

# Crowdsourced Clustering via Active Querying: Practical Algorithm with Theoretical Guarantees

Yi Chen<sup>1</sup>, Ramya Korlakai Vinayak<sup>1</sup>, Babak Hassibi<sup>2</sup>

<sup>1</sup> University of Wisconsin-Madison

<sup>2</sup> California Institute of Technology

yi.chen@wisc.edu, ramya@ece.wisc.edu, hassibi@systems.caltech.edu

## Abstract

We consider the problem of clustering  $n$  items into  $K$  disjoint clusters using noisy answers from crowdsourced workers to pairwise queries of the type: “Are items  $i$  and  $j$  from the same cluster?” We propose a novel, practical, simple, and computationally efficient active querying algorithm for crowdsourced clustering. Furthermore, our algorithm does not require knowledge of unknown problem parameters. We show that our algorithm succeeds in recovering the clusters when the crowdworkers provide answers with an error probability less than  $1/2$  and provide sample complexity bounds on the number of queries made by our algorithm to guarantee successful clustering. While the bounds depend on the error probabilities, the algorithm itself does not require this knowledge. In addition to the theoretical guarantee, we implement and deploy the proposed algorithm on a real crowdsourcing platform to characterize its performance in real-world settings. Based on both the theoretical and the empirical results, we observe that while the total number of queries made by the active clustering algorithm is order-wise better than random querying, the advantage applies most conspicuously when the datasets have small clusters. For datasets with large enough clusters, passive querying can often be more efficient in practice. Our observations and practically implementable active clustering algorithm can inform and aid the design of real-world crowdsourced clustering systems. We make the dataset collected through this work publicly available (and the code to run such experiments).

## 1 Introduction

Crowdsourcing, which refers to using a crowd of potentially non-expert humans to obtain information useful for downstream tasks, has become one of the most popular ways of collecting labelled datasets for supervised learning tasks (Sorokin and Forsyth 2008; Raykar et al. 2010). There is an abundant amount of data, e.g., billions of images and texts, that can be readily scraped from the internet. However, most of these datasets are unlabeled, and it is unclear what structures might exist in them. Crowdsourcing can be a very useful resource for exploring structure in such data (Welinder et al. 2010).

We consider the problem of *crowdsourced clustering* – finding clusters in a dataset with unlabeled items by query-

ing pairs of items for similarity: “Are items  $i$  and  $j$  from the same cluster?” Viewing the items to be clustered as nodes in a graph whose edges have not been observed, we have a problem of clustering a graph with access to a noisy oracle that can answer pairwise similarity queries.

A passive strategy for clustering in this scenario is to query all the pairs (i.e., edges), a random subset of pairs (Vinayak and Hassibi 2016), or a specifically constructed subset of pairs (Gomes et al. 2011; Ibrahim and Fu 2021) and then perform graph clustering. A major hiccup of such passive strategies is that they can only recover relatively large clusters. This is because the existing polynomial time graph clustering algorithms can only recover clusters at least  $\Omega(\sqrt{n})$  in size. This is related to the well-known hidden clique problem, where the goal is to find a hidden clique of a certain size in a random graph of size  $n$ . It is currently an open conjecture that there is no polynomial time passive algorithm that can recover a hidden clique of size smaller than  $\sqrt{n}$ . Active clustering, on the other hand, can potentially transcend this barrier. In this paper, we study how to cluster a set of items using these crowdsourced pairwise comparison queries in an active manner that overcomes the issue of recovering small clusters.

**Our Contributions:** We propose an *active crowdsourced clustering algorithm* that *does not rely on any unknown problem parameters*. It is computationally efficient, simple to implement, and capable of recovering clusters regardless of their sizes. We also provide an analysis of the proposed algorithm and sample complexity bound that guarantees the algorithm’s success in recovering all the clusters with high probability (with failure probability decaying as  $1/n$ ). A key observation is that when the crowdworkers are better than random guessers (i.e., the error probability is less than  $1/2$ ), the problem of deciding whether two items,  $i$  and  $j$ , belong to the same cluster can be recast as a problem of inferring if the true parameter of a Bernoulli random variable is above or below  $1/2$ . Our algorithm is inspired by the finite law of iterated logarithms (LIL) for multi-arm bandits (Jamieson et al. 2014; Heckel et al. 2019).

