

Fleet Supervisor Allocation: A Submodular Maximization Approach

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** In real world scenarios, the data collected by robots in diverse and un-
2 predictable environments is crucial for enhancing their models and policies. This
3 data is predominantly collected under human supervision, particularly through im-
4 itation learning (IL), where robots learn complex tasks by observing human super-
5 visors. However, the deployment of multiple robots and supervisors to accelerate
6 the learning process often leads to data redundancy and inefficiencies, especially
7 as the scale of robot fleets increases. Moreover, the reliance on teleoperation for
8 supervision introduces additional challenges due to potential network connectiv-
9 ity issues. To address these inefficiencies and the reliability concerns of network-
10 dependent supervision, we introduce an adaptive submodular maximization-based
11 policy designed for efficient human supervision allocation within multi-robot sys-
12 tems under uncertain connectivity. Our approach significantly reduces data redun-
13 dancy by balancing the informativeness and diversity of data collection, and is
14 capable of accommodating connectivity variances. We evaluated the effectiveness
15 of ASA in a simulation environment with 100 robots across four different environ-
16 ments and various network settings, including a real-world teleoperation scenario
17 over a 5G network. We trained and tested both our and the state-of-the-art policies
18 utilizing NVIDIA's Isaac Gym, and our results show that ASA enhances the return
19 on human effort by up to $5.95\times$, outperforming current baselines in all simulated
20 scenarios and providing robustness against connectivity disruptions.

21 **Keywords:** Imitation Learning, Submodular Maximization, Fleet Learning

22 1 Introduction

23 Today, diverse industries deploy robotic fleets for tasks ranging from autonomous driving [1, 2] to
24 healthcare [3] and package delivery [4]. These robots are often deployed with policies trained on a
25 dataset that is primarily based on simulations, along with a small amount of data collected through
26 real-world interactions. While effective within their training contexts, these models often fail to
27 adapt to new or evolving real-world scenarios [5], making data collection critical for the success of
28 the robotics applications [6, 7].

29 A popular approach to collecting such data is through human supervision, where humans directly
30 guide the robots to perform the tasks. These data are then used to train the robots via Imitation
31 Learning (IL), where the robots are trained to perform tasks by observing the human demonstra-
32 tions [8]. Imitation Learning (IL) has been effective in many robotics applications, ranging from
33 autonomous driving [9] to robotic manipulation [10, 11]. However, the breadth of scenarios neces-
34 sary for effective IL emphasizes the need for continual data collection [12], commonly done with
35 numerous robots in parallel. Usually, the number of humans is less than that of deployed robots. For
36 instance, a recent autonomous delivery company, Starship Technologies, operates 1700 autonomous
37 robots while teleoperating only 1% of this robotic fleet [13, 14]. The scarcity of human supervisors
38 necessitates the selection of informative robots for supervision [15, 16].

39 Human supervision is often provided through real-time teleoperation over a network, especially
40 when supervising fleets of robots distributed across the globe. For example, various companies, in-
41 cluding Cruise, utilize human supervisors located in their control centers to teleoperate autonomous

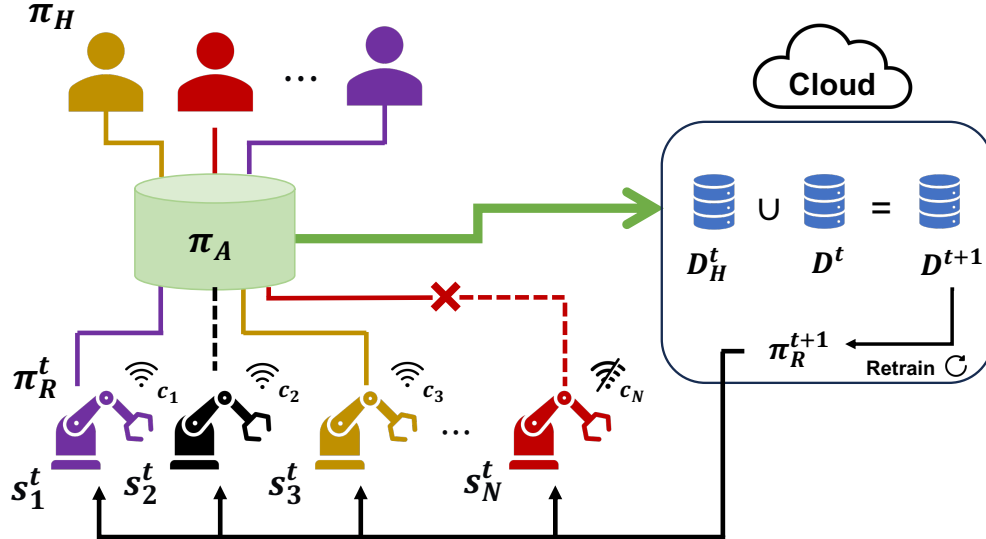


Figure 1: **Supervisor Allocation Problem:** In each time step t , the human supervisors with policy π_H can be allocated to the robots with policy π_R^t based on the allocation policy π_A . Each robot i has been operating in different states s_i^t , and the human supervisors are allocated to the robots based on the uncertainty of the robots and the similarity between the robots. Additionally, the supervision is provided through teleoperation with probability c_i ; meaning there is a chance that the connection to robot i might fail. At the end of each time step t , the data collected by human supervisors D_H^t is added to the dataset D^t to create an updated dataset D^{t+1} which is then used to train the robot policy π_R^{t+1} .

42 vehicles deployed across the world [17, 18]. However, these networks might be susceptible to con-
 43 nection failures [19], and it is important to be robust against network uncertainties. Combining these
 44 challenges with selecting informative robots, we formulate the Supervisor Allocation Problem (Fig.
 45 1), which involves managing limited human resources to maximize data diversity and quality under
 46 uncertain network connectivity. Our problem is an extension of the Interactive Fleet Learning (IFL)
 47 setting introduced by [15]. We extend the IFL setting to account for network elements that play an
 48 important role in real-world teleoperation scenarios [20].

49 We then introduce a novel human supervisor allocation policy called adaptive submodular alloca-
 50 tion (ASA). ASA distributes human supervisory capacity across a fleet, ensuring a balance of data
 51 informativeness and diversity to minimize redundancy in data collection. Our allocation policy is
 52 shown to be robust against network instabilities and is able to adapt to the dynamic nature of data
 53 collection, which we demonstrate through extensive simulations in diverse network environments,
 54 including real-world 5G scenarios. We show that ASA improves human supervision efficacy metric
 55 Return on Human Effort (RoHE) [15] by up to $5.95\times$ compared to existing benchmarks.

56 2 Related Work

57 Data collection is a critical problem in robotics and machine learning that is essential for continually
 58 improving the performance of robots [21–25]. It is closely related to active learning [26–30], where
 59 the goal is to select the most informative samples to label. Although the goal of data collection
 60 is similar to active learning, the focus is on collecting data samples that are the most informative
 61 for training the models. In our case, however, the aim is to select the robots that provide the most
 62 informative data for human supervisors.

