# Cost-Effective Algorithms for Average-Case Interactive Graph Search

Qianhao Cong
Dept. of Ind. Syst. Engg. & Mgmt.
National University of Singapore
cong_qianhao@u.nus.edu

Jing Tang*
Data Science and Analytics Thrust
The Hong Kong Uni. of Sci. and Tech.
jingtang@ust.hk

Yuming Huang
Dept. of Ind. Syst. Engg. & Mgmt.
National University of Singapore
huangyuming@u.nus.edu

Lei Chen
Data Science and Analytics Thrust
The Hong Kong Uni. of Sci. and Tech.
leichen@ust.hk

Yeow Meng Chee
Dept. of Ind. Syst. Engg. & Mgmt.
National University of Singapore
ymchee@nus.edu.sg

*Abstract*—**Interactive graph search (IGS) uses human intelligence to locate the target node in hierarchy, which can be applied for image classification, product categorization and searching a database. Specifically, IGS aims to categorize an object from a given category hierarchy via several rounds of interactive queries. In each round of query, the search algorithm picks a category and receives a boolean answer on whether the object is under the chosen category. The main efficiency goal asks for the minimum number of queries to identify the correct hierarchical category for the object. In this paper, we study the *average-case interactive graph search (AIGS)* problem that aims to minimize the expected number of queries when the objects follow a probability distribution. We propose a greedy search policy that splits the candidate categories as evenly as possible with respect to the probability weights, which offers an approximation guarantee of $O(\log n)$ for AIGS given the category hierarchy is a directed acyclic graph (DAG), where $n$ is the total number of categories. Meanwhile, if the input hierarchy is a tree, we show that a constant approximation factor of $(1+\sqrt{5})/2$ can be achieved. Furthermore, we present efficient implementations of the greedy policy, namely GreedyTree and GreedyDAG, that can quickly categorize the object in practice. Extensive experiments in real-world scenarios are carried out to demonstrate the superiority of our proposed methods.**

## I. Introduction

Crowdsourcing services, such as Amazon's Mechanical Turk and CrowdFlower, allow users to generate tasks for crowd workers to complete in exchange for rewards. The main goal of crowdsourcing is to solve problems that are hard for computers. During the past decades, crowdsourcing services have emerged as a powerful mechanism to process data for many tasks, such as object categorization [31, 39], data-labeling and dataset creation [14, 32], entity resolution [48, 49], data filtering [40, 41], and SQL-like query processing [13, 19, 23, 29, 47].

Recently, Tao et al. [46] proposed the interactive graph search (IGS) problem that aims to find an initially unknown target node on a direct acyclic graph (DAG) by asking questions like "is the target node reachable from a given node $u$?". In the beginning, all nodes on the DAG are considered as

A short version of the paper will appear in the 38th IEEE International Conference on Data Engineering (ICDE '22), May 9–12, 2022.
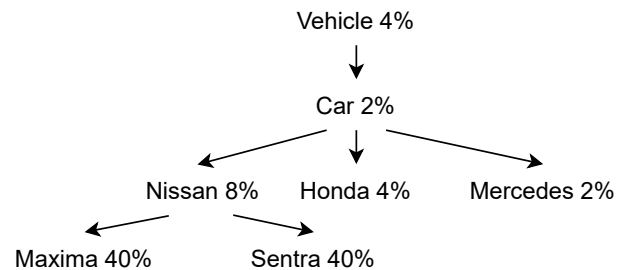*Corresponding author: Jing Tang.



Fig. 1. Image categorization [39].

candidates. By sequentially asking questions, we can gradually narrow down the candidates and finally identify the target node. In real-world scenarios, these questions may be answered through a crowdsourcing platform. Meanwhile, crowdsourcing platforms usually charge a flat fee for each question, i.e., a unit cost per question. Therefore, a natural aim is to reduce the number of queries needed for finding the target node. There are many real-world applications [39] that can be characterized by the IGS problem, such as image classification [46], product categorization [31], debugging of workflows [10, 21], etc. As an immediate illustration, the following example shows the task of labeling an image by crowdsourcing [39].

**Example 1** (Image Categorization). *In (supervised) machine learning, algorithms usually require labeled data to train the models. Image category identification via crowdsourcing is an important channel to obtain such labeled data. In particular, humans are hired to label images according to a given hierarchy by answering a sequence of questions. Consider the scenario that we attempt to label a vehicle image given a hierarchy shown in Fig. 1, where the proportions of different vehicle categories are different. Note that the target can be either a leaf node or a non-leaf node in the hierarchy. For example, with a probability of 40%, the image shows a Sentra, while with another probability of 8%, the image is a Nissan but neither a Maxima nor a Sentra.*

*We may first ask a question like "is this a car?". If the answer is no, we can immediately label the image as a non-car*

*vehicle. If the answer is yes, we continue the search process in the tree rooted at the node of car. Asking sequential questions in a similar way, we can finally find the most suitable label to describe the image. Suppose that we receive a yes-answer to both questions of "is this a car?" and "is this a Honda?". Then, a label of Honda is placed on the image.* □

A straightforward solution is a TopDown method that asks questions from the root of the hierarchy. Specifically, TopDown first queries every child of the root until it gets a *yes*-answer. If no such a child exists, the root will be returned as the target node. Otherwise, such a child will be set as the new root to repeat the process. Take Fig. 1 as an illustration and assume that *Sentra* is the target node. The TopDown algorithm first asks "is this a *car*?" and receives a *yes*-answer. Then, it asks "is this a *Nissan*?" and again gets a *yes*-answer. It further asks "is this a *Maxima*?" and this time a *no*-answer is returned. Finally, it asks "is this a *Sentra*" with a *yes*-answer so that the target node is identified successfully. This strategy can correctly find the target node but unfortunately is cost-ineffective. Since we know that most of the images are *Maxima* or *Sentra*, a smarter strategy can query these two nodes first so that with a probability of 80%, asking at most two questions can identify the target node. This observation motivates us to develop cost-effective algorithms to minimize the cost.

Tao et al. [46] studied the *worst-case* interactive graph search (WIGS) that aims to minimize the *maximum* number of queries required to identify every possible target node and proposed a series of heavy-path-based binary searches with near-optimal theoretical guarantees. However, for the task of image categorization by crowdsourcing, the data owner hires humans to label the images with a desire of minimizing the total cost of labeling all images, rather than minimizing the maximum cost of labeling one image.

**Example 2** (Query Cost). *Consider that there are* 100 *images with proportions as given in Fig. 1, e.g.,* 4 *images are non-car vehicles and* 40 *images are Maxima. An optimal solution to WIGS sequentially queries Nissan, Car, Honda, and Mercedes with answers of no, yes, no, and yes returned, assuming that the image is a Mercedes. In fact, this is the worst case for such a solution requiring* 4 *queries. Similarly, the numbers of queries asked are* 2*,* 4*,* 3*,* 3*,* 2 *and* 3 *if the correct categories are Vehicle, Car, Honda, Nissan, Maxima and Sentra, respectively. As a result, the total cost is* 260*.*

*On the other hand, consider an alternative solution that sequentially queries Maxima, Sentra, Nissan, Car, Honda, and Mercedes with a cost of* 6*, if the image is again a Mercedes. Meanwhile, if the correct categories are Vehicle, Car, Honda, Nissan, Maxima and Sentra, the costs become* 4*,* 6*,* 5*,* 3*,* 1 *and* 2*, respectively. Hence, the total cost is* 204*. Apparently, in terms of the worst-case cost, this solution requiring* 6 *queries is inferior to the optimal solution to WIGS requiring* 4 *queries only, whereas in terms of the average-case cost, the former (with a cost of* 2.04 *in average) remarkably outperforms the latter (with a cost of* 2.6 *in average).* □

Motivated by the deficiency of existing solutions, in this paper, we focus on the *average-case* interactive graph search (AIGS) problem that attempts to minimize the *expected* number of queries when the target node follows a-priori known probability distribution. It is trivial to see that minimizing the total cost of categorizing a batch of objects is factually equivalent to AIGS, assuming that we know the distribution of the target nodes in this batch of tasks. Note that in some real applications, it may be impossible to get the true data distribution. For such scenarios, we may apply a simple online learning approach that dynamically adjusts the empirical probability distribution on the fly upon obtaining the category result of each object. Usually, the empirical statistics can well characterize the real data distribution as long as a sufficiently large number of labels are obtained.

We find that obtaining the optimal solution for AIGS is computationally intractable. To address AIGS, we propose a greedy search policy that picks a query node to split the candidate nodes as evenly as possible, taking into account the probability weight of each node. With a rigorous and thorough analysis, we show that the greedy policy can achieve strong theoretical guarantees, especially a constant factor given the input hierarchy is a tree. In addition, a naive implementation of the greedy policy is to evaluate the splitting size of every candidate node in each round, which is time-consuming. To tackle this issue, we develop efficient algorithms for variant scenarios, namely GreedyTree and GreedyDAG, to accelerate the query node selection according to the greedy policy.

In summary, we make the following contributions.

- We show that computing the optimal policy for AIGS is NP-hard, even if the hierarchy has a tree structure.
- We propose a cost-effective approach for AIGS. We show that our greedy policy offers an approximation ratio of $O(\log n)$ for any input DAG with $n$ nodes and a constant factor of $(1 + \sqrt{5})/2$ when the input hierarchy is a tree (Section III).
- We devise efficient implementations of the greedy policy, incorporating several acceleration techniques. In particular, our GreedyTree and GreedyDAG algorithms run in $O(nhd)$ and $O(nm)$ time for tree and DAG hierarchies respectively, improving the naive $O(n^2m)$ time algorithm significantly, where $n$ and $m$ are the number of nodes and edges, $h$ is the length of longest path, and $d$ is the maximum out-degree of nodes in the hierarchy. (Section IV).
- We conduct extensive experiments with real-world datasets to evaluate our proposed methods, and the experimental results strongly corroborate the effectiveness and efficiency of our approach (Section V).