We implement and deploy the proposed algorithm on a real crowdsourcing platform and evaluate its performance in real-world settings. Based on both the theoretical and the empirical results, we observe that the total number of queries

made by the active clustering algorithm is order-wise better than random querying. However, the advantage of our algorithm is most conspicuous when the datasets have small clusters, which is a hard scenario for passive clustering algorithms. For datasets with large clusters, which are easier settings, a passive querying strategy of randomly querying a subset of edges followed by graph clustering can often be query efficient in practice. To the best of our knowledge, this is the first demonstration of an active clustering algorithm working in practice (beyond simulations). We make our dataset and codebase publicly available to enable further development and deployment of such systems <sup>1</sup>.

**Related Literature:** Many prior works that consider the problem of crowdsourced clustering using pairwise similarity queries employ a passive strategy with either a deterministic pattern fixed a priori (Gomes et al. 2011; Ibrahim and Fu 2021) or randomly chosen queries (Vinayak, Oymak, and Hassibi 2014; Vinayak and Hassibi 2016). Another related line of work is entity resolution in databases, where the goal is to find data records that represent the same real-world entities. There is a rich line of work in this area (see (Wang et al. 2012; Vesdapunt, Bellare, and Dalvi 2014; Verroios and Garcia-Molina 2015) and the references therein) that use heuristics-based crowdsourcing algorithms to resolve entities. Most of these works assume that there is a machine-generated similarity matrix between different data records and use this information to decide which data records to query. (Mazumdar and Saha 2017c,b) provide analysis for some of the popular heuristics and algorithms when side information is present.

A closely related work is (Yun and Proutiere 2014), which focuses on the setting with a fixed number of clusters of large sizes, i.e.,  $\Theta(n)$ , which is an easier setting for clustering. They also assume that the number of clusters is known a priori. The authors propose spectral clustering-based algorithms and theoretically analyze both passive and adaptive querying strategies.

Another closely related work is (Mazumdar and Saha 2017a), which also considers active clustering by crowdsourcing. The key differences from our setting are that they forbid repeated querying of a pair of items, and they assume that the algorithm is aware of the error probability  $p$ . Two algorithms are proposed in (Mazumdar and Saha 2017a), one that achieves a near-optimal query complexity but is computationally hard, while the other is computationally efficient but with sub-optimal query complexity. In particular, the query complexity of the computationally efficient algorithm grows quadratically in the number of clusters  $K$ , which is very costly when there are many small clusters. Both algorithms require the cluster sizes to be at least  $\Omega(\log n)$ . Furthermore, both the algorithms in (Mazumdar and Saha 2017a) require knowledge of the error probability  $p$ , which makes it difficult to deploy in practical crowdsourcing setups. Under similar assumptions of forbidden repeated queries and assuming the knowledge of error probability, another recent work (Mukherjee, Peng, and Zhang 2022)

provides efficient algorithms to recover clusters of size at  $\Omega(k \log n)$  for a fixed error probability.

In contrast, we consider the setting where the error probabilities are unknown and repeated querying of the same pair of items to different crowdworkers is allowed. Our algorithm is simple to implement, computationally efficient, capable of recovering clusters regardless of their sizes, and agnostic of the number of clusters and error probabilities while achieving near-optimal (up to logarithmic factors) query complexity. *One of the key contributions of our work is addressing how to deal with unknown parameters, which is essential for making the algorithm practical in real-world settings.* We demonstrate that repeating queries is, in fact, practical and effective in deciding cluster memberships by deploying our algorithm on a popular crowdsourcing platform and running crowdsourced clustering tasks with real crowdworkers. The goal of repeated querying in our setting is not to drive the empirical error to 0 but instead to guarantee that either the lower confidence bound on the unknown true parameter is above or the upper confidence is below  $1/2$ .

## 2 Problem Setup

In this section, we describe the problem setup, the model, and the assumptions. Consider  $n$  items that belong to  $K$  disjoint clusters. Consider a pool of crowdworkers who provide noisy answers to pairwise queries of the type: “Are items  $i$  and  $j$  from the same cluster?” Let  $\text{Query}(i, j)$  denote such a pairwise query. Let  $X_{ij}(s)$  denote the answer provided by crowdworker  $s$  to  $\text{Query}(i, j)$ . In particular,  $X_{ij}(s) = 1$  if the answer to  $\text{Query}(i, j)$  by worker  $s$  is “yes” and  $X_{ij}(s) = 0$  if the answer is “no”. For any pair of items  $i$  and  $j$ , and any positive integer  $m$ , let  $\bar{X}_{ij}(m)$  denote the average of  $m$  independent answers to the  $\text{Query}(i, j)$ , i.e.,

$$\bar{X}_{ij}(m) := \frac{1}{m} \sum_{s=1}^m X_{ij}(s). \quad (1)$$

For any item  $j$ , we use the notation  $\text{cluster}(j)$  to denote the cluster that contains item  $j$ .