63 IL is a popular approach in robotic learning, where robots learn policies from human demonstrations
 64 [31–34]. Despite its potential, the reliance on purely offline data introduces several challenges such
 65 as distribution shifts [35], which occurs when robots encounter states that were not previously experi-
 66 enced by humans. These issues can be alleviated through online data collection methods, including
 67 Dataset Aggregation (DAgger) [35] and various forms of interactive IL [36, 37]. Most interactive IL
 68 methods rely on human supervision to decide on when to intervene in the robot’s learning process.
 69 This presents scalability challenges, especially when applied to extensive robot networks [38] or
 70 during prolonged learning phases [39]. Robot-initiated interactive IL strategies like EnsembleDAg-

71 ger [40] and ThriftyDagger [41] have been proposed to mitigate these constraints, enabling robots
 72 to request human input under specific conditions. However, these methods are designed for single-
 73 robot task allocation scenarios and do not consider multi-robot scenarios. Closest to our work,
 74 Fleet-Dagger [42] has been proposed to address the supervisor allocation problem in a multi-robot
 75 scenario. However, Fleet-Dagger does not consider operational constraints that might limit the al-
 76 location of human supervisors, such as network connectivity and the potential redundancy from em-
 77 ploying multiple human supervisors in similar environments. Our work, on the other hand, focuses
 78 on learning an allocation policy that is adaptable to the operational constraints while minimizing the
 79 redundancy in the data collection, which is crucial for the system’s scalability [27].

80 One popular approach to mitigate redundancy in data collection is using submodular maximization.
 81 Submodularity refers to the property of the marginal gain of adding an item to a small set being
 82 higher than adding the same item to a large set. As submodularity is a common trend in data
 83 collection, it has been widely used in machine learning tasks such as sensor placement [43], active
 84 learning [27, 30, 44], and summarization [45]. Submodular maximization has also been extended to
 85 stochastic settings [46, 47], where the goal is to select a subset of items to maximize the expected
 86 value of a submodular function. Despite its wide use in machine learning, stochastic submodular
 87 maximization has not been used in the context of IL and multi-robot data collection scenarios. Our
 88 work is the first to use stochastic submodular maximization in the context of human supervision and
 89 multi-robot scenarios to address the supervisor allocation problem.

90 3 Problem Formulation

91 Consider a geo-distributed system of N_{robot} robots, $\mathbf{I} = \{1, \dots, N_{\text{robot}}\}$. Each robot i operates
 92 in parallel within an independent Markov Decision Process (MDP) with a different initial state.
 93 However, all robots operate within the same state and action spaces \mathbf{S} and \mathbf{A} , respectively. Each
 94 robot i observes the state of the environment $s_i^t \in \mathbf{S}$ at time t and selects an action $a_i^t \in \mathbf{A}$ based on
 95 a policy $\pi_{\mathbf{R}}^t : \mathbf{S} \rightarrow \mathbf{A}$. The robots share the same policy $\pi_{\mathbf{R}}^t$ that has been trained using the collective
 96 data \mathbf{D}^t accumulated up to time step t . We define the collection of states and actions for all robots as
 97 $\mathbf{s}^t = (s_1^t, \dots, s_{N_{\text{robot}}}^t) \in \mathbf{S}^{N_{\text{robot}}}$ and $\mathbf{a}^t = (a_1^t, \dots, a_{N_{\text{robot}}}^t) \in \mathbf{A}^{N_{\text{robot}}}$. These robots can be supervised
 98 by N_{human} human supervisors with an oracle policy $\pi_{\mathbf{H}} : \mathbf{S} \rightarrow \mathbf{A}_{\mathbf{H}}$, respectively. In addition to the
 99 robot action space \mathbf{A} , the human action space $\mathbf{A}_{\mathbf{H}}$ includes a reset action R , which can return the
 100 robot to a safe state.

101 **Supervisor Allocation and Connectivity:** In each time step t , N_{human} human supervisors can be
 102 assigned to the robots for assistance. However, the connections to the robots are unreliable, with $C =$
 103 $\{c_1, \dots, c_{N_{\text{robot}}}\} \in \mathbb{R}^{N_{\text{robot}}}$ denoting independent random variables associated with the connection
 104 reliability of the robots. $c_i \in \{0, 1\}$ indicates whether a successful connection with robot i can be
 105 established ($c_i = 1$) or not ($c_i = 0$). Under this uncertain connectivity, we are interested in finding
 106 an allocation policy $\pi_{\mathbf{A}} : \mathbb{R}^{N_{\text{robot}}} \times \mathbf{S}^{N_{\text{robot}}} \times \mathbf{A}^{N_{\text{robot}}} \times \mathbf{I} \rightarrow X$ that selects robots to be supervised
 107 $X \subseteq \mathbf{I}$ based on connection reliability C , collection of states \mathbf{s}^t and action spaces \mathbf{a}^t .

108 **Data Collection and Policy Retraining:** Upon allocation, human supervisors contribute data only
 109 from successful connections, forming the human supervision data $\mathbf{D}_{\mathbf{H}}^t$. The robot policy $\pi_{\mathbf{R}}^t$ is then
 110 updated by integrating this new data into the current dataset and retraining:

$$\mathbf{D}^{t+1} = \mathbf{D}^t \cup \mathbf{D}_{\mathbf{H}}^t, \quad \mathbf{D}_{\mathbf{H}}^t = \{(s_i, \pi_{\mathbf{H}}(s_i)) : i \in X \text{ and } c_i = 1\}, \quad (1)$$

$$\pi_{\mathbf{R}}^{t+1} = g(\pi_{\mathbf{R}}^t, \mathbf{D}^{t+1}). \quad (2)$$

111 **Objective:** Our objective is to develop an allocation policy $\pi_{\mathbf{A}}$ that maximizes the expected Return
 112 on Human Effort (RoHE) over the connectivity C . RoHE metric was introduced along with Interac-
 113 tive Fleet Learning setup [15] to set a benchmark in Fleet Learning settings. It is a ratio of the total
 114 reward obtained by the fleet to the total number of human actions:

$$\max_{\pi_{\mathbf{A}} \in \Omega} \mathbb{E}_C \left[\frac{N_{\text{human}} \sum_{i \in \mathbf{I}} \sum_{t=0}^T r(s_i^t, a_i^t)}{N_{\text{robot}} \left(1 + \sum_{t=0}^T \|\pi_{\mathbf{A}}(C, \mathbf{s}^t, \mathbf{a}^t, \mathbf{I})\|_F^2 \right)} \right]. \quad (3)$$

115 Here T is the time horizon covering all time steps, $r : \mathbf{S} \times \mathbf{A} \rightarrow \mathbb{R}$ is the reward function, $\|\cdot\|_F$ denotes
 116 the Frobenius norm, and Ω refers to a set of all allocation policies. Intuitively, RoHE measures the
 117 total performance of the robotic fleet normalized by the total number of human interventions.

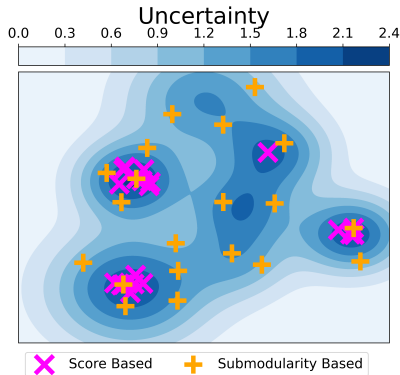


Figure 2: **Submodular Maximization Balances Uncertainty and Diversity:** This figure illustrates a toy example of our allocation problem in a 2D state space. The blue contours indicate the uncertainty levels, with darker shades representing higher uncertainty. Purple crosses (traditional score-based allocations) and yellow plus signs (our submodularity-based allocations) mark the positions of selected robots. Unlike score-based methods that often select highly uncertain but potentially overlapping states, our approach strategically picks a more diverse set of states, effectively balancing between high uncertainty and coverage, thereby reducing data redundancy and enhancing the training data’s representativeness.

118 4 A Stochastic Submodular Maximization Approach

119 We now present our novel policy, adaptive submodular allocation (ASA), for the problem outlined
 120 in Eq. 3 and its components. First, we define the stochastic submodular maximization problem,
 121 which represents the value of robot supervision, and then we define the greedy algorithm, which is
 122 used to pick the robots to supervise.

123 4.1 Submodular Maximization Problem

124 To address the optimization problem presented in Eq. 3, we use stochastic submodular maximiza-
 125 tion. Stochastic submodular maximization is particularly suited to our scenario because it leverages
 126 the diminishing returns property that naturally reflects the decrease in the marginal gain of super-
 127 vising additional robots. This aspect is vital for assessing the efficiency of human supervision.
 128 Furthermore, the method inherently discourages the selection of similar robots, thereby avoiding the
 129 assignment of humans to robots that offer overlapping information, which decreases the return on
 130 human effort. Finally, the stochastic submodular maximization accounts for the inherent uncertainty
 131 of our problem, acknowledging the non-deterministic connections to the robots, which is crucial for
 132 developing a robust solution across different connectivity patterns.

133 We first define a submodular objective function $f : 2^{N_{\text{robot}}} \times \{0, 1\}^{N_{\text{robot}}} \rightarrow \mathbb{R}$ that quantifies the
 134 value of supervising a selected set of robots X , considering the allocation reliability outcomes for
 135 these robots C . We define our objective function based on the facility location problem, a classical
 136 example of a submodular maximization objective [48], as follows:

$$f(X, C) = \sum_{i \in \mathbf{I}} \max_{j \in X} c_j M_{j,i}. \quad (4)$$

137 Here X is the set of robots selected for human supervision, and C indicates whether the connection
 138 to the robot is successful or not. $M_{j,i}$ represents the value of supervision of the robot j on the robot
 139 i , and we consider two factors: the informativeness of the robot i and the similarity between the
 140 robot j and the robot i . Additionally, our formulation is modular and can be extended to include
 141 other factors, such as prioritizing the robots that have violated the safety constraints or the robots
 142 that are in critical states. With all factors combined, we define the value of supervision $M_{j,i}$ as:

$$M_{j,i} = \mathcal{S}(j, i) * \mathcal{U}(i) + \mathcal{C}(i). \quad (5)$$

143 Here, $\mathcal{S}(j, i)$ defines the similarity between the robots i and j , and $\mathcal{U}(i)$ is the informativeness of
 144 the robot i , while $\mathcal{C}(i)$ is an indicator of whether the robot i violates the safety constraints or is in
 145 a critical state. Our definition of $M_{j,i}$ is modular, and each factor can be defined based on specific
 146 requirements. For example, the similarity function \mathcal{S} can be defined as the cosine similarity, the
 147 Euclidean distance, or any other similarity metric. The informativeness of the robot $\mathcal{U}(i)$ can be
 148 defined as the entropy of the robot’s policy or the uncertainty of the robot’s state, while the constraint
 149 function $\mathcal{C}(i)$ can be defined based on the safety constraints or the critical states for the robots.
 150 With the objective function defined, we pose the following maximization problem to optimize our
 151 allocation policy:

$$\begin{aligned} & \max_{X \subseteq \mathbf{I}} \mathbb{E}_C[f(X, C)] & (6) \\ & \text{subject to: } |X| \leq N_{\text{human}}, \end{aligned}$$

152 where the goal is to identify the subset of robots X that maximizes the expected value of f , con-
 153 strained by the number of available human supervisors N_{human} .

154 4.2 Adaptive Supervisor Allocation (ASA) Policy

155 Now, we can present our allocation policy ASA based on a greedy algorithm given in Algorithm
 156 1. Starting from an empty solution set X (line 1), ASA iteratively selects the robot with the high-
 157 est marginal gain in expectation over probabilities of connection to the robots C (line 3). Then,
 158 ASA computes the expected marginal gain of selecting the robot x^* (line 4), and if the expected
 159 marginal gain is below a certain threshold, the algorithm stops the selection process (line 5). This
 160 threshold ensures that the algorithm avoids using unnecessary human effort by stopping when the
 161 marginal gain of selecting an additional robot is low. Otherwise, the chosen robot x^* is added to
 162 the solution set X (line 8). Finally, based on the availability of the observation on whether the
 163 connection to the robot was successful or not, the connection probabilities are updated (line 9).
 164

165 Based on the availability of the observations
 166 of the connection probabilities, we define two
 167 variants of our policy: non-adaptive submod-
 168 ular allocation (n-ASA) and adaptive sub-
 169 modular allocation (ASA). In n-ASA, we are
 170 not able to observe the connection probabili-
 171 ties, and thus, the allocations are done before-
 172 hand. In ASA, on the other hand, the robots
 173 are selected iteratively; based on the success
 174 of the allocations, the connection probabili-
 175 ties are updated. To visualize the differences
 176 between n-ASA and ASA, consider the fol-
 177 lowing: in both cases, robot 1 is selected for
 178 supervision in the first iteration. While sel-
 179 ecting the robot to supervise, n-ASA con-
 180 siders both possibilities (successful and un-
 181 successful connection to robot 1) and selects
 182 the second robot that maximizes the expected
 183 marginal gain over both cases. However, in ASA, after selecting robot 1 for supervision, we observe
 184 whether the connection was successful or not and select the next robot that maximizes the expected
 185 marginal gain based on the observation.

