

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

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** The pre-train and fine-tune approach in machine learning has been
2 highly successful across various domains, enabling rapid task learning by utiliz-
3 ing existing data and pre-trained models from the internet. We seek to apply this
4 approach to robotic reinforcement learning, allowing robots to learn new tasks
5 with minimal human involvement by leveraging online resources. We introduce
6 ROBOFUME, a reset-free fine-tuning system that pre-trains a versatile manipu-
7 lation policy from diverse prior experience datasets and autonomously learns a
8 target task with minimal human input. In real-world robot manipulation tasks, our
9 method can incorporate data from an external robot dataset and improve perfor-
10 mance on a target task in as little as 3 hours of autonomous real-world experience.
11 We also evaluate our method against various baselines in simulation experiments.
12 Website: tinyurl.com/robofume

13 **Keywords:** Autonomous RL, scalable robot learning, vision language models for
14 robotics

15 1 Introduction

16 In many domains that involve machine learning, a widely successful paradigm for learning task-
17 specific models is to first pre-train a general-purpose model from an existing diverse prior dataset,
18 and then adapt the model with a small addition of task-specific data [1, 2, 3, 4, 5]. This paradigm is
19 attractive to real-world robot learning, since collecting data on a robot is expensive, and fine-tuning
20 an existing model on a small task-specific dataset could substantially improve the data efficiency
21 for learning a new task. Pre-training a policy with offline reinforcement learning and then fine-
22 tuning it with online reinforcement learning is a natural way to implement this paradigm in robotics.
23 However, numerous challenges arise when using this recipe in practice. First, off-the-shelf robot
24 datasets often use different objects, fixture placements, camera viewpoints, and lighting conditions
25 compared to the local robot platform. This causes non-trivial distribution shifts between pre-training
26 and online fine-tuning data, which makes effectively fine-tuning a robot policy difficult. Indeed,
27 most existing works only show the benefit of the pre-train and fine-tune paradigm where the robot
28 uses the same hardware instance in both pre-training and fine-tuning phases [6, 7]. Second, training
29 or fine-tuning a policy in the real world often requires extensive human supervision, which includes
30 manually resetting the environment between trials [8, 9, 10] and engineering reward functions [11,
31 12, 7]. In this work, our goal is to address these two challenges and develop a practical framework
32 that enables robot fine-tuning with minimal time and human effort.

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

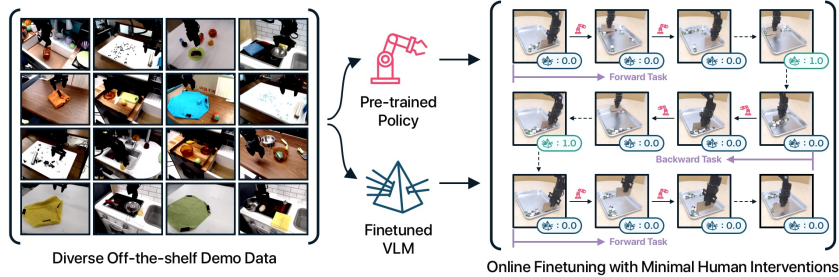


Figure 1: We propose a system that enables autonomous and efficient real-world robot learning. First, we pre-train a **multi-task policy** and fine-tune a pre-trained **Vision-Language Model (VLM)** as a **reward model** using diverse off-the-shelf demonstration datasets. Then, we fine-tune the **pre-trained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

37 ual environment resets using reset-free RL [11, 7, 13], where an agent alternates between running
 38 a task policy and a reset policy during training while updating both with online experience. How-
 39 ever, these works do not leverage diverse off-the-shelf robot datasets. Recent advances in offline
 40 RL algorithms have enabled policies to leverage diverse offline data and improve further via online
 41 fine-tuning [14, 15], but these new methods have not been integrated into a system that aims at min-
 42 imizing human supervision in the fine-tuning phase. There are also works that propose to eliminate
 43 the need for human-specified reward functions by learning reward prediction models [16, 17, 18, 13],
 44 but we found that many of these proposed models can be brittle when deployed in a real-world RL
 45 fine-tuning setup. In summary, although prior works have presented individual components that are
 46 vital to building a working system for efficient and human-free robot learning, it is not clear which
 47 components one should use to put together such a system and how.

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

69 We evaluate our framework by pre-training it on the Bridge dataset [19] and testing it on a diverse
 70 set of real-world downstream tasks: cloth folding, cloth covering, sponge pick-and-place, placing
 71 lid on a pot, and putting a pot in a sink. We find that our system provides substantial improvements
 72 over offline-only methods with as little as 3 hours of real-world training. We perform more quanti-
 73 tative experiments in a simulation setup, where we illustrate that our method outperforms imitation
 74 learning and offline RL methods that either do not perform fine-tuning online or do not incorporate
 75 diverse prior data.

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

80 2 Related Work

81 **Offline RL.** Offline RL algorithms [20, 21, 22, 23, 24, 25] provide a framework for initializing robot
82 manipulation policies from offline demonstrations or interaction datasets. Such algorithms can also
83 be extended to include an online fine-tuning phase after training a policy offline [26, 27, 28, 29,
84 30, 31, 32, 15]. Our work utilizes a recent offline RL algorithm, calibrated Q-learning (CalQL)
85 [15], a state-of-the-art method that effectively learns from offline data and continuously improves
86 the policy’s performance online by explicitly correcting the scale of the learned Q-values. We show
87 that integrating CalQL helps our framework effectively utilize diverse prior datasets that have large
88 distribution shifts from real-world online interactions.

89 **Reset-free RL.** Training an RL policy on a real robot typically requires manual environment resets.
90 To eliminate such need to manually reset environments, prior works have studied approaches to
91 learn robot policies in a ‘reset-free’ setup. Some work [33, 34, 11, 35, 36] cast the ‘reset-free’
92 learning problem as a multi-task learning problem, observing that by learning a set of tasks where
93 some of the tasks could reset others, an agent could then be trained to perform all of those tasks
94 without needing manual resets. Other works [37, 38, 39, 12, 17, 13, 7] learn both a task policy and a
95 reset policy for performing the task and resetting to the initial state distribution. Our work takes an
96 approach between the two classes of approaches, learning a language-conditioned multi-task policy
97 that can perform both the target task and the reset for the target task. Most of these prior works learn
98 from scratch rather than incorporating prior data and assume that a reward function is available.
99 ARIEL [7] combines incorporating prior data with reset-free learning but assumes a hand-crafted
100 reward function for each environment. They also collect their own prior dataset on the same robot
101 hardware set-up as their target task. MEDAL++ [13] learns a reward classifier with demonstration
102 and online interaction data via adversarial training, but does not consider incorporating diverse prior
103 data. Leveraging diverse, off-the-shelf prior demonstration datasets is desirable since these datasets
104 are readily available to use and can help a system obtain a policy initialization for efficient fine-
105 tuning on a target task. Our system offers an approach to both incorporate diverse prior data and
106 improve the autonomy of the fine-tuning phase by learning a model for predicting rewards. In
107 particular, we found out that by leveraging diverse demonstration data, our system requires only
108 about 3 hours of training in the real world compared to 10-30 hours in MEDAL++.

109 **Reward learning.** Early works have studied the problem of learning a reward or cost function in
110 imitation learning. These works leverage inverse optimal control (IOC) or inverse reinforcement
111 learning (IRL) to extract a reward function directly from expert demonstrations [40, 41, 42]. With
112 the advent of deep neural networks, more recent works have explored learning a reward model
113 for an imitation learning or RL policy [43, 44, 45, 46, 47, 17, 13]. When using classifier-based
114 reward models in reinforcement learning, RL agents can exploit the learned model by exploring
115 states unseen during classifier training, tricking the model to output incorrect rewards. To solve such
116 an exploitation issue, many works that learn reward models leverage adversarial learning, where a
117 system learns a discriminator that identifies states similar to those in demonstrations as positives and
118 those visited by the policy as negatives [44, 47, 17, 13]. However, prior work has found this training
119 objective to be sensitive to distribution shifts between offline and online setups, such as lighting and
120 camera view changes [48]. In this work, we fine-tune vision language models (VLM), pre-trained on
121 internet-scale data, to construct a reward model. Large scale pre-training can learn representations
122 that are robust to natural variations such as lighting, camera shifts and distractors [49, 16].