## II. PROBLEM DEFINITION

In this paper, we study the problem of *average-case interactive graph search (AIGS)*. That is, we sequentially ask the crowd some simple reachability questions with boolean answers, i.e., *yes* or *no*, to gradually narrow down the potential categories of the object. Meanwhile, the cost of crowdsourcing is usually

TABLE I
FREQUENTLY USED NOTATIONS.

| Notation | Description |
|----------|-------------|
| $G = (V, E)$ | a DAG $G$ with node set $V$ and edge set $E$ |
| $n$ | the number of nodes in $G$, i.e., $n = |V|$ |
| $m$ | the number of edges in $G$, i.e., $m = |E|$ |
| $z$ | target node, e.g., the category of an unlabeled object |
| $reach(u)$ | query on node $u$, i.e., whether $z$ is reachable from $u$ in $G$ |
| $G_u$ | the subgraph of $G$ rooted at node $u$ |

---

**Algorithm 1:** FrameworkIGS($G$)

**Input:** an input DAG $G$
**Output:** the target node $z$ in $G$
1 **while** $G$ *has at least two nodes* **do**
2     select a query node $u$ in $G$;
3     **if** $reach(u) = yes$ **then**
4        $G \leftarrow G_u$; // $G_u$ is the subgraph of $G$ rooted at $u$
5     **else**
6        $G \leftarrow G \setminus G_u$;

7 **return** the only node in $G$;

---

determined by the number of questions asked. Therefore, our aim is to correctly identify the hierarchical category by asking the minimum *expected* number of questions, i.e., minimizing the *expected* cost. For ease of reference, Table I summarizes the notations that are frequently used.

Specifically, AIGS is formally defined as follows. We abstract a category hierarchy as a directed acyclic graph (DAG) $G = (V, E)$ with a set $V$ of $n$ nodes and a set $E$ of $m$ edges. We assume that there is only one *root* in $G$. If there are multiple roots, we can simply add a dummy node to $G$ with an outgoing edge to every original root, which generates a new DAG with one root only. Given an unknown *target node* (e.g., an unlabeled object) $z \in V$, there is an *oracle* that can answer questions. For any *query node* $q \in V$, the oracle returns a boolean answer, denoted as $reach(q)$, as follows,

$$reach(q) = \begin{cases} yes, & \text{if there is a directed path from } q \text{ to } z, \\ no, & \text{otherwise.} \end{cases}$$

It is easy to see that $reach(q) = yes$ if and only if $q$ can *reach* $z$, which characterizes the *reachability* from $q$ to $z$. For IGS, the algorithm interactively picks a query node $u$ and receives a boolean answer $reach(u)$. The algorithm repeats the process until it determines the target node $z$.

In particular, in each round of interaction, if it gets a *yes*-answer, the target node must be reachable from the query node $u$. That is, if $reach(u) = yes$, it clearly holds that the target node $z \in G_u$ and the search graph is updated to $G_u$, where $G_u$ is the subgraph of $G$ rooted at the query node $u$. Otherwise, if $reach(u) = no$, $G$ is updated by $G \setminus G_u$ as $z \in G \setminus G_u$. The algorithm stops to return the target node when the search graph just has one node. We present the main procedure of IGS in Algorithm 1, referred to as FrameworkIGS. As can be seen, the critical task of the algorithm design for IGS lies in the choice of query node (Line 2).

We define the *cost* as the number of questions the algorithm asks, assuming that each question charges a fixed price (e.g., $1 per question). In fact, as shall be discussed in Section III-D, we can extend to a general scenario where the payment for each question is distinct to reveal the difficulties of questions, e.g., hard questions are more expensive than easy questions. A natural goal of IGS is to identify the target node with the minimum cost. In this paper, we consider that each node $v$ is associated with a probability $p(v)$ measuring the likelihood of $v$ being the target node, e.g., a-priori known data distribution. Given a set $S$ of nodes, let $p(S) = \sum_{v \in S} p(v)$ denote the probability of the target node being one node in

$S$, e.g., $p(V) = 1$. Then, we define the *expected cost* of identifying a target node as the expected number of questions asked by the algorithm, taking into account the randomness of the target node. Therefore, we aim to study the interactive question-asking strategies (also known as policies) for the *average-case interactive graph search (AIGS)* problem with the goal of minimizing the expected cost. For convenience, we abuse notation and use the terms "policy" and "algorithm" interchangeably.

**Definition 1** (Average-Case Interactive Graph Search)**.** *Given a DAG hierarchy $G = (V, E)$ and an (unknown) target node $z \in V$ following the a-priori known probability distribution $p(\cdot)$, AIGS asks for a query policy to identify the target node such that the expected cost is minimized.*

**Remark.** In the AIGS problem, we assume a-priori known probability distribution. In practice, the probability distribution can be empirically learned from the historical data, e.g., a simple statistical result of the categorized objects. Moreover, when there is a batch of objects (e.g., a batch of unlabeled images) required to be categorized, these objects generally follow the probability distribution. As a result, minimizing the total cost of categorizing a batch of objects, which is the average cost of each object multiplying the number of objects, is equivalent to the optimization problem of AIGS.

### III. GREEDY POLICY FOR AIGS

In this section, we show that finding the optimal solution for AIGS is computationally intractable and propose a natural greedy policy with provable approximation guarantees.

#### A. Impossibility Results

**Lemma 1.** *The AIGS problem is NP-hard. In fact, there cannot be any $o(\log n)$-approximate algorithm for AIGS unless* NP $\subseteq$ DTIME $\left(n^{O(\log \log n)}\right)$.

To prove Lemma 1, we begin by showing that the IGS problem is equivalent to the search in *partially ordered set (poset)* problem [12]. We first introduce the definition of poset.

**Definition 2** (Partially Ordered Set)**.** *A partially ordered set is a pair $(V, \leq_R)$, where $V$ is a set of objects and $\leq_R$ is a partially ordered relation which satisfies the following properties:*

- *Reflexivity: $a \leq_R a$.*

- *Antisymmetry: if $a \leq_R b$ and $b \leq_R a$, then $a = b$.*
- *Transitivity: if $a \leq_R b$ and $b \leq_R c$, then $a \leq_R c$.*

Reflexivity shows that every object is related to itself. Antisymmetry indicates that for any two distinct objects, they cannot relate to each other. Transitivity reflects that the relation is transitive.

**Definition 3** (Search in Poset)**.** *Given an (unknown) target object $z$ in a poset $(V, \leq_R)$, the search problem aims to locate the target objects by queries. To search the target object, each query on object $x \in V$ will return a boolean result:*

- *yes, if the target object is related to $x$, i.e., $z \leq_R x$;*
- *no, otherwise.*

We make a connection between the IGS problem and the search problem in poset.

**Lemma 2.** *The IGS problem is equivalent to the search problem in a poset.*

*Proof.* The reachability of IGS satisfies the following properties.

- Reflexivity: Each node can reach itself.
- Antisymmetry: For two nodes $u, v$, if $u$ can reach $v$ and $v$ can reach $u$, then $u$ should equal to $v$.
- Transitivity: For three nodes $u, v, w$, if $u$ can reach $v$ and $v$ can reach $w$, then $u$ can reach $w$ by taking $v$ as a relay node.

Hence, the reachability relation in the IGS problem is a partially ordered relation. Moreover, given a target node $z \in V$, the query on node $q \in V$ will get a *yes* answer if and only if $z$ is reachable from $q$, which satisfies query in a poset. Therefore, the IGS problem is a search problem in a poset.

On the other hand, given a poset, we can construct a DAG hierarchy as follows. For any two distinct objects $a$ and $b$ in the poset, if $a \leq_R b$ and there does not exist any $c \in V \setminus \{a, b\}$ such that $a \leq_R c$ and $c \leq_R b$, i.e., $a$ is directly related to $b$, then $a$ is the child of $b$ in the hierarchy. It is easy to verify that the constructed hierarchy is a DAG and for two objects $a$ and $b$ satisfying $a \leq_R b$, it must hold that $a$ is reachable from $b$ in the hierarchy. This implies that searching a poset can be represented by an IGS problem. $\square$

Now we are ready to prove Lemma 1.

*Proof of Lemma 1.* Cicalese et al. [8] proved that deriving the optimal solution for the average-case poset search is NP-hard even if the poset has a tree structure (i.e., the input hierarchy of AIGS is a tree). Moreover, Cicalese et al. [8] pointed out that for the problem of minimizing the weighted average number of queries to identify an initially unknown object for general posets, there cannot be any $o(\log n)$-approximate algorithm unless $\text{NP} \subseteq \text{DTIME}\left(n^{O(\log \log n)}\right)$, where $n$ is the number of nodes. Combining with Lemma 2, we complete the proof. $\square$

Lemma 1 implies that in general, it is impossible to construct the optimal solution for AIGS in polynomial time unless $\text{P} = \text{NP}$. Even worse, devising a good approximate algorithm for AIGS within a factor of $o(\log n)$ is also impossible. In the following, we propose a greedy algorithm that can provide provable theoretical guarantees.

### B. Greedy Policy

Recall that a key step of AIGS is to select the query node $u$ which will split the candidate DAG $G$ into two DAGs $G_u$ and $G \setminus G_u$, where $G_u$ is the subgraph of $G$ rooted at $u$ with all its descendants. Then, the target node $z$ is either in $G_u$ or in $G \setminus G_u$. Hence, to minimize the expected cost, intuitively, a "good" strategy is to filter nodes with total weight as large as possible with respect to the probability distribution. That is, the query splits the candidate set of nodes into two parts with the most equal probability weights. This strategy can avoid asking a large number of questions to identify the nodes with high probability. In this paper, we utilize such a natural greedy policy to address the AIGS problem. Specifically, the greedy policy picks the node $u$, referred to as the *middle point*, such that the difference between $p(G \setminus G_u)$ and $p(G_u)$ is minimized.

**Definition 4** (Middle Point)**.** *Given a DAG $G$ associated with a weight $p(v)$ for each node $v$, the middle point $u^*$ of $G$ is defined as*

$$u^* := \underset{u \in V}{\arg\min} \, |p(G_u) - p(G \setminus G_u)| = \underset{u \in V}{\arg\min} \, |2p(G_u) - p(G)|.$$

Note that there may exist multiple middle points in a DAG. If such a case happens, the greedy policy just arbitrarily chooses one of them.

### C. Theoretical Guarantees

In what follows, we show that the greedy policy can achieve an approximation guarantee of $O(\log n)$ for AIGS, which matches the best achievable ratio. A key step of our analysis is to map AIGS to the binary decision tree problem [5, 18].

There has been extensive research [5, 25] showing that the minimum (probability) weight of nodes will have a significant impact on the approximation ratio for the decision tree problem. Fortunately, a rounding technique [5] can be used tackle the negative impact of the minimum weight. In particular, for each node $u$, its weight will be rounded from $p(u)$ to $w(u)$ as follows,

$$w(u) = \left\lceil \frac{n^2 \cdot p(u)}{\max_{v \in V} p(v)} \right\rceil. \tag{1}$$

Then, the greedy policy is used to construct the decision tree based on the rounded weights $w(\cdot)$, namely rounded greedy policy. Chakaravarthy et al. [5] proved that the rounded greedy policy achieves an approximation ratio of $2(1 + 3 \ln n)$.[1] This approximation result also applies to AIGS.

**Theorem 1.** *The rounded greedy algorithm provides an approximation ratio of $2(1 + 3 \ln n)$ for AIGS.*

---

[1]The original analysis [5] requires the probabilities to be rational numbers, whereas we assume the weights to be real numbers. However, this mismatch does not matter a lot, since the machine uses a limited number of bits to record a real number with sufficient precision.
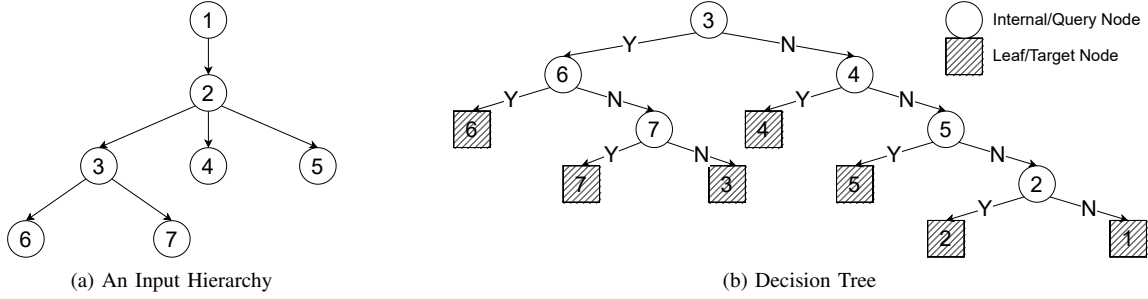
(a) An Input Hierarchy        (b) Decision Tree

Fig. 2. Decision tree for the query policy of interactive graph search.

To prove Theorem 1, we make the connection between our AIGS problem and the decision tree problem [5, 18]. In particular, a target node of AIGS corresponds to a leaf node of decision tree while a query node of AIGS corresponds to an internal node of decision tree (i.e., whether the target node is reachable from the query node). Hence, we revisit the AIGS problem from the view of decision tree. In particular, as the answer is a boolean, i.e., *yes* or *no*, a binary decision tree is considered in practice, which is formally defined as follows.

**Definition 5** (Binary Decision Tree [5]). *The binary decision tree problem take as input a table having $N$ rows and $M$ columns. Each row is called an object and columns are binary attributes of these objects. The aim of this problem is to output a binary tree where each leaf is associated with an object and each internal node is associated with a test for one attribute. If an object passes a test, it goes to the left branch; otherwise, it goes to the right branch. Each object can be uniquely identified by a path from the root to its representing leaf. The goal of a decision tree problem is to find such a tree that minimizes the weighted sum of the depths of all the leaves, considering that each object (i.e., leaf) is associated with a weight.*

**Lemma 3.** *The AIGS problem can be reduced to the binary decision problem.*

*Proof.* To reduce the AIGS problem to the binary decision tree problem, we can consider each node in the input hierarchy as an object and the reachability relations as the attributes. Specifically, we can transfer the input hierarchy with $n$ nodes to an $n \times n$ table where the attribute of $i$-th row and $j$-th column equals to 1 if node $i$ is reachable from node $j$, and equals to 0 otherwise. Then, a query on node $j$ in AIGS is exactly a test for one attribute $j$ in the decision tree. As a result, minimizing the expected cost for AIGS is equivalent to minimizing the weighted sum of the depths of all the leaves for the binary decision tree problem. This completes the proof. □

Based on this observation, we construct a binary decision tree according to the search strategy of AIGS and calculate the cost based on the constructed decision tree.

**Definition 6** (Decision Tree for AIGS). *We model the query policy of AIGS as a decision tree, where each internal node represents a question (i.e., the query node) and each leaf node represents a search result (i.e., the target node). Starting from the root of a decision tree, the algorithm asks questions following the below rules.*

- *If the answer is yes, i.e., the target node is reachable from the query node, it goes to the left child.*
- *Otherwise, if the answer is no, it goes to the right child.*

*The above process stops at a leaf node, which will be returned as the search result.*

According to the above definition, it is trivial to see that given any DAG $G$, the size of its decision tree $D$ is at most twice the size of $G$, since every node of $G$ appears as a leaf node of $D$ and meanwhile at most all nodes of $G$ serve as internal nodes of $D$. Furthermore, it is easy to observe that the number of questions asked for identifying a target node is equal to the depth of the corresponding leaf node in $D$, i.e., the length of its path from the root in $D$. Meanwhile, the weight of each leaf node $v$ is equal to the probability $p(v)$ of the node $v$ being the target node. Hence, we formally define the cost of a query policy leveraging the notion of decision tree.

**Definition 7** (Cost of Query Policy). *Given a query policy for the AIGS problem, let $D$ be the decision tree constructed according to the policy. Denote by $L(D)$ the set of all leaf nodes in $D$, and by $\ell(v)$ the depth of every leaf node $v \in L(D)$. Then, the cost of $D$, i.e., the expected cost of the query policy, is given by*

$$cost(D) = \sum_{v \in L(D)} p(v)\ell(v). \tag{2}$$

**Example 3.** *Fig. 2(a) gives the graph representation of the hierarchy given in Fig. 1. Suppose that the nodes have equal weights, i.e., $p(v) = 1/7$ for every $v \in V$. Fig. 2(b) shows the decision tree of a (greedy) policy. Assume that node 5 is the target node. Then, the algorithm will sequentially query on nodes 3, 4 and 5 with answers of no, no and yes, respectively, so that the target node is figured out with no ambiguity. Moreover, the expected cost of this policy is $\frac{1}{7} \times (2 \times 2 + 3 \times 3 + 2 \times 4) = 3$.* □

Now, we are ready to prove Theorem 1.

*Proof of Theorem 1.* Our above analysis shows that AIGS can be transferred to a special case of average-case binary decision tree. In addition, Chakaravarthy et al. [5] proved that the rounded greedy policy achieves an approximation ratio of $2(1 + 3 \ln n)$ for average-case binary decision tree. Putting it together completes the proof. □

Next, we consider two special cases of AIGS: (i) the input hierarchy is a tree, and (ii) the probability of every node being the target is identical, i.e., $p(v) = 1/n$ for every $v \in V$. We show that the approximation ratios achieved by the greedy policy under the two scenarios are much better than $O(\log n)$ leveraging the theoretical developments of poset [9] and decision tree [30].

**Theorem 2.** *If the input hierarchy is a tree, the greedy policy provides an approximation ratio of $\frac{1+\sqrt{5}}{2}$ for AIGS.*

*Proof.* Let $\widetilde{D}$ be the decision tree constructed by greedy policy and $D^*$ be the optimal decision tree, with respect to the cost function defined in (2). For such a tree-like poset, Cicalese et al. [9] proved that the cost of $\widetilde{D}$ satisfies

$$cost(\widetilde{D}) \leq \frac{1 + \sqrt{5}}{2} \cdot cost(D^*). \tag{3}$$

Furthermore, according to Lemma 2, we know that $cost(\widetilde{D})$ is exactly the expected number of questions asked by the greedy policy for AIGS. Putting it together completes the proof. $\square$

**Theorem 3.** *If every node has an equal probability to be the target node, i.e., $p(v) = \frac{1}{n}$ for every $v \in V$, the greedy policy provides an approximation ratio of $O\left(\frac{\log n}{\log \log n}\right)$ for AIGS.*

*Proof.* Li et al. [30] showed that for any instance of the decision tree problem on $n$ objects with equal probability to be the target, the greedy decision tree $\widetilde{D}$ and the optimal decision tree $D^*$ satisfy

$$cost(\widetilde{D}) \leq \frac{6 \log n}{\log cost(D^*)} \cdot cost(D^*).$$

Moreover, for any binary tree with $n$ leaf nodes, it is trivial to verify that the total depth of all leaf nodes is minimized on the complete binary tree, which is at least $n \log_2 n$. This implies the optimal solution has a cost of at least $\log_2 n$. As a result, the greedy policy returns an $O\left(\frac{\log n}{\log \log n}\right)$-approximate solution for AIGS when every node has an equal probability to be the target node. Hence, the theorem is proved. $\square$

### D. Extension to Heterogeneous Cost

In the above discussion, we consider that each query charges a fixed price. In some scenarios, different questions may have different difficulties, which asks for heterogeneous payments for different questions, e.g., \$0.5 for an easy question and \$1.5 for a hard question. We show that the greedy policy with a slight modification on the definition of *middle point*, taking the cost weight of each query into consideration, can address the problem with strong theoretical guarantees.

Consider the cost-sensitive AIGS (CAIGS) problem that the query on node $v$ charges a price of $c(v)$. It is easy to observe that the cost for identifying a target node $z$ is equal to the total cost of the query nodes on the path from the root to the leaf node $z$ in $D$. Similar to Definition 7, the cost of a query policy for CAIGS can be defined as follows.

**Definition 8** (Cost of Query Policy for CAIGS). *Given a query policy for the CAIGS problem, let $D$ be the decision*



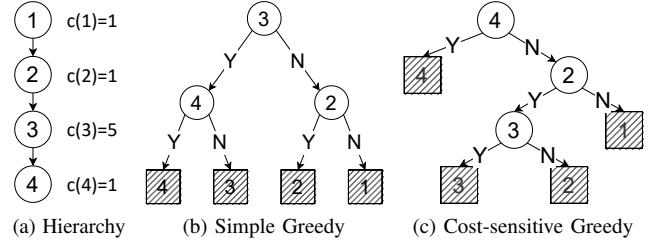(a) Hierarchy     (b) Simple Greedy     (c) Cost-sensitive Greedy

Fig. 3. Cost-sensitive greedy for CAIGS.

*tree constructed according to the policy. Denote by $L(D)$ the set of all leaf nodes in $D$, and by $\hat{\ell}(v)$ the total cost of the query nodes on the path from the root to the leaf node $v$ in $D$. Then, the cost of $D$, i.e., the expected cost of the query policy, is given by*

$$cost(D) = \sum_{v \in L(D)} p(v)\hat{\ell}(v). \tag{4}$$

To cope with the heterogeneous cost of each query node, we attempt to select the query node such that (i) it splits the category graph as evenly as possible with respect to the probability weights, and meanwhile (ii) its query cost is as cheap as possible. For goal (i), we try to find a node $u$ that maximizes $p(G_u)p(G \setminus G_u)$. That is, since $p(G_u) + p(G \setminus G_u) = p(G)$ is a constant, maximizing $p(G_u)p(G \setminus G_u)$ is to minimize the difference between $p(G_u)$ and $p(G \setminus G_u)$, i.e., $\min|p(G_u) - p(G \setminus G_u)|$. For goal (ii), we directly minimize $c(u)$. Finally, to optimize both goals simultaneously, we select the node $u$ that maximizes $\frac{p(G_u)p(G \setminus G_u)}{c(u)}$. We define such a node as the cost-sensitive middle point.

**Definition 9** (Cost-Sensitive Middle Point). *Given a DAG $G$ associated with a probability weight $p(v)$ for $v$ being the target node and a query cost $c(v)$ on node $v$, the cost-sensitive middle point $u^*$ of $G$ is defined as*

$$u^* := \arg\max_{u \in V} \frac{p(G_u)p(G \setminus G_u)}{c(u)}.$$

Note that when each query charges a unit price, the node $u^*$ maximizing $p(G_u)p(G \setminus G_u)$ also minimizing $|p(G_u) - p(G \setminus G_u)|$, which generalizes the middle point defined in Definition 4 for the case of nodes with homogeneous costs. In the following, we show that the cost-sensitive rounded greedy algorithm, which picks the cost-sensitive middle point as the query node in each round with respect to the rounded probability weights by equation (1), achieves an approximation ratio of $2(1 + 3 \ln n)$ for CAIGS.

**Theorem 4.** *The cost-sensitive rounded greedy algorithm provides an approximation ratio of $2(1 + 3 \ln n)$ for CAIGS.*

*Proof (Sketch).* Adler and Heeringa [1] showed that when the probability weight $p(v)$ is identical for every $v$, i.e., $p(v) = 1/n$, the cost-sensitive greedy policy considering heterogeneous cost for each test can obtain the same approximation ratio as homogeneous cost on the binary decision tree problem. Combing with the rounding technique [5] tailored to the

**Algorithm 2:** GreedyNaive($G$)

**Input:** an input DAG $G$
**Output:** the target node

1   $sum\_prob \leftarrow 1$;
2   **while** $|G| > 1$ **do**
3     $min\_prob \leftarrow +\infty$;
4     **foreach** *node* $v \in G$ **do**
5       $reach\_prob \leftarrow$ GetReachableSetWeight($G, v$);
6       **if** $|2 \cdot reach\_prob - sum\_prob| < min\_prob$ **then**
7         $q \leftarrow v$;
8         $q\_prob \leftarrow reach\_prob$;
9         $min\_prob \leftarrow |2 \cdot reach\_prob - sum\_prob|$;
10    **if** $reach(q) = yes$ **then**
11      $G \leftarrow G_q$;
12      $sum\_prob \leftarrow q\_prob$;
13    **else**
14      $G \leftarrow G \setminus G_q$;
15      $sum\_prob \leftarrow sum\_prob - q\_prob$;
16   **return** *the only node in $G$;*

---

**Algorithm 3:** GetReachableSetWeight($G, v$)

**Input:** a DAG $G$ and a node $v$
**Output:** the total probability of all $v$'s reachable nodes in $G$

1   $prob \leftarrow 0$;
2   initialize a queue $S$ and insert the node $v$ into $S$;
3   mark every node in $G$ as *unvisited*;
4   **while** $S$ *is not empty* **do**
5     push out a node from $S$ as $u$;
6     **foreach** *child* $v$ *of* $u$ **do**
7       **if** $v$ *is not visited* **then**
8         $prob \leftarrow prob + p(v)$;
9         insert $v$ into $S$;
10        mark $v$ as *visited*;
11   **return** *prob;*

---

weighted probability $p(v)$ yields that the cost-sensitive rounded greedy algorithm has an approximation ratio of $2(1 + 3\ln n)$ for the CAIGS problem. $\qquad \square$

The following example demonstrates the effectiveness of the cost-sensitive rounded greedy algorithm.

**Example 4.** *As an example, consider a simple category hierarchy given in Fig. 3(a). Without considering the query cost of each node (i.e., homogeneous costs), as shown in Fig. 3(b), the greedy strategy first selects node 3; then chooses node 4 if a yes-answer is returned to further distinguish nodes 3 and 4, and chooses node 2 otherwise to distinguish nodes 1 and 2. Suppose that $c(1) = c(2) = c(4) = 1$ and $c(3) = 5$. Then, the expected cost of such a greedy strategy is $5 + 1 \times 0.5 + 1 \times 0.5 = 6$. On the other hand, as shown in Fig. 3(c), the cost-sensitive greedy algorithm will first select node 4 rather than node 3, since $\frac{p(G_4)p(G \setminus G_4)}{c(4)} = \frac{0.25 \times 0.75}{1} = 0.1875$ while $\frac{p(G_3)p(G \setminus G_3)}{c(3)} = \frac{0.5 \times 0.5}{5} = 0.05$. If a no-answer is returned (no further action is required if a yes-answer is returned), it selects node 2. Only if a yes-answer is further returned, it chooses node 3. Thus, the expected cost of the cost-sensitive greedy algorithm is $1 + 1 \times 0.75 + 5 \times 0.5 = 4.25$, which is significantly smaller than $6$ by the simple greedy algorithm without considering the query cost of each node.* $\qquad \square$

### E. Additional Discussions

In the above analysis, we prove the approximation guarantees in the most general case. However, there may be special cases where the theoretical and experimental results are significantly better. For example, Nowak [37] showed that the greedy algorithm can achieve a better guarantee given that the hierarchy has some geometric conditions. However, such conditions are hard to meet in practice. Deriving a succinct condition is certainly an interesting topic.

In this paper, we only consider a purely sequential query policy. Apparently, queries may be asked in a batch to reduce interactions. Designing effective batched algorithms is an important extension and is rather challenging. For AIGS on a tree, we can ask a batch of $k$ questions simultaneously leveraging the $k$-partition scheme [26] to ensure provable guarantees. However, for AIGS on a general DAG, it remains an open problem for devising effective algorithms in a batched version with bounded guarantees.

Furthermore, our algorithms are designed for regular crowds. Some companies may employ in-house experts in practice. In such an expert environment, we can ask some more complicated queries which are difficult to the regular crowd but tractable to their in-house experts. There is a trade-off between the difficulty and the number of queries. Constructing a flexible strategy that properly balances the trade-off is an interesting problem, and we leave it as the future work.

Note also that the AIGS problem is motivated by an application of image categorization, but our algorithms can be also applied to other content formats, e.g., textual tasks. Some existing experiment results [22] revealed that textual understanding often takes more time than visual. Therefore, the cost is likely to be higher if directly applying our approach to textual tasks.

### IV. INSTANTIATIONS OF GREEDY POLICY

A key challenge for instantiating the greedy policy is to find the *middle point* of the candidate hierarchy $G$. In this section, we first present a naive instantiation by independently calculating the total weight of all reachable nodes for every candidate node $v \in G$. To make use of the graph structure of the candidate hierarchy, we then design efficient algorithms to instantiate the greedy policy.

### A. A Naive Instantiation

Algorithm 2 gives the pseudo-code for a naive implementation of the greedy policy, referred to as GreedyNaive. It

---

**Algorithm 4:** GreedyTree($T$)

**Input:** a tree $T$
**Output:** the target node
1   $r \leftarrow$ the root of $T$;
2   initialize $\tilde{p}(v) \leftarrow p(T_v)$ and $size(v) \leftarrow |T_v|$ for each $v \in T$ by performing SetWeightDFS($T, r$) (Algorithm 5);
3   **while** $size(r) > 1$ **do**
4      $v \leftarrow r$;
5      **while** $2\tilde{p}(v) > \tilde{p}(r)$ **and** $v$ *is not a leaf node* **do**
6          $u \leftarrow v$;
7          $v \leftarrow$ the child of $u$ with the largest $\tilde{p}(v)$;
8      **if** $|2\tilde{p}(u) - \tilde{p}(r)| \leq |2\tilde{p}(v) - \tilde{p}(r)|$ **then** $q \leftarrow u$;
9      **else** $q \leftarrow v$;
10     **if** $reach(q) = yes$ **then** $r \leftarrow q$;
11     **else**
12        **foreach** *node $v$ on the path from $r$ to $q$* **do**
13           $\tilde{p}(v) \leftarrow \tilde{p}(v) - \tilde{p}(q)$;
14           $size(v) \leftarrow size(v) - size(q)$;

15   **return** *node $r$*;

---

simply enumerates all candidate nodes and computes the total probability of all nodes reachable from every node independently in each round of query (Line 3–9). In particular, the total probability of the subgraph rooted at $v$ is calculated via a subroutine GetReachableSetWeight (Line 5) given by Algorithm 3. Given a graph $G$ and a query node $v$, the GetReachableSetWeight algorithm performs a breadth first search (BFS) that starts from $v$. However, this naive instantiation can introduce heavy computational overhead.

**Time Complexity.** During the search process, as each round eliminates at least one candidate node, there will be at most $n$ rounds of finding the middle point. In each round, it will enumerate at most $n$ nodes in the candidate set and take $O(m)$ time for computing the total probability of the subgraph rooted at each candidate node, where $m$ is the number of edges in $G$. After finding the middle point, the algorithm needs to update the graph by performing a BFS according to the answer with $O(m)$ time. Hence, the total time complexity of GreedyNaive is $O(n^2 m)$.

**Discussion.** The aforementioned instantiation of the greedy policy is straightforward and intuitive, but it is far from optimized in terms of its efficiency. In fact, such a naive implementation may not handle large-scale hierarchies due to its relatively high time complexity. To tackle the efficiency issue, there are two key challenges to be solved—(i) how to find the middle point efficiently and (ii) how to update the graph efficiently after getting the query result.

*B. Efficient Instantiation on Tree*

We first look at a simple scenario when the input hierarchy is a tree and denote such a hierarchy as $T$ (instead of $G$). In the GreedyNaive algorithm, it may be quite time-consuming

---

**Algorithm 5:** SetWeightDFS($T, u$)

**Input:** a tree $T$ and a node $u$
**Output:** set the probability $\tilde{p}(u) = p(T_u)$ of subtree $T_u$ rooted at $u$ and its size $size(u) = |T_u|$
1   initialize $\tilde{p}(u) \leftarrow p(u)$ and $size(u) \leftarrow 1$;
2   **if** *node $u$ is not a leaf* **then**
3      **foreach** *child $v$ of $u$* **do**
4          SetWeightDFS($T, v$);
5          update $\tilde{p}(u) \leftarrow \tilde{p}(u) + \tilde{p}(v)$;
6          update $size(u) \leftarrow size(u) + size(v)$;

---

by performing BFS to find the middle point. We address this challenge utilizing the notion of *weighted heavy path* which extends the concept of heavy-path [44] taking into account the probability of each node. In the following, we formally define the weighted heavy path.

**Definition 10** (Weighted Heavy Path). *Given an internal node $u$ of $T$, let $v$ be a child of $u$ with the largest subtree weight[2] (with ties broken arbitrarily). Then, the edge between $u$ and $v$ is said to be heavy, and the other out-going edges of $u$ are said to be light. A weighted heavy path is a maximal path by concatenating heavy edges, i.e., the path cannot be extended with another heavy edge.*

According to the above definition, it is easy to see that for each node $u$ in $T$, there is at most one in-coming/out-going edge of $u$ is heavy. This implies that every node appears in one and exactly one weighted heavy path. We denote by $H(T, u)$ the weighted heavy path of $T$ that contains the node $u$. We show that $H(T, r)$ contains a middle point of $T$, where $r$ is the root of $T$.

**Theorem 5.** *Given a tree $T$, let $H(T, r)$ be the weighted heavy path containing the root $r$ of $T$, with respect to the node probability. Then, the node $u \in H(T, r)$ that minimizes $|2 \cdot p(T_u) - p(T)|$ is the middle point of $T$, where $T_u$ denotes the subtree of $T$ rooted at $u$.*

*Proof.* Consider any internal node $u \in T$, let $v$ be the child of $u$ such that the edge $(u, v)$ is heavy and $x$ be any other child of $u$. Clearly, it holds that $p(T_u) \geq 2p(T_x)$, since $p(T_v) \geq p(T_x)$ and $p(T_u) \geq p(T_v) + p(T_x) + p(u)$. Thus,

$$|2p(T_x) - p(T)| = p(T) - 2p(T_x).$$

Moreover, for any descendant $y$ of $x$, we have $p(T_x) \geq p(T_y)$, which indicates that

$$|2p(T_y) - p(T)| = p(T) - 2p(T_y) \geq p(T) - 2p(T_x).$$

Now, we consider the node $v$. There are two cases: (i) $p(T) \geq 2p(T_v)$ and (ii) $p(T) < 2p(T_v)$. For case (i), obviously,

$$|2p(T_v) - p(T)| = p(T) - 2p(T_v) \leq p(T) - 2p(T_x).$$

---

[2]The weight of a subtree is the total probability of nodes therein.

**Algorithm 6:** GreedyDAG($G$)

**Input:** an input DAG $G$
**Output:** the target node

1 round $w(v) \leftarrow \left\lceil \frac{n^2 \cdot p(v)}{\max_{v \in V} p(v)} \right\rceil$ for each $v$;
2 initialize $r \leftarrow$ root of $G$, and $\tilde{w}(v) \leftarrow w(G_v)$ for each $v$;
3 **while** $r$ *is not a leaf* **do**
4      initialize $q \leftarrow r$, queue $S \leftarrow \{r\}$, and $min\_w \leftarrow +\infty$;
5      **while** $S$ *is not empty* **do**
6          pop out a node from $S$ as $u$;
7          **foreach** *child $v$ of $u$* **do**
8              **if** $|2\tilde{w}(v) - \tilde{w}(r)| < min\_w$ **then**
9                  $q \leftarrow v$;
10                  $min\_w \leftarrow |2\tilde{w}(v) - \tilde{w}(r)|$;
11              **if** $2\tilde{w}(v) > \tilde{w}(r)$ **then** insert $v$ into $S$;
12      **if** $query(q) = yes$ **then** $r \leftarrow q$;
13      **else**
14          **foreach** $v \in G_q$ **do** AdjustWeight($G, v$);
15          delete every node in $G_q$ from $G$;
16 **return** *node $r$;*

---

**Algorithm 7:** AdjustWeight($G, v$)

**Input:** a DAG $G$ and the node $v$ to be deleted

1 initialize a queue $S$ and insert the node $v$ into $S$;
2 mark every node in $G$ as *unvisited*;
3 **while** $S$ *is not empty* **do**
4      pop out a node from $S$ as $u$;
5      **foreach** *each parent $v$ of $u$* **do**
6          **if** $v$ *is not visited* **then**
7              $\tilde{w}(v) \leftarrow \tilde{w}(v) - w(v)$;
8              insert $v$ into $S$ and mark $v$ as *visited*;

---

For case (ii), as $p(T_v) + p(T_x) \leq p(T)$, it also holds that

$$|2p(T_v) - p(T)| = 2p(T_v) - p(T) \leq p(T) - 2p(T_x).$$

Putting it together, we have

$$|2p(T_v) - p(T)| \leq |2p(T_x) - p(T)| \leq |2p(T_y) - p(T)|.$$

This implies that $v$ dominates all other children of $u$ as well as all their descendants. Hence, starting with $u = r$, we just keep the subtree rooted at the child $v$ of $u$ such that the edge $(u, v)$ is heavy, and meanwhile remove all subtrees rooted at the other children of $u$, since all nodes in latter are dominated by $v$. We then set $u = v$ and repeat the process until $v$ is a leaf node. This procedure actually generates the weighted heavy path $H(T, r)$ containing the root $r$ of $T$. Therefore, by definition, the node $u \in H(T, r)$ that minimizes $|2 \cdot p(T_u) - p(T)|$ is the middle point of $T$. □

Based on the Theorem 5, we can find the middle point by considering the nodes on the weighted heavy path to avoid enumerating all nodes in the whole tree. Algorithm 4 presents our improved algorithm for the instantiation of the greedy policy for AIGS on trees, namely GreedyTree. Initially, it invokes SetWeightDFS (Algorithm 5) by performing one round depth-first search (DFS) starting from the root $r$ of $T$ to compute $\tilde{p}(v) = p(T_v)$ and $size(v) = |T_v|$ for every node $v \in T$ (Line 2). To find the middle point, GreedyTree traverses the weighted heavy path using a top-down strategy starting from the root of tree (Lines 4–9). To update the candidate tree after we get the answer, instead of being fully recomputed, these weights are incrementally updated only if the answer to the query is *no* (Lines 10–14). That is, for every node $v$ on the path from the root $r$ to the query node $q$, its weight $\tilde{p}(v)$ will be decreased by a value of $\tilde{p}(q)$ after removing subtree $T_q$,

while the weight $\tilde{p}(u)$ remains the same for any $u$ not on this path. On the other hand, if the answer is *yes*, the algorithm just simply sets the query node $q$ as the new root $r$.