Suppose the workers were perfect; then, with  $\Theta(nK)$  queries, we could assign all the items to the correct clusters. However, the workers on crowdsourcing platforms are not experts and hence make errors.

**Model:** We consider the following two-coin model for worker errors. When two items  $i$  and  $j$  are from the same cluster, for all workers  $s$ ,

$$X_{ij}(s) = \begin{cases} 1 & \text{with probability } p, \\ 0 & \text{with probability } 1 - p. \end{cases} \quad (2)$$

When  $i$  and  $j$  are not from the same cluster, for all workers  $s$ ,

$$X_{ij}(s) = \begin{cases} 1 & \text{with probability } q, \\ 0 & \text{with probability } 1 - q. \end{cases} \quad (3)$$

We note that this is similar to the Stochastic Block Model (SBM) used in analyzing graph clustering or community detection problems (Holland, Laskey, and Leinhardt 1983; Condon and Karp 2001), where  $X_{ij} = 1$  denotes an edge and  $X_{ij} = 0$  denotes no edge between two nodes  $i$  and  $j$ .

<sup>1</sup><https://github.com/kitkatdaru/crowd-active-clustering>

**Assumptions:** We assume that the answers given by different workers are independent, i.e.,  $X_{ij}(s)$  and  $X_{ij}(s')$  are independent when  $s \neq s'$ . We also assume that while the workers make errors, they are better than random guessers, i.e.,  $1 \geq p > \frac{1}{2} > q \geq 0$ .

### 3 Active Clustering Algorithm

In this section, we present the active querying algorithm for crowdsourced clustering. The algorithm proceeds by building clusters from scratch. Initially, we have a set of items to be clustered. A randomly chosen item is set as its own cluster at the beginning. We start by picking an item  $i$  that is yet to be clustered and query it with existing clusters to decide its membership. To decide the membership of item  $i$  with cluster( $j$ ), an item is picked at random from cluster( $j$ ) and the Query( $i, j$ ) is repeated with different crowdworkers until the membership of  $i$  can be established with confidence. If it does not belong to any of the existing clusters, then it starts a new cluster. This process continues until all the items are clustered. *The key challenge is to decide the cluster memberships with guaranteed confidence when the error probabilities are unknown.* We propose using the finite law of iterated logarithms (Jamieson et al. 2014) to obtain time-varying confidence bounds, which are monotonically decreasing in time  $t$  and are valid for all  $t$ . The detailed pseudocode for the algorithm is given in Algorithm 1.

### 4 Performance Guarantees

We analyze Algorithm 1 under the error model inspired by the SBM and the assumptions described in Section 2. Let  $\Delta = \frac{1}{2} \min\{p - \frac{1}{2}, \frac{1}{2} - q\}$ . We can guarantee the following performance for Algorithm 1 under the assumptions on our model:

**Theorem 1** (Main Theorem). *Algorithm 1 succeeds in recovering all the clusters exactly with at most  $\mathcal{O}\left(\frac{nK}{\Delta^2} \log\left(n \log \frac{1}{\Delta}\right)\right)$  queries overall, with high probability.*

Note that high probability here refers to an upper bound on the probability of failure that decays at the rate upper bounded by  $1/n$ . While the bound on the query complexity is a function of the problem parameters,  $K$ ,  $p$ , and  $q$ , the algorithm itself does not need to know these parameters.

If  $p < 1$  and  $q > 0$ ,  $\Omega\left(\frac{1}{\Delta^2}\right)$  repetitions per query are needed and hence  $\Omega\left(\frac{nK}{\Delta^2}\right)$  queries are necessary for Algorithm 1 to succeed with probability at least  $3/4$ . Hence, the upper bound on the number of queries is optimal up to log factors. Furthermore, the extra price that Algorithm 1 pays for not knowing  $\Delta$  (as compared to the sample complexity if we knew  $p$  and  $q$ ) is a  $\log\left(\log \frac{b_2}{\Delta}\right)$  term. Comparing our bounds with the asymptotic lower bounds in (Yun and Proutiere 2014) and the lower bounds in (Mazumdar and Saha 2017a), we note that our bounds are within log factors of optimal query complexity. In particular, for  $p = 1 - q$ , the lower bound is of the order  $\Omega\left(\frac{nK}{\Delta^2}\right)$ .

