRR*, vol. abs/2006.09359, 2020.
- [27] R. Rafailov, T. Yu, A. Rajeswaran, and C. Finn, “Offline reinforcement learning from images with latent space models,” in *Learning for Dynamics and Control*. PMLR, 2021, pp. 1154–1168.
- [28] J. Lyu, X. Ma, X. Li, and Z. Lu, “Mildly conservative q-learning for offline reinforcement learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 1711–1724, 2022.
- [29] A. Beeson and G. Montana, “Improving td3-bc: Relaxed policy constraint for offline learning and stable online fine-tuning,” *arXiv preprint arXiv:2211.11802*, 2022.
- [30] J. Wu, H. Wu, Z. Qiu, J. Wang, and M. Long, “Supported policy optimization for offline reinforcement learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 31 278–31 291, 2022.
- [31] S. Lee, Y. Seo, K. Lee, P. Abbeel, and J. Shin, “Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble,” in *Conference on Robot Learning*. PMLR, 2022, pp. 1702–1712.
- [32] M. S. Mark, A. Ghadirzadeh, X. Chen, and C. Finn, “Fine-tuning offline policies with optimistic action selection,” in *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.
- [33] K. Lu, A. Grover, P. Abbeel, and I. Mordatch, “Reset-free lifelong learning with skill-space planning,” *arXiv preprint arXiv:2012.03548*, 2020.
- [34] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan, “Learning to walk in the real world with minimal human effort,” *arXiv preprint arXiv:2002.08550*, 2020.
- [35] A. Gupta, C. Lynch, B. Kinman, G. Peake, S. Levine, and K. Hausman, “Demonstration-bootstrapped autonomous practicing via multi-task reinforcement learning,” *arXiv preprint arXiv:2203.15755*, vol. 1, 2022.
- [36] K. Xu, Z. Hu, R. Doshi, A. Rovinsky, V. Kumar, A. Gupta, and S. Levine, “Dexterous manipulation from images: Autonomous real-world rl via substep guidance,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 5938–5945.
- [37] B. Eysenbach, S. Gu, J. Ibarz, and S. Levine, “Leave no trace: Learning to reset for safe and autonomous reinforcement learning,” *arXiv preprint arXiv:1711.06782*, 2017.
- [38] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, V. Kumar, and S. Levine, “The ingredients of real-world robotic reinforcement learning,” *arXiv preprint arXiv:2004.12570*, 2020.
- [39] A. Sharma, A. Gupta, S. Levine, K. Hausman, and C. Finn, “Autonomous reinforcement learning via subgoal curricula,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 18 474–18 486, 2021.
- [40] A. Y. Ng, S. Russell, *et al.*, “Algorithms for inverse reinforcement learning,” in *Icml*, vol. 1, 2000, p. 2.
- [41] P. Abbeel and A. Y. Ng, “Apprenticeship learning via inverse reinforcement learning,” in *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 1.
- [42] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, *et al.*, “Maximum entropy inverse reinforcement learning,” in *Aaai*, vol. 8. Chicago, IL, USA, 2008, pp. 1433–1438.
- [43] J. Ho, J. Gupta, and S. Ermon, “Model-free imitation learning with policy optimization,” in *International conference on machine learning*. PMLR, 2016, pp. 2760–2769.
- [44] J. Ho and S. Ermon, “Generative adversarial imitation learning,” *Advances in neural information processing systems*, vol. 29, 2016.
- [45] C. Finn, S. Levine, and P. Abbeel, “Guided cost learning: Deep inverse optimal control via policy optimization,” in *International conference on machine learning*. PMLR, 2016, pp. 49–58.
- [46] J. Fu, K. Luo, and S. Levine, “Learning robust rewards with adversarial inverse reinforcement learning,” *arXiv preprint arXiv:1710.11248*, 2017.

- [47] J. Fu, A. Singh, D. Ghosh, L. Yang, and S. Levine, "Variational inverse control with events: A general framework for data-driven reward definition," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, ser. NIPS'18. Red Hook, NY, USA: Curran Associates Inc., 2018, p. 8547–8556.
- [48] A. Singh, L. Yang, K. Hartikainen, C. Finn, and S. Levine, "End-to-end robotic reinforcement learning without reward engineering," *arXiv preprint arXiv:1904.07854*, 2019.
- [49] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta, "R3m: A universal visual representation for robot manipulation," *arXiv preprint arXiv:2203.12601*, 2022.
- [50] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain, "Time-contrastive networks: Self-supervised learning from video," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 1134–1141.
- [51] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. G. Arenas, H.-T. L. Chiang, T. Erez, L. Hasenclever, J. Humplik, *et al.*, "Language to rewards for robotic skill synthesis," *arXiv preprint arXiv:2306.08647*, 2023.
- [52] Y. Du, K. Konyushkova, M. Denil, A. Raju, J. Landon, F. Hill, N. de Freitas, and S. Cabi, "Vision-language models as success detectors," *arXiv preprint arXiv:2303.07280*, 2023.
- [53] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "Minigt-4: Enhancing vision-language understanding with advanced large language models," *arXiv preprint arXiv:2304.10592*, 2023.
- [54] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch, Q. Vuong, A. He, *et al.*, "Bridgedata v2: A dataset for robot learning at scale," *arXiv preprint arXiv:2308.12952*, 2023.