186 When the marginal threshold parameter is set to zero, ASA is equivalent to the greedy algorithm for
 187 submodular maximization, which is proven to approximate the optimal solution for the submodular
 188 maximization problem in Eq. 6 with a factor of $1 - 1/e$ [46]. Additionally, n-ASA approximates the
 189 optimal adaptive policy with a factor of $(1 - 1/e)^2$ [46]. To compute selected robots faster, we use
 190 the lazy greedy algorithm [47]; this has the same time complexity as Algorithm 1, but has a better
 191 empirical performance.

192 5 Experiments

193 We consider a fleet of $N_{\text{robot}} = 100$ robots that can be supervised by $N_{\text{human}} = 5$ human supervisors.
 194 The human supervisors are implemented as reinforcement learning agents using the Proximal Policy
 195 Optimization (PPO) algorithm [49]. We utilize the behavior cloning algorithm to initialize the robot
 196 policies based on an offline dataset of 5000 state-action pairs and use our allocation policy to collect
 197 data from the robots and update the models. In all of our experiments, when the robots violate the
 198 constraints, we perform a hard reset to bring the robots back to a safe state. We set the hard reset
 199 time $t_R = 5$ timesteps, the minimum intervention time $t_T = 5$ timesteps, and the fleet operation
 200 time $T = 10,000$ timesteps. We average the results over 3 random seeds for each task and network
 201 configuration. We have chosen these parameters to align with the settings used in the environments
 202 of the benchmark algorithms [42] for a fair comparison.

203 **Environments:** We consider four different environments in our experiments: (1) Humanoid, where
 204 the robots focus on bipedal locomotion; (2) ANYmal, where the robots focus on quadruped locomo-

Algorithm 1 ASA Policy

Input: connectivities of robots C , set of all robots \mathbf{I}
Output: robots selected for supervision X

```

1: Initialize  $X \leftarrow \emptyset$ 
2: for  $k = 1$  to  $N_{\text{human}}$  do
3:    $x^* \leftarrow \operatorname{argmax}_{x \in \mathbf{I} \setminus X} \mathbb{E}_C[f(X, C)]$ 
4:   Compute  $\Delta \leftarrow \mathbb{E}_C[f(X \cup \{x^*\}, C) - f(X, C)]$ 
5:   if  $\Delta < \text{threshold}$  then
6:     break
7:   end if
8:    $X \leftarrow X \cup x^*$ 
9:   If possible, observe whether connection to
   robot  $x^*$  is successful and update connectivities
    $C$ 
10: end for

```

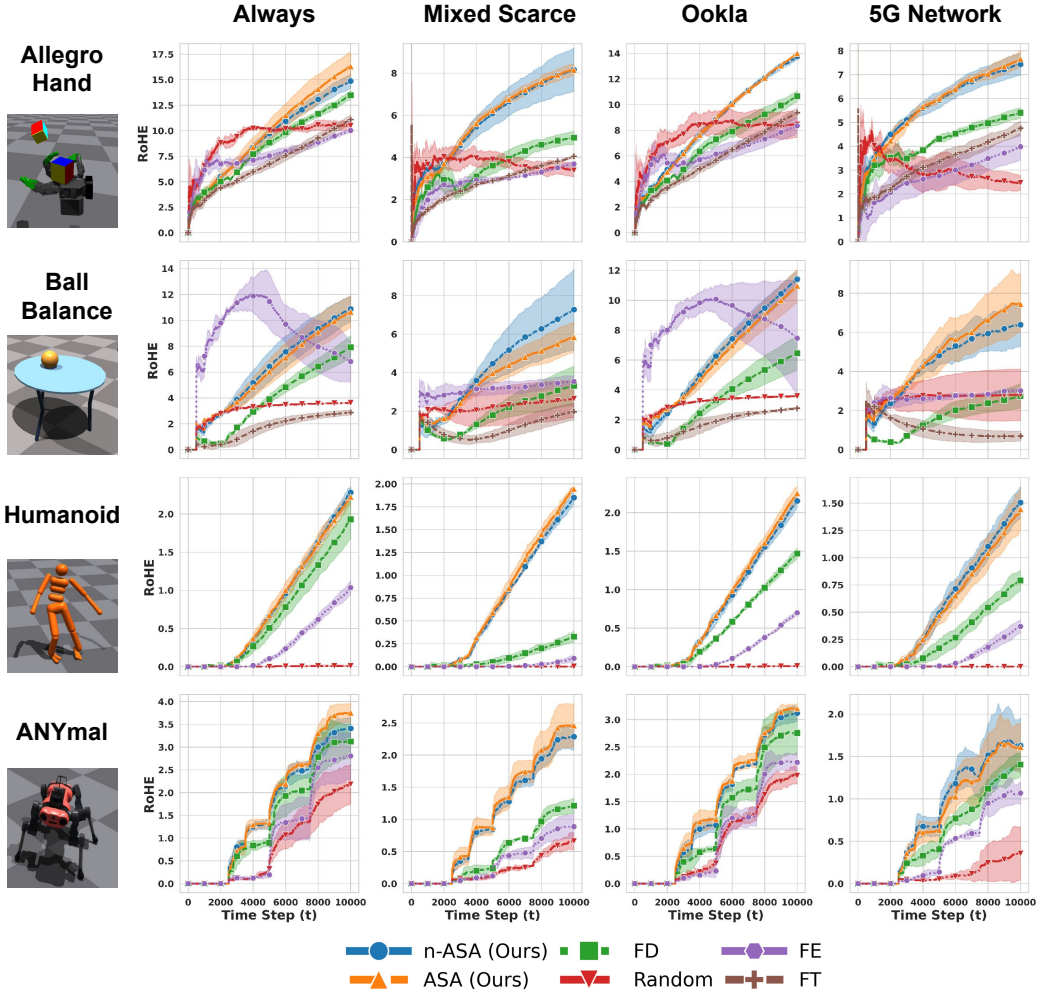


Figure 3: **Our ASA and n-ASA policies outperform other benchmarks across all environments and network combinations.** Here, each row represents a different environment, and each column corresponds to a different network configuration. ASA and n-ASA performance is affected least by changes in the network configurations because of their stochastic submodular maximization-based policies that can incorporate network uncertainties. Additionally, the submodular maximization objective improves the performance when there are no network uncertainties (column 1) due to its ability to cover diverse and informative scenarios.

205 tion with the ANYmal robot; (3) Allegro Hand, where the robots focus on dexterous manipulation
 206 tasks; and (4) Ball Balance where the robots focus on balancing a ball on a plate. Each environment
 207 defines constraint violations specifically similar to [50] requiring human intervention to reset the
 208 robots to a safe state.