The following corollary provides the general version of the main theorem,

---

Algorithm 1: Active Crowdclustering by Crowdsourcing; without the knowledge of  $p$  and  $q$

---

```

1: Input: set of items to be clustered  $V$ ,  $\zeta \in (0, 1)$ ,  $\delta \in (0, \log(1 + \zeta)/e)$ 
2: Output:  $\mathcal{C}$ , clusters
3: Pick  $i \in V$  randomly
4: Initialize  $\mathcal{C} = \{\mathcal{C}_1 := \{i\}\}$ 
5:  $V \leftarrow V \setminus \{i\}$ 
6: while  $V \neq \emptyset$  do
7:   Pick  $v \in V$  randomly
8:   for  $k \in [|\mathcal{C}|]$  do
9:     Pick  $u \in \mathcal{C}_k$  randomly
10:     $\bar{X}_{vu}(0) \leftarrow 0$ 
11:    for each time step  $t$  do
12:       $X_{vu}(t) \leftarrow \text{Query}(v, u)$  {Query to a distinct crowdworker}
13:       $\bar{X}_{vu}(t) \leftarrow \frac{t-1}{t} \bar{X}_{vu}(t-1) + \frac{1}{t} X_{vu}(t)$ 
        {Cumulative empirical average of the answers}
14:       $\psi(t) \leftarrow (1 + \sqrt{\zeta}) \sqrt{\frac{1+\zeta}{2t} \log\left(\frac{(1+\zeta)t}{\delta}\right)}$ 
        {Confidence interval}
15:      if  $\bar{X}_{uv}(t) - \psi(t) > \frac{1}{2}$  then
16:         $\mathcal{C}_k \leftarrow \mathcal{C}_k \cup \{v\}$  {Assign  $v$  to  $\mathcal{C}_k$ }
17:         $V \leftarrow V \setminus \{v\}$ 
18:        goto Line 6
19:      end if
20:      if  $\bar{X}_{uv}(t) + \psi(t) < \frac{1}{2}$  then
21:        goto Line 8 with  $k$  increments by 1 {Move to next cluster}
22:      end if
23:    end for
24:  end for
25:  if  $v$  is not assigned to  $\mathcal{C}_k, \forall k$  then
26:     $\mathcal{C} \leftarrow \mathcal{C} \cup \{v\}$  {Start a new cluster with  $v$ }
27:     $V \leftarrow V \setminus \{v\}$ 
28:  end if
29: end while
30: return  $\mathcal{C}$ 

```

---

**Corollary 2.** *For any  $\zeta \in (0, 1)$ ,  $c \geq 3$ ,  $\delta = \frac{\delta'}{nc} \in (0, \log(1 + \zeta)/e)$ , then with probability at least  $1 - 1/n$ , Algorithm 1 succeeds in recovering all the clusters exactly and the total number of queries made is upper bounded by  $\mathcal{O}\left(nK \frac{b_1}{\Delta^2} \log\left(\frac{nc}{b_3 \delta'} \log \frac{b_2}{\Delta}\right)\right)$ , where  $b_1 = 3$ ,  $b_2 = (1 + \zeta)^2$ ,  $b_3 = \frac{1}{(2(1 + \sqrt{\zeta}))^3}$ .*

Note that  $c$ ,  $\delta$  and  $\zeta$  can be chosen such that the failure probability decays as  $1/\text{poly}(n)$  and this does not require the knowledge of error probabilities  $p$  and  $q$ . We further note that the choice of  $\delta$  and  $\zeta$  also affects the size of confidence interval  $\psi(t)$  and hence the number of queries made by Algorithm 1. The bound presented in Theorem 1 is obtained by choosing  $c = 4$  and  $\zeta = 0.1151$ .

Our proof is based on constructing monotonically decreasing confidence intervals around the cumulative empirical mean of the answers for a pair  $(u, v)$  as a function of time, i.e., the number of queries so far,  $\bar{X}_{uv}(t)$ , and the

unknown true mean, which are simultaneously valid for all times with high probability. We construct these confidence intervals,  $\psi(t)$ , using the finite law of iterated logarithms (LIL) ((Jamieson et al. 2014)). With the help of the LIL-based confidence intervals,  $\psi(t)$ , our algorithm decides whether to assign  $u$  to the same cluster as  $v$  or decides  $u$  and  $v$  are not in the same cluster. This is done based on whether the lower confidence bound goes above  $1/2$  or the upper confidence bound goes below  $1/2$  without knowing the true mean a priori. We analyze the *stopping time*, i.e., the number of repeated queries until either the lower confidence bound goes above  $1/2$  or the upper confidence interval goes below  $1/2$ , for each queried pair to upper bound the number of repeated queries taken per pair. Combining with all the pairs queried, which is at most  $\Theta(nK)$ , provides the overall query complexity bound.

## Discussion

In this section, we reflect on Algorithm 1, discuss extensions, and compare with passive querying.

**General bound with confusion matrix** While the main result is presented with a simple two-coin error model where the probabilities  $p$  and  $q$  capture the intra- and inter-cluster error probabilities, respectively, the analysis can be extended to a more general setting. Let  $P \in [0, 1]^{n \times n}$  be the *confusion matrix* associated with the  $n$  items being clustered, where each entry  $P_{ij}$  is the probability of the answer to Query( $i, j$ ) is 1. The assumption that the workers are better than random guessers in this general case implies that  $P_{ij} > 1/2$  when  $i$  and  $j$  are from the same cluster and  $P_{ij} < 1/2$  otherwise. Define  $\Delta_{ij} := |P_{ij} - 1/2|$ . The proof of Theorem 1 can be modified to obtain the following upper bound on the total number of queries made by Algorithm 1 in the general case for successfully recovering clusters with high probability,

$$\sum_{i,j:\{i,j\} \in \Omega} \frac{b_1}{\Delta_{ij}^2} \log \left( \frac{n^c}{b_3 \delta} \log \frac{b_2}{\Delta_{ij}} \right),$$

where  $\Omega$  is the set of queries made and  $|\Omega| \leq nK$ .

**Modifications in querying** In Algorithm 1, to decide if item  $i$  belongs to cluster( $j$ ), a random item  $j$  is picked as a representative from that cluster and Query( $i, j$ ) is repeated until a decision can be made about the membership of  $i$ . Instead of repeating Query( $i, j$ ) with the same representative item, we could pick a random element from cluster( $j$ ) for each repetition. For the assumed model, there is no statistical change from our assumptions and hence the guarantee provided by Theorem 1 still holds. For the general confusion matrix case described above, careful bookkeeping is needed as instead of  $\mathbb{E}(\bar{X}_{ij}(t)) = P_{ij}$ , we will have  $\mathbb{E}(\bar{X}_{ij}(t)) = \frac{1}{t} \sum_{s=1}^t P_{ij}(s)$ . In practice, switching to different representative elements from a cluster could help avoid being stuck with a bad representative picked by chance in the beginning from cluster( $j$ ). This is illustrated in Section 6 with experiments that start with initialized clusters and we compare random versus non-random choices of candidate items for repeated querying while deciding cluster membership (Table 2). While the edge error rates are similar, there is a slight

advantage in total queries for random candidates for repetition, which arises from occasional bad choices for representative candidates that the non-random alternative could get stuck in for some queries.

**Active vs. Passive Queries** Here, we discuss the pros and cons of active querying for crowdsourced clustering when compared to using passive queries. Crowdsourced clustering using passive queries has been previously approached with a two-step process (Gomes et al. 2011; Vinayak and Hassibi 2016; Ibrahim and Fu 2021). In the first step, a random or carefully designed pre-determined subset of the  $\binom{n}{2}$  pairs of items, say  $\lceil r \binom{n}{2} \rceil$  with  $r \in (0, 1]$  are queried to partially fill a noisy adjacency matrix. In the second step, a graph clustering algorithm runs on it. We focus on computationally efficient (polynomial time) clustering algorithms for this discussion.

- **Active querying succeeds regardless of cluster sizes:** Computationally efficient graph clustering algorithms, e.g., spectral clustering (McSherry 2001; Rohe, Chatterjee, and Yu 2011), convex clustering algorithms (Chen et al. 2014; Vinayak, Oymak, and Hassibi 2014; Jalali et al. 2015), have a bottleneck in terms of the size of the smallest cluster that can be recovered. In particular, the smallest cluster has to be sufficiently large, i.e., at least  $\Omega(\sqrt{n})$ , to be recovered. This bottleneck of on the minimum cluster size is conjectured to also be necessary for any polynomial-time graph clustering algorithm (related to the *hidden clique conjecture*). Therefore, using any known computationally efficient graph clustering algorithms with passive querying can only recover clusters of size at least  $\Omega(\sqrt{n})$ . On the contrary, the sufficient condition for the exact recovery of the clusters (Theorem 1) using Algorithm 1, which is computationally efficient, holds regardless of the cluster sizes. This is also illustrated by our experiments on real crowdsourcing platform (see Section 6, