

# Multi-Task Interactive Robot Fleet Learning with Visual World Models

**Anonymous Author(s)**

Affiliation

Address

email

1           **Abstract:** Recent advancements in large multi-task robot agents hold promise of  
2           robot fleets operating in household and industrial settings, capable of performing  
3           diverse tasks across various environments. Despite these advances, these robots  
4           often struggle with generalization and robustness in real-world, unstructured set-  
5           tings, limiting their practical deployment. To address these issues, we propose  
6           FIREWORK, a framework for multi-task interactive robot fleet learning, supervis-  
7           ing the continual policy update throughout deployment with efficient multi-task  
8           runtime monitoring and humans in the loop. We realize this with a visual world  
9           model to predict future task outcomes, which serves as the backbone for learned  
10          representation across downstream error prediction tasks. We train dynamic error  
11          predictors that auto-adjust their prediction criterias to adapt to the robot's evolv-  
12          ing autonomy, reducing human workload over time. Evaluations on large-scale  
13          benchmarks demonstrate our framework's effectiveness in improving multi-task  
14          policy performance and monitoring accuracy. We demonstrate the performance  
15          of FIREWORK on robocasa in simulation and mutex in real world, two large, di-  
16          verse multi-task environment benchmarks. More information at project website  
17          <https://firework-corl2024.github.io/>

18           **Keywords:** Robot Manipulation, Interactive Imitation Learning, Fleet Learning

## 19          1 Introduction

20          In recent years, great progress has been made in developing generally capable robot agents for  
21          performing a wide range of tasks [1, 2, 3]. The rapid progress in generalist robots holds promise  
22          for a future of robot fleets [4, 5, 6] operating in households and industrial settings on a generalist  
23          multi-task policy. Despite the progress in research development, the robots can still fall short of gen-  
24          eralization and robustness when being deployed in the real world, where environments are diverse  
25          and unstructured; this undermines the safety and reliability of the system and limits the model's  
26          applicability to real-world deployment.

27          To mitigate these challenges, numerous works in the field of interactive imitation learning (IIL)  
28          [7, 8, 9, 10, 11, 12, 13] and interactive fleet learning (IFL) [14, 5, 6] have been proposed. Prior  
29          works in human-in-the-loop learning [7, 13, 15, 10] have proposed to use human monitoring and  
30          human correction on the job to ensure trustworthy deployment; however, they require constant hu-  
31          man monitoring and incur high human workload. To reduce the amount of human workload, the  
32          runtime monitoring approaches [16, 17, 18, 19, 12, 20, 21] propose using error detectors to automati-  
33          cally monitor robot performance, detect system anomalies and ask for human control when needed.  
34          Prior works has proposed out-of-distribution (OOD) detection [20, 21, 18] or failure detection [12]  
35          methods to detect error cases of robot execution. However, these works have primarily been ap-  
36          plied for task-specific policies within a controlled single-task setting, limiting the efficiency and  
37          effectiveness in the multi-task fleet deployment.

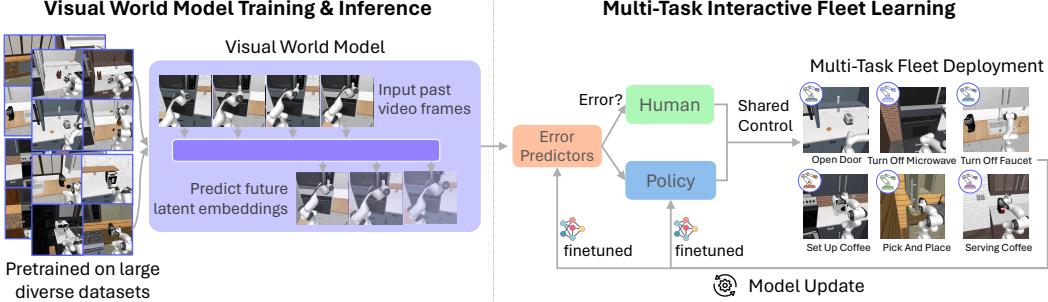


Figure 1: **Multi-Task Fleet Learning.** We first train a visual world model and then use it for multi-task fleet learning and deployment.

38 We introduce FIREWORK<sup>1</sup>, a framework for multi-task interactive robot fleet learning, where multi-  
 39 task policies can be deployed to massive robot fleets with humans in the loop and efficient run-  
 40 time monitoring. The key challenges are 1) to enable a multi-task policy to be continually updated  
 41 throughout deployment with the diverse, massive-scale deployment data, enabling sharing of knowl-  
 42 edge across tasks; 2) to build a runtime monitoring mechanism which efficiently monitor different  
 43 tasks simultaneously across different diverse scenes throughout the deployment; To tackle the first  
 44 challenge, we train a multi-task policy that is able to be constantly finetuned on deployment data  
 45 learning from the data on the job; To tackle the second challenge, we train a visual world model  
 46 as the backbone of the runtime monitoring, whose learned representation can be shared and used  
 47 across downstream error predictors for various tasks.

48 Recent advancement of world models [22, 23, 24, 25] enables simulation into the future, predicting  
 49 future progress outcomes. Inspired by such progress, we build a visual world model that is trained  
 50 on large diverse experience of robot trajectory frames. This enables future prediction of robot task  
 51 progress, allowing us to see future outcome and preempt potential failures from happening. The  
 52 learned representation from the model are then and shared in downstream error prediction across  
 53 different tasks. Building on the frozen representation from the world model, we train error predictors  
 54 based on failure and OOD detection which can be continuously finetuned throughout deployment.  
 55 Rather than having a fixed detection threshold as in prior works, we auto-adjust the error predictor’s  
 56 prediction criteria based on the current task performance using the human interaction trace as a  
 57 proxy. We find that the auto threshold adjustment can customize to the evolving level of robot  
 58 autonomy, giving better runtime monitoring results.