209 **Network Configurations:** We use 4 different network configurations based on various connectivity
 210 probabilities: (1) **Always**, where the robots can always be supervised by the human supervisors; (2)
 211 **Mixed-Scarce**, where some of the robots have a high probability of connection while others have a
 212 low; (3) **Ookla**, where the robots have a varying probability of connection based on cellular network
 213 performance metrics [51] (4) **5G**, which is connectivity data that we collected over a real-world 5G
 214 network in a university robotics lab.

215 **Metrics:** We evaluate the performance of the allocation policies based on the following metrics:
 216 (1) Return on Human Effort (RoHE), which was given in Eq. 3 and (2) the cumulative number of
 217 successfully completed tasks by the entire fleet, which we will refer to as cumulative success. To
 218 simplify, RoHE measures the fleet performance per human intervention, while cumulative success
 219 only considers the total successful task completion without considering the number of interventions.
 220 For example, simply allocating all human supervisors would improve cumulative success but decrease
 221 RoHE. An ideal allocation policy should balance the two, as an ideal system would require a
 222 high total success while using humans as efficiently as possible.

223 **Baselines:** We compare the following baselines: (1) **Random**, which randomly selects the robots to
 224 be supervised by the human supervisors at each time step; (2) **Fleet-EnsembleDagger (FE)**, which
 225 utilizes variance for uncertainty estimation, combining it with constraint-based prioritization [40];
 226 (3) **Fleet-ThriftyDagger (FT)**, which merges uncertainty and goal-oriented prioritization, adapt-
 227 ing ThriftyDagger for fleet setting [41], for environments with a defined goal; (4) **Fleet-Dagger**
 228 **(FD)**, which prioritizes the robots violating constraints and selects the robots with the highest un-
 229 certainty and risk of failure for fleet supervision [42] (5) **Non-Adaptive Submodular Allocation**
 230 **(n-ASA)** and (6) **Adaptive Submodular Allocation (ASA)**, which are the two variants of our pro-
 231 posed method based on submodular maximization in the absence and presence of the observation of
 232 the connection to the robots described in Section 4; please see the Appendix for the exact similarity,
 233 uncertainty and constraint functions we have used in the submodular maximization objective given
 234 in Eq. 4.

235 How do ASA and n-ASA perform under different network configurations?

236 We evaluate the performance of our policies, ASA and n-ASA, under different network configura-
 237 tions for each environment. The RoHE metric for each time step has been shown in Fig. 3, and
 238 cumulative success values at the final time step have been presented in Fig. 4. In both metrics, we
 239 can see that our ASA and n-ASA allocation policies are able to outperform other benchmarks. This
 240 is because our policies can incorporate network uncertainty information into their allocation policy
 241 through stochastic submodularity. On the other hand, other benchmark policies are not designed to
 242 incorporate such network uncertainty. We can see that our allocation policy outperforms other
 243 baselines in terms of the RoHE metric by up to $5.95\times$, $2.03\times$, $1.65\times$, and $2.47\times$ in Humanoid,
 244 ANYmal, Allegro Hand, and Ball Balance environments, respectively.

245 How do ASA and n-ASA compare when the 246 network connectivity is stable?

247 To test whether our RoHE gains are only due to adaptability to different network configura-
 248 tions, we have also simulated a network where all robots are always reachable. In column 1
 249 of Fig. 3, we can see that our ASA and n-ASA policies still outperform other allocation bench-
 250 marks thanks to their ability to diversify the selected robots to cover more states. As we have
 251 shown in the 2D toy example in Fig. 2, rather than only focusing on the states with high un-
 252 certainty, ASA and n-ASA consider the whole state space to collect combined data that is more
 253 informative.
 254
 255
 256
 257
 258
 259

260 How does the availability of observation af- 261fect our allocation policies?

262 Although we know that ASA approximates the optimal solution with a stricter bound than n-
 263 ASA, in practice, these allocation policies perform very similarly. Both policies are supe-
 264 rior to other benchmarks in all configurations and have similar RoHE and cumulative success
 265 metrics. This performance similarity between ASA and n-ASA further proves that our perfor-
 266 mance gains are mainly a result of our stochastic submodular maximization approach rather
 267 than the observation of whether the connections with the previous robot are successful or not.
 268 This flexibility enables our policies to be applied in various real-world scenarios where the
 269 observations might be impossible.
 270
 271
 272
 273
 274
 275
 276

277 Can ASA and n-ASA improve cumulative success and RoHE metrics at the same time?

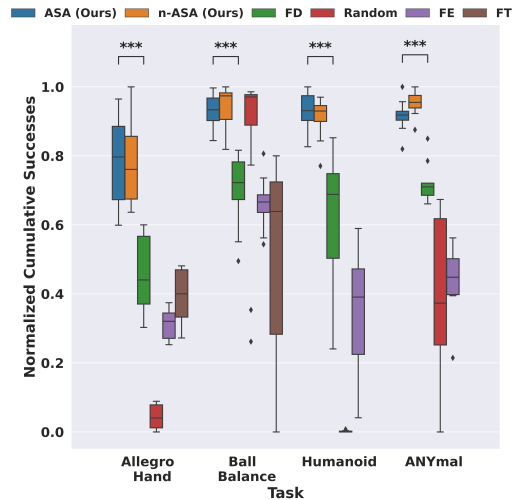


Figure 4: **ASA and n-ASA outperform all benchmarks in cumulative successes metric.** We present a box plot of the normalized cumulative success metric at the final time step in each environment and allocation policy. ASA and n-ASA can achieve higher cumulative success values with less standard deviation thanks to their robustness against network uncertainties and superior robot selection. Interestingly, Random policy achieves comparable results in the Ball Balance environment simply by allocating more humans to the process, sacrificing the return on human metric 3. On the other hand, ASA and n-ASA achieve higher performance by not allocating more humans but through a selection of better representative robots.

ALLOCATION POLICY	ALLEGROHAND		ANYMAL		BALLBALANCE		HUMANOID	
	ROHE	CUMULATIVE SUCCESS	ROHE	CUMULATIVE SUCCESS	ROHE	CUMULATIVE SUCCESS	ROHE	CUMULATIVE SUCCESS
RANDOM	2.47	228.67	0.36	47.33	2.81	1401.33	0	0
FT	4.75	2131.66	-	-	0.69	344.66	-	-
FE	3.98	1789.67	1.07	137	3.01	1315	0.37	169.33
FD	5.4	2213	1.40	200	2.75	1209	0.79	319.33
N-ASSA (OURS)	7.66	3705.33	1.61	246.33	7.45	1758.33	1.44	444.33
ASSA (OURS)	7.43	3658	1.63	239.33	6.39	1703.33	1.51	470

Table 1: Our proposed n-ASA and ASA policies outperform other baselines in all environments (columns) in real-world 5G network data in terms of cumulative success and return on human effort (RoHE). The results are consistent across all tasks, showing the adaptability of our method to different environments. Additionally, ASA and n-ASA outperform other baselines in both cumulative success and RoHE, meaning our allocation policy is both efficient in human effort and achieves higher cumulative success.