**Time Complexity.** Computing the weight of each subtree via SetWeightDFS (Algorithm 5) takes $O(n)$ time by performing one round DFS. Denote by $d$ the maximum out-degree of the nodes and by $h$ the length of the longest path in $T$. In each round of query, GreedyTree starts from the root and goes down along the weighted heavy path for at most $h$ layers and the node in each layer has at most $d$ children. Thus, the time complexity for the query node selection in each round is $O(hd)$. After obtaining the query response, it takes $O(h)$ time to update the weights for all ancestors of the query node. Therefore, the total time complexity of GreedyTree is $O(nhd)$.[3]

### C. Efficient Instantiation on DAG

We then move to the general case when the hierarchy is a DAG. Algorithm 6 instantiates the rounded greedy policy with $2(1 + 3 \ln n)$-approximation (see Theorem 1), namely GreedyDAG. It first initializes $\tilde{w}(v) = w(G_v)$ by GetReachableSetWeight (Algorithm 3) with respect to the rounded weight $w(\cdot)$ (Lines 1–2). Then, it goes down from the root to search for the middle point by BFS (Lines 4–11). Specifically, if the total weight $\tilde{w}(v)$ of reachable nodes from $v$ satisfies $2\tilde{w}(v) \leq \tilde{w}(r)$, where $r$ is the root of current candidate DAG, then $v$ dominates all its descendants. As a result, all descendants of such a node $v$ are skipped. Upon receiving the query answer, the graph updates the weights of $v$'s ancestors for each $v \in G_q$ only if the answer is *no* (Lines 12–15). In particular, it calls AdjustWeight (Algorithm 7) to update the weights by performing a reverse BFS from each $v \in G_q$.

**Time Complexity.** The total time complexity for computing $\tilde{w}(v)$ for every node is $O(nm)$. For each query, it takes $O(m)$ time to identify the middle point. Since there will be at most $n$ rounds of finding the middle point, the total time used for identifying query nodes is $O(nm)$. Finally, to remove a node from the graph after receiving the query result, it takes $O(m)$ time to perform a reverse BFS to update the graph via AdjustWeight (Algorithm 7). Since it at most removes $n$ nodes, the total update time is $O(nm)$. Therefore, the total

---

[3]If a max-heap is used to store the weights of the children of each node, one can verify that GreedyTree just takes $O(nh \log d)$ time.

## TABLE II
### STATISTICS OF DATASETS.

| Dataset | #nodes | Height | Max Deg. | Type | #objects |
|---------|--------|--------|----------|------|----------|
| Amazon | 29,240 | 10 | 225 | Tree | 13,886,889 |
| ImageNet | 27,714 | 13 | 402 | DAG | 12,656,970 |

## TABLE III
### COST UNDER REAL DATA DISTRIBUTION.

| Dataset | TopDown | MIGS | WIGS | **GreedyTree/GreedyDAG** |
|---------|---------|------|------|--------------------------|
| Amazon | 92.23 | 89.19 | 37.35 | **21.02** |
| ImageNet | 101.18 | 96.28 | 30.18 | **22.29** |

time complexity of GreedyDAG is $O(nm)$, which significantly improves the GreedyNaive algorithm of $O(n^2m)$.

## V. EXPERIMENTS

In this section, we evaluate the performance of our proposed approach with extensive experiments on two real-world datasets. The key metrics used for comparison are cost (i.e., the number of queries) and running time. All experiments are conducted on a machine with an Intel i7-7700 CPU and 32GB RAM. All the algorithms are implemented in Python.

### A. Experimental Setting

**Datasets.** Similar to previous work [31, 46], we test our algorithms using the following two real-world datasets, including Amazon and ImageNet. Table II summarizes some important attributes of the tested datasets.

- Amazon [20]. This is a dataset including the product hierarchy at Amazon with a tree structure. This dataset contains information of products sold at Amazon, like reviews and product metadata. Specifically, the record has a field named *categories*, and we can consider this field as a path starting from the root of the hierarchy to this product category. By combining these paths together, we can get a tree hierarchy with 29,240 nodes.

- ImageNet [14]. This is a large-scale hierarchical image dataset using the structure of WordNet [35]. Each category is represented as a tag called as *synset* in the XML file (https://www.imagenet.org/api/xml/structure_released.xml), and the subcategory is given inside each tag explicitly. The attribute *wnid* of each tag assigns an unique id to each category. We note that the category with *wnid* = "fa11misc" contains miscellaneous images that do not conform to WordNet. We extract all categories except the one with *wnid* = "fa11misc" and get a DAG with 27,714 nodes.

**Metrics.** The primary metric evaluated in this section is the *cost*, i.e., the expected number of queries. Moreover, to evaluate the efficiency of proposed instantiations of the greedy policy, a comparison of running time will be included as well.

**Competing Algorithms.** We compare our GreedyTree (on Amazon) and GreedyDAG (on ImageNet) algorithms against three baselines, including the native TopDown method, the heavy-path-based binary search method for the WIGS problem proposed by Tao et al. [46], referred to as WIGS, and the search method using multiple-choice query proposed by Li et al. [31], referred to as MIGS. For MIGS [31], we consider

the number of choices read by the crowd as the cost, since a $k$-choice query can be decomposed to $k$ binary queries.

### B. Experimental Results

*1) Comparison of Cost:* We first compare the cost of different algorithms. In the experiments, we count the number of objects in each category to obtain the real probability distribution. However, the real probability distribution may not be known in practice. To tackle this issue, we apply a simple online learning approach that dynamically adjusts the empirical probability distribution on the fly upon obtaining the category result of each object. Furthermore, to demonstrate the robustness of our algorithms, we evaluate several synthetic probability distributions, including both homogeneous and heterogeneous probability settings.

**Real Data Distribution.** Table III shows the average number of queries asked by different algorithms for locating the target node given the a-priori known real data distribution. The result shows that our GreedyTree and GreedyDAG methods significantly outperform all the three baselines. Specifically, compared with the naive TopDown method, our GreedyTree and GreedyDAG algorithms save $77.21\%$ and $77.97\%$ cost on the Amazon and ImageNet datasets, respectively, while compared with MIGS, our GreedyTree and GreedyDAG algorithms save $76.43\%$ and $76.85\%$ cost on the two datasets. Meanwhile, we observe that, although WIGS outperforms both TopDown and MIGS, our GreedyTree and GreedyDAG algorithms still achieve remarkable improvements of $43.72\%$ and $26.14\%$ on the two datasets. These results demonstrate the superiority of our proposed algorithms.

Interestingly, our experiment results show that TopDown and MIGS incur comparable cost. This is because the goal of Li et al. [30] is to minimize the expected number of multiple-choice queries. However, the workload to answer a multiple-choice query may vary significantly. For example, if MIGS queries on the root node of ImageNet, a question with around 100 choices is asked, while some other queries may only have a few choices. Hence, if we take into account the number of choices read by the crowd, we will find MIGS can only reduce the crowd workload slightly by around $5\%$ compared with TopDown.

**Learning Distribution on the Fly.** Instead of assuming the a-priori known real data distribution, we learn the empirical distribution on the fly. That is, when we label the $i$-th object, we use the statistics of the first $(i-1)$ labeled objects as the input probability distribution. At the very beginning when no object is labeled, we assume that all categories occur with an equal probability. We record the average cost for every 10 thousand
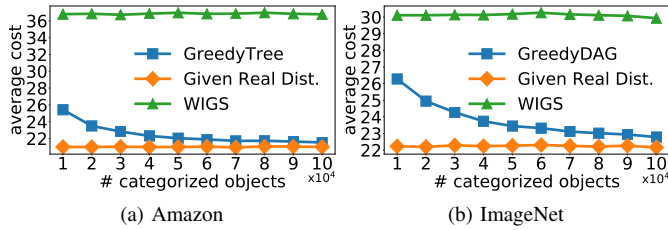
Fig. 4. Average cost vs. the number of categorized objects.

objects. Note that when the objects arrive in different sequences, our algorithms will incur different costs. We generate 20 traces by randomly shuffling the objects and report the average results. For the sake of visualization, we include WIGS, GreedyTree and GreedyDAG using the offline distribution extracted from the real data as baselines for comparison, since the costs of TopDown, MIGS are significantly higher. Fig. 4 shows the results. We observe that the average costs of the baselines, including WIGS, GreedyTree and GreedyDAG using the offline data distribution as an input, are almost the same as the values in Table III. The reason is that the baseline algorithms do not rely on the online learned distribution and the distribution of every 10 thousand objects roughly remains the same. As a contrary, the average cost of GreedyTree and GreedyDAG using the online learned distribution decreases along with the number of objects labeled, and gradually converges to that using the real data distribution. Intuitively, the more categorization results we get, the more accurately we can predict the distribution of the future objects that will be labeled, which will eventually help reduce the cost. Interestingly, we find that although both datasets have more than 10 million objects, our algorithms equipped with a practical online learning approach perform very close (with a difference less than 3%) to those using the real data distribution just after obtaining the categories of 50 thousand objects.

**Synthetic Data Distribution.** To further evaluate the impact of data distribution, we use several representative synthetic probability distributions, including the unweighted setting and three weighted settings. For the unweighted setting, each node in the category hierarchy will be the target node with an equal probability, i.e., $p(v) = 1/n$ for each $v$. For the weighted settings, we first randomly assign a value $x_v$ to each node $v$ according to certain distributions, including uniform distribution, exponential distribution, and Zipf distribution [52], and then normalize $x_v$ as $v$'s occurring probability, i.e., $p(v) = \frac{x_v}{\sum_v x_v}$. Specifically, for uniform distribution, we first generate a random number $x_v$ uniformly distributed in the range of $(0, 1)$ and then assign a node with the probability $p(v) = \frac{x_v}{\sum_v x_v}$. Similarly, for exponential distribution, the random number is generated from the distribution $\text{Exp}(1)$, while for Zipf distribution, the probability density function is $f(x; a) = \frac{x^{-a}}{\zeta(a)}$, where $\zeta$ is the Riemann Zeta function and $a$ is the parameter which is set as $a = 2$ by default. Note that Zipf distribution is a long-tail probability distribution. Since each node is assigned with a random probability and the probability will influence the

TABLE IV
COST UNDER SEVERAL PROBABILITY SETTINGS ON AMAZON.

| Distribution | TopDown | MIGS | WIGS | **GreedyTree** |
|---|---|---|---|---|
| Equal | 81.17 | 80.81 | 27.42 | **25.35** |
| Uniform | 81.28 | 81.19 | 27.47 | **23.68** |
| Exponential | 82.42 | 81.65 | 27.37 | **22.70** |
| Zipf | 82.09 | 81.94 | 27.55 | **14.03** |

TABLE V
COST UNDER SEVERAL PROBABILITY SETTINGS ON IMAGENET.

| Distribution | TopDown | MIGS | WIGS | **GreedyDAG** |
|---|---|---|---|---|
| Equal | 123.31 | 126.12 | 34.56 | **31.48** |
| Uniform | 125.82 | 124.66 | 34.55 | **28.66** |
| Exponential | 125.41 | 127.39 | 34.57 | **27.00** |
| Zipf | 125.24 | 133.48 | 34.74 | **14.41** |

performance, we repeat the experiments 20 times and report the average result.

Table IV and Table V show the expected costs of different algorithms with various probability settings on the Amazon and ImageNet datasets, respectively. We observe that our GreedyTree and GreedyDAG algorithms consistently and significantly outperform TopDown, MIGS and WIGS for all the four probability settings tested. Moreover, we find that the costs of TopDown and WIGS almost remain the same for different distributions. The reason is that these baseline algorithms do not take into account the probability distribution and the expected probability associated with each node is the same as the generated probability is an independent and identically distributed random variable. On the contrary, the costs of our GreedyTree and GreedyDAG methods are lower if the distribution of probability is more skewed. That is, the costs of GreedyTree and GreedyDAG under the Zipf distribution are substantially less than those under the exponential distribution, and are in turns smaller than those under the uniform distribution, and are in turns further less than those under the unweighted setting. In particular, under the Zipf distribution where a few nodes have extremely large probabilities, the improvement percentage of GreedyTree and GreedyDAG over WIGS reaches 59.07% and 58.52% on the Amazon and ImageNet datasets, respectively, whereas for the unweighted setting, the reduction on cost is around 10%. These results demonstrate the superiority of our cost-effective design for the AIGS problem.

The experiment results in Table IV and Table V shows that a skewed probability distribution is in favor of our approach. We further confirm this observation by testing different parameters $a$ for the Zipf distribution. That is, the smaller the value parameter $a$ is, the more biased the distribution of probability is. Fig. 5 gives the result. It clearly shows that the costs of GreedyTree and GreedyDAG increase along with the parameter $a$ of the Zipf distribution, and finally approaches the costs when each node has an identical probability. This is because if a node is associated with a high probability, our
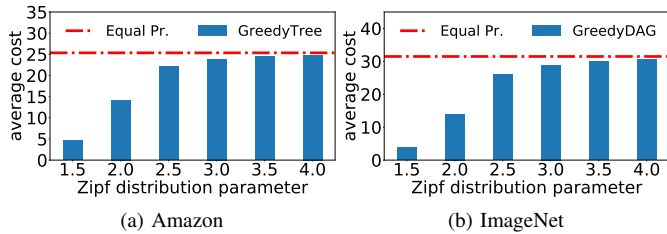
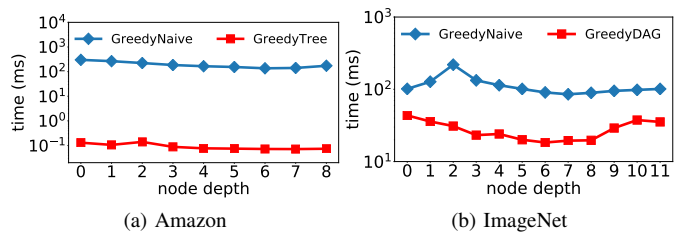Fig. 5. Cost vs. parameter of Zipf distribution.



Fig. 6. Running time.

approach will ask the question that can directly lead to this node, which can avoid many unnecessary queries.

*2) Comparison of Running Time:* In this section, we evaluate the efficiency of our GreedyTree and GreedyDAG algorithms and consider the GreedyNaive algorithm as a baseline for comparison. We randomly select 1,000 nodes in each depth as the target nodes and compute the average running time. Note that one node may be chosen multiple times. For example, in depth 0, there is only one node, i.e., the root, and the root will be selected 1,000 times as the target node. Fig. 6 shows the average running time varying along with the node depth. On the Amazon dataset with a tree hierarchy, our GreedyTree algorithm runs three orders of magnitudes faster than the GreedyNaive algorithm. On the ImageNet dataset with a general DAG hierarchy, our GreedyDAG algorithm is also noticeably faster than the GreedyNaive algorithm. These results demonstrate the efficiency of our proposed approach.

## VI. RELATED WORK

### A. Human-Based Computation

Our work is related to human-based computation, which aims to solve tasks that are more challenging for computers but are easier to handle for humans. One idea is to ask for some information from humans and tackle the challenge with human intelligence. Crowdsourcing has emerged as a major technique that enables programmers to employ humans to solve these problems. There have been many works in the database area to study how to use human assistance to process data [17, 33, 34], such as SQL-like query processing [13, 23, 29, 42], retrieving the maximum item from a set [47], and finding the highest-ranked object [19]. Moreover, crowdsourcing is also applied in some other applications like data integration and data cleaning [49], best path selection in geo-positioning services [50], and filtering data based on a set of properties [41]. Some similar tasks also occur in the machine learning area. In (supervised) machine learning, some algorithms require labeled data for training. However, the labeled data may be difficult or expensive to obtain. In active learning, the algorithm can choose the data from which it learns knowledge and requests the corresponding labels from human [11, 43]. This task is similar to our work in the sense that they all incorporate humans into the computation.

### B. Interactive Graph Search

Interactive graph search (IGS) aims to find the target node on DAG with the assist of humans. The representative application of IGS is object categorization [6, 31], hierarchy creation [7, 45], and faceted search [2]. This problem is originally proposed by Tao et al. [46] and they devised algorithms using the heavy-path-based binary search technique with near-optimal theoretical bounds considering the *worst-case* cost, referred to as WIGS which aims to minimize the cost in the worst-case. Prior to this work, Parameswaran et al. [39] proposed the human-assisted graph search problem (HumanGS), which studied a graph search problem similar to IGS under the offline setting. Recently, Li et al. [31] studied a variant IGS problem by asking multiple-choice questions and assuming the input hierarchy is a tree, namely MIGS. Zhu et al. [51] later studied the budget constrained interactive graph search for multiple targets, which considers that there are more than one target nodes on the hierarchy and the budget may not allow us to find the target nodes exactly. The techniques developed for these problems are ineffective for addressing our AIGS problem, due to the differences in problem definitions.

### C. Partially Ordered Set

From an algorithmic perspective, the IGS problem has a close relationship with the search in a partially ordered set (poset) problem. This is an extension of the well-studied problem of search in a fully ordered set [24]. Most of the researchers focus on the tree-like poset and reduce the number of comparisons in the *worst-case*, which is equivalent to the WIGS [46] problem given the input hierarchy is a tree. It has been proved that this problem is equivalent to finding the edge ranking of a simple tree corresponding to the Hasse diagram [15]. There has been an efficient algorithm in linear time to solve this problem [3, 36, 38]. However, optimizing the worst-case cost for the general poset, i.e., WIGS [46] given a DAG hierarchy, is shown to be NP-hard, and there exists an $O\left(\frac{\log n}{\log \log n}\right)$-approximate polynomial-time algorithm [15]. In contrast to the problem of optimizing the worst-case cost, minimizing the average-case cost for search in poset is more challenging, which is NP-hard even if the poset has a tree structure [8]. In fact, it has been proved that there is no $o(\log n)$-approximate algorithm unless $NP \subseteq DTIME\left(n^{O(\log \log n)}\right)$ [8]. However, surprisingly, the simple greedy method can provide a constant approximate guarantee for the tree-like poset [9, 27].

In addition, Dereniowski et al. [16] studied a problem similar to MIGS and proposed a quasi-polynomial time approximation scheme, while they tried to minimize the cost in the worst-case. Bose et al. [4] further gave an online $O(\log \log n)$-competitive search tree data structure to such a MIGS-like problem.

## D. Decision Tree

Moreover, as discussed in Section III, AIGS is also highly relevant to decision tree. The decision tree problem [5, 28] is also known as the split tree problem [25] and binary identification problem [18] in the literature. The task is to find a scheme to unambiguously identify objects with the minimized average number of queries by a set of binary tests. The decision tree problem is a classic formulation of variant problems such as active learning, entity identification and diagnosis [11]. We can consider the AIGS problem as a special case of the binary decision tree problem, and the search strategy of AIGS can be easily represented as a binary decision tree. It has been proven that computing the optimal solution for the decision tree problem is NP-hard [28]. Dasgupta [11] and Kosaraju et al. [25] showed that the nature greedy algorithm has an approximation ratio of $O(\log(\frac{1}{p_{min}}))$, where $p_{min}$ is the minimum occurring probability of objects. Through some simple modifications, the rounded greedy algorithm [5, 25] can easily achieve an approximation ratio of $O(\log n)$ regardless of the occurring probability.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we study the problem of average-case interactive graph search (AIGS). We show that the AIGS problem is NP-hard and propose cost-effective greedy algorithms with provable theoretical guarantees, e.g., $(1 + \sqrt{5})/2$ if the input hierarchy is a tree and $O(\log n)$ in general. Moreover, we devise efficient instantiations to accelerate the algorithms. With extensive experiments on real data, we show that our solutions considerably outperform the state of the art for AIGS.