59 We test our multi-task interactive robot fleet learning framework on large multi-task benchmarks  
 60 both in simulation and on real robots. Our key findings are 1) our runtime monitoring can effectively  
 61 supervise diverse multi-task scenarios with high performances; 2) the multi-task policy is able to  
 62 continually improve over time with the data collected at fleet deployment; 3) our error predictors for  
 63 runtime monitoring surpass baselines in their accuracy and use of human workload. In summary,  
 64 our contributions are as follows:

- 65 1. A framework of multi-task interactive robot fleet learning where the multi-task robot policy  
 66 efficiently improves over deployment with runtime monitoring and human in the loop;
- 67 2. A runtime monitoring mechanism that is based on visual world model backbone and comes  
 68 with task-adaptive error prediction thresholds;
- 69 3. The showcase of high performance of such multi-task fleet learning system both in simula-  
 70 tion and on real robots.

<sup>1</sup>FIREWORK: Fleet-based Interactive Robot LEarning with Visual WORld Models for Multi-tasK Setting.

71 **2 Related Work**

72 **Multi-task Robot Learning.** There has been significant progress in the development of multi-task  
73 robotics agents [2, 26, 27, 1, 28, 29, 30, 31], driven by advancement in modern advanced policy  
74 architectures like transformer [1, 32] and diffusion models [33, 34], collection of larger robotics  
75 datasets [35, 36, 37, 38] and advancing robotics benchmarks [32, 39]. Instead of performing on  
76 task-specific domains, these models are generally capable to perform a diverse range of tasks, where  
77 the tasks can be specified commonly from language [29, 1] or multimodal prompts [27, 26].

78 **Robot Fleet Learning.** There has been growing interest in robot fleet learning [40, 4, 5, 41, 42, 6]  
79 as to scale up real-world robot deployment and data collection [43, 44, 45]. Prior literature has  
80 studied the different problems in fleet learning. Hoque et al. [14] proposed the interactive fleet  
81 learning (IFL) framework and studies the *resource allocation* problem [40]. Works in *decentralized*  
82 and *federated learning* [5, 41, 42] studies policy merging and data selection from individual fleets.  
83 Other works studies the *system* aspect of managing a robot fleet [6]. In this work, we studies the  
84 learning aspect of learning from the multi-task deployment data aggregated from a fleet of robot  
85 performing diverse tasks [46, 43].

86 **Interactive Imitation Learning.** To ensure trustworthy deployment and help robot learn on the  
87 job, human-in-the-loop methods [10, 9, 7, 13, 15, 8, 11] has been proposed to use human mon-  
88 itoring and intervention in robot task execution processes. To further reduce human workload in  
89 monitoring, runtime monitoring and error detection techniques [12, 17, 47, 16, 48, 49] seek to de-  
90 tect failures and anomalies during robot task execution. Two main bodies of work are unsupervised  
91 out-of-distribution (OOD) detection [18, 50, 19, 51, 20, 21], and failure detection either with binary  
92 classification [52, 53, 54] or learning risk function estimate from trajectory data [12, 14]. Driven by  
93 the recent development of foundation models, there has been works that uses LLMs / VLMs as error  
94 detectors for LLM- / VLM-based policies [55, 56, 57], or for object detection [58]. These works  
95 have yet to show applicability to robot motor action manipulation policies in close-loop due to the  
96 need to run runtime monitoring with fast, consistent inference speed on real robots [1].

97 **3 FIREWORK: Multi-Task Interactive Robot Fleet Learning**

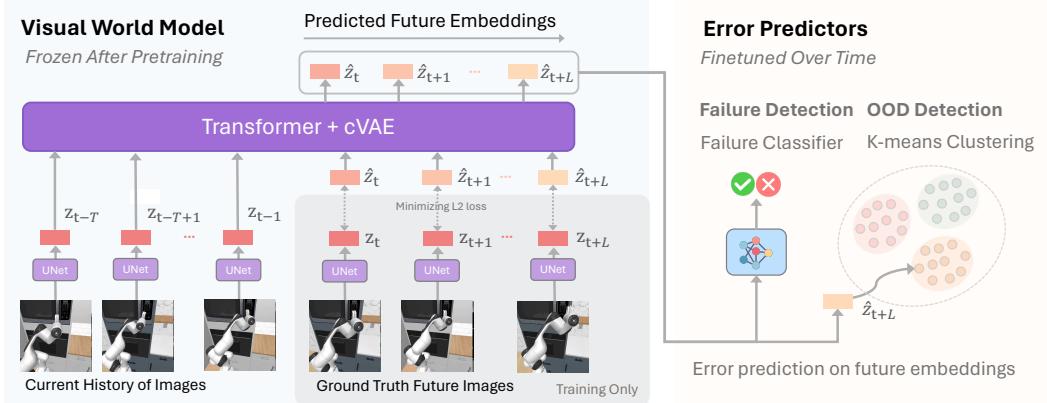
98 **3.1 Background**

99 **3.1.1 Problem Formulation**

100 We formulate multi-task interactive robot fleet learning as finite-horizon Markov Decision Process  
101 (MDP), where  $N$  robots operates in  $N$  independent MDPs. The robot  $i$  operate in the  $i$ -th MDP  
102 specified by  $\mathcal{M}_i = (S, A, \mathcal{T}, H_i, \mu_i^0, R_i)$ , where  $S$  and  $A$  are the state and action spaces of the  
103 robot,  $\mathcal{T} : S \times A \rightarrow S$  is the transition dynamics,  $H_i$  is the horizon length,  $\mu_i^0$  is the initial state  
104 distribution, and  $R_i : S \times A \rightarrow \mathbb{R}$  is the reward function. In a sparse-reward setting,  $R_i$  is replaced  
105 with a goal predicate  $g_i : S \rightarrow \{0, 1\}$ . The collection of  $\{\mathcal{M}_i\}_{i=1}^N$  can be reformulated as one single  
106 MDP, with shared  $S, A, T$ . The data from all robots are aggregated to learn one unified multi-task  
107 policy  $\pi(a | s, g_i)$  that maximizes the expected return:  $\max_{\pi} J(\pi) = \mathbb{E}_{s_t, a_t \sim \pi, \mu_0} \left[ \sum_{t=1}^H g_i(s_t) \right]$ .

108 **3.1.2 Multi-Task Interactive Fleet Deployment**

109 We consider an interactive learning framework [7, 13, 17] for a fleet of robots [4], where learning  
110 and deployment happens iteratively with a human  $\pi_H$  in the loop. The robot bootstrap its initial  
111 policy  $\pi_1$  by doing behavioral cloning from a set of human demonstrations,  $\mathcal{D}^0$ . During each round  
112 deployment round  $i$ , each robot perform task execution with  $\pi_i$  with runtime monitoring, supervised  
113 by the error predictors,  $E_i$ . During policy execution,  $E_i$  determines whether current state can po-  
114 tentially lead to error states. Upon detecting an potential error, it indicates to the human and ask  
115 humans to monitor the process. While monitoring, the human can chooses to actively intervene [59]  
116 and take over the control if necessary, letting  $\pi_H(s, g)$  overriding  $\pi(s, g)$ . We add the trajectories



**Figure 3: Model Architecture.** The visual world model is trained with Unet encoder and decoder, as well a cVAE + Transformer based prediction model. The world model can be used to predict the future embeddings from the current state. The learned representations are used for error predictions, including failure and OOD detection.

117  $\tau = (s_t, a_t, r_t, c_t)$  to the data buffer of the current deployment round,  $\mathcal{D}'$ , where  $c_t$  indicates whether  
 118 timestep  $t$  for trajectory  $\tau$  is from human action. We update the data buffer  $\mathcal{D}^{i+1} = \mathcal{D}^i \cup \mathcal{D}'$ , which  
 119 will be used to update the policy  $\pi_{i+1}$  and error predictors  $E_{i+1}$  for the next round.

### 120 3.2 Runtime Monitoring for Multi-Task Fleet Deployment

121 In this section, we discuss FIREWORK’s run-  
 122 time monitoring mechanism, which monitors  
 123 different tasks simultaneously across different  
 124 diverse scenes throughout the deployment. The  
 125 key challenge is to efficiently monitor different  
 126 tasks simultaneously.

127 We envision such general runtime monitor-  
 128 ing to meet two key design requirements: 1)  
 129 generally-capable: the error predictor models  
 130 should not be task-specific and should be effi-  
 131 ciently trained and used for a variety of tasks;  
 132 2) task-adaptive: it should adapt to different  
 133 tasks and their evolving progress during de-  
 134 ployment. Additionally, as a runtime monitor-  
 135 ing algorithm, it should be failure-preemptive,  
 136 predicting the errors before they occur to ensure  
 137 prompt detection and intervention.

138 Motivated by such design ideas, we build a  
 139 multi-task runtime monitoring model that con-  
 140 tinuously adapts and improves for individual  
 141 tasks during deployment. Our key idea is to  
 142 train a visual world model backbone as a sim-  
 143 ulator for general task progress understanding,  
 144 and use the backbone for downstream error predictions across different tasks. We present the model  
 145 architecture in Figure 3.

#### 146 3.2.1 Training the Backbone Visual World Model

147 Recent advancement of world models [22, 23, 24, 25] enables simulation into the future. Inspired by  
 148 such progress, we train a visual world model on large diverse experience of robot trajectory frames.

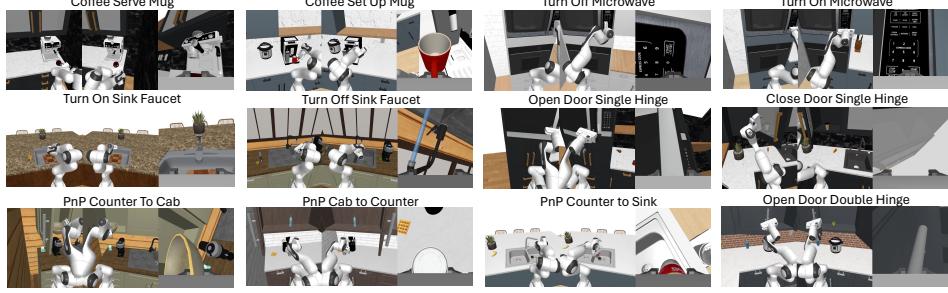


Figure 4: **Evaluation Tasks.** We show here all 12 simulation tasks from Robocasa; we will show the real-world Mutex tasks in the Appendix.

149 This enables future prediction of robot task progress, allowing us to see future outcome and preempt  
 150 potential failures from happening.

151 FIREWORK trains an autoregressive visual world model  $\mathcal{W}$  as the backbone for downstream error  
 152 prediction. The world model  $\mathcal{W} = (E_\gamma, D_\lambda, T_\psi)$  consists of an encoder  $E_\gamma$ , a decoder  $D_\lambda$ , and a  
 153 conditional next state prediction model  $T_\psi$ . In training,  $E_\gamma$  encodes image observation  $x_t$  at timestep  
 154  $t$  into latent embedding  $z_t$ . The decoder  $D_\lambda$  reconstructs  $z_t$  to image  $\hat{x}_t$ , and minimizes the image  
 155 reconstruction L2 loss.  $T_\psi$  inputs the history of  $T$  timesteps of embeddings  $(z_{t-T}, z_{t-T+1}, \dots, z_t)$ ,  
 156 and outputs  $z_{t+1}^\hat{+}$ . It autoregressively predicts  $(z_{t+1}^\hat{+}, z_{t+2}^\hat{+}, \dots, z_{t+L}^\hat{+})$  for  $L$  steps into the future using  
 157 the last  $T$  timesteps of embeddings and reconstructed embeddings, and minimizes the embedding  
 158 reconstruction loss between  $(z_{t+1}^\hat{+}, z_{t+2}^\hat{+}, \dots, z_{t+L}^\hat{+})$  and  $(z_{t+1}, z_{t+2}, \dots, z_{t+L})$ .  $E_\gamma$  and  $D_\lambda$  are imple-  
 159 mented with UNet [60, 33], and  $T_\psi$  uses conditional Variational Autoencoder (cVAE) [61].  $T_\psi$  is  
 160 jointly trained with  $E_\gamma, D_\lambda$  on the same latent space.

161 We use a stochastic latent space rather than a deterministic one, since the stochastic latent space  
 162 of cVAE supports multiple future sampling and facilitates better prediction [62, 17]. We use trans-  
 163 former architecture [63] for the encoder, decoder and prior network for the cVAE in  $T_\psi$ .

### 164 3.2.2 Building the Downstream Error Predictors

165 The visual world model gives information of how the world states change and how tasks progress  
 166 over time. The learned representation from the world model can be useful for downstream error  
 167 prediction tasks, where we predict into the future and see if the future state is an error. Since the  
 168 individual tasks have different conditions, and can have different performance, we will train different  
 169 error predictors on top of the frozen representation that are customized to the different tasks. We  
 170 predict two kinds of errors: failure and out-of-distribution.