278 We can clearly see in Fig. 3 and Fig. 4 that our ASA and n-ASA policies achieve both the highest
279 RoHE and cumulative success in all network configurations. Our policies are able to balance these
280 two metrics thanks to their threshold criteria, preventing the allocation of humans to uninformative
281 robots. For example, we can see that in the Ball Balance environment, the Random allocation policy
282 is able to achieve comparable cumulative success in Fig. 4, but it fails to reach comparable RoHE
283 values (see row 2 in Fig. 3). This suggests that Random policy achieved high cumulative success
284 values by simply allocating more humans but failed to optimize human efficiency.

285 5.1 Physical 5G Network Connectivity Data

286 In addition to the simulated network connectivity data, we also evaluate our allocation policies on
287 real-world 5G network connectivity data collected in the field. To create such a dataset, we utilize
288 a local 5G network dedicated to testing the real-time teleoperation of the robots over a period of
289 24 hours. Then, we divide the geographic area into 100 different regions with the same number
290 of users (robots) in each region and calculate the average latency and throughput of the network for
291 each region. We use this data to create network connectivity where the robots with higher latency and
292 lower throughput have a lower probability of establishing a successful connection with the human
293 supervisors. Please refer to the Appendix for further details on 5G network data collection and the
294 exact setup we used.

295 **Results:** We present our results on the real-world 5G network data in Table 1. The results show
296 that our proposed method outperforms other baselines in terms of the RoHE and cumulative success
297 metrics by up to $2.47\times$ and $1.67\times$, respectively under the 5G network configuration.

298 **Limitations:** Our work has several limitations. First, it uses only real-world data collected from 5G
299 field trials without hardware robotics experiments. Additionally, it assumes that network connectiv-
300 ity and robot states and policies are independent across robots, while in real-world scenarios, robots
301 might share the same network or physical location, meaning their policies might affect each other.

302 6 Conclusion and Future Work

303 We present novel supervisor allocation policies, ASA and n-ASA, for assigning human supervi-
304 sors to the robotic fleet for data collection. ASA and n-ASA are based on stochastic submodular
305 maximization, providing a modular approach to incorporate different allocation objectives, informa-
306 tiveness metrics, as well as re-training methods. Our allocation policies beat current benchmarks in
307 terms of performance metrics such as RoHE and cumulative success in all environments and net-
308 work configurations. These performance gains are thanks to its stochastic submodular maximization
309 objective, which incorporates network connectivity in the allocation process while balancing the di-
310 versity and informativeness of selected robots. Finally, we collect real-world 5G network data from
311 a field dedicated to teleoperated robots and show the applicability of our allocation policy in real-
312 world scenarios as well.

313 In a future project, we plan to extend our work to include hardware robotics experiments, including
314 teleoperation over a 5G network possible in an application such as autonomous driving. We also
315 plan to investigate the impact of different imitation learning methods to test the generalizability of
316 our allocation policies.

References

- [1] Kara Carlson. You now can ride in a driverless car in austin, as gm-owned cruise expands rideshare services, Dec 2022. URL <https://www.statesman.com/story/business/technology/2022/12/21/cruise-car-company-launches-austin-driverless-rideshare-service/69743913007/>.
- [2] Anthony James. Waymo expands us public operations with deployments in los angeles and austin, Mar 2024. URL <https://www.autonomousvehicleinternational.com/news/robotaxis/waymo-expands-us-operations-with-deployments-in-los-angeles-and-austin.html>.
- [3] Robert Valner, Houman Masnavi, Igor Rybalskii, Rauno Pölluäär, Erik Kõiv, Alvo Aabloo, Karl Kruusamäe, and Arun Singh. Scalable and heterogenous mobile robot fleet-based task automation in crowded hospital environments—a field test. *Frontiers in Robotics and AI*, 9, 08 2022. doi:10.3389/frobt.2022.922835.
- [4] Joseph Quinlivan. How amazon deploys collaborative robots in its operations to benefit employees and customers, Jun 2023. URL <https://www.aboutamazon.com/news/operations/how-amazon-deploys-robots-in-its-operations-facilities>.
- [5] Cosmin Paduraru, Daniel Jaymin Mankowitz, Gabriel Dulac-Arnold, Jerry Li, Nir Levine, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, 110:2419 – 2468, 2021. URL <https://api.semanticscholar.org/CorpusID:234868359>.
- [6] Sergey Levine, Aviral Kumar, G. Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *ArXiv*, abs/2005.01643, 2020. URL <https://api.semanticscholar.org/CorpusID:218486979>.
- [7] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pages 991–1002. PMLR, 08–11 Nov 2022. URL <https://proceedings.mlr.press/v164/jang22a.html>.
- [8] Takayuki Osa, Joni Pajarinen, Gerhard Neumann, J. Andrew Bagnell, Pieter Abbeel, and Jan Peters. An algorithmic perspective on imitation learning. *Foundations and Trends® in Robotics*, 7(1-2):1–179, 2018. ISSN 1935-8253. doi:10.1561/23000000053. URL <http://dx.doi.org/10.1561/23000000053>.
- [9] Eli Bronstein, Mark Palatucci, Dominik Notz, Brandyn Allen White, Alex Kuefler, Yiren Lu, Supratik Paul, Payam Nikdel, Paul Mougín, Hongge Chen, Justin Fu, Austin Abrams, Punit Shah, Evan Raca, Benjamin Frenkel, Shimon Whiteson, and Drago Anguelov. Hierarchical model-based imitation learning for planning in autonomous driving. *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8652–8659, 2022. URL <https://api.semanticscholar.org/CorpusID:252968066>.
- [10] Bin Fang, Shi-Dong Jia, Di Guo, Muhua Xu, Shuhuan Wen, and Fuchun Sun. Survey of imitation learning for robotic manipulation. *International Journal of Intelligent Robotics and Applications*, 3:362 – 369, 2019. URL <https://api.semanticscholar.org/CorpusID:202733441>.
- [11] Zipeng Fu, Tony Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *ArXiv*, abs/2401.02117, 2024. URL <https://api.semanticscholar.org/CorpusID:266755740>.
- [12] Carlos Celemin, Rodrigo Pérez-Dattari, Eugenio Chisari, Giovanni Franzese, Leandro de Souza Rosa, Ravi Prakash, Zlatan Ajanović, Marta Ferraz, Abhinav Valada, and Jens Kober. Interactive imitation learning in robotics: A survey. *Foundations and Trends® in Robotics*, 10(1-2):1–197, 2022. ISSN 1935-8253. doi:10.1561/23000000072. URL <http://dx.doi.org/10.1561/23000000072>.
- [13] Kevin Jost. Starship goes deep on college food deliveries, Dec 2022. URL <https://insideautonomousvehicles.com/starship-goes-deep-on-college-food-deliveries/>.

- 372 [14] Mark DiPietro. Starship technologies paving way for robotics in food delivery, Oct 2022.
373 URL <https://blogs.shu.edu/stillmanexchange/2022/10/15/starship-technologies-paving-way-for-robotics-in-food-delivery/#:~:text=Founded%20by%20two%20Skype%20co,over%2010%2C000%20orders%20per%20day.>
374
375
376
- 377 [15] Ryan Hoque, Lawrence Yunliang Chen, Satvik Sharma, K Dharmarajan, Brijen Thananjeyan,
378 P. Abbeel, and Ken Goldberg. Fleet-dagger: Interactive robot fleet learning with scalable
379 human supervision. In Conference on Robot Learning, 2022.
- 380 [16] Shivin Dass, Karl Pertsch, Hejia Zhang, Youngwoon Lee, Joseph J. Lim, and Stefanos Niko-
381 laidis. Pato: Policy assisted teleoperation for scalable robot data collection, 2023.
- 382 [17] View All Posts by Mario Herger. First documented intervention by a cruise teleoperator, Febru-
383 ary 2023. URL <https://thelastdriverlicenseholder.com/2023/02/23/first-documented-intervention-of-a-cruise-teleoperator/>.
384
- 385 [18] View All Posts by Mario Herger. Remote-controlled driving as demonstrated by phantom auto,
386 February 2024. URL <https://thelastdriverlicenseholder.com/2024/02/01/remote-controlled-driving-as-demonstrated-by-phantom-car/>.
387
- 388 [19] Abhinav Dahiya. Route Planning and Operator Allocation in Robot Fleets. PhD thesis, Uni-
389 versity of Waterloo, 2023.
- 390 [20] Adriana Noguera Cundar, Reza Fotouhi, Zachary Ochitwa, and Haron Obaid. Quantifying the
391 effects of network latency for a teleoperated robot. Sensors, 23:8438, 10 2023. doi:10.3390/
392 s23208438.
- 393 [21] Sandeep Chinchali, E. Pergament, M. Nakanoya, E. Cidon, E. Zhang, D. Bharadia, M. Pavone,
394 and S. Katti. Harvestnet: Mining valuable training data from high-volume robot sensory
395 streams. In 2020 International Symposium on Experimental Robotics (ISER), Valetta, Malta,
396 2020.
- 397 [22] Sandeep Chinchali, E. Pergament, M. Nakanoya, E. Cidon, E. Zhang, D. Bharadia, M. Pavone,
398 and S. Katti. Sampling training data for continual learning between robots and the cloud. In
399 2020 International Symposium on Experimental Robotics (ISER), Valetta, Malta, 2020.
- 400 [23] Yuchong Geng, Dongyue Zhang, Po-han Li, Oguzhan Akcin, Ao Tang, and Sandeep P Chin-
401 chali. Decentralized sharing and valuation of fleet robotic data. In 5th Annual Conference on
402 Robot Learning, Blue Sky Submission Track, 2021.
- 403 [24] Oguzhan Akcin, Po-han Li, Shubhankar Agarwal, and Sandeep P. Chinchali. Decentralized
404 data collection for robotic fleet learning: A game-theoretic approach. In Karen Liu, Dana
405 Kulic, and Jeff Ichnowski, editors, Proceedings of The 6th Conference on Robot Learning,
406 volume 205 of Proceedings of Machine Learning Research, pages 978–988. PMLR, 14–18
407 Dec 2023. URL <https://proceedings.mlr.press/v205/akcin23a.html>.
- 408 [25] Crystal Chao, Maya Cakmak, and Andrea L. Thomaz. Transparent active learning for robots.
409 In 2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI), pages
410 317–324, 2010. doi:10.1109/HRI.2010.5453178.
- 411 [26] Burr Settles. Active learning literature survey. Computer Sciences Technical Report 1648,
412 University of Wisconsin–Madison, 2009.
- 413 [27] Oguzhan Akcin, Orhan Unuvar, Onat Ure, and Sandeep P. Chinchali. Fleet active learning:
414 A submodular maximization approach. In Jie Tan, Marc Toussaint, and Kourosh Darvish,
415 editors, Proceedings of The 7th Conference on Robot Learning, volume 229 of Proceedings
416 of Machine Learning Research, pages 1378–1399. PMLR, 06–09 Nov 2023. URL <https://proceedings.mlr.press/v229/akcin23a.html>.
417
- 418 [28] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep Bayesian active learning with im-
419 age data. In Doina Precup and Yee Whye Teh, editors, Proceedings of the 34th International
420 Conference on Machine Learning, volume 70 of Proceedings of Machine Learning Research,
421 pages 1183–1192. PMLR, 06–11 Aug 2017. URL <http://proceedings.mlr.press/v70/gall17a.html>.
422
- 423 [29] David A. Cohn, Zoubin Ghahramani, and Michael I. Jordan. Active learning with statistical
424 models. JOURNAL OF ARTIFICIAL INTELLIGENCE RESEARCH, 4:129–145, 1996.

