

Fleet-Dagger: Interactive Robot Fleet Learning with Scalable Human Supervision

Anonymous Author(s)

Affiliation

Address

email

Abstract:

Commercial and industrial deployments of robot fleets often fall back on remote human teleoperators during execution when robots are at risk or unable to make task progress. With continual learning, interventions from the remote pool of humans can also be used to improve the robot fleet control policy over time. A central question is how to effectively allocate limited human attention to individual robots. Prior work addresses this in the single-robot, single-human setting. We formalize the Interactive Fleet Learning (IFL) setting, in which multiple robots interactively query and learn from multiple human supervisors. We present a fully implemented open-source IFL benchmark suite of GPU-accelerated Isaac Gym environments for the evaluation of IFL algorithms. We propose Fleet-Dagger, a family of IFL algorithms, and compare a novel Fleet-Dagger algorithm to 4 baselines in simulation. We also perform 1000 trials of a physical block-pushing experiment with 4 ABB YuMi robot arms. Experiments suggest that the allocation of humans to robots significantly affects the performance of the fleet, and that our algorithm achieves up to $8.8\times$ higher return on human effort than baselines. See <https://sites.google.com/view/fleet-dagger> for supplemental material.

Keywords: Fleet Learning, Interactive Learning, Human Robot Interaction

1 Introduction

Waymo, Zoox, Amazon, Nimble Robotics, and Plus One Robotics use remote human supervision of robot fleets in applications ranging from self-driving taxis to automated warehouse fulfillment [1, 2, 3, 4]. These robots intermittently cede control during task execution to remote human supervisors for corrective interventions. The interventions take place either during learning, when they are used to improve the robot policy, or at execution time, when the policy is no longer updated but robots can still request human assistance when needed to improve reliability. In the *continual learning* setting, these occur simultaneously: the robot policy has been deployed but continues to be updated indefinitely with additional intervention data. Furthermore, any individual robot can share its intervention data with the rest of the fleet. As opposed to robot swarms that must coordinate with each other to achieve a common objective, a robot *fleet* is a set of independent robots simultaneously executing the same control policy for the same task in parallel environments. We refer to the setting of a robot fleet learning via interactive requests for human supervision as *Interactive Fleet Learning (IFL)*.

Of central importance in IFL is the supervisor allocation problem: how should limited human supervision be allocated to robots in a manner that maximizes the throughput of the fleet? Prior work studies this in the single-robot, single-human case. A variety of interactive learning algorithms have been proposed that estimate quantities such as uncertainty [5], novelty [6, 7, 8], risk [7, 8], and predicted action discrepancy [9, 10]. However, it remains unclear which algorithms are the most effective when generalized to the multi-human, multi-robot case.

38 To this end, we formalize the IFL problem and present the IFL Benchmark (IFLB), a new open-
39 source Python toolkit and benchmark for developing and evaluating human-to-robot allocation
40 algorithms. The IFLB includes environments from Isaac Gym [11], which enabled efficient simulation
41 of thousands of learning robots for the first time in 2021. This paper makes the following contributions:
42 (1) a formalism for interactive fleet learning, (2) the IFLB, an open-source software benchmark and
43 toolkit for IFL algorithms with 3 Isaac Gym environments for complex robotic tasks, (3) Fleet-
44 DAgger, a novel family of IFL algorithms, (4) results from large-scale simulation experiments with a
45 fleet of 100 robots, and (5) results from 1000 physical trials with 4 physical robot arms and 2 human
46 supervisors providing teleoperation remotely over the Internet.

47 **2 Related Work**

48 **2.1 Allocating Human Supervisors to Robots at Execution Time**

49 For human-robot teams, deciding when to transfer control between robots and humans during
50 execution is a widely studied topic in the literature of both sliding autonomy [12, 13, 14] and Human-
51 Robot Interaction (HRI). In sliding autonomy, also known as adjustable autonomy [15, 16] or adaptive
52 automation [17], humans and robots dynamically adjust their level of autonomy and transfer control
53 to each other during execution [14, 17]. Since identifying which robot to assist in a large robot fleet
54 can be overwhelming for a human operator [18, 19, 20, 21, 22], several strategies, such as using
55 a cost-benefit analysis to decide whether to request operator assistance [12] and using an advising
56 agent to filter robot requests [23], have been proposed to improve the performance of human-robot
57 teams [23, 21, 24] and increase the number of robots that can be controlled [25], a quantity known
58 as “fan-out” [26]. Other examples include user modeling [12, 13, 21, 24] and studying interaction
59 modes [27] for better system and interface design [28, 29]. Zheng et al. [30] propose to compute the
60 time until stopping for each robot based on its estimated risk and prioritize the robots accordingly. Ji
61 et al. [31] consider the setting where physical assistance is required to resume tasks for navigation
62 robots and formalize single-human, multi-robot allocation as graph traversal. Dahiya et al. [32]
63 formulate the problem of multi-human, multi-robot allocation during robot execution as a Restless
64 Multi-Armed Bandit problem. In addition, allocating humans to address robot requests has also been
65 studied from the perspective of queueing theory and scheduling theory [33, 34, 22, 35, 36, 37]. The
66 vast majority of the human-robot teaming and queueing theory work, however, does not involve
67 learning; the robot control policies are assumed to be fixed. In contrast, we study supervisor allocation
68 during robot learning, where the effectiveness of allocation policies affects not only human burden
69 and the performance of human-robot team but also the efficiency of policy learning.

70 **2.2 Single-Human, Single-Robot Interactive Learning**

71 Imitation learning (IL) is a paradigm of robot learning in which a robot uses demonstrations from a
72 human to initialize and/or improve its policy [38, 39, 40, 41, 42, 43]. However, learning from purely
73 offline data often suffers from distribution shift [44, 45], as compounding approximation error leads
74 to states that were not visited by the human. This can be mitigated with online data collection with
75 algorithms such as Dataset Aggregation (DAgger) [44] and interactive imitation learning [46, 47, 48].
76 Human-gated interactive IL algorithms [49, 50, 51] require the human to monitor the robot learning
77 process and decide when to take and cede control of the system. While intuitive, these approaches
78 are not scalable to large fleets of robots or the long periods of time involved in continual learning,
79 as humans cannot effectively focus on many robots simultaneously [20, 21, 22] and are prone to
80 fatigue [52]. To reduce the burden on the supervisor, several robot-gated interactive IL algorithms
81 such as SafeDAgger [9], EnsembleDAgger [5], LazyDAgger [10], and ThriftyDAgger [8] have been
82 proposed, in which the robot actively solicits human interventions when certain criteria are met.
83 Interactive reinforcement learning (RL) [53, 54, 55, 56] is another active area of research in which
84 robots learn from both online human feedback and their own experience. However, these interactive
85 learning algorithms are designed for and studied in the single-human, single-robot setting. Other
86 works related to single-robot interactive learning include task allocation [57]; in contrast, we focus
87 on efficient robot learning of a single task.

88 2.3 Multi-Robot Interactive Learning

89 In this paper, we study allocation policies for multiple humans and multiple robots. While many
 90 existing works [58, 59, 60, 61, 62] have leveraged NVIDIA’s Isaac Gym’s [11] capability of parallel
 91 simulation to accelerate reinforcement learning with multiple robots, these approaches consider
 92 robots in isolation without interactive human supervision. The work that is closest to ours is by
 93 Swamy et al. [63], who study the multi-robot, single-human problem of allocating the attention of
 94 one human operator during robot fleet learning. They propose to learn an internal model of human
 95 preferences as a human supervises a small fleet of 4 robots and use this model to assist the human in
 96 supervising a larger fleet of 12 robots. While this approach mitigates the scaling issue in human-gated
 97 interactive IL, a small fleet of robots can be difficult for a single human supervisor to optimally
 98 control. In contrast, we consider robot-gated algorithms for requesting and allocating supervision
 99 from multiple humans.

100 To the best of our knowledge, this work is the first to formalize and study multi-human, multi-robot
 101 interactive learning. This problem setting poses unique challenges, especially as the size of the fleet
 102 grows large relative to the number of humans, as each human allocation affects both the robot that
 103 receives supervision and the robots that do not receive human attention.

104 3 Interactive Fleet Learning Problem Formulation

105 We consider a fleet of N robots operating in parallel as a set of N independent and identical Markov
 106 decision processes (MDPs) $\{\mathcal{M}_i\}_{i=1}^N$ specified by the tuple $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$ with the same state space \mathcal{S} ,
 107 action space \mathcal{A} , unknown transition dynamics $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$,
 108 discount factor $\gamma \in [0, 1)$, and initial state distribution p_0 . We assume the MDPs have an identical
 109 indicator function $c(s) : \mathcal{S} \rightarrow \{0, 1\}$ that identifies which states $s \in \mathcal{S}$ violate a *constraint* in the
 110 MDP. States that violate MDP constraints are fault states from which the robot cannot make further
 111 progress. For instance, the robot may be stuck on the side of the road or have incurred hardware
 112 damage. We assume that the timesteps are synchronized across all robots and that they share the
 113 same non-stationary policy $\pi_{\theta_t} : \mathcal{S} \rightarrow \mathcal{A}$, parameterized by θ_t at each timestep t , where θ_t is updated
 114 over time with continual learning.

115 The collection of $\{\mathcal{M}_i\}_{i=1}^N$ can be reformulated as a single MDP $\mathcal{M} = (\mathcal{S}^N, \mathcal{A}^N, \bar{p}, \bar{r}, \gamma)$, composed
 116 of vectorized states and actions of all robots in the fleet (denoted by bold font) and joint transition
 117 dynamics. In particular, $\mathbf{s} = (s_1, \dots, s_N) \in \mathcal{S}^N$, $\mathbf{a} = (a_1, \dots, a_N) \in \mathcal{A}^N$, $\bar{p}(\mathbf{s}^{t+1} | \mathbf{s}^t, \mathbf{a}^t) =$
 118 $\prod_{i=1}^N p(s_i^{t+1} | s_i^t, a_i^t)$, and $\bar{r}(\mathbf{s}, \mathbf{a}) = \sum_{i=1}^N r(s_i, a_i)$.

119 We assume that robots can query a set of $M \ll N$ human supervisors for assistance interactively
 120 (i.e., during execution of π_{θ_t}). We assume that each human can help only one robot at a time and that
 121 all humans have the same policy $\pi_H : \mathcal{S} \rightarrow \mathcal{A}_H$, where $\mathcal{A}_H = \mathcal{A} \cup \{R\}$ and R is a *hard reset*, an
 122 action that resets the MDP to the initial state distribution $s^0 \sim p_0$. As opposed to a *soft reset* that
 123 can be performed autonomously by the robot via a reset action $r \in \mathcal{A}$ (e.g., a new bin arrives in an
 124 assembly line), a *hard reset* requires human intervention due to constraint violation (i.e., entering
 125 some s where $c(s) = 1$). A human assigned to a robot either performs hard reset R (if $c(s) = 1$)
 126 or teleoperates the robot system with policy π_H (if $c(s) = 0$). A hard reset R takes t_R timesteps to
 127 perform, and all other actions take 1 timestep.

128 Supervisor allocation (i.e., the assignment of humans to robots) is determined by an allocation policy

$$\omega : (\mathbf{s}^t, \pi_{\theta_t}, \boldsymbol{\alpha}^{t-1}, \mathbf{x}^t) \mapsto \boldsymbol{\alpha}^t \in \{0, 1\}^{M \times N} \quad \text{s.t.} \quad \sum_{j=1}^N \alpha_{ij}^t \leq 1 \quad \text{and} \quad \sum_{i=1}^M \alpha_{ij}^t \leq 1 \quad \forall i, j, \quad (1)$$

129 where \mathbf{s}^t are the current states for each of the robots, $\boldsymbol{\alpha}^t$ is an $M \times N$ binary matrix that indicates
 130 which robots will receive assistance from which human at the current timestep t , and \mathbf{x}^t is an
 131 augmented state containing any auxiliary information for each robot, such as the type and duration
 132 of an ongoing intervention. Unlike Dahiya et al. [32] which studies execution-time allocation, the
 133 allocation policy ω here depends on the current robot policy π_{θ_t} , which in turn affects the speed of
 134 the policy learning. While there are a variety of potential objectives to consider, e.g., minimizing

135 constraint violations in a safety-critical environment, we define the IFL objective as *return on human*
 136 *effort (ROHE)*:

$$\max_{\omega \in \Omega} \mathbb{E}_{\tau \sim p_{\omega, \theta_0}(\tau)} \left[\frac{\sum_{t=0}^T \gamma^t \bar{r}(\mathbf{s}^t, \mathbf{a}^t)}{\sum_{t=0}^T \|\omega(\mathbf{s}^t, \pi_{\theta_t}, \boldsymbol{\alpha}^{t-1}, \mathbf{x}^t)\|_F^2} \right], \quad (2)$$

137 where Ω is the set of allocation policies, T is the total amount of time the fleet operates (rather than
 138 the time horizon of an individual task execution), θ_0 are the initial parameters of the robot policy, and
 139 $\|\cdot\|_F$ is the Frobenius norm. The objective is the expected ratio of the cumulative reward across all
 140 timesteps and all robots to the total amount of human time spent helping robots with allocation policy
 141 ω . Intuitively, the ROHE measures the performance of the robot fleet normalized by the total human
 142 effort required to achieve this performance. We provide a more thorough derivation of the ROHE
 143 objective in Appendix 8.1.

144 Since human teleoperation with π_H provides additional online data, this data can be used to update
 145 the robot policy π_{θ_t} :

$$\begin{cases} D^{t+1} \leftarrow D^t \cup D_H^t \text{ where } D_H^t := \{(s_i^t, \pi_H(s_i^t)) : \pi_H(s_i^t) \neq R \text{ and } \sum_{j=1}^N \alpha_{ij}^t = 1\} \\ \pi_{\theta_{t+1}} \leftarrow f(\pi_{\theta_t}, D^{t+1}) \end{cases}, \quad (3)$$

146 where f is a policy update function (e.g., gradient descent).

147 4 Interactive Fleet Learning Algorithms

148 4.1 Fleet-Dagger

149 Given the problem formulation above, we propose Fleet-Dagger, a family of IFL algorithms, where
 150 an *IFL algorithm* is an algorithm for supervisor allocation for IFL (i.e., it specifies an $\omega \in \Omega$ as defined
 151 in Section 3). Fleet-Dagger seeks to maximize the IFL objective in the context of interactive imitation
 152 learning with dataset aggregation: its policy update function f is supervised learning on D^t , which
 153 consists of all human data collected so far (Section 3). In addition, each Fleet-Dagger algorithm
 154 defines a unique *priority function* $\hat{p} : (s, \pi_{\theta_t}) \rightarrow [0, \infty)$ that indicates how much priority score to
 155 assign to each robot based on its state s and the current policy π_{θ_t} , where, similar to scheduling
 156 theory, a higher value indicates a higher priority robot. To reduce thrashing [10, 8], Fleet-Dagger
 157 algorithms also specify t_T , the minimum time a human supervisor must spend teleoperating a robot.

158 Fleet-Dagger uses priority function \hat{p} and t_T to define an allocation ω as follows. At each timestep t ,
 159 Fleet-Dagger first scores all robots with \hat{p} and sorts the robots by their priority values. If a human
 160 supervisor is currently performing hard reset action R and t_R timesteps have not elapsed, that human
 161 continues to help that robot. If a human is currently teleoperating a robot and the minimum t_T
 162 timesteps have not elapsed, that human continues to teleoperate the robot. If a robot with a human
 163 supervisor continues to be high priority after the minimum intervention time (t_R for a hard reset or
 164 t_T for teleoperation) has elapsed, that human remains assigned to the robot. If a human is available
 165 to help a robot, the human is reassigned to the robot with the highest priority value that is currently
 166 unassisted. Finally, if a robot has priority $\hat{p}(\cdot) = 0$, it does not receive assistance even if a human is
 167 available. We include the pseudocode of Fleet-Dagger in the supplemental material.

168 4.2 Fleet-Dagger Algorithms

169 All algorithms below specify a unique priority function \hat{p} , which is synthesized with Fleet-Dagger as
 170 described in Section 4.1 to specify an allocation ω .

171 **Behavior Cloning:** At all timesteps t , this baseline gives priority $\hat{p}(\cdot) = 1$ for robots that have
 172 violated a constraint (i.e., $c(s_i^t) = 1$) and require a hard reset, and $\hat{p}(\cdot) = 0$ for all other robots. We
 173 refer to this as C -prioritization for Constraint. Thus, the robot fleet can only receive hard resets from
 174 human supervisors (i.e., no human teleoperation). Without C -prioritization, robots that require hard
 175 resets would remain indefinitely idle and not learn.

176 **Random:** This baseline simply assigns a random priority for each robot at each timestep. Additionally,
 177 to control the total amount of human supervision, we introduce a threshold hyperparameter such that
 178 if a robot’s priority value is below the threshold, its priority is set to zero and it will not request help.

179 **Fleet-EnsembleDagger** This baseline adapts EnsembleDagger [5] to the IFL setting. EnsembleDagger
180 uses the output variance among an ensemble of neural networks bootstrapped on subsets of the
181 training data as an estimate of epistemic uncertainty; accordingly, we define the robot priority for
182 Fleet-EnsembleDagger as ensemble variance. Since ensemble variance is designed for continuous
183 action spaces, for environments with discrete action spaces we instead estimate the uncertainty with
184 the Shannon entropy [64] among the outputs of a single classifier network. We refer to this priority
185 function as U -prioritization for Uncertainty. Finally, since EnsembleDagger was not designed for
186 environments with constraint violations and idle robots will negatively affect the ROHE, we add
187 C -prioritization for a more fair comparison. Specifically, given an uncertainty threshold value,
188 robots with uncertainty above threshold are prioritized first in order of their uncertainty, followed by
189 constraint-violating robots.

190 **Fleet-ThriftyDagger:** This baseline adapts the ThriftyDagger algorithm [8] to the IFL setting.
191 ThriftyDagger uses a synthesis of uncertainty (which we refer to as the U -prioritization value)
192 and the probability of task failure (estimated with a goal critic Q-function) to query a human for
193 supervision. Since Fleet-Dagger requires a single metric by which to compare different robots, we
194 adapt ThriftyDagger to the fleet setting by calculating a linear combination of the U -prioritization
195 value and the probability of task failure after normalizing each value with running estimates of their
196 means and standard deviations. As in [8], we pretrain the goal critic on an offline dataset of human
197 and robot task execution. Similar to Fleet-EnsembleDagger, we first prioritize by the combined
198 uncertainty-goal value above a parameterized threshold, followed by C -prioritization.

199 **Constraint-Uncertainty-Risk (C.U.R.):** Here we propose a novel Fleet-Dagger algorithm. As the
200 name suggests, C.U.R. does C -prioritization (prioritize all constraint-violating robots), followed by U -
201 prioritization (prioritize uncertain robots above a minimum threshold), followed by R -prioritization,
202 where R stands for risk, which we define as the probability of constraint violation. Intuitively, idle
203 robots should be reset in order to continue making progress, uncertain robots should receive more
204 human supervision in areas with little to no reference behavior to imitate, and robots at risk should
205 request human teleoperation to safety before an expensive hard reset. As in [65], we estimate the
206 probability of constraint violation with a safety critic Q-function. We initialize the safety critic on an
207 offline dataset of constraint violations. In addition, we implement an initial warmup period during
208 which constraint violations are assigned *zero priority* rather than high priority. Here, the intuition is
209 that rather than attending to hard resets for an initially low-performing policy, human intervention
210 should instead be spent on valuable teleoperation data that can improve the robot policy. Hence,
211 during the warmup period, constraint-violating robots remain idle and human attention is allocated to
212 the teleoperation of a smaller number of robots.

213 5 Interactive Fleet Learning Benchmark

214 While many algorithms have been proposed for interactive learning [10, 5, 8, 9], to our knowledge
215 there exists no unified benchmark for evaluating them. To facilitate reproducibility and standardized
216 evaluation for IFL algorithms, we introduce the Interactive Fleet Learning Benchmark (IFLB). The
217 IFLB is an open-source Python implementation of IFL with a suite of simulation environments and a
218 modular software architecture for rapid prototyping and evaluation of new IFL algorithms.

219 5.1 Environments

220 The IFLB is built on top of NVIDIA Isaac Gym [11], a highly optimized software platform for end-
221 to-end GPU-accelerated robot learning released in 2021, without which the simulation of hundreds of
222 learning robots would be computationally intractable. The IFLB currently supports the following 3
223 Isaac Gym environments with high-dimensional continuous state and action spaces: (1) **Humanoid**, a
224 bipedal legged locomotion task from OpenAI Gym [66], (2) **Anymal**, a quadruped legged locomotion
225 task with the ANYmal robot by ANYbotics, and (3) **AllegroHand**, a task involving dexterous
226 manipulation of a cube with a 4-finger Allegro Hand by Wonik Robotics. Constraint violation is
227 defined as (1) the humanoid falling down, (2) the ANYmal falling down on its torso or knees, and (3)
228 dropping the cube from the hand, respectively. See Figure 1 for images of each of these tasks. While
229 these three tasks are the current fully supported environments in the IFLB, adding new Isaac Gym

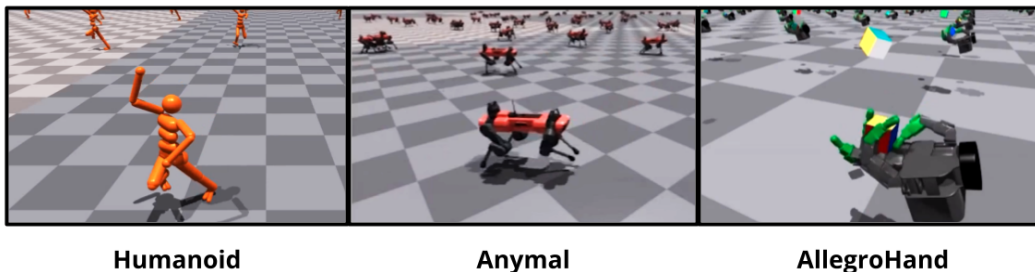


Figure 1: Isaac Gym benchmark environments in the IFLB.

230 environments to the IFLB is straightforward and will be done by the maintainers on a rolling basis
 231 (and can also easily be done by end users).

232 5.2 Software Architecture

233 The IFLB defines 3 interfaces for the development of IFL algorithms: (1) agents, (2) supervisors, and
 234 (3) allocations. An *agent* is an implementation of the robot fleet policy π_{θ_i} (Section 3), such as an
 235 IL or RL agent. A *supervisor* is an implementation of the supervisor policy π_H (Section 3), such
 236 as a fully trained RL agent, a model-based planner, or a teleoperation interface for remote human
 237 supervisors. Lastly, an *allocation* is an implementation of the priority function \hat{p} (Section 4), such as
 238 C.U.R. priority or ThriftyDagger priority. For reference, the IFLB includes an imitation learning
 239 agent, a fully trained RL supervisor using Isaac Gym’s reference PPO [67] implementation, and
 240 all allocations from Section 4, which we use in our experiments. Users of the IFLB can flexibly
 241 implement their own IFL algorithms by defining new agents, supervisors, and allocations.

242 Given an agent, supervisor, allocation, and environment, the IFLB runs Fleet-Dagger as described in
 243 Section 4.1. IFLB allows flexible command line configuration of all parameters of the experiment
 244 (e.g., t_T , t_R , N , M) as well as the parameters of the agent, supervisor, and allocation. If desired, the
 245 code can also be modified to support families of IFL algorithms other than Fleet-Dagger. A code
 246 sample for reproducing simulation results is available in the supplement and a public codebase will
 247 be released on Github after the anonymous review period.

248 6 Experiments

249 6.1 Metrics

250 Throughout online training, we measure four metrics at each timestep t : (1) the cumulative number
 251 of successful task completions across the fleet and up to time t ; (2) cumulative hard resets (i.e.,
 252 constraint violations); (3) cumulative idle time, i.e., how long robots spend idle in constraint-violating
 253 states waiting for hard resets; and (4) the return on human effort (ROHE), approximated as the
 254 ratio of cumulative task successes to the cumulative amount of human time (both hard resets and
 255 teleoperation) in hundreds of timesteps for simulation experiments and tens of timesteps for physical
 256 experiments. We do not measure ROHE for Behavior Cloning as it is an offline algorithm. For the
 257 Humanoid and Anymal locomotion environments, success is defined as reaching the episode horizon
 258 without constraint violation and with reward of at least 95% of that of the supervisor policy. For
 259 AllegroHand, a goal-conditioned task, success is defined by reaching a goal state. For the physical
 260 block-pushing experiment (Section 6.3), a robot achieves 1 success for a push in the correct direction
 261 and 5 successes for reaching the goal.

262 6.2 IFLB Simulation Experiments

263 **Experimental Setup:** We evaluate all Fleet-Dagger algorithms in the 3 benchmark simulation
 264 environments: Humanoid, Anymal, and AllegroHand. We use reinforcement learning agents fully
 265 trained with PPO [67] as the algorithmic supervisor π_H . We initialize the robot policy π_{θ_0} with
 266 behavior cloning on an offline dataset of 5000 state-action pairs. For a fair comparison, the Behavior
 267 Cloning baseline is given additional offline data equal to the average amount of human time solicited
 268 by C.U.R. by operation time boundary T . The Random baseline’s priority threshold is set such that

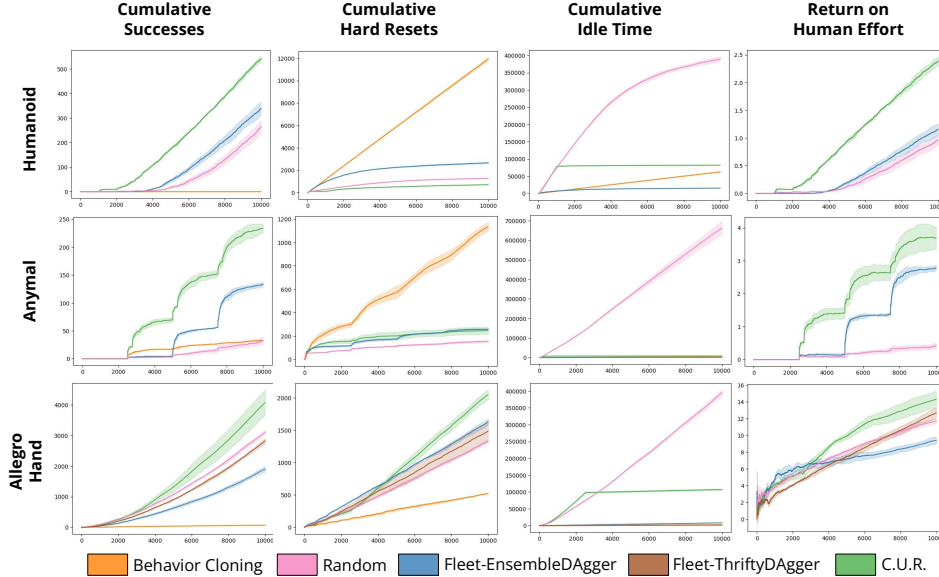


Figure 2: Simulation results in the Isaac Gym benchmark tasks with $N = 100$ robots and $M = 10$ human supervisors, where the x -axis is timesteps from 0 to $T = 10,000$. We plot the metrics described in 6.1. The C.U.R. algorithm outperforms all baselines on all environments in terms of ROHE and cumulative successes and is competitive with baselines for cumulative hard resets and idle time.

269 in expectation, it reaches the average amount of human time solicited by C.U.R. by time T . Since
 270 Fleet-ThriftyDagger requires a goal-conditioned task, it is only evaluated on AllegroHand. All
 271 training runs are executed with $N = 100$ robots, $M = 10$ humans, $t_T = 5$, $t_R = 5$, and operation
 272 time $T = 10000$, and are averaged over 3 random seeds. In the appendix, we provide hyperparameter
 273 details for each algorithm, ablation studies on each component of the C.U.R. algorithm, and an
 274 analysis of hyperparameter sensitivity to the number of humans M , minimum intervention time t_T ,
 275 and hard reset time t_R . The code supplement also provides instructions for reproducing results.

276 **Results:** We plot results in Figure 2. First, we observe that the choice of IFL algorithm has a
 277 significant impact on all metrics in all environments, indicating that allocation matters in the IFL
 278 setting. We also observe that the robot fleet achieves a higher throughput (number of cumulative task
 279 successes) with C.U.R. allocation than baselines in all environments at all times. C.U.R. also attains
 280 a higher ROHE, indicating more efficient use of human supervision. An increase in ROHE over
 281 time signifies that the improvement in the robot policy π_{θ_t} outpaces cumulative human supervision,
 282 indicating that the IFL algorithms learn not only where to assign available humans but also when
 283 to *stop* requesting unnecessary supervision. C.U.R. also incurs fewer hard resets than baselines,
 284 especially Behavior Cloning, which must constantly hard reset robots with a low-performing offline
 285 policy. For AllegroHand, however, C.U.R. incurs higher hard resets and a smaller ROHE margin
 286 over baselines. We hypothesize that since the task is too challenging to execute without human
 287 supervision in the given fleet operation time, prioritizing hard resets ironically only gives the robots
 288 additional opportunities to violate constraints. We also see that C -prioritization effectively eliminates
 289 cumulative idle time; C.U.R. idle time flattens out after the initial warmup.

290 6.3 Physical Block-Pushing Experiment

291 **Experimental Setup:** Finally, we evaluate Fleet-DAGger in a physical block-pushing experiment
 292 with $N = 4$ ABB YuMi robot arms and $M = 2$ human supervisors. See Figure 3 for the hardware
 293 setup. The objective of each robot is to reach a goal position is randomly sampled from the allowable
 294 region of the workspace. At each timestep, the robot chooses one of four discrete pushing actions
 295 corresponding to pushing each of the four vertical faces of the cube orthogonally by a fixed distance.
 296 The robot policy takes an overhead image observation of the cube in the workspace with the goal
 297 programatically generated in the image. Hard resets are physical adjustments of the cube, while
 298 teleoperation is performed over the Internet by a remote human supervisor, who specifies one of

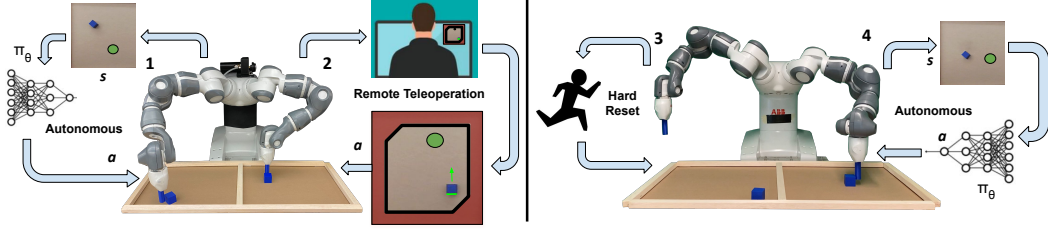


Figure 3: **Physical Task Setup:** an example timestep t in the physical experiment with 2 humans and 4 independent identical robot arms each executing the block pushing task. **Robot 1** queries robot policy π_{θ_t} for an action given an overhead image of the workspace and executes it in the workspace. **Robot 2** is teleoperated by a remote Human 1, where the human views the overhead image and specifies a pushing action through a user interface. The red region at the edges of the workspace are constraint violation regions. Human 2 is performing a physical hard reset for **Robot 3**, which has violated a constraint in a previous timestep. **Robot 4** autonomously executes the same robot policy as that of Robot 1 on its own state.

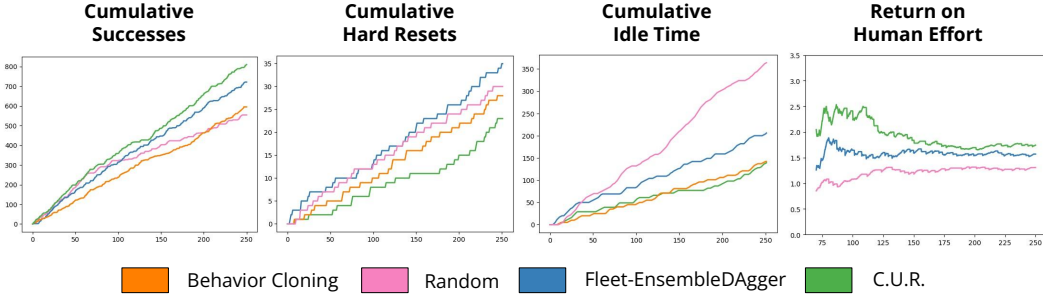


Figure 4: Physical results for the block-pushing task with 4 robots and 2 humans, where the x -axis is timesteps. C.U.R. achieves higher ROHE and cumulative successes as well as lower cumulative hard resets and idle time. However, ROHE does not noticeably improve, likely due to the low fleet operation time T .

299 the 4 pushing actions via a keyboard interface. We set $t_T = 3$, $t_R = 5$, and $T = 250$ for a total of
 300 4×250 robot actions per algorithm. All algorithms are initialized with an offline dataset of 5000
 301 image-action pairs (500 teleoperated actions with $10\times$ data augmentation).

302 **Results:** We plot results in Figure 4. We observe that the C.U.R. algorithm achieves a higher ROHE,
 303 higher cumulative successes, lower hard resets, and lower idle time than baselines, albeit by a small
 304 margin. However, the ROHE for all algorithms remains flat over time instead of improving. Results
 305 suggest that (1) the fleet operation time of $T = 250$ is too short to observe improvement for a
 306 relatively high-performing initial policy π_{θ_0} and (2) U -prioritization in its current form is less suitable
 307 for real-world multimodal human supervisors than it is for deterministic algorithmic supervisors.
 308 Since a human may arbitrarily choose one of multiple equally suitable actions, high robot uncertainty
 309 over these actions does not necessarily translate to a need for human supervision.

310 7 Limitations and Future Work

311 The IFL formulation has a number of modeling assumptions that limit its generality. (1) The human
 312 supervisors are assumed to be homogeneous, (2) all robots operate in the same state space with
 313 the same distribution, (3) all robots are independent and we do not consider robot coordination,
 314 (4) humans have perfect situational awareness [22] and can move to different robots without any
 315 switching latency, and (5) we assume a constant hard reset time. In terms of experiments, the
 316 simulations have algorithmic rather than human supervision, and the physical task is relatively
 317 straightforward with discrete planar actions.

318 In future work, we will run more physical experiments to evaluate the applicability of the IFL
 319 benchmark results to real world tasks. We will also study reinforcement learning algorithms for IFL
 320 and extend the number of features and supported environments in the IFL benchmark suite. We hope
 321 that other robotics researchers will develop their own IFL algorithms and evaluate them using the
 322 benchmark toolkit to accelerate progress.

References

- 323
- 324 [1] J. Smith. Robotic Arms Are Using Machine Learning to Reach Deeper Into Distribution. *Wall*
325 *Street Journal*, Jan. 2022. ISSN 0099-9660. URL <https://tinyurl.com/2p89zb35>.
- 326 [2] Logistics Automation with Plus One Robotics. *Parcel Moni-*
327 *tor*, May 2022. URL [https://www.parcelmonitor.com/blog/](https://www.parcelmonitor.com/blog/tech-spotlight-logistics-automation-with-plus-one-automation/)
328 [tech-spotlight-logistics-automation-with-plus-one-automation/](https://www.parcelmonitor.com/blog/tech-spotlight-logistics-automation-with-plus-one-automation/).
- 329 [3] E. Chu. How Zoox Builds Autonomous Vehicles from the Wheels Up-Blog.
330 *AI Exchange*, Apr. 2022. URL [https://exchange.scale.com/public/blogs/](https://exchange.scale.com/public/blogs/how-zoox-builds-autonomous-vehicles-from-the-wheels-up)
331 [how-zoox-builds-autonomous-vehicles-from-the-wheels-up](https://exchange.scale.com/public/blogs/how-zoox-builds-autonomous-vehicles-from-the-wheels-up).
- 332 [4] A. C. Madrigal. Waymo’s Robot Cars, and the Humans Who Tend to Them. *The Atlantic*,
333 Aug. 2018. URL [https://www.theatlantic.com/technology/archive/2018/08/](https://www.theatlantic.com/technology/archive/2018/08/waymos-robot-cars-and-the-humans-who-tend-to-them/568051/)
334 [waymos-robot-cars-and-the-humans-who-tend-to-them/568051/](https://www.theatlantic.com/technology/archive/2018/08/waymos-robot-cars-and-the-humans-who-tend-to-them/568051/). Section: Tech-
335 nology.
- 336 [5] K. Menda, K. Driggs-Campbell, and M. J. Kochenderfer. EnsembleDagger: A Bayesian
337 Approach to Safe Imitation Learning. In *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and*
338 *Systems (IROS)*, 2019.
- 339 [6] B. Kim and J. Pineau. Maximum mean discrepancy imitation learning. In *Robotics: Science*
340 *and systems*, 2013.
- 341 [7] M. Laskey, S. Staszak, W. Hsieh, J. Mahler, F. Pokorny, A. Dragan, and K. Goldberg. SHIV:
342 Reducing supervisor burden using support vectors for efficient learning from demonstrations in
343 high dimensional state spaces. In *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2016.
- 344 [8] R. Hoque, A. Balakrishna, E. Novoseller, A. Wilcox, D. S. Brown, and K. Goldberg. Thriftydag-
345 ger: Budget-aware novelty and risk gating for interactive imitation learning. *Conference on*
346 *Robot Learning (CoRL)*, 2021.
- 347 [9] J. Zhang and K. Cho. Query-efficient imitation learning for end-to-end autonomous driving. In
348 *Association for the Advancement of Artificial Intelligence (AAAI)*, 2017.
- 349 [10] R. Hoque, A. Balakrishna, C. Putterman, M. Luo, D. S. Brown, D. Seita, B. Thananjeyan,
350 E. Novoseller, and K. Goldberg. LazyDagger: Reducing context switching in interactive
351 imitation learning. In *International Conference on Automation Sciences and Engineering*
352 *(CASE)*, 2021.
- 353 [11] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,
354 A. Allshire, A. Handa, and G. State. Isaac gym: High performance gpu-based physics simulation
355 for robot learning. *ArXiv*, abs/2108.10470, 2021.
- 356 [12] B. Sellner, R. Simmons, and S. Singh. User modelling for principled sliding autonomy in
357 human-robot teams. In *Multi-Robot Systems. From Swarms to Intelligent Automata Volume III*,
358 pages 197–208. Springer, 2005.
- 359 [13] B. Sellner, F. W. Heger, L. M. Hiatt, R. Simmons, and S. Singh. Coordinated multiagent teams
360 and sliding autonomy for large-scale assembly. *Proceedings of the IEEE*, 94(7):1425–1444,
361 2006.
- 362 [14] M. B. Dias, B. Kannan, B. Browning, E. Jones, B. Argall, M. F. Dias, M. Zinck, M. Veloso,
363 and A. Stentz. Sliding autonomy for peer-to-peer human-robot teams. In *Proceedings of the*
364 *international conference on intelligent autonomous systems*, pages 332–341, 2008.
- 365 [15] P. Scerri, D. V. Pynadath, and M. Tambe. Towards adjustable autonomy for the real world.
366 *Journal of Artificial Intelligence Research*, 17:171–228, 2002.