There are some possible future research directions. For AIGS, the questions are answered by employees via a crowdsourcing platform. The employees may make mistakes during answering questions, which will introduce some noise to the search process. Existing experiment [31, 46] has shown that some noise is even persistent resulting from incomplete or questionable ground truth in the dataset or the subjective judgment from employees. Dealing with the negative influence of noise, especially persistent noise, is a challenge. On the theoretical side, it is interesting to derive better approximation ratios leveraging the structure of category hierarchy, such as height and maximum degree.

## REFERENCES

[1] M. Adler and B. Heeringa. Approximating optimal binary decision trees. *Algorithmica*, 62:1112–1121, 2012.

[2] S. Basu Roy, H. Wang, G. Das, U. Nambiar, and M. Mohania. Minimum-effort driven dynamic faceted search in structured databases. In *Proc. ACM CIKM*, pages 13–22, 2008.

[3] Y. Ben-Asher, E. Farchi, and I. Newman. Optimal search in trees. *SIAM Journal on Computing*, 28(6):2090–2102, 1999.

[4] P. Bose, J. Cardinal, J. Iacono, G. Koumoutsos, and S. Langerman. Competitive online search trees on trees. In *Proc. SODA*, pages 1878–1891, 2020.

[5] V. T. Chakaravarthy, V. Pandit, S. Roy, P. Awasthi, and M. Mohania. Decision trees for entity identification: Approximation algorithms and hardness results. In *Proc. PODS*, pages 53–62, 2007.

[6] K. Chakrabarti, S. Chaudhuri, and S.-w. Hwang. Automatic categorization of query results. In *Proc. ACM SIGMOD*, pages 755–766, 2004.

[7] L. B. Chilton, G. Little, D. Edge, D. S. Weld, and J. A. Landay. Cascade: Crowdsourcing taxonomy creation. In *Proc. ACM CHI*, pages 1999–2008, 2013.

[8] F. Cicalese, T. Jacobs, E. Laber, and M. Molinaro. On the complexity of searching in trees and partially ordered structures. *Theoretical Computer Science*, 412(50):6879–6896, 2011.

[9] F. Cicalese, T. Jacobs, E. Laber, and M. Molinaro. Improved approximation algorithms for the average-case tree searching problem. *Algorithmica*, 68(4):1045–1074, 2014.

[10] S. Cohen, S. Cohen-Boulakia, and S. Davidson. Towards a model of provenance and user views in scientific workflows. In *Proc. DILS*, pages 264–279, 2006.

[11] S. Dasgupta. Analysis of a greedy active learning strategy. In *Proc. NeurIPS*, pages 337–344, 2004.

[12] C. Daskalakis, R. M. Karp, E. Mossel, S. J. Riesenfeld, and E. Verbin. Sorting and selection in posets. *SIAM Journal on Computing*, 40(3):597–622, 2011.

[13] S. B. Davidson, S. Khanna, T. Milo, and S. Roy. Using the crowd for top-k and group-by queries. In *Proc. ICDT*, pages 225–236, 2013.

[14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE CVPR*, pages 248–255, 2009.

[15] D. Dereniowski. Edge ranking and searching in partial orders. *Discrete Applied Mathematics*, 156(13):2493–2500, 2008.

[16] D. Dereniowski, A. Kosowski, P. Uznanski, and M. Zou. Approximation strategies for generalized binary search in weighted trees. In *Proc. ICALP*, pages 84:1–84:14, 2017.

[17] M. J. Franklin, D. Kossmann, T. Kraska, S. Ramesh, and R. Xin. CrowdDB: Answering queries with crowdsourcing. In *Proc. ACM SIGMOD*, pages 61–72, 2011.

[18] M. R. Garey. Optimal binary identification procedures. *SIAM Journal on Applied Mathematics*, 23(2):173–186, 1972.

[19] S. Guo, A. Parameswaran, and H. Garcia-Molina. So who won? dynamic max discovery with the crowd. In *Proc. ACM SIGMOD*, pages 385–396, 2012.

[20] R. He and J. McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proc. WWW*, pages 507–517, 2016.

[21] T. Heinis and G. Alonso. Efficient lineage tracking for scientific workflows. In *Proc. ACM SIGMOD*, pages 1007–1018, 2008.

[22] A. Jain, A. D. Sarma, A. Parameswaran, and J. Widom. Understanding workers, developing effective tasks, and enhancing marketplace dynamics: A study of a large crowdsourcing marketplace. *Proc. VLDB Endowment*, 10 (7):829–840, 2017.

[23] A. M. E. W. D. Karger and S. M. R. Miller. Human-powered sorts and joins. *Proc. VLDB Endowment*, 5(1): 13–24, 2011.

[24] D. E. Knuth. *The Art of Computer Programming, Volume 3: Sorting and Searching*. Addison Wesley Longman Publishing Co., Inc., 1998.

[25] S. R. Kosaraju, T. M. Przytycka, and R. Borgstrom. On an optimal split tree problem. In *Proc. WADS*, pages 157–168, 1999.

[26] S. Kundu and J. Misra. A linear tree partitioning algorithm. *SIAM Journal on Computing*, 6(1):151–154, 1977.

[27] E. Laber and M. Molinaro. An approximation algorithm for binary searching in trees. *Algorithmica*, 59(4):601–620, 2011.

[28] H. Laurent and R. L. Rivest. Constructing optimal binary decision trees is np-complete. *Information Processing Letters*, 5(1):15–17, 1976.

[29] G. Li, C. Chai, J. Fan, X. Weng, J. Li, Y. Zheng, Y. Li, X. Yu, X. Zhang, and H. Yuan. CDB: Optimizing queries with crowd-based selections and joins. In *Proc. ACM SIGMOD*, pages 1463–1478, 2017.

[30] R. Li, P. Liang, and S. Mussmann. A tight analysis of greedy yields subexponential time approximation for uniform decision tree. In *Proc. SODA*, pages 102–121, 2020.

[31] Y. Li, X. Wu, Y. Jin, J. Li, and G. Li. Efficient algorithms for crowd-aided categorization. *Proc. VLDB Endowment*, 13(8):1221–1233, 2020.

[32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *Proc. ECCV*, pages 740–755, 2014.

[33] X. Liu, M. Lu, B. C. Ooi, Y. Shen, S. Wu, and M. Zhang. Cdas: A crowdsourcing data analytics system. *Proc. VLDB Endowment*, 5(10):1040–1051, 2012.

[34] A. Marcus, E. Wu, D. R. Karger, S. Madden, and R. C. Miller. Crowdsourced databases: Query processing with people. In *Proc. CIDR*, pages 211–214, 2011.

[35] G. A. Miller. *WordNet: An Electronic Lexical Database*. MIT Press, 1998.

[36] S. Mozes, K. Onak, and O. Weimann. Finding an optimal tree searching strategy in linear time. In *Proc. SODA*, pages 1096–1105, 2008.

[37] R. Nowak. Generalized binary search. In *Proc. Allerton*, pages 568–574, 2008.

[38] K. Onak and P. Parys. Generalization of binary search: Searching in trees and forest-like partial orders. In *Proc. IEEE FOCS*, pages 379–388, 2006.

[39] A. Parameswaran, A. D. Sarma, H. Garcia-Molina, N. Polyzotis, and J. Widom. Human-assisted graph search: It's okay to ask questions. *Proc. VLDB Endowment*, 4 (5):267–278, 2011.

[40] A. Parameswaran, S. Boyd, H. Garcia-Molina, A. Gupta, N. Polyzotis, and J. Widom. Optimal crowd-powered rating and filtering algorithms. *Proc. VLDB Endowment*, 7(9):685–696, 2014.

[41] A. G. Parameswaran, H. Garcia-Molina, H. Park, N. Polyzotis, A. Ramesh, and J. Widom. Crowdscreen: Algorithms for filtering data with humans. In *Proc. ACM SIGMOD*, pages 361–372, 2012.

[42] H. Park, H. Garcia-Molina, R. Pang, N. Polyzotis, A. Parameswaran, and J. Widom. Deco: A system for declarative crowdsourcing. *Proc. VLDB Endowment*, 5 (12):1990–1993, 2012.

[43] B. Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 2009.

[44] D. D. Sleator and R. E. Tarjan. A data structure for dynamic trees. *Journal of Computer and System Sciences*, 26(3):362–391, 1983.

[45] Y. Sun, A. Singla, D. Fox, and A. Krause. Building hierarchies of concepts via crowdsourcing. In *Proc. IJCAI*, pages 844–851, 2015.

[46] Y. Tao, Y. Li, and G. Li. Interactive graph search. In *Proc. ACM SIGMOD*, pages 1393–1410, 2019.

[47] P. Venetis, H. Garcia-Molina, K. Huang, and N. Polyzotis. Max algorithms in crowdsourcing environments. In *Proc. WWW*, pages 989–998, 2012.

[48] N. Vesdapunt, K. Bellare, and N. Dalvi. Crowdsourcing algorithms for entity resolution. *Proc. VLDB Endowment*, 7(12):1071–1082, 2014.

[49] J. Wang, T. Kraska, M. J. Franklin, and J. Feng. CrowdER: Crowdsourcing entity resolution. *Proc. VLDB Endowment*, 5(11):1483–1494, 2012.

[50] C. J. Zhang, Y. Tong, and L. Chen. Where to: Crowd-aided path selection. *Proc. VLDB Endowment*, 7(14): 2005–2016, 2014.

[51] X. Zhu, X. Huang, B. Choi, J. Jiang, Z. Zou, and J. Xu. Budget constrained interactive search for multiple targets. *Proc. VLDB Endowment*, 14(6):890–902, 2021.

[52] G. K. Zipf. *Selected Studies of the Principle of Relative Frequency in Language*. Harvard University Press, 1932.