171 **Failure Prediction.** We train a failure classifier  $F_\sigma$  on the frozen latent space for each specific  
 172 task. The failure labels are derived from the trajectories  $\tau = (s_t, a_t, r_t, c_t)$  from the last round  
 173  $\mathcal{D}'$  as detailed in Section 3.1.2, where  $c_t$  can indicate moments when human interventions occur.  
 174 The segment of trajectory  $(s_t, a_t, r_t, c_t)_{i-M}^{i-1}$  before the beginning of each human intervention  $s$   
 175  $c_i = \text{human}$  are used as failure labels [13, 17]. A three-class classification is used to distinguish  
 176 between failure states, normal rollouts, and intervention states, so as to distinguish the nuanced  
 177 difference between failures and interventions that often occur at similar states. The classifier is  
 178 trained using a cross-entropy loss function  $\mathcal{L}_F = -\sum_{i=1}^n y_i \log(\hat{y}_i)$  with balanced sampling, where  
 179  $y_i \in \{\text{normal, failure, human}\}$ .

180 **Out-of-Distribution (OOD) Prediction.** We determine whether a state is OOD using k-means clus-  
 181 tering. We use the frozen visual world model  $\mathcal{W}$  to generate embeddings for a subset of trajectories  
 182 sampling from the data buffer  $\mathcal{D}^i$  at the deployment round  $i$ . We then reduce the dimensions of this  
 183 latent space to  $l$  using Principal Component Analysis (PCA) and calculate  $c$  k-means centroids for  
 184 each task. To perform OOD prediction for an embedding  $z$ , we first reduce its dimensions to  $l$  using  
 185 PCA and find the nearest centroid of the last predicted latent space, and calculate their L2 distance  
 186  $d$ . The state is classified as OOD if  $d$  exceeds the task threshold,  $\alpha_g$ . We determine this threshold

187 by calculating  $d_g^\theta$ , the distance of the top  $\theta_g$  percentile of the distances to the nearest centroids from  
188 the validation latent embeddings; we will explain how we determine  $\theta$  in the following section. We  
189 observed that using k-means instead of the nearest neighbor approach [17] yields better space and  
190 time efficiency.

191 **Adaptive Decision Boundary Auto-Adjustment.** One challenge for multi-task runtime monitoring  
192 is that task performance varies across different tasks, and moreover task performance changes along  
193 deployment updates. We make our error predictors adapt to the varying performance of different  
194 tasks in the multi-task deployment. Our intuition is to assign looser decision boundaries to tasks  
195 with higher performance, which means the policy needs less supervision; and tighter boundaries to  
196 tasks with lower performance, allowing the error predictor to be more conservative and ask for more  
197 human monitoring to supervise the policy execution. We realize this idea for the failure predictor  
198 and OOD predictor respectively as follows. For the failure classifier, we fine-tune the classifier  
199 model on the deployment data from the most recent round,  $\mathcal{D}'$ , with the newest human intervention  
200 labels that reflects the most up-to-date risk perception of the robot on this task.

201 For the OOD detector, our intuition is that 1) we should enlarge the boundary for a state to be  
202 “out-of-distribution” since the robot has seen more states at deployment, and hence more states  
203 would be in-distribution; 2) the classification criteria changes as the robot policy improves over a  
204 task. When selecting embeddings to perform k-means clustering, we sample from  $D^i$ , which the  
205 past deployment data aggregated across previous rounds. As discussed in the previous section, to  
206 do OOD classification for an embedding  $z$ , we use a threshold  $\alpha_g$  for its distance  $d$  to its nearest  
207 centroid, which depends on  $\theta_g$ . Ideally,  $\theta_g$  should depend on the policy performance on this task  
208  $g$ . We use the human intervention ratio from the most recent round,  $p_H$ , as a proxy for the policy  
209 performance, and set  $\theta_g$  to be a function of the human intervention ratio from the last round. We  
210 model this function as a exponential decay function,  $\theta_g = a + b e^{-c p_H}$ , and obtain the parameters  
211 by fitting the curve to a few test trajectories from different tasks to capture the relationship between  
212 the two values. We found this function to be robust across rounds or tasks, and we use the same  
213 function throughout the simulation and real world experiments. More discussion can be found in the  
214 Appendix.

### 215 3.2.3 FIREWORK in Operation

216 Here we discuss the design choices of training and running FIREWORK in operation.

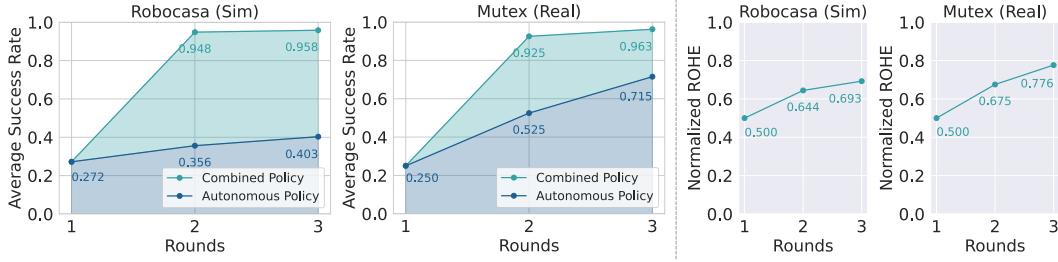
217 **Continual Model Improvement.** The visual world model  $\mathcal{W}$  is trained once and frozen afterwards,  
218 serving as a backbone for the general representation. Meanwhile, we want our system to contin-  
219 ually improve over time. Therefore, the policy and the error predictors are being finetuned over  
220 deployment: the policy is continually finetuned on the aggregation of previous rounds of data with  
221 weighted sampling, the failure classifier is finetuned over the last round of intervention labels while  
222 the OOD classifier increases its latent space coverage every round.

223 **Error Predictors in Operation.** For failure and OOD predictions, the visual world model is used  
224 for inference to predict the latent embeddings  $L$  steps ahead of the current frame. We do the  $L$ -step  
225 prediction  $N$  times to sample  $N$  possible futures by drawing samples from the cVAE latent space.  
226 Each predicted future embedding is individually assessed by the error predictors. For the failure  
227 prediction, we average the failure label for each future embedding; for OOD prediction, we average  
228 the distances from each of the future states to their nearest clustering centroid, and compare the final  
229 average distance to the OOD threshold.

230 **Training the Multi-Task Policy.** We train a Transformer-based [63] multi-task policy shown in  
231 Figure 2, which inputs image observation, robot proprioceptive states, and task language goal, and  
232 output robot motor actions, following the design from robomimic [32] and robocasa [39].

## 233 4 Experiments

234 In our experiments, we seek to answer the following questions: How well does FIREWORK continue  
235 to improve its policy over deployment with effective runtime monitoring? 2) How does FIREWORK’s



**Figure 5: System Performance Results.** The autonomous multi-task policy continues to improve over deployment. Left: Combined and Autonomous Policy Performance. Right: ROHE. Our runtime monitoring is able to ensure high combined policy performance as well as ROHE.

236 runtime monitoring prediction compare with baselines? 3) Are error prediction results predicted at  
237 the right timings qualitatively?

#### 238 4.1 Evaluation Setup

239 To measure the effectiveness of the multi-task fleet learning, we want to see how the policy and  
240 runtime monitoring perform at the deployment and how the system performance evolves over time.  
241 Therefore, we evaluate the system performance in a human-in-the-loop setting of an iterative de-  
242 ployment loop, following prior works [64, 13, 17].

243 **Evaluation Setting.** Our evaluation process consists of 3 rounds of updates, with each round involv-  
244 ing runtime-monitoring. The initial round requires full human supervision as there is no intervention  
245 data available to train error predictors at this time. For the following 2 rounds, human supervision is  
246 only needed when an error is predicted. For each round, the deployment data and intervention labels  
247 are collected to train the policy and error predictors for the next round.

248 **Evaluation Metrics.** To evaluate the policy performance, we utilize **Autonomous Performance**  
249 to measure how effectively the policy can achieve its goals without human intervention. For the  
250 effectiveness of runtime-monitoring, we use Combined Policy Performance and Return of Human  
251 Effort (ROHE) [12, 17]. **Combined Policy Performance** measures the overall system performance.  
252 **Return of Human Effort (ROHE)** evaluates how well the system utilizes human effort.

253 **Evaluation Environments.** We evaluate our system both in the simulator and real-world environ-  
254 ment. For simulation, we use robocasa [39] as our simulator benchmark, a challenging benchmark  
255 due to its visually diverse tasks and environments with a wide range of objects, layouts, and scenes.  
256 The visual world model is trained on all task data generated from mimicgen [38] (5k trajectories  $\times$   
257 20 task suites). The policy is trained on 12 task suites with human demonstration data (50 trajec-  
258 tories per task). The policy takes two workspace camera images (left and right view), and the wrist  
259 camera image. For real world, we use mutex [27] as our real world benchmark. The visual world  
260 model is trained on 50 tasks from the dataset, and the policy is trained on 10 tasks from 5 task suites.  
261 The policy takes in one workspace camera image and one wrist camera image.

	MoMaRT	PATO	ThriftyDagger	FIREWORK
Combined Policy Performance	0.567	0.873	0.733	<b>0.958</b>
ROHE	0.434	0.522	0.518	<b>0.693</b>

**Table 1: Baseline Results Comparison.** FIREWORK surpassed the baselines both in combined policy performance and ROHE.

#### 262 4.2 Evaluation Results

263 *System Performance Over Time.* We present our system-level performance results in Figure 5. Our  
264 first finding is that the autonomous policy consistently improves over the deployment rounds, with  
265 13% increase in simulation and 45% in the real world. Second, FIREWORK’s combined policy  
266 performance is consistently high, getting higher than 95% in both simulation and in real world.  
267 This shows that the runtime monitoring mechanism is effective at monitoring the deployment and

268 requesting human help when necessary. Third, FIREWORK also displays better ROHE performance  
 269 over time. With finetuning from the deployment data, and a adaptive threshold, the error predictors  
 270 At the same time, the improving policy also fails less, which makes the error predictors predict  
 271 fewer errors. This shows the interplay of better policy trained from the deployment data and more  
 272 confident error predictors.

273 *Baseline Comparison.* How does FIREWORK’s multi-  
 274 task runtime monitoring module with visual world model  
 275 compare with prior methods? We compare FIREWORK  
 276 with three prior works in runtime monitoring: **MoMart**  
 277 [20], **PATO** [21] and **ThriftyDagger** [12]. **MoMart** per-  
 278 form OOD detection using the VAE reconstruction loss  
 279 of current input images. **PATO** performs OOD detec-  
 280 tion by combining the variance of ensemble policy and  
 281 the variance of VAE reconstruction of future image goals.  
 282 **ThriftyDagger** combines OOD detection and failure de-  
 283 tection, by measuring the the variance of ensemble policy  
 284 and using a risky Q function to estimate the risky value  
 285 of the current state. Since the baselines are originally de-  
 286 signed for single task environments, we evaluate with 5  
 287 tasks in robocasa, each chosen from the 5 different task  
 288 categories. For each baseline, we train a separate model  
 289 for each of the tasks. We use the same multi-task policy  
 290 as FIREWORK, and compare using the runtime monitor-  
 291 ing metrics for one deployment round. We implement  
 292 these baselines following their original designs. Specifi-  
 293 cally, **MoMart** and **PATO** have a defined fixed threshold  
 294 for their detection criteria. **ThriftyDagger**’s threshold is  
 295 defined by the desired intervention ratio [12]. It calculate the risk value threshold based on the  
 296  $(1 - K)$ -th percentile of all risk values in the data, where  $K$  is the desired human intervention ra-  
 297 dio. For a fair comparison, we use the same human intervention ratio as FIREWORK in the same  
 298 deployment round.

299 Table 1 shows the comparison of runtime monitoring metrics of FIREWORK with the baselines.  
 300 FIREWORK surpasses the baselines both in combined policy performance and ROHE. One of our  
 301 observation is that one fixed threshold value does not work for all the tasks. Due to the difference  
 302 in task distribution and performance, the same threshold might be predicting all intervention for one  
 303 task and zero intervention for the other, incurring high human workload for one and hurting task  
 304 performance in another.

305 *Analysis and Discussion.* Can FIREWORK catch important errors and give meaningful intervention  
 306 timings? We conduct qualitative analysis of where our error predictors predict error during robot  
 307 execution, using active human monitoring. Specifically, we have a human operator fully supervise  
 308 the robot execution of a policy and can intervene whenever an unsafe state is observed. We then  
 309 apply the learned error predictors to the collected trajectories, and compare failure states identified  
 310 by the error predictors with those intervened by the human operator. Figure 6 shows examples of  
 311 the human intervention region and OOD / Failure errors predicted by FIREWORK, corresponding to  
 312 the human risk assessment. We present more model analysis and discussion in the Appendix.

## 313 5 Conclusion

314 We presented FIREWORK a framework of multi-task interactive robot fleet learning. This consists of  
 315 a multi-task robot policy and a runtime monitoring mechanism with visual world model backbone  
 316 and task-adaptive error prediction thresholds. We showed the high performance of this multi-task  
 317 fleet learning system both in simulation and on real robots.

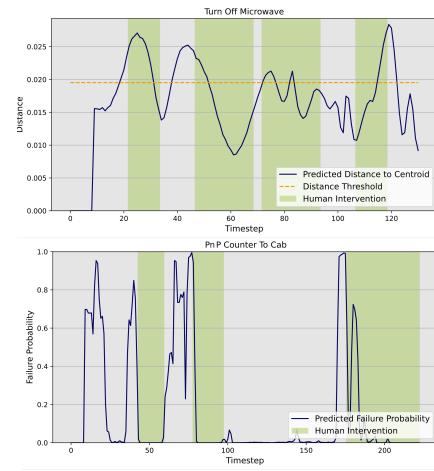


Figure 6: **Visualization of the error predictor timings compared with human intervention.** Green area indicates actual human intervention. Top: OOD Prediction. Bottom: Failure Prediction.

318 **References**

- 319 [1] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-  
320 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi,  
321 R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manju-  
322 nath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao,  
323 M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran,  
324 V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-1:  
325 Robotics transformer for real-world control at scale, 2022.
- 326 [2] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez,  
327 Y. Sulsky, J. Kay, J. T. Springenberg, T. Eccles, J. Bruce, A. Razavi, A. Edwards, N. Heess,  
328 Y. Chen, R. Hadsell, O. Vinyals, M. Bordbar, and N. de Freitas. A generalist agent, 2022.
- 329 [3] D. Shah, A. Sridhar, N. Dashora, K. Stachowicz, K. Black, N. Hirose, and S. Levine. Vint: A  
330 foundation model for visual navigation, 2023.
- 331 [4] R. Hoque, L. Y. Chen, S. Sharma, K. Dharmarajan, B. Thananjeyan, P. Abbeel, and K. Gold-  
332 berg. Fleet-dagger: Interactive robot fleet learning with scalable human supervision. In *Con-  
333 ference on Robot Learning*, pages 368–380. PMLR, 2023.
- 334 [5] L. Wang, K. Zhang, A. Zhou, M. Simchowitz, and R. Tedrake. Robot fleet learning via policy  
335 merging, 2024.
- 336 [6] M. Müller, S. Brahmbhatt, A. Deka, Q. Leboutet, D. Hafner, and V. Koltun. Openbot-fleet: A  
337 system for collective learning with real robots, 2024.
- 338 [7] A. Mandlekar, D. Xu, R. Martín-Martín, Y. Zhu, L. Fei-Fei, and S. Savarese. Human-in-the-  
339 loop imitation learning using remote teleoperation. In *arXiv preprint arXiv:2012.06733*, 2020.
- 340 [8] E. Chisari, T. Welschehold, J. Boedecker, W. Burgard, and A. Valada. Correct me if i am  
341 wrong: Interactive learning for robotic manipulation. In *RAL*, volume 7, pages 3695–3702,  
342 2021.
- 343 [9] J. Spencer, S. Choudhury, M. Barnes, M. Schmittle, M. Chiang, P. Ramadge, and S. Srinivasa.  
344 Learning from interventions: Human-robot interaction as both explicit and implicit feedback.  
345 In *16th Robotics: Science and Systems, RSS 2020*. MIT Press Journals, 2020.
- 346 [10] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. Hg-dagger: Interac-  
347 tive imitation learning with human experts. *2019 International Conference on Robotics and  
348 Automation (ICRA)*, May 2019. doi:10.1109/icra.2019.8793698. URL <http://dx.doi.org/10.1109/ICRA.2019.8793698>.
- 349 [11] R. Hoque, A. Balakrishna, C. Puterman, M. Luo, D. S. Brown, D. Seita, B. Thananjeyan,  
350 E. Novoseller, and K. Goldberg. Lazydagger: Reducing context switching in interactive im-  
351 itation learning. In *2021 IEEE 17th International Conference on Automation Science and  
352 Engineering (CASE)*, pages 502–509. IEEE, 2021.
- 353 [12] R. Hoque, A. Balakrishna, E. Novoseller, A. Wilcox, D. S. Brown, and K. Goldberg.  
354 Thriftydagger: Budget-aware novelty and risk gating for interactive imitation learning. *arXiv  
355 preprint arXiv:2109.08273*, 2021.
- 356 [13] H. Liu, S. Nasiriany, L. Zhang, Z. Bao, and Y. Zhu. Robot learning on the job: Human-in-  
357 the-loop autonomy and learning during deployment. In *Robotics: Science and Systems (RSS)*,  
358 2023.
- 359 [14] R. Hoque, L. Y. Chen, S. Sharma, K. Dharmarajan, B. Thananjeyan, P. Abbeel, and K. Gold-  
360 berg. Fleet-dagger: Interactive robot fleet learning with scalable human supervision, 2022.

- 362 [15] Q. Li, Z. Peng, and B. Zhou. Efficient learning of safe driving policy via human-ai copilot  
363 optimization, 2022.
- 364 [16] E. Yel and N. Bezzo. Fast run-time monitoring, replanning, and recovery for safe autonomous  
365 system operations. In *2019 IEEE/RSJ International Conference on Intelligent Robots and*  
366 *Systems (IROS)*, pages 1661–1667, 2019. doi:10.1109/IROS40897.2019.8968498.
- 367 [17] H. Liu, S. Dass, R. Martín-Martín, and Y. Zhu. Model-based runtime monitoring with interac-  
368 tive imitation learning. In *IEEE International Conference on Robotics and Automation (ICRA)*,  
369 2024.
- 370 [18] M. Salehi, H. Mirzaei, D. Hendrycks, Y. Li, M. H. Rohban, and M. Sabokrou. A unified survey  
371 on anomaly, novelty, open-set, and out-of-distribution detection: Solutions and future chal-  
372 lenges. *CoRR*, abs/2110.14051, 2021. URL <https://arxiv.org/abs/2110.14051>.
- 373 [19] C. Richter and N. Roy. Safe visual navigation via deep learning and novelty detection.  
374 *Robotics: Science and Systems*, 07 2017.
- 375 [20] J. Wong, A. Tung, A. Kurenkov, A. Mandlekar, L. Fei-Fei, S. Savarese, and R. Martín-Martín.  
376 Error-aware imitation learning from teleoperation data for mobile manipulation. In *Proceed-  
377 ings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learn-  
378 ing Research*, pages 1367–1378, 2022.
- 379 [21] S. Dass, K. Pertsch, H. Zhang, Y. Lee, J. J. Lim, and S. Nikolaidis. Pato: Policy assisted  
380 teleoperation for scalable robot data collection, 2022.
- 381 [22] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman,  
382 E. Luhman, C. Ng, R. Wang, and A. Ramesh. Video generation models as world simulators.  
383 2024.
- 384 [23] J. Bruce, M. Dennis, A. Edwards, J. Parker-Holder, Y. Shi, E. Hughes, M. Lai, A. Mavalankar,  
385 R. Steigerwald, C. Apps, Y. Aytar, S. Bechtle, F. Behbahani, S. Chan, N. Heess, L. Gonzalez,  
386 S. Osindero, S. Ozair, S. Reed, J. Zhang, K. Zolna, J. Clune, N. de Freitas, S. Singh, and  
387 T. Rocktäschel. Genie: Generative interactive environments, 2024.
- 388 [24] J. Xiang, G. Liu, Y. Gu, Q. Gao, Y. Ning, Y. Zha, Z. Feng, T. Tao, S. Hao, Y. Shi, Z. Liu, E. P.  
389 Xing, and Z. Hu. Pandora: Towards general world model with natural language actions and  
390 video states. 2024.
- 391 [25] A. Hu, L. Russell, H. Yeo, Z. Murez, G. Fedoseev, A. Kendall, J. Shotton, and G. Corrado.  
392 Gaia-1: A generative world model for autonomous driving, 2023.
- 393 [26] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu,  
394 and L. Fan. Vima: General robot manipulation with multimodal prompts. In *Fortieth Interna-  
395 tional Conference on Machine Learning*, 2023.
- 396 [27] R. Shah, R. Martín-Martín, and Y. Zhu. Mutex: Learning unified policies from multimodal  
397 task specifications. In *7th Annual Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=PwqiqaEzJ>.
- 398 [28] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,  
399 A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman,  
400 A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal,  
401 L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao,  
402 K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut,  
403 H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao,  
404 P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web  
405 knowledge to robotic control, 2023.

- 407 [29] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. Bc-z:  
 408 Zero-shot task generalization with robotic imitation learning, 2022.
- 409 [30] K. Bousmalis, G. Vezzani, D. Rao, C. Devin, A. X. Lee, M. Bauza, T. Davchev, Y. Zhou,  
 410 A. Gupta, A. Raju, A. Laurens, C. Fantacci, V. Dalibard, M. Zambelli, M. Martins, R. Pevce-  
 411 viciute, M. Blokzijl, M. Denil, N. Batchelor, T. Lampe, E. Parisotto, K. Źołna, S. Reed, S. G.  
 412 Colmenarejo, J. Scholz, A. Abdolmaleki, O. Groth, J.-B. Regli, O. Sushkov, T. Rothörl, J. E.  
 413 Chen, Y. Aytar, D. Barker, J. Ortiz, M. Riedmiller, J. T. Springenberg, R. Hadsell, F. Nori, and  
 414 N. Heess. Robocat: A self-improving generalist agent for robotic manipulation, 2023.
- 415 [31] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,  
 416 T. Kreiman, C. Xu, J. Luo, Y. L. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh,  
 417 C. Finn, and S. Levine. Octo: An open-source generalist robot policy, 2024.
- 418 [32] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese,  
 419 Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations  
 420 for robot manipulation, 2021.
- 421 [33] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models, 2020.
- 422 [34] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion  
 423 policy: Visuomotor policy learning via action diffusion, 2024.
- 424 [35] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and  
 425 S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets,  
 426 2021.
- 427 [36] Open X-Embodiment Collaboration et al. Open X-Embodiment: Robotic learning datasets  
 428 and RT-X models. In *Proceedings of the IEEE International Conference on Robotics and*  
 429 *Automation (ICRA)*, Yokohama, Japan, 2024.
- 430 [37] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,  
 431 M. K. Srirama, L. Y. Chen, et al. Droid: A large-scale in-the-wild robot manipulation dataset,  
 432 2024.
- 433 [38] A. Mandlekar, S. Nasiriany, B. Wen, I. Akinola, Y. Narang, L. Fan, Y. Zhu, and D. Fox.  
 434 Mimicgen: A data generation system for scalable robot learning using human demonstrations.  
 435 In *7th Annual Conference on Robot Learning*, 2023.
- 436 [39] S. Nasiriany, A. Maddukuri, L. Zhang, A. Parikh, A. Lo, A. Joshi, A. Mandlekar, and Y. Zhu.  
 437 Robocasa: Large-scale simulation of everyday tasks for generalist robots. In *Robotics: Science*  
 438 and *Systems (RSS)*, 2024.
- 439 [40] G. Swamy, S. Reddy, S. Levine, and A. D. Dragan. Scaled autonomy: Enabling human opera-  
 440 tors to control robot fleets, 2020.
- 441 [41] O. Akcin, P.-h. Li, S. Agarwal, and S. P. Chinchali. Decentralized data collection for robotic  
 442 fleet learning: A game-theoretic approach. In K. Liu, D. Kulic, and J. Ichnowski, edi-  
 443 tors, *Proceedings of The 6th Conference on Robot Learning*, volume 205 of *Proceedings*  
 444 of *Machine Learning Research*, pages 978–988. PMLR, 14–18 Dec 2023. URL <https://proceedings.mlr.press/v205/akcin23a.html>.
- 446 [42] O. Akcin, O. Unuvar, O. Ure, and S. P. Chinchali. Fleet active learning: A submodular max-  
 447 imization approach. In J. Tan, M. Toussaint, and K. Darvish, editors, *Proceedings of The*  
 448 *7th Conference on Robot Learning*, volume 229 of *Proceedings of Machine Learning Re-*  
 449 *search*, pages 1378–1399. PMLR, 06–09 Nov 2023. URL <https://proceedings.mlr.press/v229/akcin23a.html>.

- 451 [43] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly,  
452 M. Kalakrishnan, V. Vanhoucke, and S. Levine. QT-Opt: Scalable deep reinforcement learning  
453 for vision-based robotic manipulation. In *CoRL*, 2018.
- 454 [44] A. . Team, J. Aldaco, T. Armstrong, R. Baruch, J. Bingham, S. Chan, K. Draper, D. Dwibedi,  
455 C. Finn, P. Florence, S. Goodrich, W. Gramlich, T. Hage, A. Herzog, J. Hoech, T. Nguyen,  
456 I. Storz, B. Tabanpour, L. Takayama, J. Tompson, A. Wahid, T. Wahrburg, S. Xu,  
457 S. Yaroshenko, K. Zakka, and T. Z. Zhao. Aloha 2: An enhanced low-cost hardware for  
458 bimanual teleoperation, 2024.
- 459 [45] M. Ahn, D. Dwibedi, C. Finn, M. G. Arenas, K. Gopalakrishnan, K. Hausman, B. Ichter,  
460 A. Irpan, N. Joshi, R. Julian, S. Kirmani, I. Leal, E. Lee, S. Levine, Y. Lu, I. Leal, S. Maddineni,  
461 K. Rao, D. Sadigh, P. Sanketi, P. Sermanet, Q. Vuong, S. Welker, F. Xia, T. Xiao, P. Xu, S. Xu,  
462 and Z. Xu. Autort: Embodied foundation models for large scale orchestration of robotic agents,  
463 2024.
- 464 [46] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic  
465 grasping with deep learning and large-scale data collection, 2016.
- 466 [47] K.-C. Hsu, H. Hu, and J. F. Fisac. The safety filter: A unified view of safety-critical con-  
467 trol in autonomous systems. 2023. URL <https://api.semanticscholar.org/CorpusID:261697421>.
- 468 [48] R. Sinha, E. Schmerling, and M. Pavone. Closing the loop on runtime monitors with fallback-  
469 safe mpc. In *Proc. IEEE Conf. on Decision and Control*, 2023.
- 470 [49] K. Nakamura, R. Tian, and A. V. Bajcsy. A general calibrated regret metric for detecting and  
471 mitigating human-robot interaction failures. *ArXiv*, abs/2403.04745, 2024. URL <https://api.semanticscholar.org/CorpusID:268264177>.
- 472 [50] R. Sinha, A. Sharma, S. Banerjee, T. Lew, R. Luo, S. M. Richards, Y. Sun, E. Schmerling, and  
473 M. Pavone. A system-level view on out-of-distribution data in robotics, 2022.
- 474 [51] K. Menda, K. Driggs-Campbell, and M. J. Kochenderfer. EnsembleDagger: A bayesian ap-  
475 proach to safe imitation learning. In *2019 IEEE/RSJ International Conference on Intelligent  
476 Robots and Systems (IROS)*, pages 5041–5048. IEEE, 2019.
- 477 [52] A. Diryag, M. Mitić, and Z. Miljković. Neural networks for prediction of robot failures. *Pro-  
478 ceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering  
479 Science*, 228(8):1444–1458, 2014.
- 480 [53] C. Gokmen, D. Ho, and M. Khansari. Asking for help: Failure prediction in behavioral cloning  
481 through value approximation. *ArXiv*, abs/2302.04334, 2023.
- 482 [54] A. Xie, F. Tajwar, A. Sharma, and C. Finn. When to ask for help: Proactive interventions in  
483 autonomous reinforcement learning, 2022.
- 484 [55] A. Z. Ren, A. Dixit, A. Bodrova, S. Singh, S. Tu, N. Brown, P. Xu, L. Takayama, F. Xia,  
485 J. Varley, Z. Xu, D. Sadigh, A. Zeng, and A. Majumdar. Robots that ask for help: Uncertainty  
486 alignment for large language model planners, 2023.
- 487 [56] Z. Liu, A. Bahety, and S. Song. Reflect: Summarizing robot experiences for failure explanation  
488 and correction, 2023.
- 489 [57] L. Guan, Y. Zhou, D. Liu, Y. Zha, H. B. Amor, and S. Kambhampati. "task success" is not  
490 enough: Investigating the use of video-language models as behavior critics for catching unde-  
491 sirable agent behaviors, 2024.
- 492 [58] A. Elhafsi, R. Sinha, C. Agia, E. Schmerling, I. Nesnas, and M. Pavone. Semantic anomaly  
493 detection with large language models, 2023.

- 496 [59] A. Hu, G. Corrado, N. Griffiths, Z. Murez, C. Gurau, H. Yeo, A. Kendall, R. Cipolla, and  
497 J. Shotton. Model-based imitation learning for urban driving, 2022.
- 498 [60] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image  
499 segmentation, 2015.
- 500 [61] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional  
501 generative models. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, edi-  
502 tors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.,  
503 2015. URL [https://proceedings.neurips.cc/paper\\_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/8d55a249e6baa5c06772297520da2051-Paper.pdf).
- 504 [62] P. Wu, A. Escontrela, D. Hafner, P. Abbeel, and K. Goldberg. Daydreamer: World models for  
505 physical robot learning. In *Conference on Robot Learning*, pages 2226–2240. PMLR, 2023.
- 506 [63] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polo-  
507 sukhin. Attention is all you need, 2023.
- 508 [64] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta,  
509 E. Orbay, S. Savarese, and L. Fei-Fei. Roboturk: A crowdsourcing platform for robotic skill  
510 learning through imitation. In *CoRL*, pages 879–893, 2018.