- 367 [16] D. Kortenkamp, D. Keirn-Schreckenghost, and R. P. Bonasso. Adjustable control autonomy for
368 manned space flight. In *2000 IEEE Aerospace Conference. Proceedings (Cat. No. 00TH8484)*,
369 volume 7, pages 629–640. IEEE, 2000.
- 370 [17] T. B. Sheridan. Adaptive automation, level of automation, allocation authority, supervisory
371 control, and adaptive control: Distinctions and modes of adaptation. *IEEE Transactions on*
372 *Systems, Man, and Cybernetics-Part A: Systems and Humans*, 41(4):662–667, 2011.
- 373 [18] J. Y. Chen and M. J. Barnes. Supervisory control of multiple robots: Effects of imperfect
374 automation and individual differences. *Human Factors*, 54(2):157–174, 2012.
- 375 [19] M. Lewis. Human interaction with multiple remote robots. *Reviews of Human Factors and*
376 *Ergonomics*, 9(1):131–174, 2013.
- 377 [20] S.-Y. Chien, M. Lewis, S. Mehrotra, and K. Sycara. Imperfect automation in scheduling operator
378 attention on control of multi-robots. In *Proceedings of the Human Factors and Ergonomics*
379 *Society Annual Meeting*, volume 57, pages 1169–1173. SAGE Publications Sage CA: Los
380 Angeles, CA, 2013.
- 381 [21] J. R. Peters, V. Srivastava, G. S. Taylor, A. Surana, M. P. Eckstein, and F. Bullo. Human
382 supervisory control of robotic teams: Integrating cognitive modeling with engineering design.
383 *IEEE Control Systems Magazine*, 35(6):57–80, 2015.
- 384 [22] S.-Y. Chien, Y.-L. Lin, P.-J. Lee, S. Han, M. Lewis, and K. Sycara. Attention allocation
385 for human multi-robot control: Cognitive analysis based on behavior data and hidden states.
386 *International Journal of Human-Computer Studies*, 117:30–44, 2018.
- 387 [23] A. Rosenfeld, N. Agmon, O. Maksimov, and S. Kraus. Intelligent agent supporting human–
388 multi-robot team collaboration. *Artificial Intelligence*, 252:211–231, 2017.
- 389 [24] J. W. Crandall, M. L. Cummings, M. Della Penna, and P. M. De Jong. Computing the effects
390 of operator attention allocation in human control of multiple robots. *IEEE Transactions on*
391 *Systems, Man, and Cybernetics-Part A: Systems and Humans*, 41(3):385–397, 2010.
- 392 [25] S. A. Zanlongo, P. Dirksmeier, P. Long, T. Padir, and L. Bobadilla. Scheduling and path-planning
393 for operator oversight of multiple robots. *Robotics*, 10(2):57, 2021.
- 394 [26] D. R. Olsen Jr and S. B. Wood. Fan-out: Measuring human control of multiple robots. In
395 *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 231–238,
396 2004.
- 397 [27] M. Cakmak, C. Chao, and A. L. Thomaz. Designing interactions for robot active learners. *IEEE*
398 *Transactions on Autonomous Mental Development*, 2(2):108–118, 2010.
- 399 [28] S. Amershi, M. Cakmak, W. B. Knox, and T. Kulesza. Power to the people: The role of humans
400 in interactive machine learning. *Ai Magazine*, 35(4):105–120, 2014.
- 401 [29] J. Y. Chen and M. J. Barnes. Human–agent teaming for multirobot control: A review of human
402 factors issues. *IEEE Transactions on Human-Machine Systems*, 44(1):13–29, 2014.
- 403 [30] K. Zheng, D. F. Glas, T. Kanda, H. Ishiguro, and N. Hagita. Supervisory control of multiple
404 social robots for navigation. In *2013 8th ACM/IEEE International Conference on Human-Robot*
405 *Interaction (HRI)*, pages 17–24. IEEE, 2013.
- 406 [31] T. Ji, R. Dong, and K. Driggs-Campbell. Traversing supervisor problem: An approximately
407 optimal approach to multi-robot assistance. *arXiv preprint arXiv:2205.01768*, 2022.
- 408 [32] A. Dahiya, N. Akbarzadeh, A. Mahajan, and S. L. Smith. Scalable operator allocation for
409 multi-robot assistance: A restless bandit approach. *IEEE Transactions on Control of Network*
410 *Systems*, 2022.

- 411 [33] R. B. Cooper. Queueing theory. In *Proceedings of the ACM'81 conference*, pages 119–122,
412 1981.
- 413 [34] L. Sha, T. Abdelzaher, A. Cervin, T. Baker, A. Burns, G. Buttazzo, M. Caccamo, J. Lehoczky,
414 A. K. Mok, et al. Real time scheduling theory: A historical perspective. *Real-time systems*, 28
415 (2):101–155, 2004.
- 416 [35] M. Gombolay, R. Wilcox, and J. Shah. Fast scheduling of multi-robot teams with temporospatial
417 constraints. 2013.
- 418 [36] W. B. Rouse. Adaptive allocation of decision making responsibility between supervisor and
419 computer. In *Monitoring behavior and supervisory control*, pages 295–306. Springer, 1976.
- 420 [37] A. Daw, R. C. Hampshire, and J. Pender. How to staff when customers arrive in batches. *arXiv*
421 *e-prints*, pages arXiv–1907, 2019.
- 422 [38] J. A. Bagnell. An invitation to imitation. Technical Report CMU-RI-TR-15-08, Carnegie
423 Mellon University, Pittsburgh, PA, 3 2015.
- 424 [39] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *Proceedings of the*
425 *27th International Joint Conference on Artificial Intelligence (IJCAI)*, 7 2018.
- 426 [40] F. Codevilla, E. Santana, L. A. M., and A. Gaidon. Exploring the limitations of behavior cloning
427 for autonomous driving. *International Conference on Computer Vision*, 2019.
- 428 [41] D. A. Pomerleau. Efficient training of artificial neural networks for autonomous navigation.
429 *Neural Computation*, 3(1), 1991.
- 430 [42] J. Ho and S. Ermon. Generative adversarial imitation learning. In *Advances in Neural Informa-*
431 *tion Processing Systems*, 2016.
- 432 [43] B. D. Argall, S. Chernova, M. Veloso, and B. Browning. A survey of robot learning from
433 demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.
- 434 [44] S. Ross, G. J. Gordon, and J. A. Bagnell. A reduction of imitation learning and structured
435 prediction to no-regret online learning. In *International Conference on Artificial Intelligence*
436 *and Statistics (AISTATS)*, 2011.
- 437 [45] M. Laskey, J. Lee, R. Fox, A. Dragan, and K. Goldberg. DART: Noise injection for robust
438 imitation learning. In *Conf. on Robot Learning (CoRL)*, 2017.
- 439 [46] S. Chernova and M. Veloso. Interactive policy learning through confidence-based autonomy.
440 *Journal of Artificial Intelligence Research*, 34:1–25, 2009.
- 441 [47] S. Jauhri, C. Celemin, and J. Kober. Interactive imitation learning in state-space. In *Conference*
442 *on Robot Learning (CoRL)*. PMLR, 2020.
- 443 [48] M. Rigter, B. Lacerda, and N. Hawes. A framework for learning from demonstration with
444 minimal human effort. *IEEE Robotics and Automation Letters*, 5(2):2023–2030, 2020.
- 445 [49] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. HG-Dagger: Interactive
446 imitation learning with human experts. In *Proc. IEEE Int. Conf. Robotics and Automation*
447 *(ICRA)*, 2019.
- 448 [50] J. Spencer, S. Choudhury, M. Barnes, M. Schmittle, M. Chiang, P. Ramadge, and S. Srinivasa.
449 Learning from interventions: Human-robot interaction as both explicit and implicit feedback.
450 In *Proc. Robotics: Science and Systems (RSS)*, 2020.
- 451 [51] A. Mandlekar, D. Xu, R. Martín-Martín, Y. Zhu, L. Fei-Fei, and S. Savarese. Human-in-the-loop
452 imitation learning using remote teleoperation, 2020.

- 453 [52] R. R. Murphy and E. Rogers. Cooperative assistance for remote robot supervision. *Presence: Teleoperators & Virtual Environments*, 5(2):224–240, 1996.
- 454
- 455 [53] L. Xie, S. Wang, S. Rosa, A. Markham, and N. Trigoni. Learning with training wheels: Speeding up training with a simple controller for deep reinforcement learning. In *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2018.
- 456
- 457
- 458 [54] A. Kurenkov, A. Mandlekar, R. Martin-Martin, S. Savarese, and A. Garg. Ac-teach: A bayesian actor-critic method for policy learning with an ensemble of suboptimal teachers. In *Conf. on Robot Learning (CoRL)*, 2019.
- 459
- 460
- 461 [55] F. Wang, B. Zhou, K. Chen, T. Fan, X. Zhang, J. Li, H. Tian, and J. Pan. Intervention aided reinforcement learning for safe and practical policy optimization in navigation. In *Conf. on Robot Learning (CoRL)*, 2018.
- 462
- 463
- 464 [56] G. Kahn, P. Abbeel, and S. Levine. LaND: Learning to Navigate from Disengagements. In *arXiv preprint arXiv:2010.04689*, 2020.
- 465
- 466 [57] S. Vats, O. Kroemer, and M. Likhachev. Synergistic scheduling of learning and allocation of tasks in human-robot teams. *arXiv preprint arXiv:2203.07478*, 2022.
- 467
- 468 [58] G. B. Margolis, G. Yang, K. Paigwar, T. Chen, and P. Agrawal. Rapid locomotion via reinforcement learning. *arXiv preprint arXiv:2205.02824*, 2022.
- 469
- 470 [59] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control. *arXiv preprint arXiv:2203.06173*, 2022.
- 471
- 472 [60] A. Escontrela, X. B. Peng, W. Yu, T. Zhang, A. Iscen, K. Goldberg, and P. Abbeel. Adversarial motion priors make good substitutes for complex reward functions. *arXiv preprint arXiv:2203.15103*, 2022.
- 473
- 474
- 475 [61] X. B. Peng, Y. Guo, L. Halper, S. Levine, and S. Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *arXiv preprint arXiv:2205.01906*, 2022.
- 476
- 477 [62] N. Rudin, D. Hoeller, P. Reist, and M. Hutter. Learning to walk in minutes using massively parallel deep reinforcement learning. In *Conference on Robot Learning*, pages 91–100. PMLR, 2022.
- 478
- 479
- 480 [63] G. Swamy, S. Reddy, S. Levine, and A. D. Dragan. Scaled autonomy: Enabling human operators to control robot fleets. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5942–5948. IEEE, 2020.
- 481
- 482
- 483 [64] C. Shannon. A mathematical theory of communication. *Bell System Technical Journal (1948)*, 1948.
- 484
- 485 [65] B. Thananjeyan*, A. Balakrishna*, S. Nair, M. Luo, K. Srinivasan, M. Hwang, and J. E. Gonzalez. Recovery rl: Safe reinforcement learning with learned recovery zones. 2020.
- 486
- 487 [66] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *ArXiv*, abs/1606.01540, 2016.
- 488
- 489 [67] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 490