# An Active Perception Game for Robust Autonomous Exploration

Siming He, Yuezhan Tao, Igor Spasojevic, Vijay Kumar and Pratik Chaudhari
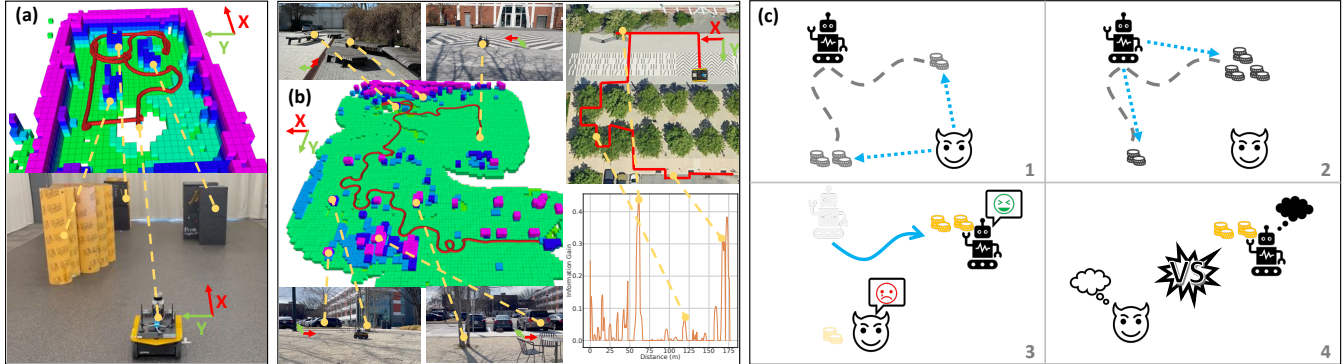


**Fig. 1: (a) and (b): A Jackal ground robot exploring the environments with the proposed algorithm.** The robot actively plans paths that maximize the predicted information gain in indoor (a) and outdoor (b) environments. **(c): The evolution of an active perception game.** The environment owns a certain amount of information. While the environment aims to prevent the agent from collecting the information, the agent seeks to maximize its information gain. For each episode, the environment tries to mislead the agent by revealing erroneous information gain of two action sequences in Step 1. If the agent makes decisions based on the erroneous information gain, the agent could be misled to take actions with low (true) information gain. Therefore, before making decisions, the agent has to predict the true information gain by fixing the erroneous information gain in Step 2. If the agent makes a good prediction, it will collect high (true) information gain in Step 3. In Step 4, the agent and the environment think about their winning strategies for the next episode.

*Abstract*— We formulate active perception for an autonomous agent that explores an unknown environment as a two-player zero-sum game: the agent aims to maximize information gained from the environment while the environment aims to minimize the information gained by the agent. In each episode, the environment reveals a set of actions with their potentially erroneous information gain. In order to select the best action, the robot needs to recover the true information gain from the erroneous one. The robot does so by minimizing the discrepancy between its estimate of information gain and the true information gain it observes after taking the action. We propose an online convex optimization algorithm that achieves sub-linear expected regret $O(T^{3/4})$ for estimating the information gain. We also provide a bound on the regret of active perception performed by any (near-)optimal prediction and trajectory selection algorithms. We evaluate this approach using semantic neural radiance fields (NeRFs) in simulated realistic 3D environments to show that the robot can discover up to 12% more objects using the improved estimate of the information gain. On the M3ED dataset, the proposed algorithm reduced the error of information gain prediction in occupancy map by over 67%. In real-world experiments using occupancy maps on a Jackal ground robot, we show that this approach can calculate complicated trajectories that efficiently explore all occluded regions.

## I. INTRODUCTION

As robots transition from controlled laboratory environments into real-world settings, their ability to comprehend their surroundings becomes increasingly crucial. This ability

General Robotics, Automation, Sensing and Perception (GRASP) Laboratory, University of Pennsylvania. Email: {siminghe, yztao, igorspas, kumar, pratikac}@seas.upenn.edu

will enable robots to plan and make decisions in complex environments adaptively. There is an increase of research interest in active perception [1], [2]. Active perception algorithms can be crucial in search and rescue, planetary exploration, environmental monitoring, structure inspection, etc.

Active perception is commonly formulated as finding a control sequence that maximizes the information gained from the environment. This requires the agent to predict the information gain of putative actions. Typically, this is done on the basis of some map that the robot builds using past observations. The "true" information gained is the incremental information from the actual realized future observations. Of course, the robot could never correctly estimate this true information gain *a priori*. But it can compare the estimated information gain before each observation to the true information gained after the observation, to build a better estimate for the next episode. This paper develops an approach to do so, and, as a consequence, obtain a robust active perception pipeline that ekes out more from potentially poor past exploratory actions. We formulate the active perception problem as a two-player zero-sum game between the agent and the environment as illustrated and explained in Fig. 1. An introductory video along with the experiments can be accessed at: https://youtu.be/p7Elon583q4.

The contributions of our paper are as follows:

- We formulate robust active perception as a two-player zero-sum game between the agent and the environment

in Section II. Then, we prove the regret bound on active perception performance given any (near-)optimal prediction and trajectory selection algorithms in Section IV-C.

- We present a robust active perception algorithm to make and improve information gain prediction online in Section III. We prove this algorithm has a sub-linear expected regret $O(T^{3/4})$ for prediction in Section IV-B.
- As shown in Section V, we integrate the proposed algorithms into two different active perception systems that use semantic NeRF and occupancy maps respectively. Then, we conduct extensive experiments 1) in a photo-realistic simulator, 2) on the M3ED dataset, and 3) on a customized Clearpath Jackal platform. The experiments demonstrate the effectiveness and generalizability of the proposed method.

### A. Related Work

Related works have studied active perception systems based on different scene representations. Each representation has its unique application and challenges in the context of active perception.

*1) Gaussian Processes (GP):* Active sensing in GP has been well studied in [3]. Finding optimal sensor placement is an NP-complete problem. The work obtains near-optimal online bounds for greedy algorithms based on the submodularity of mutual information (MI). The analysis assumes GP has a known covariance function which is usually unrealistic in real-world applications. An extension [4] uses a near-optimal exploration-exploitation approach for stationary GP with an unknown covariance function. For non-stationary GP, a single estimation of covariance function is not enough. The paper [4] proposed manually dividing the whole space into smaller areas with a stationary process. Then, a similar analysis also gives an efficient and near-optimal algorithm for informative path planning in GP [5], [6], [7].

These works provide a near-optimal performance guarantee when the covariance function and MI prediction are easy and accurate. However, many robotics tasks require more complicated representations such as non-stationary GP, occupancy map, point cloud, or semantic map. In those representations, different regions could have very different and diverse characteristics. It's unrealistic to model every small area one by one. In addition, the observation model could be unclear in semantic and implicit representations. As a result, the mutual information prediction could be inaccurate, and the performance guarantee would not hold. To bridge this gap, our paper firstly establishes the relationship between mutual information prediction loss and performance guarantee of active perception. Then, we provide an online algorithm with a expected regret guarantee on the prediction loss.

*2) Occupancy Map:* Occupancy map is commonly used in robotics tasks. Due to the complicated nature of such representation, most active perception algorithms with occupancy maps don't provide performance guarantee analysis. Next-best-views (NBV) methods samples and selects viewpoint that maximizes some utility functions. Utilities based on the amount of observed volume is used in [8], [9], [10]. The following works [11], [12] used surface reconstruction utility in addition to volumetric utility. Information-theoretic methods [13], [14], [15] design efficient algorithms to calculate mutual information between range sensor data and occupancy maps. Since the occupancy of unseen areas is unknown, NBV and information-theoretic methods could underestimate or overestimate the utility. To address this issue, several works [16], [17], [18] propose map completion. The methods train neural networks to predict the occupancy of unseen areas based on existing datasets. Those methods work well in environments similar to the datasets but would struggle in out-of-distribution environments. Hence, the methods could not guarantee a general algorithm performance at test time. In contrast, based on online learning algorithm, our method provides performance guarantees at test time even in out-of-distribution or adversarial environments.

*3) Semantic and Implicit representation:* Semantic scene representations become popular in recent years and have been used for loop-closure, navigation, and decision making. Uncertainty in semantics [19], [20], [21] is used as a utility for exploration. Our theory and algorithm can be applied to the methods and provide performance guarantees. Reinforcement learning methods [22], [23], [24] are also used for active perception by learning to act based on past data. Unlike those methods, our method will not struggle in out-of-distribution environments.

## II. PROBLEM FORMULATION

We will first discuss the formal description of the active perception game. We consider a sequential game between Nature and the robot that evolves over $T$ episodes as follows. At the start, Nature selects a scene $\xi$ that is unknown to the robot but is kept fixed throughout the game. In each episode, $i \in \{1, \ldots, T\}$, Nature selects and reveals to the robot $N$ vantage points $X_i$ and the erroneous information gain $\tilde{I}_i$ that these viewpoints provide about the unknown scene $\xi$. Let $Y_i$ be the set of all measures that the robot could obtain from each of the locations $X_i$ respectively. The robot selects a subset $x_i$ of points in $X_i$ subject to some budget constraints and collects $n$ measurements

$$y_{\text{future}}(i) = y_{ni+1}^{n(i+1)} \in Y_i$$

from each of these viewpoints. We use $y_{ni+1}^{n(i+1)}$ to denote the set of $n$ observations in episode $i$, i.e., from the $(ni+1)$-th to the $(n(i+1))$-th observation in all collected observations. Using these observations, the robot can calculate the true information gain $I_i$ of these locations $x_i$. The game then

transitions to the subsequent episode with Nature revealing a new set of candidate vantage points $X_{i+1}$ and their erroneous information gain, and so on... The game ends after $T$ episodes. The history of measurements of the agent in episode $i$ is then

$$y_{\text{past}}(i) \doteq y_1^{ni}.$$

The protocol of the game is summarized in Definition 2.1. The objective of the agent is to maximize the amount of true information gain collected about the environment. In the following sections, we describe in more detailed notions such as accuracy, map measurements, and a suitable notion of (approximate) optimality. We use $\xi$ to denote the hidden quantity of interest, such as the map of the environment. If we did not have any observations, the scene is drawn from a probability distribution $p(\xi)$. If we have past observations $y_{\text{past}}(i)$, it tells us that the scene will be drawn from a more concentrated distribution $p(\xi_{\text{past}}(i)) = p(\xi \mid y_{\text{past}}(i))$.

*Definition 2.1 (Active Perception Game):* For episodes $i = 1, \ldots, T$:

0) Nature selects $X_i$, together with $I(y_{\text{past}}(i); y_{\text{future}}(i))$ and a erroneous version $\tilde{I}(\xi_{\text{past}}(i); y_{\text{future}}(i))$ for $y_{\text{future}}(i) \in \mathcal{P}(Y_i)$, where $\mathcal{P}(Y_i)$ is the power set of $Y_i$.
1) With budget constraints, robot selects a subset $y_{\text{future}}(i)$ of observations based on context $c_i$ which consists of $y_{\text{past}}(i)$, $\{\tilde{I}(\xi_{\text{past}}(s); y_{\text{future}}(s)) \mid y_{\text{future}}(j) \in \mathcal{P}(Y_s), s \le i\}$, and $\{I(\xi_{\text{past}}(s); y_{\text{future}}(s))\}_{s<i}$.
2) Robot receives measurements $y_{\text{future}}(i)$ and $I(\xi_{\text{past}}(i); y_{\text{future}}(i))$ after visiting points $x_i$.

The information gain is defined as the mutual information between past and future observations:

$$I(\xi_{\text{past}}(i); y_{\text{future}}(i)) = H(\xi_{\text{past}}(i)) + H(y_{\text{future}}(i)) - H(\xi_{\text{past}}(i), y_{\text{future}}(i))$$

where $H(\xi_{\text{past}}(i)) = -\int dp(\xi_{\text{past}}(i)) \log p(\xi_{\text{past}}(i))$. The task of the agent is to develop a policy $\pi_{1 \le i \le T}$ to minimize

$$H(\xi \mid c_T) \tag{1}$$

which is equivalent to maximize

$$\sum_{i=1}^{T} H(\xi_{\text{past}}(i)) - H(\xi_{\text{past}}(i) \mid y_{\text{future}}(i)) \tag{2}$$

where each $\pi_i$ can be regarded as a mapping of the following form $\pi_i : c_i \mapsto y_{\text{future}}(i)$.

## III. ACTIVE PERCEPTION ALGORITHM

Active perception aims at obtaining information of the scene through control of an agent. We can formalize it as selecting control sequences $u_{\text{future}}$ to maximize the mutual information in each episode $i$

$$\hat{u}_{\text{future}} \in \underset{u_{\text{future}}}{\arg\max} \, I(y_{\text{past}}(i); y_{\text{future}}(i)). \tag{3}$$

The mutual information objective depends on selected future observations $y_{\text{future}}(i)$ which is determined by the control $u_{\text{future}}$. Since we don't know true information gain
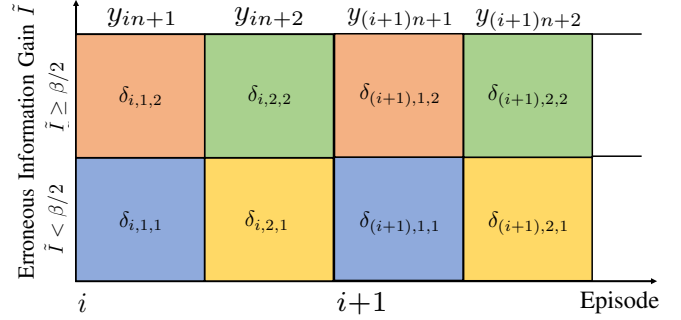


Fig. 2: **Illustration of Subgroup (Definition 3.1) and Improvement (Definition 3.2) Functions.** Assume there are $b = 2$ bins and $n = 2$ observations per episode. Four subgroups are colored differently. Observation $y_{in+j}$ denotes the $j$-th observation in episode $i$. Each observation can be assigned into one of the four subgroups based on its erroneous gain $\tilde{I}$ and its location $j$ in each episode. For example, the second observation ($j = 2$) in episode $i$ is $y_{in+2}$. If $\tilde{I}$ of $y_{in+2}$ is less than $\beta/2$, the algorithm refines the erroneous gain by $f(\tilde{I}) = \tilde{I} + \delta_{i,2,1}$. At the end of episode $i$, $\delta_{i,2,1}$ is refined to $\delta_{(i+1),2,1}$ based on true information gain $I$.

$I(y_{\text{past}}(i); y_{\text{future}}(i))$ beforehand, the agent makes an estimate $\hat{I}(\xi_{\text{past}}(i); y_{\text{future}}(i))$ of the true information gain based on the erroneous $\tilde{I}(\xi_{\text{past}}(i); y_{\text{future}}(i))$. The estimate is based on the improvement function

$$f_i : \tilde{I}(\xi_{\text{past}}(i); y_{\text{future}}(i)) \mapsto \hat{I}(\xi_{\text{past}}(i); y_{\text{future}}(i)). \tag{4}$$

Then, the agent maximizes $\hat{I}(\xi_{\text{past}}(i); y_{\text{future}}(i))$.

*Remark 3.1:* The erroneous information gain could be any estimate of the actual information gain. Therefore, our method is a general way that can be applied to numerous existing active perception systems that predict $\tilde{I}(\xi_{\text{past}}; y_{\text{future}})$. For many active perception systems, the predictions are usually inaccurate due to (a) complexity of scene representation; (b) aleatoric uncertainty from unmodeled sensor noise; (c) assumption of independent observations.

We show in Section IV that our algorithms have *no-regret* predictions of information gain. This leads to a bounded regret for active perception performance.

### A. Information Gain Prediction

An agent can make prediction based on erroneous information gain by IPO (Improving Prediction Online) in Algorithm 1. At the end of each episode, the agent also learns to make better prediction by LIP (Learning to Improve Prediction) in Algorithm 2.

We firstly explain the notations we used in the algorithms. For episode $i$ and $j \in \{1, \ldots, n\}$, let

$$I_{in+j} \equiv I(\xi_{\text{past}}; y_{in+j} \mid y_{in}^{in+j-1}). \tag{5}$$

By chain rule, we have $I(\xi_{\text{past}}; y_{\text{future}}(i)) = \sum_{j=1}^{n} I_{in+j}$. Similarly, we define $\tilde{I}_{in+j}$ and $\hat{I}_{in+j}$.

The subgroup and improvement function could be any reasonable function. Considering the generalization and the sample complexity of our algorithms, we choose classes of functions that are simple but effective. We illustrate our choices in Fig. 2.

**Algorithm 1** Improving Prediction Online (IPO)

   **Input** $\{\tilde{I}\}_{in+j}$, improvement function $f_i$
   **Output** improved prediction $\{\hat{I}\}_{in+j}$
1: **begin**
2:    $A \leftarrow \{\}$
3:    **for** $\tilde{I} \in \{\tilde{I}\}_{in+j}$ **do**
4:       $A \leftarrow A \cup \{f_i(\tilde{I})\}$
5:    **end for**
6:    **return** $A$
7: **end**

---

***Definition 3.1 (Subgroup Function $\mathbb{1}$):*** Assuming $\tilde{I}_{in+j}$ is in $[0, \beta]$, we discretize $[0, \beta]$ to $b$ bins. The initial prediction $\tilde{I}_{in+j}$ is in the $k$-th bin if $\tilde{I}_{in+j} \in [(k-1)\beta/b, k\beta/b)$. For $k = 1, \dots, b$, we define subgroup function $\mathbb{1}_k : \mathbb{R} \to \{0, 1\}$ such that

$$\mathbb{1}_k(x) = \begin{cases} 1, & \text{if } x \in \left[\frac{(k-1)\beta}{b}, \frac{k\beta}{b}\right) \\ 0, & \text{otherwise} \end{cases}. \quad (6)$$

***Definition 3.2 (Improvement Function $f$):*** The function $f : \mathbb{R} \to \mathbb{R}$ is defined as

$$\hat{I}_{in+j} = f(\tilde{I}_{in+j}) = \tilde{I}_{in+j} + \sum_{k=1}^{b} \mathbb{1}_k(\tilde{I}_{in+j})\delta_{i,j,k}. \quad (7)$$

$\delta_{i,j,k}$ is a learned constant to improve the predictions at $k$-th bin at episode $i$ for subgroup $(j, k)$.

We use value of initial prediction as context because it is closely related to the notion of multicalibration which guarantees low errors even when conditioned on initial predictions [25]. Initial prediction usually assumes the independence of observations $I(\xi_{\text{past}}; y_{\text{future}}(i)) = \sum_{j=1}^{n} I(\xi_{\text{past}}; y_{in+j})$. To reduce prediction errors due to this assumption, the improvement $\delta_{i,j,k}$ also depends on $j$. We could choose different subgroup definition for prediction improvement. The analysis and implementation can be easily extended to other definitions of $\mathbb{1}$ and $f$.

Since we have a set of pixel-wise mutual information feedback at each time in real applications, we use the notation $\{I\}_{in+j}$ to denote a set of mutual information feedback at time $in + j$. We are improving the mutual information prediction of each pixel. We also use the notation $\mathbf{c}(i), \eta, d$ in LIP. $\mathbf{c}(i) \in \mathbb{N}^{n \times b}$ keeps track of the number of past updates in each subgroup. $\mathbf{c}_{j,k}(i)$ is the number of times that predictions are in subgroup (j,k) up to and including episode $i$. $\eta$ is the learning rate and $d$ is the gradient of loss.

In IPO, we want to return a set $A$ of mutual information predictions at each pixel. The algorithm directly applies our $f$ to get each prediction from erroneous information gain.

In LIP, GRADIENT() returns the (sub)gradient of loss $|I - \hat{I}|$. From line 12 to line 17, the algorithm refines the improvement function by updating $\delta_{i,j,k}$.

---

**Algorithm 2** Learning to Improve Prediction (LIP)

   **Input** $\{(\{\tilde{I}\}_{in+j}, \{\hat{I}\}_{in+j}, \{I\}_{in+j})\}_{j=1}^{n}$,
   improvement function $f_i$ , number of updates $\mathbf{c}_i$
   **Output** $f_{i+1}, \mathbf{c}(i+1)$
1: **procedure** GRADIENT($\hat{I}$, $I$)
2:    **if** $\hat{I} \geq I$ **then**
3:       **return** -1
4:    **else**
5:       **return** 1
6:    **end if**
7: **end procedure**
8:
9: **begin**
10:    **for** $j \leftarrow 1$ to $n$ **do**
11:       **for** $(\tilde{I}, \hat{I}, I) \in (\{\tilde{I}\}_{in+j}, \{\hat{I}\}_{in+j}, \{I\}_{in+j})$ **do**
12:          $d \leftarrow$ GRADIENT($\hat{I}$, $I$)
13:          $k \leftarrow \text{argmax}_k \mathbb{1}_k(\tilde{I})$
14:          $\eta_{\text{previous}} \leftarrow \beta/2\sqrt{1/\mathbf{c}_{j,k}(i)}$
15:          $\mathbf{c}_{j,k}(i) \leftarrow \mathbf{c}_{j,k}(i) + 1$
16:          $\eta_{\text{current}} \leftarrow \beta/2\sqrt{1/\mathbf{c}_{j,k}(i)}$
17:          $\delta_{i,j,k} \leftarrow \delta_{i,j,k} \cdot \eta_{\text{current}}/\eta_{\text{previous}} + d \cdot \eta_{\text{current}}/2$
18:       **end for**
19:    **end for**
20:    **return** $f_i, \mathbf{c}(i)$
21: **end**

---

**Algorithm 3** Informative Path Planning (IPP)

   **Input** prediction mechanism $\hat{I}$, optimization oracle $\pi$, randomization level $\tau$
   **Output** the selected future observations $y_{\text{future}}$
1: **begin**
2:    $y \leftarrow \pi(\hat{I})$
3:    **for** $j \in \{1, 2, \dots, n\}$ **do**
4:       $z \sim \text{Uniform}(0, 1)$
5:       **if** $z < \tau$ **then** $y_j \leftarrow$ randomly selected $y' \in Y_i$
6:       **end if**
7:    **end for**
8:    **return** $y$
9: **end**

---

*B. Informative Path Planning*

Optimizing over the objective (3) is known to be NP-complete and many approximation algorithms have been proposed to achieve near-optimality with tractable runtime [3].

In episode $i$, we denote $\pi^*$ the policy to find the $y_{\text{future}}^*(i)$ that maximizes predicted information gain given prediction mechanism $\hat{I}$, i.e.,

$$\pi^*(\hat{I}) = \underset{y_{\text{future}}(i)}{\text{argmax}} \, \hat{I}(\xi_{\text{past}}; y_{\text{future}}(i)).$$

The information gain given this policy is

$$v(\pi^*(\hat{I}), I) = I(\xi_{\text{past}}; y_{\text{future}}^*(i)). \quad (8)$$

Using $\pi^*$ is feasible for a small set of sampled observations since we can find $y_{\text{future}}^*(i)$ by iterating through all possible

---

**Algorithm 4** Combined Algorithm

---

1: **begin**
2:     $c \leftarrow \mathbf{1} \in \mathbb{R}^{n \times b}$
3:     **for** $i \in \{1, \ldots, T\}$ **do**
4:         Nature reveals $\tilde{I}(\xi_{\text{past}}(i); y_{\text{future}}), y_{\text{future}} \in \mathcal{P}(Y_i)$
5:         $\hat{I}(\cdot) \leftarrow \text{IPO}(\cdot, f_i, \tau)$
6:         $y_{\text{future}}(i) \leftarrow \text{IPP}(\hat{I}(\cdot), \pi)$
7:         Agent receives $y_{\text{future}}(i)$ and $I(\xi_{\text{past}}(i), y_{\text{future}}(i))$
8:         $\mathcal{I} \leftarrow (\tilde{I}, \hat{I}, I)(\xi_{\text{past}}(i), y_{\text{future}}(i))$
9:         $(f_{i+1}, c(i+1)) \leftarrow \text{LIP}(\mathcal{I}, f_i, c(i))$
10:     **end for**
11: **end**

---

$y_{\text{future}}(i)$. For larger set of sampled observations, approximation algorithm for the policy has to be used. For example, the eSIP algorithm $\pi_e$ in [5] is proved to satisfy

$$\sum_{t=1}^{T} v(\pi_e(\hat{I}), I) \geq (1 - 1/e)/(1 + \log_2 k) \sum_{t=1}^{T} v(\pi^*(\hat{I}), I).$$

In general, our algorithm and analysis work for any $\pi$ such that

$$\sum_{t=1}^{T} v(\pi(\hat{I}), I) \geq \gamma v(\pi^*(\hat{I}), I), \gamma \in (0, 1].$$

IPP in Algorithm 3 takes in an optimization oracle $\pi$ and the prediction mechanism $\hat{I}$. Line 2 stores the selected path. For each observation on the path, we sample $z$ from a uniform distribution in line 3. We denote $y_j$ the j-th observation. In line 5, with probability $\tau$, we set $y_j$ to be a randomly selected observation from $Y_i$. $\tau \in (0, 1)$ controls the amount of randomization (exploration) by our algorithm.

### C. Combined Algorithm

In Algorithm 4, we integrate IPO, IPP, LIP into the active perception system. Given the context from Nature (line 3), the agent uses information gain prediction to select measurements (line 4 and 5) to observe. After observing the measurements and receiving true information gain (line 6), the agent learns to improve prediction (line 8).

## IV. ANALYSIS

Though our algorithms are based on pixel-wise mutual information feedback $\{I\}_{in+j}$, we assume single pixel observation and mutual information feedback $I_{in+j}$ at each time step in this analysis section for conciseness. The analysis and conclusion would be the same up to constant terms.

### A. Regret for Online Prediction with Full Information

Our problem is in bandit-feedback setting (i.e., only receive true information gain for observed views). We firstly analyze our algorithm in full-information setting (i.e., receive true information gain for all views in $Y_i$). Then, we show a reduction from bandit-feedback setting to full-information setting. In full-information setting, we set $\tau = 0$.

Since each subgroup is improved independently, we start by analyzing a single subgroup $(j, k)$. The number of predictions in subgroup $(j, k)$ in episode $i$ is

$$c_{j,k}(i) = \sum_{y_{in+j} \in Y_i} \mathbb{1}_k(\tilde{I}_{in+j})$$

and the total predictions in the subgroup is

$$c_{j,k} = \sum_{i=1}^{T} c_{j,k}(i)$$

for full-information setting. We denote

$$\{I_a\}_{a=1}^{c_{j,k}(i)}$$

the set of predictions that are in subgroup $(j, k)$ in episode $i$. This is different from the notation $I_{in+j}$ in previous section, but note that

$$\{\{I_a\}_{a=1}^{c_{j,k}} \mid j = 1, \ldots, n; k = 1, \ldots, b\}$$

is a partition of all predictions

$$\{I_{in+j} \mid y_{in+j} \in Y_i; i = 1, \ldots, T; j = 1, \ldots, n\}.$$

Online prediction in full-information setting is an online convex optimization problem: for each episode $i$ and subgroup $(j, k)$, Nature chooses a convex loss by selecting all $I_a - \tilde{I}_a$, then the agent selects predictions $\delta_{i,j,k}$ to minimize the loss. The convex loss is

$$l_{i,j,k} = \sum_{a=1}^{c_{j,k}(i)} |I_a - \hat{I}_a| = \sum_{a=1}^{c_{j,k}(i)} |I_a - \tilde{I}_a - \delta_{i,j,k}|. \quad (9)$$

We define regret as the performance measure. When the agent finishes episode $i$ and looks back, it realizes $\bar{\delta}_{j,k}$ is the prediction that it should have made. Regret is the additional loss due to predicting $\delta_{i,j,k}$ instead of $\bar{\delta}_{j,k}$.

***Definition 4.1 (Subgroup Prediction Regret):*** For $k = 1, \ldots, b$ and $j = 1, \ldots, n$, we have subgroup prediction regret

$$R(j, k) = \sum_{i=1}^{T} \left[ l_{i,j,k} - \bar{l}_{i,j,k} \right] \quad (10)$$

where $\bar{l}_{i,j,k}$ is the loss when

$$\bar{\delta}_{j,k} = \arg\min_{\delta} \sum_{i=1}^{T} \sum_{a=1}^{c_{j,k}(i)} \left| I_a - \tilde{I}_a - \delta \right|. \quad (11)$$

The overall regret is

$$R = \sum_{j=1}^{n} \sum_{k=1}^{b} R(j, k). \quad (12)$$

An algorithm is said to be *no-regret* if the average regret $R/T \to 0$ as $T \to \infty$. *Follow the leader* algorithms solve online convex optimizations by selecting the best past prediction as the next prediction [26]. Regularization is usually added to the algorithms to avoid oscillatory prediction. We firstly show that our LIP algorithm is a special case of *follow the regularized leader* in Lemma 6.1. Then, we provide the regret bound of this special case in Lemma 6.2.

By choosing a suitable scheduler of the learning rate, we can obtain the sub-linear regret of the proposed algorithm

***Theorem 4.1:*** LIP uses learning rate $\eta_{i,j,k} = \beta\sqrt{1/i}$. The

regret bound with this learning rate is
$$R(j,k) \leq \beta(N\sqrt{T}+1).$$
And the bound on R is in worst case
$$nb\beta(N\sqrt{T}+1).$$

*Proof:* Plug in the specified $\eta_{i,j,k}$ to the bound in Lemma 6.2. $R(j,k)$ is bounded due to Holder's inequality $\sum_{i=1}^{n}\sqrt{1/i} \leq 2\sqrt{n}$. R is bounded by Jensen's inequality. ∎

***Remark 4.1:*** Our online prediction algorithm with IPO and LIP has *no-regret* since $\beta(N\sqrt{T}+1)/T$ and $nb\beta(N\sqrt{T}+1)/T$ goes to 0 as $T \to \infty$.

### B. Regret for Online Prediction with Bandit Feedback

We follow the similar process in [27] to reduce bandit-feedback setting to full-information setting. With bandit feedback, the agent can only observe part of the loss $|I_a - \hat{I}_a|$ that is corresponding to $y_a \in y_{\text{future}}(i)$ where $(I_a, y_a)$ is the $a$-th observed (information gain, observation) in subgroup $(j,k)$ in episode $i$. For $a = 1, \ldots, c_{j,k}(i)$, $y_a$ is observed with probability $p(y_a \in y_{\text{future}}(i))$. To use the full-information algorithm in Section IV-A for bandit-feedback problem, we firstly need to hallucinate the loss. Let $\hat{l}_{i,j,k}(a)$ be the hallucination for the $a$-th loss in episode $i$ for subgroup $(j,k)$. For $a = 1, \ldots, c_{j,k}(i)$,

$$\hat{l}_{i,j,k}(a) = \begin{cases} \frac{|I_a - \hat{I}_a|}{p(y_a \in y_{\text{future}}(i))} & y_a \in y_{\text{future}}(i) \\ 0 & \text{otherwise} \end{cases}. \quad (13)$$

We show this filled loss is unbiased Lemma 6.3 which is important for analyzing regret in expectation. Additionally, to make this filled loss well-defined, we need $p(y_a \in y_{\text{future}}(i))$ to be strictly positive for all $y_a \in y_{\text{future}}(i)$. We achieve this by creating a randomized $\check{y}_{\text{future}}(i)$ of $y_{\text{future}}(i)$. Let $\tau \in (0, \frac{1}{2})$ in LIP. The $j$-th element of $\check{y}_{\text{future}}(i)$ equals to the $j$-th element of $y_{\text{future}}(i)$ with probability $1 - \tau$. With probability $\tau$, we let the $j$-th element of $\check{y}_{\text{future}}(i)$ be randomly sampled from $Y_i$. This randomization gives strictly positive probability for observing each observation as shown in Lemma 6.4. Intuitively, in bandit-feedback setting where the agent cannot collect true information gain for all observations in $Y_i$, the agent needs to explore (by randomization) to understand the true information gain of different observations.

After randomizing LIP, we prove the expected regret in bandit-feedback setting in Lemma 6.5. With a suitable randomization level, the algorithms have a sublinear regret.

***Theorem 4.2:*** Let $\tau = T^{-1/4}$. The upper bound of expected regret $\mathbb{E}(R_{\text{bandit}})$ is
$$N\beta(NT^{3/4} + T^{1/4}) + \beta T^{3/4} = O(T^{3/4}).$$

*Proof:* Plug in the specified $\tau$ into the bound in Lemma 6.5. ∎

***Remark 4.2:*** The randomized algorithm is *no-regret* in bandit-feedback setting.

### C. Regret Bound for Robust Active Perception

We want to understand the relationship among quality of information gain predictions, optimality of path planning algorithms, and performance of active perception in this section. Let's denote $\bar{I}(\xi_{\text{past}}; y_{\text{future}})$ the prediction of competitor where each $\bar{I}_i, i = 1, \ldots, c_{j,k}$ is based on (11). Assume that 1) the overall regret R is upper-bounded by $\alpha$; 2) the performance of competitor $\sum_{j=1}^{n}\sum_{k=1}^{b}\sum_{i=1}^{T}\bar{l}_{i,j,k}$ is upper-bounded by $\alpha'$. For concise notations, we will use $I, \hat{I}, \tilde{I}, \bar{I}$ without $(\xi_{\text{past}}, y_{\text{future}})$ in the following section.

We define the regret of active perception as the extra information gain if the agent makes decision based on $\bar{I}$) instead of $\hat{I}$. We firstly show that improved predictions of information gain lead to a tighter bound on the regret of active perception.

***Theorem 4.3:*** The bound on the regret of active perception is

$$\mathbb{E}\left[\sum_{i=1}^{T} v(\pi^*(\bar{I}), I) - v(\pi^*(\hat{I}), I)\right] \leq 2(\alpha + \alpha') \quad (14)$$

See *Proof* in Section VI-B.

The bound $\alpha$ is shown in Theorem 4.2 to be sub-linear in episodes $T$. The bound $\alpha'$ depends on the capacity of the competitor class. Since the competitor class can choose predictions in hindsight, the bound $\alpha'$ should be linear in $T$ with a *small* factor given a suitable competitor class. Additionally, we want to take the optimality of path planning algorithms into account.

***Theorem 4.4 (Active Perception Performance):*** Assume we use a policy $\pi$ in IPP such that
$$v(\pi(\hat{I}), \hat{I}) \geq \gamma v(\pi^*(\hat{I}), \hat{I})$$
where $\gamma$ is described in Section IV-C. The regret of active perception with IPP is

$$\mathbb{E}\left[\sum_{i=1}^{T} v(\pi^*(\bar{I}), I) - v(\pi(\hat{I}), I)\right]. \quad (15)$$

This regret is upper bounded by
$$4(\alpha + \alpha') + (1 - \gamma)v(\pi^*(\hat{I}), \hat{I}).$$

See *Proof* in Section VI-B.

***Remark 4.3:*** If we are using the brute force algorithm $\pi^*$, the regret guarantee is
$$4(N\beta(NT^{3/4} + T^{1/4}) + \beta T^{3/4} + \alpha').$$
And if we are using the eSIP algorithm $\pi_e$, the regret guarantee is
$$4(N\beta(NT^{3/4} + T^{1/4}) + \beta T^{3/4} + \alpha')$$
$$+ (\log_2 k + 1/e)/(1 + \log_2 k)v(\pi^*(\hat{I}), \hat{I})$$
by plugging in the $\gamma$ for eSIP.

## V. RESULTS

To evaluate the effectiveness and generalizability of our method, we conducted three sets of experiments: (1) we

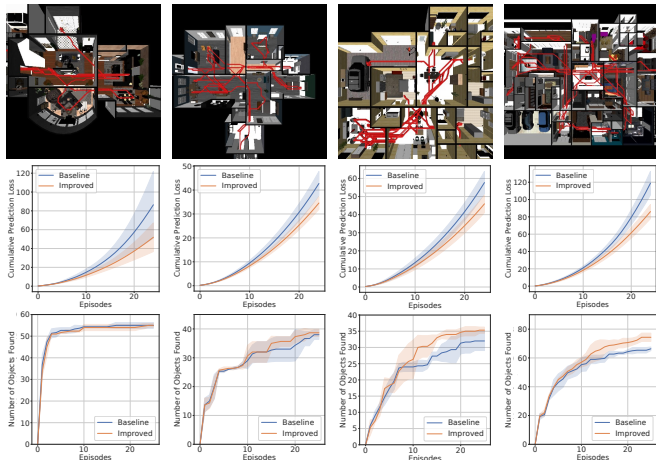Scene 102344529  Scene 102344280  Scene 102344250  Scene 102816036

**Fig. 3: Visualization of simulation experiment performance. Top**: Aerial views of the simulation environments and the trajectory of the proposed algorithm in each environment. **Middle**: We compare the baseline [20] and the proposed methods based on the mean and standard deviation of cumulative prediction loss. **Bottom**: We compare the baseline and the proposed methods based on the number of objects found. We calculate the means and standard deviations after three repetitions of each experiment with the same initialization.

evaluated our method in a photorealistic simulator; (2) we studied the effectiveness of our method with noisy sensor measurements using a real-world dataset; (3) we carried out real-world experiments with our algorithms running fully onboard the robot.

### A. Simulation Experiment

*1) Experiment Setup:* The 3D photorealistic Habitat-Sim [28], [29], [30] is used for our simulation experiments. RGBD images and semantic segmentations are used to train a semantic NeRF during exploration [20]. In all of the simulation experiments, the robot started at a random location in the map and collects an initial set of observations to build the map. Then, we uniformly sampled 20 goals on the x-y plane and generated the paths with Dijkstra's algorithm. We discretized the path with $k$ steps and evaluated the information at each step. The path with the highest aggregated information was passed into Rotorpy [31] to generate a dynamically feasible trajectory for execution.

We use a deep ensemble of NeRFs as $\xi$ in Section II. The $y$ includes RGB images, depth measurements, semantic segmentation, and transmittance. For RGB and depth, we model each ray $r$ in $y$ as a Gaussian distribution $r|\xi_{past} \sim \mathcal{N}(\mu, \sigma)$ where $\mu$ and $\sigma$ are estimated from the ensemble [20]. Hence, $H(r) = \frac{1}{2}\log(2\pi\sigma^2) + \frac{1}{2}$. Since the entropy depends on $\sigma$ which could be arbitrarily large, we set an upper limit to the entropy in practice. In all NeRF experiments, we set the upper bound of the mutual information to be 5. Semantic segmentation generates a distribution over the categories for entropy calculation. We also set the upper bound of the mutual information of segmentation to be 5. Transmittance of a ray is a Bernoulli random variable of

whether the ray hits obstacles. Hence, we set the upper bound of the mutual information of transmittance to be 1. In this setting, we have six categories of information gain $I_R, I_G, I_B, I_D, I_S, I_O$ corresponding to RGBD, segmentation, and transmittance. We improve the information prediction in each of the category independently using our algorithms. We choose the number of bins $b$ in Definition 3.1 to be 100 and the number of measurements $n$ in Section II to be 40. Our algorithm generates the aggregated information prediction $\hat{I}$ for each sampled path and executes the one that maximizes it.

After the execution of a path, we update the map based on collected observations $y_{future}$. Since we assume the absolute certainty in the observation, observed $y_{future}$ is a Dirac-delta distribution which is non-zero only at the observations' value. The feedback is computed as $H(\xi_{past}) - H(\xi_{past}|y_{future})$ using collected observations from the episode. Given the erroneous information gain, agent's improved predictions, and the feedback, we can refine the improvement function following LIP. The level of randomization $\tau$ is set to 0 in practice since randomized path is costly for SWaP-constrained agents.

*2) Experiment Results:* We tested our method in four different simulation environments of varying sizes and complexities. For each experiment, we calculated the cumulative loss of information gain prediction, which is defined as the absolute difference between the information prediction and the feedback. Besides, the number of objects observed during the experiment was recorded. As a result, the proposed method was capable of generating more accurate information predictions and selecting paths that provide observations of more semantic objects, as illustrated in Fig. 3 and Table I. From the statistics of the experiments, the proposed method reduced the error of information gain prediction by 18.6% - 40.0% on average and enabled the robot to observe up to 12% more objects. We envision the improvements will be more significant with the increase in the scale and complexity of the environments.

### B. M3ED Dataset Experiment

*1) Experiment Setup:* We carried out experiments with a high-quality synchronized dataset M3ED [32] to study the performance of our method on information gain prediction under unmodeled sensor noise. In these experiments, we treated the raw lidar data and the associated odometry as the ground truth. We added Gaussian noises and Impulse noises to the ground truth measurements separately to simulate unmodeled sensor noises. Specifically, we applied Gaussian noise with zero mean and standard deviation that equals $1/8$ of the raw depth measurements. For Impulse noise, we replaced half of the measurements with new measurements that are uniformly sampled between $1/3$ and $5/3$ of the original value.

|  | Scene 102344529 | | Scene 102344280 | | Scene 102344250 | | Scene 102816036 | |
|---|---|---|---|---|---|---|---|---|
|  | Loss[1] | #Objects[2] | Loss | #Objects | Loss | #Objects | Loss | #Objects |
| **Baseline** | 86 ± 35.4 | 55 ±0.8 | 43 ± 5.0 | 38±1.6 | 57 ± 6.5 | 32 ± 2.9 | 118 ±13.5 | 66.3 ± 1.2 |
| **Improved** | 52 ± 15.1 | 55 ± 0.8 | 35 ± 2.6 | 38.7 ± 1.2 | 46 ± 5.0 | 35.3 ± 0.9 | 86 ± 8.1 | 74.3 ± 2.9 |

[1] The loss of cumulative information gain prediction.    [2] The number of objects found.

**TABLE I: Exploration statistics for simulation environments.** We conducted experiments with the baseline [20] and proposed methods in each of the Habitat scene in Fig. 3. We obtain the mean and standard deviation of Loss and #Objects by repeating every experiment three times with the same initialization.

|  | Indoor: building_loop | | | Outdoor: penno_short_loop | | |
|---|---|---|---|---|---|---|
| Noise Type → | None | Gaussian | Impulse | None | Gaussian | Impulse |
| **Baseline** | 0.475 | 0.484 | 0.490 | 0.474 | 0.484 | 0.488 |
| **Improved** | 0.163 | 0.160 | 0.162 | 0.162 | 0.163 | 0.172 |

**TABLE II: Prediction improvement for M3ED dataset.** We compare the mean prediction loss of the baseline and proposed methods under different types of depth measurement noise described in Section V-B.

We represented the environment with 3D voxels and assumed they were independent of each other. All voxels were initialized with a probability of occupancy to be $0.5$. For each measurement and its associated odometry, ray-castings were conducted, followed by a log-odds-based map update. To predict the information gain, rays were uniformly sampled at each position. For each ray $r$ that intersects with voxels $m = \{m_1, \ldots, m_n\}$ in $\xi_{\text{past}}$, we calculated $I(\xi_{\text{past}}; r) = H(m) - H(m|r)$. Since each ray maximally intersects with $u = (max\ ray\ length/voxel\ size)$ voxels, we normalized mutual information $H(m)/u - H(n|r)/u$ to better fit this with our algorithms. We summed up the information of all sampled rays and treated this as the information gain of the observation $I(\xi_{\text{past}}, y_{\text{future}})$.

After every observation, we calculated the feedback information gain and the mean prediction loss between the prediction and the feedback. Subsequently, we followed the process described in LIP to improve future information predictions. We set $\beta = 1$, $b = 100$, and $n = 1$ here for our algorithms in Section III-A.

*2) Experiment Results:* Results in Table II demonstrate that our proposed algorithm significantly improved the accuracy of the prediction of information gain given unmodeled sensor noises. In particular, the prediction loss was reduced by over 67%. This set of experiments demonstrates that the proposed method is capable of providing more accurate information gain prediction even with unmodeled sensor noises.

### C. Real World Experiment

*1) Experiment Setup:* To validate the effectiveness of our proposed method, we carried out real-world experiments by deploying it on a customized Clearpath Jackal platform. Details of our platform can be found in [33]. The Jackal UGV is equipped with an AMD Ryzen 3600 CPU and an NVIDIA GTX 1650 GPU. In addition, it carries an Ouster OS1-64 LiDAR. We used Faster-LIO [34] for state estimation and move-base [35] for global and local planning. In all our real-world experiments, we used the same map representations and the definition of information gain as described in Section V-B. To reduce computational overhead and enable real-time onboard operations, we set steps $k$ to be $5$. In each episode, we sampled an $11 \times 11$ grid centered around the robot with cell length $2$ meters. Then we predicted the information gain at each cell lattice. The path with the maximum predicted information gain is searched with DFS on the grid with a max length of $10$ meters. Similar to simulation experiments, we set $\tau$ to $0$.

*2) Experiment Results:* We carried out experiments in both an indoor environment and an urban outdoor environment. One representative result in each environment is shown in Fig. 1. In the indoor environment Fig. 1(a), the robot actively navigated to the regions that are occluded from its start position to perceive more information. Similarly, as shown in outdoor environment in Fig. 1(b), the robot actively navigated to occluded areas that have high actual information gain. Areas that correspond to high information gain from the plot are highlighted. With the proposed algorithm running online, the robot was able to actively plan paths to maximize the gathered information, resulting in full coverage of areas with high information content.

## VI. DISCUSSION

Active exploration of unknown environments is a challenging task. Fundamentally, this is because the robot does not have an accurate model of the environment to explain its observations, and as a consequence, its estimated information gain is incorrect. We developed an online learning approach to incrementally refine the estimated information gain. We showed, across simulation experiments, evaluations on real-world datasets, and experiments on real robotic platforms, that this approach not only improves the estimate of information gain, but also enables robust active perception policies that can explore new environments effectively. These approaches are essential good real-world performance in applications such as search and rescue missions in unknown environments.

### REFERENCES

[1] J. A. Placed, J. Strader, H. Carrillo, N. Atanasov, V. Indelman, L. Carlone, and J. A. Castellanos, "A survey on active simultaneous localization and mapping: State of the art and new frontiers," 2023.

[2] R. Zeng, Y. Wen, W. Zhao, and Y.-J. Liu, "View planning in robot active vision: A survey of systems, algorithms, and applications," 2020.

[3] A. Krause, A. Singh, and C. Guestrin, "Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies," *J. Mach. Learn. Res.*, vol. 9, 2008.

[4] A. Krause and C. Guestrin, "Nonmyopic active learning of gaussian processes: an exploration-exploitation approach," in *Proceedings of the 24th international conference on Machine learning*, 2007.

[5] A. Singh, A. Krause, C. Guestrin, and W. J. Kaiser, "Efficient informative sensing using multiple robots," *J. Artif. Int. Res.*, vol. 34, no. 1, 2009.

[6] J. Binney and G. S. Sukhatme, "Branch and bound for informative path planning," in *IEEE International Conference on Robotics and Automation*, 2012.

[7] A. Viseras, D. Shutin, and L. Merino, "Robotic active information gathering for spatial field reconstruction with rapidly-exploring random trees and online learning of gaussian processes," *Sensors*, vol. 19, no. 5, 2019.

[8] C. Connolly, "The determination of next best views," in *IEEE International Conference on Robotics and Automation*, vol. 2, 1985.

[9] A. Bircher, M. Kamel, K. Alexis, H. Oleynikova, and R. Siegwart, "Receding horizon "next-best-view" planner for 3d exploration," in *IEEE International Conference on Robotics and Automation*, 2016.

[10] L. Yoder and S. Scherer, *Autonomous Exploration for Infrastructure Modeling with a Micro Aerial Vehicle*. Cham: Springer International Publishing, 2016.

[11] R. Border, J. D. Gammell, and P. Newman, "Surface edge explorer (see): Planning next best views directly from 3d observations," in *IEEE International Conference on Robotics and Automation*, 2018.

[12] L. Schmid, M. Pantic, R. Khanna, L. Ott, R. Siegwart, and J. Nieto, "An efficient sampling-based method for online informative path planning in unknown environments," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, 2020.

[13] Z. Zhang, T. Henderson, V. Sze, and S. Karaman, "Fsmi: Fast computation of shannon mutual information for information-theoretic mapping," in *International Conference on Robotics and Automation*, 2019.

[14] B. Charrow, S. Liu, V. Kumar, and N. Michael, "Information-theoretic mapping using cauchy-schwarz quadratic mutual information," in *IEEE International Conference on Robotics and Automation*, 2015.

[15] T. Henderson, V. Sze, and S. Karaman, "An efficient and continuous approach to information-theoretic exploration," in *IEEE International Conference on Robotics and Automation*, 2020.

[16] R. Shrestha, F.-P. Tian, W. Feng, P. Tan, and R. Vaughan, "Learned map prediction for enhanced mobile robot exploration," in *International Conference on Robotics and Automation*, 2019.

[17] Y. Tao, E. Iceland, B. Li, E. Zwecher, U. Heinemann, A. Cohen, A. Avni, O. Gal, A. Barel, and V. Kumar, "Learning to explore indoor environments using autonomous micro aerial vehicles," 2023.

[18] Y. Tao, Y. Wu, B. Li, F. Cladera, A. Zhou, D. Thakur, and V. Kumar, "Seer: Safe efficient exploration for aerial robots using learning to predict information gain," in *IEEE International Conference on Robotics and Automation*, 2023.

[19] G. Georgakis, B. Bucher, K. Schmeckpeper, S. Singh, and K. Daniilidis, "Learning to map for active semantic goal navigation," 2022.

[20] S. He, C. D. Hsu, D. Ong, Y. S. Shao, and P. Chaudhari, "Active perception using neural radiance fields," 2023.

[21] Y. Tao, X. Liu, I. Spasojevic, S. Agarwal, and V. Kumar, "3d active metric-semantic slam," *IEEE robotics & automation letters.*, vol. 9, no. 3, 2024-3.

[22] A. Dosovitskiy and V. Koltun, "Learning to act by predicting the future," 2017.

[23] X. Sun, Y. Wu, S. Bhattacharya, and V. Kumar, "Multi-agent exploration of an unknown sparse landmark complex via deep reinforcement learning," 2022.

[24] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov, "Learning to explore using active neural slam," 2020.

[25] I. Globus-Harris, D. Harrison, M. Kearns, A. Roth, and J. Sorrell, "Multicalibration as boosting for regression," 2023.

[26] A. Roth, "Learning in games," 2023. [Online]. Available: https://www.cis.upenn.edu/~aaroth/GamesInLearning.pdf

[27] A. Slivkins, "Lecture notes: Bandits, experts and games (lecture 8)," 2016. [Online]. Available: https://www.cs.umd.edu/~slivkins/CMSC858G-fall16/lecture8-both.pdf

[28] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, D. Parikh, and D. Batra, "Habitat: A Platform for Embodied AI Research," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.

[29] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. Chaplot, O. Maksymets, A. Gokaslan, V. Vondrus, S. Dharur, F. Meier, W. Galuba, A. Chang, Z. Kira, V. Koltun, J. Malik, M. Savva, and D. Batra, "Habitat 2.0: Training home assistants to rearrange their habitat," in *Advances in Neural Information Processing Systems*, 2021.

[30] X. Puig, E. Undersander, A. Szot, M. D. Cote, T.-Y. Yang, R. Partsey, R. Desai, A. W. Clegg, M. Hlavac, S. Y. Min, *et al.*, "Habitat 3.0: A co-habitat for humans, avatars and robots," 2023.

[31] S. Folk, J. Paulos, and V. Kumar, "Rotorpy: A python-based multirotor simulator with aerodynamics for education and research," 2023.

[32] K. Chaney, F. Cladera, Z. Wang, A. Bisulco, M. A. Hsieh, C. Korpela, V. Kumar, C. J. Taylor, and K. Daniilidis, "M3ed: Multi-robot, multi-sensor, multi-environment event dataset," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2023.

[33] I. D. Miller, F. Cladera, T. Smith, C. J. Taylor, and V. Kumar, "Stronger together: Air-ground robotic collaboration using semantics," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, 2022.

[34] C. Bai, T. Xiao, Y. Chen, H. Wang, F. Zhang, and X. Gao, "Faster-lio: Lightweight tightly coupled lidar-inertial odometry using parallel sparse incremental voxels," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, 2022.

[35] E. Marder-Eppstein. (2024) ROS move_base package.

## APPENDIX

### A. Lemmas

**Lemma 6.1:** Follow the regularized leader with regularization term of $\delta_a^2/\eta_a$ results in the update rule of LIP:

$$\delta_{a+1} = \frac{\eta_a}{\eta_{a-1}}\delta_a + \frac{\eta_a}{2}d_a \qquad (16)$$

where $\eta_a, \eta_{a-1}$ are the learning rate for previous and current updates. Gradient $d_a$ is the output of GRADIENT$(\cdot,\cdot)$.

*Proof:* In *follow the regularized leader*, $\delta_a$ equals to

$$\arg\min_\delta \left(\sum_{s=1}^a \delta(-d_s) + \frac{\delta^2}{\eta_a}\right).$$

Since the objective is convex in $\delta$, we can find the minimizer by solving $\sum_{s=1}^a(-d_s) + 2\delta/\eta_a = 0$ which gives

$$\delta_{a+1} = \eta_a\sum_{s=1}^a \frac{d_s}{2}.$$

Given $\delta_a = \eta_{a-1}\sum_{s=1}^{a-1} d_s/2$, we get

$$\delta_{a+1} = \eta_a\frac{\delta_a}{\eta_{a-1}} + \eta_a\frac{d_a}{2}.$$

∎

**Lemma 6.2:** *Follow the regularized leader* with regularization $\delta_a^2/\eta_a$ has the following regret bound:

$$\mathrm{R}(j,k) \le \sum_{i=1}^T \frac{N\eta_{i,j,k}}{2} + \frac{\beta^2}{\eta_1} \qquad (17)$$

*Proof:* Firstly, $\mathrm{R}(j,k) \le \sum_{i=1}^T(\delta_{i,j,k} - \bar{\delta})d_{i,j,k}$ by the first-order condition of convexity. We can decompose it as

$$\sum_{i=1}^T \delta_{(i+1),j,k}d_{i,j,k} - \bar{\delta}d_{i,j,k} + \delta_{i,j,k}d_{i,j,k} - \delta_{(i+1),j,k}d_{i,j,k}.$$

Based on Lemma 1.4.1 and Lemma 1.4.2 in [26], we have

$$\sum_{i=1}^{T} \delta_{(i+1),j,k} d_{i,j,k} \le \sum_{i=1}^{T} \bar{\delta} d_{i,j,k}$$

and

$$\sum_{i=1}^{T} \delta_{(i+1),j,k} d_{i,j,k} - \bar{\delta} d_{i,j,k} \le \frac{\beta^2}{\eta_1}.$$

Secondly, we have

$$\sum_{i=1}^{T} \delta_{i,j,k} d_{i,j,k} - \delta_{(i+1),j,k} d_{i,j,k}$$

$$= \sum_{i=1}^{T} (\delta_{i,j,k} - \delta_{(i+1),j,k}) d_{i,j,k} \le \sum_{i=1}^{T} |\delta_{i,j,k} - \delta_{(i+1),j,k}| |d_{i,j,k}|$$

$$\le \sum_{i=1}^{T} |\delta_{i,j,k} - \delta_{(i+1),j,k}| N = \sum_{i=1}^{T} |N \eta_{i,j,k} \frac{d_{i,j,k}}{2}|$$

$$\le \sum_{i=1}^{T} N \frac{\eta_{i,j,k}}{2}.$$

∎

**Lemma 6.3:** The filled loss $\hat{l}_{i,j,k}(a)$ is an unbiased estimator of the actual loss, i.e., $\mathbb{E}_{y_a \in y_{\text{future}}(i)}(\hat{l}_{i,j,k}(a)) = |I_a - \hat{I}_a|$.

*Proof:* $\mathbb{E}(\hat{l}_{i,j,k}(a)) = p(y_a \in y_{\text{future}}(i))|I_a - \hat{I}_a|/p(y_a \in y_{\text{future}}(i)) + 0 = |I_a - \hat{I}_a|$. ∎

**Lemma 6.4:** After the randomization, $p(y_a \in y_{\text{future}}(i)) \ge \tau/N > 0$.

*Proof:* With probability $\tau$, any $y \in Y_i$ is selected with probability $1/N$, i.e., $p(y_a \in y_{\text{future}}(i)) \ge \tau/N > 0$. ∎

**Lemma 6.5:** We have the following bound on expected regret in bandit-feedback setting:

$$\mathbb{E}(R_{\text{bandit}}) \le N\beta(N\sqrt{T} + 1)/\tau + \tau\beta T.$$

*Proof:* With probability $(1 - \tau)$, we are using the full-information algorithm with filled loss which gives us the expected bound $N\beta/\tau(N\sqrt{T} + 1)$ by Theorem 4.1 and Lemma 6.3. This bound uses $N\beta/\tau$ instead of the original $\beta$ since the filled loss is bounded by $N\beta/\tau$ according to Lemma 6.4. With probability $\tau$, we are doing uniform sampling. The regret of this part is upper bounded by $\beta T$. Therefore, the overall regret for bandit feedback setting is

$$\mathbb{E}(R_{\text{bandit}}) \le (1 - \tau)\mathbb{E}(\hat{R}) + \tau\beta T$$
$$\le (1 - \tau)N\beta/\tau(N\sqrt{T} + 1) + \tau\beta T$$
$$\le N\beta(N\sqrt{T} + 1)/\tau + \tau\beta T.$$

∎

### B. Proof of Theorems

*Proof* of Theorem 4.3: Let the regret at each round be

$$\alpha_i = \sum_{j=1}^{n} \sum_{k=1}^{b} \sum_{a=1}^{c_{j,k}(i)} |\hat{I}_a - I_a| + |I_a - \bar{I}_a|.$$

We have

$$|\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\hat{I}), \hat{I}) - v(\pi^*(\hat{I}), \bar{I}))|$$

$$= |\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\hat{I}), \hat{I} - \bar{I}))|$$

$$\le |\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\hat{I}), \alpha_i))| = |\mathbb{E}(\sum_{i=1}^{T})\alpha_i| \le \alpha.$$

The first equality follows from linearity of $v$ in the second argument. The second equality holds from the definition of $v$. Since $\hat{I} - \bar{I} = \sum_{j=1}^{n} \sum_{k=1}^{b} \sum_{a=1}^{c_{j,k}(i)} \hat{I}_a - \bar{I}_a = \sum_{j=1}^{n} \sum_{k=1}^{b} \sum_{a=1}^{c_{j,k}(i)} \hat{I}_a - I_a + I_a - \bar{I}_a \le \alpha_i$, inequalities follow from the bound on regret in (12). Similarly,

$$|\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\bar{I}), \hat{I}) - v(\pi^*(\bar{I}), \bar{I}))| \le \alpha$$

$$|\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\bar{I}), I) - v(\pi^*(\bar{I}), \bar{I}))| \le \alpha'$$

$$|\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\hat{I}), \bar{I}) - v(\pi^*(\hat{I}), I))| \le \alpha'$$

In addition, we have

$$\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\hat{I}), \hat{I})) \ge \mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\bar{I}), \hat{I}))$$

since $\pi^*(\hat{I})$ is the optimal solution for $\hat{I}$. Then

$$\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\bar{I}), I) - v(\pi^*(\hat{I}), I))$$

$$= \mathbb{E}\Big[\sum_{i=1}^{T} \Big(v(\pi^*(\bar{I}), I) - v(\pi^*(\bar{I}), \bar{I})\Big) + \Big(v(\pi^*(\bar{I}), \bar{I})$$

$$- v(\pi^*(\bar{I}), \hat{I})\Big) + \Big(v(\pi^*(\bar{I}), \hat{I}) - v(\pi^*(\hat{I}), \hat{I})\Big)$$

$$+ \Big(v(\pi^*(\hat{I}), \hat{I}) - v(\pi^*(\hat{I}), \bar{I})\Big) + \Big(v(\pi^*(\hat{I}), \bar{I})$$

$$- v(\pi^*(\hat{I}), I)\Big)\Big] \le 2(\alpha + \alpha')$$

We have shown that each pair of terms in the equality is bounded which leads to the final inequality. ∎

*Proof* of Theorem 4.4: We decompose the regret to obtain the bound:

$$\mathbb{E}(\sum_{i=1}^{T} v(\pi^*(\bar{I}), I) - v(\pi(\hat{I}), I))$$

$$= \mathbb{E}\Big[\sum_{i=1}^{T} \Big(v(\pi^*(\bar{I}), I) - v(\pi^*(\hat{I}), I)\Big) + \Big(v(\pi^*(\hat{I}), I)$$

$$- v(\pi^*(\hat{I}), \hat{I})\Big) + \Big(v(\pi^*(\hat{I}), \hat{I}) - v(\pi(\hat{I}), \hat{I})\Big)$$

$$+ \Big(v(\pi(\hat{I}), \hat{I}) - v(\pi(\hat{I}), I)\Big)\Big]$$

$$\le 2(\alpha + \alpha') + 2(\alpha + \alpha') + (1 - \gamma)v(\pi^*(\hat{I}), \hat{I})$$

$$= 4(\alpha + \alpha') + (1 - \gamma)v(\pi^*(\hat{I}), \hat{I}).$$

The first term in the inequality follows from Theorem 4.3. The second term in the inequality follows from the inequality in the proof of Theorem 4.3. The third term in the inequality follows from the bound for $v$. ∎

## C. Additional Visualizations of Simulation Experiments
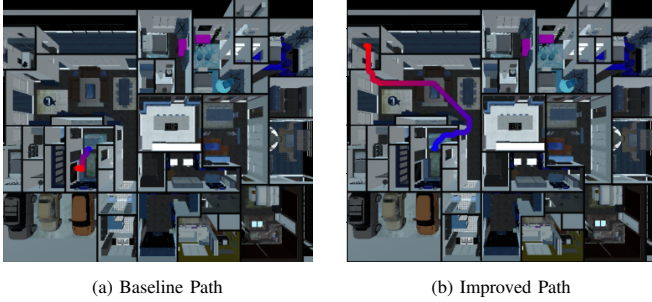


(a) Baseline Path

(b) Improved Path

**Fig. 4: Comparison of selected goals by baseline and improved methods in Scene 102816036.** The robot samples twenty candidate goals. If the information gain prediction is not improved, the robot chooses a sub-optimal goal as shown in **(a)**. After improving information gain prediction by the proposed method, the robot chooses a goal in an unexplored room as shown in **(b)**.
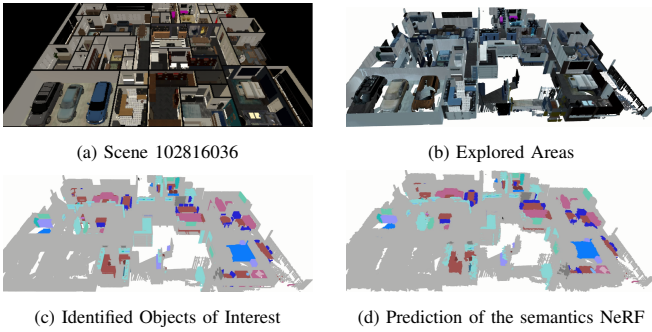


(a) Scene 102816036

(b) Explored Areas

(c) Identified Objects of Interest

(d) Prediction of the semantics NeRF

**Fig. 5: Map of Scene 102816036 after autonomous exploration. (a)** shows Scene 102816036. **(b)** shows the explored areas by the proposed algorithm after 25 episodes. **(c)** is the ground-truth point cloud of explored areas. Different colors correspond to ground-truth object categories. **(d)** is the point cloud extracted from the semantic NeRF after 25 episodes. Different colors correspond to different object categories predicted by the NeRF.
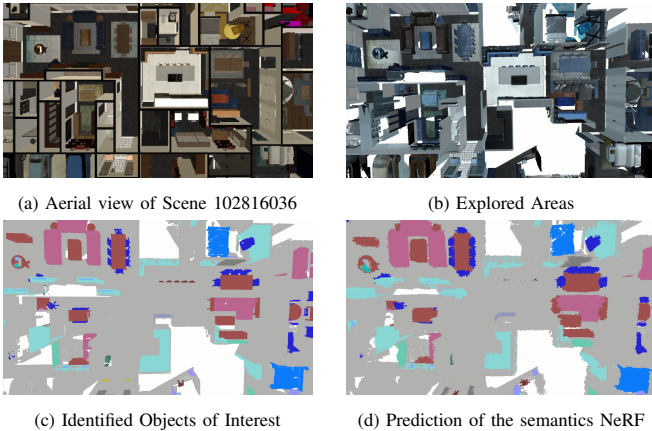


(a) Aerial view of Scene 102816036

(b) Explored Areas

(c) Identified Objects of Interest

(d) Prediction of the semantics NeRF

**Fig. 6: Aerial map of Scene 102816036 after autonomous exploration. (a)** shows aerial view of Scene 102816036. **(b)** shows the explored areas by the proposed algorithm after 25 episodes. **(c)** is the ground-truth point cloud of explored areas. Different colors correspond to ground-truth object categories. **(d)** is the point cloud extracted from the semantic NeRF after 25 episodes. Different colors correspond to different object categories predicted by the NeRF.
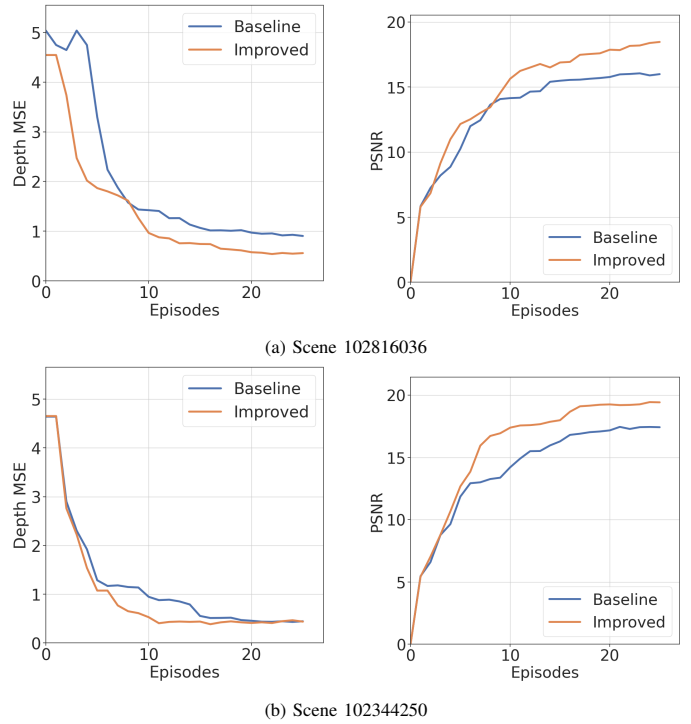


(a) Scene 102816036



(b) Scene 102344250

**Fig. 7: Evaluation of trained NeRFs.** For each scene, NeRFs are trained online with the supervision of collected measurements. For each room in each scene, we choose four test viewpoints with yaw angles $0, \pi/2, \pi, 3\pi/2$ to cover the room. We evaluate the trained NeRFs based on the test viewpoints after each episode. **(a)** and **(b)** are the plots of evaluation for Scene 102816036 and Scene 102344250 respectively.