123 **Leveraging pre-trained representations as reward predictors.** Several recent works have shown
124 positive results in utilizing pre-trained vision models [50, 16], large language models (LLMs) [51]
125 or vision language models (VLMs) [52] as reward predictors. We tried VIP [16], a method that

126 pre-trains a visual representation for generating dense reward functions for novel robotic tasks, and
127 found it insufficient for the real-world robot fine-tuning setup. In this work, we fine-tune a pre-
128 trained VLM [53] and find that it performs most effectively as a reward model. Our proposed system
129 is flexible and can easily be adapted to use other pre-trained visual representations and VLMs.

130 3 Preliminaries

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

137 To facilitate learning on the downstream task, we also assume the availability of a small set of target
138 task demos \mathcal{D}_f , target task reset demos \mathcal{D}_b , and target task failure states \mathcal{D}_{\ominus} . The reset demos \mathcal{D}_b
139 come from the reset task τ_b which resets the environment from an end state of τ_f to the initial state
140 distribution of τ_f . The failure states \mathcal{D}_{\ominus} consist entirely of image observations that correspond to
141 unsuccessful states and are collected to aid with the VLM reward learning. In addition to all the
142 given data ($\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\ominus}$), each task τ is also accompanied with a language description l .

143 4 ROBOFUME

144 Our work focuses on designing an efficient and scalable framework for pre-training on a diverse
145 set of prior demonstrations and autonomously fine-tuning on target tasks. Our system consists of
146 an offline pre-training phase and an online fine-tuning phase. In Section 4.1, we discuss how we
147 pre-train a language-conditioned multi-task policy on diverse data that can be fine-tuned online effi-
148 ciently. Online fine-tuning requires a reward function to label successes and failures. In Section 4.2,
149 we introduce a VLM-based classifier for providing a reward signal to the policy in the fine-tuning
150 phase. Finally, in Section 4.3, we describe how to autonomously adapt the pre-trained policy in
151 the fine-tuning phase by utilizing the VLM-based reward classifier as a reward signal and chaining
152 forward and backward behaviors to practice the task with minimal human interventions.

153 4.1 Pre-Training a Multi-Task Policy on Diverse Prior Data

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

161 CalQL requires the training of an actor and a critic. Since we use image observations, we addition-
162 ally train an encoder $\phi(s_{\text{img}})$ that projects the images into a lower-dimensional space before giving
163 them as inputs to the actor and critic. The encoder ϕ is a 4-layer CNN, and is optimized exclu-
164 sively against the critic loss. To best utilize the multi-task data, we encode task descriptions l using
165 pre-trained CLIP embeddings, resulting in an embedding $z = \text{CLIP}(l)$ which is used as the task
166 representation. The policy then takes as inputs a concatenation of the encoded image observation
167 $\phi(s_{\text{img}})$, task representation z , and proprioceptive information s_p , processes the concatenated vector
168 through an MLP, and produces the output action a .

169 In addition to updating the policy using CalQL, we regularize policy learning with a behavior cloning
170 (BC) loss, which encourages the behaviors to stay close to the seen demonstrations. Not only does
171 this regularization improve performance of the offline pre-training, but we find that it also makes it

172 less likely for the autonomous fine-tuning procedure to exploit false positive rewards from the VLM
173 reward model. The weight of the BC regularization term is chosen such that the scales of the RL
174 loss and the BC loss are similar throughout the pre-training phase. We train the policy π and the
175 critic Q with datasets $\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b$. After the offline learning phase, the policy and critic contain
176 knowledge of all tasks in the prior data and the target task.

177 4.2 Fine-Tuning A Vision-Language Model for Rewards

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

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

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

204 4.3 Autonomous Online Fine-tuning

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

209 Since we aim for a fully autonomous setup, we roll out the policy in a reset-free manner, alternating
210 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
211 fine-tuned VLM from the previous subsection as the sparse reward function for the RL algorithm.
212 When the VLM predicts the task has been completed successfully, we terminate the episode and
213 switch the language instruction for the policy to complete the other task. In addition to switching
214 tasks upon completion as predicted by the VLM, we switch after a fixed number of timesteps (150) to
215 ensure the robot does not become stuck in bad states. As mentioned in Section 4.1, we fine-tune the
216 policy using CalQL with an additional BC regularization term on the critic. We find that without a
217 BC regularization term, behaviors degrade over the course of training. By constraining the policy to
218 stay close to the expert demonstrations from the target and reset tasks, the agent becomes less likely
219 to exploit false-positives from the VLM reward model. We use the same fixed BC regularization
220 weight throughout fine-tuning as we did on the offline pre-training phase. Our fine-tuning pipeline

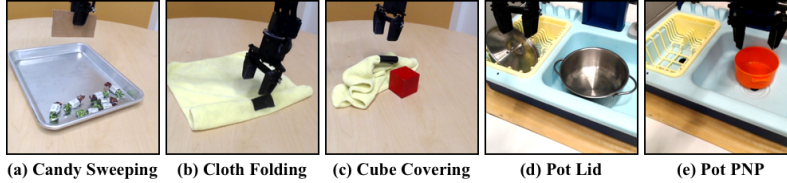


Figure 2: **Illustrations of the five real-world evaluation tasks.** (a) Sweep candies to the top of the tray. (b) fold the yellow cloth. (c) cover a red wooden cube using the cloth. (d) place the lid on top of the metallic pot. (e) move the orange pot from the sink to the drying rack.

221 is implemented on top of the implementation of MEDAL++ [13]. Please refer to this work for more
 222 details on our training procedure.

223 5 Experiments

224 We design our experiments to answer the following questions: Is our method able to improve its
 225 performance through near autonomous online interactions? How does our proposed VLM reward
 226 function mechanism compare to existing alternatives? And, how does each component of ROBO-
 227 FUME or data affect the performance of our method?

228 5.1 Real Robot Experiments

229 **Setup.** We evaluate ROBOFUME on five different real-robot manipulation tasks. We use a WidowX
 230 250 robotic arm with a single third-person camera (Logitech C920, resizing images to 100x100
 231 pixels). Figure 2 shows the five tasks we fine-tune and evaluate on. Our method runs autonomously
 232 executing back and forth the target task and the reset target task for a fixed number of steps or
 233 until the VLM predicts success. For tasks involving deformable objects (the two cloth tasks) we
 234 manually reset the object to the initial forward pose every 15-25 episodes, and for the rest of the tasks
 235 we reset every 30-35 episodes. Tasks that use the kitchen-sink environment (*pot lid* and *pot pnp*)
 236 frequently experience episode interruptions when the robot arm applies more than the maximum
 237 allowed torque, for example, when close to the sink borders. All tasks use 50 forward and 50
 238 backward demos for the target task, and fewer than 20 combined trajectories of failures. We use
 239 demos from the BridgeDataV2 [19, 54] for pre-training our language-conditioned policy, selecting
 240 approximately 1,000 trajectories with relevant behaviors per task.

241 **Results.** Table 1 shows the results of our method after pretraining (labeled OFFLINE) and after au-
 242 tonomous fine-tuning (labeled FT 30K STEPS), comparing with a behavior cloning (BC) baseline.
 243 BC trains a language-conditioned policy on all the prior and target data. After 30k steps of au-
 244 tonomous online interaction, our method shows relative improvement of 51% upon the pre-trained
 245 performance, and outperforms BC by 58% on an average. For pick and place tasks (*pot lid* and
 246 *pot pnp*), the fine-tuned policy was more likely to retry the action if it initially failed to grasp the
 247 object. For candy sweeping, BC and the pre-trained policy were prone to overshooting and pushing
 248 on the border of the tray after the first sweep, whereas fine-tuning the policy enabled the policy to
 249 chain multiple sweeping attempts for higher success. Additionally, we find that policies learned by
 250 ROBOFUME (both offline and after fine-tuning) to be more robust to scene distractors on the *candy*
 251 *sweeping* task, as reported in Table 3. The policies were trained without any distractors, but multiple
 252 objects not seen during training were placed in the background during evaluation. ROBOFUME
 253 policies retained 68% of its original performance, compared to BC which retained only 10% of its
 254 original performance. We hypothesize that BC might be more sensitive to spurious features, whereas
 255 ROBOFUME learns from more predictive features, leading to more robust policies.

256 5.2 Simulation Experiments and Ablations

257 We use a suite of simulated robotic manipulation environments to ablate contributions of different
 258 components of our algorithm. We test on three simulated environments used in [6]. We consider

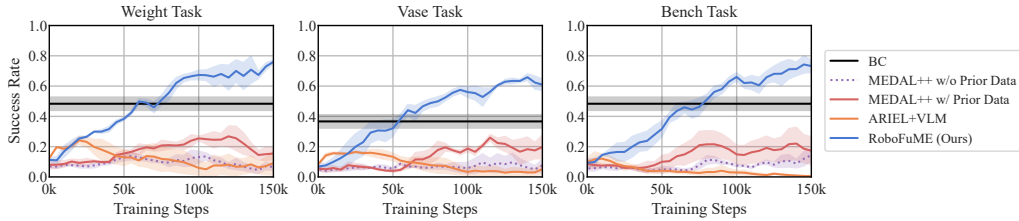


Figure 3: **Performance of our method on three simulated environments.** We report the success rate over the course of training, averaged over three seeds. Our method ROBOFUME outperforms BC, ARIEL+VLM [7], and MEDAL++ [13] consistently on all three domains.

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 1: **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.

259 three bin-sorting tasks in which different objects (a vase, a tiny bench, and a dumbbell weight) have
 260 to be placed on the correct bin based on the language instruction, given only a sparse binary reward.
 261 We provide 10 forward and reset demonstrations for each task, 30 failure demos, and 10 demos each
 262 for 20 prior tasks that show picking and placing diverse objects on the same environment. For all
 263 methods that require online experience, we reset the environment every 1,000 environment steps, i.e.
 264 every 25 episodes of interactions. We compare our method against the following baselines: (1) *BC*
 265 behavior clones on all prior and target data; (2) *MEDAL++* learns separate forward and backward
 266 policies from target forward and backward task demonstrations and performs reset-free learning using
 267 an adversarially trained classifier as a reward signal; (3) *MEDAL++ with prior data* modifies
 268 *MEDAL++* to a single language-conditioned multi-task policy and adds all prior demonstration data
 269 into the replay buffer; (4) *ARIEL+VLM* modifies ARIEL [7] to use our VLM reward models as re-
 270 ward signal, instead of a handcrafted ground-truth reward. The results of our simulation experiments
 271 are presented in Figure 3. In all simulation tasks, our method ROBOFUME consistently outperforms

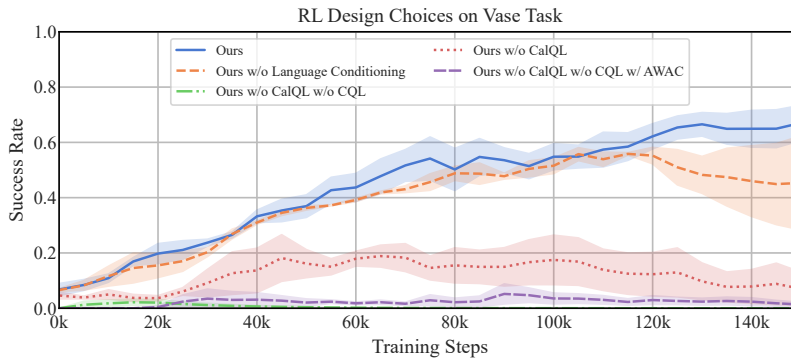


Figure 4: **Performance of our method on the Vase simulated task with different actor-critic update objectives.** Fine-tuning with CalQL is critical to obtain stable improvements on this task, as training with CQL, AWAC, or SAC yields poor performance. We also find that language-conditioned policies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

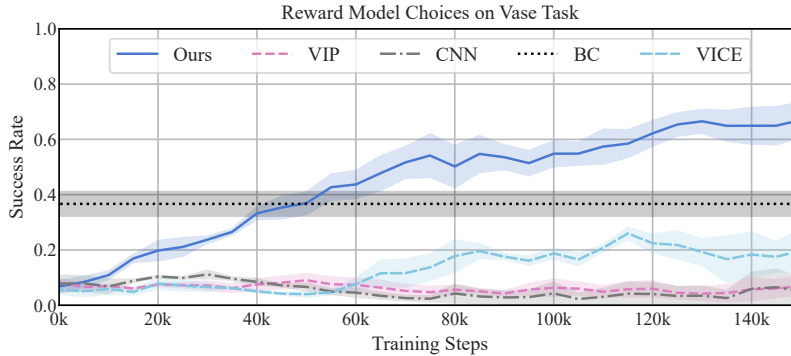


Figure 5: **Performance of our method on the simulated Vase task using different reward functions.** Our method uses a fine-tuned VLM reward function and outperforms VICE rewards, whereas CNN and VIP rewards fail to improve online.

272 prior methods, achieving success rates at least 20% higher than all baselines within 200k steps of
 273 online fine-tuning.

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

279 **Ablations on Reward Models.** We compare our VLM reward function against other choices of
 280 automatic reward functions on the Vase simulated task in Figure 5. VICE [47] adversarially trains
 281 a binary classifier using positive samples from successful demonstrations, and labeling online expe-
 282 rience as negative. We find that offline pre-training sufficiently limits the exploitation of the frozen
 283 VLM reward, outperforming VICE and thus, bypassing the need for adversarially trained reward
 284 functions. Such adversarial training can often learn to discriminate based on spurious shifts in the
 285 real world, such as lighting or scene changes, leading to instability in training outside simulation.
 286 VIP [16] trains a representation function such that the distance in representation space between the
 287 current observation and a goal image can be used to construct a dense reward function. We find that
 288 in the Vase simulated task, VIP fails to obtain good behaviors. Qualitatively, we observe VIP to
 289 be prone to false positives, which are exploited by the RL algorithm. To test the importance of the
 290 VLM large-scale pre-training compared to our fine-tuning procedure, we train a CNN classifier from
 291 scratch using the same data as we used to fine-tune the VLM, leading to unsatisfactory performance
 292 compared to fine-tuning a VLM.

293 5.3 Additional Analysis

294 In order to understand the performance of our method in the real robot experiments better, we per-
 295 formed additional analysis to examine various aspects of our framework. Please find these experi-
 296 ments in the Appendix.

297 6 Conclusion and Future Work

298 We introduced an autonomous framework that leverages existing diverse prior robot demonstration
 299 datasets and improves performance in a new robot manipulation skill by finetuning online. By
 300 combining state-of-the-art offline-to-online RL algorithms, reset-free RL, and VLM-based reward
 301 models, our framework can fine-tune efficiently and nearly autonomously. Integrating this work
 302 with new VLM models that can exhibit robust zero-shot performance on unseen manipulation tasks
 303 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*, 1(8):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*, 33:1877–1901, 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*, pages 16000–16009, 2022.
- [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*, pages 1652–1662. PMLR, 2023.
- [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)*, pages 378–383. IEEE, 2016.
- [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)*, pages 2351–2358. IEEE, 2017.
- [10] K. Ploeger, M. Lutter, and J. Peters. High acceleration reinforcement learning for real-world juggling with binary rewards. In *Conference on Robot Learning*, pages 642–653. PMLR, 2021.
- [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)*, pages 6664–6671. IEEE, 2021.
- [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*, pages 308–319. PMLR, 2022.
- [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.

- 347 [17] A. Sharma, R. Ahmad, and C. Finn. A state-distribution matching approach to non-episodic
348 reinforcement learning. *arXiv preprint arXiv:2205.05212*, 2022.
- 349 [18] Y. J. Ma, W. Liang, V. Som, V. Kumar, A. Zhang, O. Bastani, and D. Jayaraman.
350 Liv: Language-image representations and rewards for robotic control. *arXiv preprint*
351 *arXiv:2306.00958*, 2023.
- 352 [19] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and
353 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets.
354 *arXiv preprint arXiv:2109.13396*, 2021.
- 355 [20] Y. Wu, G. Tucker, and O. Nachum. Behavior regularized offline reinforcement learning. *arXiv*
356 *preprint arXiv:1911.11361*, 2019.
- 357 [21] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and
358 scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 359 [22] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review,
360 and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 361 [23] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma. Mopo: Model-
362 based offline policy optimization. *Advances in Neural Information Processing Systems*, 33:
363 14129–14142, 2020.
- 364 [24] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforce-
365 ment learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- 366 [25] W. Zhou, S. Bajracharya, and D. Held. Plas: Latent action space for offline reinforcement
367 learning. In *Conference on Robot Learning*, pages 1719–1735. PMLR, 2021.
- 368 [26] N. Ashvin, D. Murtaza, G. Abhishek, and L. Sergey. Accelerating online reinforcement learn-
369 ing with offline datasets. *CoRR*, vol. *abs/2006.09359*, 2020.
- 370 [27] R. Rafailov, T. Yu, A. Rajeswaran, and C. Finn. Offline reinforcement learning from images
371 with latent space models. In *Learning for Dynamics and Control*, pages 1154–1168. PMLR,
372 2021.
- 373 [28] J. Lyu, X. Ma, X. Li, and Z. Lu. Mildly conservative q-learning for offline reinforcement
374 learning. *Advances in Neural Information Processing Systems*, 35:1711–1724, 2022.
- 375 [29] A. Beeson and G. Montana. Improving td3-bc: Relaxed policy constraint for offline learning
376 and stable online fine-tuning. *arXiv preprint arXiv:2211.11802*, 2022.
- 377 [30] J. Wu, H. Wu, Z. Qiu, J. Wang, and M. Long. Supported policy optimization for offline re-
378 inforcement learning. *Advances in Neural Information Processing Systems*, 35:31278–31291,
379 2022.
- 380 [31] S. Lee, Y. Seo, K. Lee, P. Abbeel, and J. Shin. Offline-to-online reinforcement learning via
381 balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, pages 1702–
382 1712. PMLR, 2022.
- 383 [32] M. S. Mark, A. Ghadirzadeh, X. Chen, and C. Finn. Fine-tuning offline policies with optimistic
384 action selection. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.
- 385 [33] K. Lu, A. Grover, P. Abbeel, and I. Mordatch. Reset-free lifelong learning with skill-space
386 planning. *arXiv preprint arXiv:2012.03548*, 2020.
- 387 [34] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan. Learning to walk in the real world with minimal
388 human effort. *arXiv preprint arXiv:2002.08550*, 2020.

- 389 [35] A. Gupta, C. Lynch, B. Kinman, G. Peake, S. Levine, and K. Hausman. Demonstration-
390 bootstrapped autonomous practicing via multi-task reinforcement learning. *arXiv preprint*
391 *arXiv:2203.15755*, 1, 2022.
- 392 [36] K. Xu, Z. Hu, R. Doshi, A. Rovinsky, V. Kumar, A. Gupta, and S. Levine. Dexterous manipulation
393 from images: Autonomous real-world rl via substep guidance. In *2023 IEEE International*
394 *Conference on Robotics and Automation (ICRA)*, pages 5938–5945. IEEE, 2023.
- 395 [37] B. Eysenbach, S. Gu, J. Ibarz, and S. Levine. Leave no trace: Learning to reset for safe and
396 autonomous reinforcement learning. *arXiv preprint arXiv:1711.06782*, 2017.
- 397 [38] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, V. Kumar, and S. Levine. The
398 ingredients of real-world robotic reinforcement learning. *arXiv preprint arXiv:2004.12570*,
399 2020.
- 400 [39] A. Sharma, A. Gupta, S. Levine, K. Hausman, and C. Finn. Autonomous reinforcement learn-
401 ing via subgoal curricula. *Advances in Neural Information Processing Systems*, 34:18474–
402 18486, 2021.
- 403 [40] A. Y. Ng, S. Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1,
404 page 2, 2000.
- 405 [41] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Pro-*
406 *ceedings of the twenty-first international conference on Machine learning*, page 1, 2004.
- 407 [42] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey, et al. Maximum entropy inverse reinforce-
408 ment learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- 409 [43] J. Ho, J. Gupta, and S. Ermon. Model-free imitation learning with policy optimization. In
410 *International conference on machine learning*, pages 2760–2769. PMLR, 2016.
- 411 [44] J. Ho and S. Ermon. Generative adversarial imitation learning. *Advances in neural information*
412 *processing systems*, 29, 2016.
- 413 [45] C. Finn, S. Levine, and P. Abbeel. Guided cost learning: Deep inverse optimal control via
414 policy optimization. In *International conference on machine learning*, pages 49–58. PMLR,
415 2016.
- 416 [46] J. Fu, K. Luo, and S. Levine. Learning robust rewards with adversarial inverse reinforcement
417 learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 418 [47] J. Fu, A. Singh, D. Ghosh, L. Yang, and S. Levine. Variational inverse control with events: A
419 general framework for data-driven reward definition. In *Proceedings of the 32nd International*
420 *Conference on Neural Information Processing Systems, NIPS’18*, page 8547–8556, Red Hook,
421 NY, USA, 2018. Curran Associates Inc.
- 422 [48] A. Singh, L. Yang, K. Hartikainen, C. Finn, and S. Levine. End-to-end robotic reinforcement
423 learning without reward engineering. *arXiv preprint arXiv:1904.07854*, 2019.
- 424 [49] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual represen-
425 tation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 426 [50] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain.
427 Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE international*
428 *conference on robotics and automation (ICRA)*, pages 1134–1141. IEEE, 2018.
- 429 [51] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. G. Arenas, H.-T. L. Chiang, T. Erez,
430 L. Hasenclever, J. Humplik, et al. Language to rewards for robotic skill synthesis. *arXiv*
431 *preprint arXiv:2306.08647*, 2023.

- 432 [52] Y. Du, K. Konyushkova, M. Denil, A. Raju, J. Landon, F. Hill, N. de Freitas, and S. Cabi.
433 Vision-language models as success detectors. *arXiv preprint arXiv:2303.07280*, 2023.
- 434 [53] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language
435 understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- 436 [54] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch,
437 Q. Vuong, A. He, et al. Bridgedata v2: A dataset for robot learning at scale. *arXiv preprint*
438 *arXiv:2308.12952*, 2023.

439 **7 Appendix**

440 **7.1 Method Details**

441 We present an overview of our system 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.
 $\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\otimes} \leftarrow \text{load\_data}()$ .
 $\hat{r} \leftarrow \text{finetune\_vlm}(\hat{r}, \{\mathcal{D}_{\text{prior}}, \mathcal{D}_f, \mathcal{D}_b, \mathcal{D}_{\otimes}\})$ .
// 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  $N_{\text{upd\_ratio}}$  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
    | |  $s \leftarrow s'$ .

```

442 **7.2 Additional Analysis in Real Robot Experiments**

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 2: **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 3: **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.

443 **How Accurate is the VLM Reward?** We analyze the performance of the VLM reward over the
444 course of fine-tuning for real-robot experiments. In Table 2, we report the false positive rate, false
445 negative rate, accuracy, and precision metrics for the VLM reward. The metrics are computed on
446 the data collected during fine-tuning against a hand-engineered ground truth reward. We observe
447 that while false negative rates are high, false positive rates are low across all tasks. This asymmetry

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

Table 4: **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.

448 is crucial for successful RL fine-tuning, as RL policies can learn poor behaviors by exploiting false
 449 positives, but labeling some successful rollouts as negatives does not necessarily impede learning.

450 **How Important is Diverse Prior Data and Language Conditioning?** We ablate the contribution
 451 of diverse prior data and language-conditioned policies to ROBOFuME by evaluating the offline
 452 performance on the *candy sweeping* task, reported in Table 4. When pre-training without using prior
 453 data, that is, exclusively using target data, our method is able to sweep less than half the amount
 454 of candies on average. Similarly, we find that one-hot task encodings perform substantially worse
 455 than language-conditioned policies, as the prior dataset used in real-robot training is larger and more
 456 diverse compared to the simulation experiments.