- 425 [30] Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch
426 acquisition for deep bayesian active learning. In H. Wallach, H. Larochelle, A. Beygelzimer,
427 F. d'Alché-Buc, E. Fox, and R. Garnett, editors, Advances in Neural Information Processing
428 Systems, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neur](https://proceedings.neurips.cc/paper_files/paper/2019/file/95323660ed2124450caaac2c46b5ed90-Paper.pdf)
429 [ips.cc/paper_files/paper/2019/file/95323660ed2124450caaac2c46b](https://proceedings.neurips.cc/paper_files/paper/2019/file/95323660ed2124450caaac2c46b5ed90-Paper.pdf)
430 [5ed90-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/95323660ed2124450caaac2c46b5ed90-Paper.pdf).
- 431 [31] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In
432 Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence,
433 IJCAI-18, pages 4950–4957. International Joint Conferences on Artificial Intelligence Organi-
434 zation, 7 2018. doi:10.24963/ijcai.2018/687. URL [https://doi.org/10.24963/ijc](https://doi.org/10.24963/ijcai.2018/687)
435 [ai.2018/687](https://doi.org/10.24963/ijcai.2018/687).
- 436 [32] Felipe Codevilla, Eder Santana, Antonio M. López, and Adrien Gaidon. Exploring the limita-
437 tions of behavior cloning for autonomous driving. 2019 IEEE/CVF International Conference
438 on Computer Vision (ICCV), pages 9328–9337, 2019. URL [https://api.semanticsc](https://api.semanticscholar.org/CorpusID:125953399)
439 [holar.org/CorpusID:125953399](https://api.semanticscholar.org/CorpusID:125953399).
- 440 [33] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In Proceedings of
441 the 30th International Conference on Neural Information Processing Systems, NIPS'16, page
442 4572–4580, Red Hook, NY, USA, 2016. Curran Associates Inc. ISBN 9781510838819.
- 443 [34] Brenna Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot
444 learning from demonstration. Robotics and Autonomous Systems, 57:469–483, 05 2009. doi:
445 [10.1016/j.robot.2008.10.024](https://doi.org/10.1016/j.robot.2008.10.024).
- 446 [35] Stephane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and
447 structured prediction to no-regret online learning. In Geoffrey Gordon, David Dunson, and
448 Miroslav Dudík, editors, Proceedings of the Fourteenth International Conference on Artificial
449 Intelligence and Statistics, volume 15 of Proceedings of Machine Learning Research, pages
450 627–635, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL [https://proceedi](https://proceedings.mlr.press/v15/ross11a.html)
451 [ngs.mlr.press/v15/ross11a.html](https://proceedings.mlr.press/v15/ross11a.html).
- 452 [36] Sonia Chernova and Manuela Veloso. Interactive policy learning through confidence-based
453 autonomy. J. Artif. Int. Res., 34(1):1–25, jan 2009. ISSN 1076-9757.
- 454 [37] Snehal Jauhri, Carlos Celemin, and Jens Kober. Interactive imitation learning in state-space. In
455 Jens Kober, Fabio Ramos, and Claire Tomlin, editors, Proceedings of the 2020 Conference on
456 Robot Learning, volume 155 of Proceedings of Machine Learning Research, pages 682–692.
457 PMLR, 16–18 Nov 2021. URL [https://proceedings.mlr.press/v155/jauhri](https://proceedings.mlr.press/v155/jauhri21a.html)
458 [i21a.html](https://proceedings.mlr.press/v155/jauhri21a.html).
- 459 [38] Shih-Yi Chien, Yi-Ling Lin, Pei-Ju Lee, Shuguang Han, Michael Lewis, and Katia Sycara.
460 Attention allocation for human multi-robot control: Cognitive analysis based on behavior data
461 and hidden states. International Journal of Human-Computer Studies, 117:30–44, 2018. ISSN
462 1071-5819. doi:<https://doi.org/10.1016/j.ijhcs.2018.03.005>. URL [https://www.sc](https://www.sciencedirect.com/science/article/pii/S107158191830096X)
463 [iencedirect.com/science/article/pii/S107158191830096X](https://www.sciencedirect.com/science/article/pii/S107158191830096X). Cognitive
464 Assistants.
- 465 [39] Robin R. Murphy and Erika Rogers. Cooperative Assistance for Remote Robot Supervision.
466 Presence: Teleoperators and Virtual Environments, 5(2):224–240, 08 1996. doi:10.1162/pres
467 [.1996.5.2.224](https://doi.org/10.1162/pres.1996.5.2.224). URL <https://doi.org/10.1162/pres.1996.5.2.224>.
- 468 [40] Kunal Menda, Katherine Rose Driggs-Campbell, and Mykel J. Kochenderfer. Ensembledag-
469 ger: A bayesian approach to safe imitation learning. In 2019 IEEE/RSJ International
470 Conference on Intelligent Robots and Systems, IROS 2019, Macau, SAR, China, November
471 3-8, 2019, pages 5041–5048. IEEE, 2019. doi:10.1109/IROS40897.2019.8968287. URL
472 <https://doi.org/10.1109/IROS40897.2019.8968287>.
- 473 [41] Ryan Hoque, Ashwin Balakrishna, Ellen Novoseller, Albert Wilcox, Daniel S. Brown, and
474 Ken Goldberg. Thriftydagger: Budget-aware novelty and risk gating for interactive imitation
475 learning. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, Proceedings of the
476 5th Conference on Robot Learning, volume 164 of Proceedings of Machine Learning Research,
477 pages 598–608. PMLR, 08–11 Nov 2022. URL [https://proceedings.mlr.press/v](https://proceedings.mlr.press/v164/hoque22a.html)
478 [164/hoque22a.html](https://proceedings.mlr.press/v164/hoque22a.html).

- 479 [42] Ryan Hoque, Lawrence Yunliang Chen, Satvik Sharma, Karthik Dharmarajan, Brijen Thanan-
480 jeyan, Pieter Abbeel, and Ken Goldberg. Fleet-dagger: Interactive robot fleet learning
481 with scalable human supervision. In Karen Liu, Dana Kulic, and Jeff Ichnowski, edi-
482 tors, Proceedings of The 6th Conference on Robot Learning, volume 205 of Proceedings
483 of Machine Learning Research, pages 368–380. PMLR, 14–18 Dec 2023. URL [https://](https://proceedings.mlr.press/v205/hoque23a.html)
484 proceedings.mlr.press/v205/hoque23a.html.
- 485 [43] Andreas Krause and Daniel Golovin. Submodular function maximization. In Tractability,
486 2014.
- 487 [44] Baharan Mirzasoleiman, Amin Karbasi, Rik Sarkar, and Andreas Krause. Distributed sub-
488 modular maximization. Journal of Machine Learning Research, 17(235):1–44, 2016. URL
489 <http://jmlr.org/papers/v17/mirzasoleiman16a.html>.
- 490 [45] Hui Lin and Jeff Bilmes. A class of submodular functions for document summarization. In
491 Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics:
492 Human Language Technologies - Volume 1, HLT ’11, pages 510–520, USA, 2011. Association
493 for Computational Linguistics. ISBN 9781932432879.
- 494 [46] Arash Asadpour, Hamid Nazerzadeh, and Amin Saberi. Stochastic submodular maximization.
495 In Christos Papadimitriou and Shuzhong Zhang, editors, Internet and Network Economics,
496 pages 477–489, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.
- 497 [47] Daniel Golovin and Andreas Krause. Adaptive submodularity: Theory and applications in
498 active learning and stochastic optimization. Journal of Artificial Intelligence Research, 42, 03
499 2010. doi:10.1613/jair.3278.
- 500 [48] A. M. Frieze. A cost function property for plant location problems. Math. Program., 7(1):
501 245–248, dec 1974. ISSN 0025-5610. doi:10.1007/BF01585521. URL [https://doi.or](https://doi.org/10.1007/BF01585521)
502 [g/10.1007/BF01585521](https://doi.org/10.1007/BF01585521).
- 503 [49] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
504 policy optimization algorithms, 2017.
- 505 [50] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
506 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac
507 gym: High performance gpu-based physics simulation for robot learning, 2021.
- 508 [51] Ookla. Internet speed dataset. [https://www.kaggle.com/datasets/dhruvildav](https://www.kaggle.com/datasets/dhruvildave/ookla-internet-speed-dataset)
509 [e/ookla-internet-speed-dataset](https://www.kaggle.com/datasets/dhruvildave/ookla-internet-speed-dataset), 2022. [Online; accessed 15-February-2024].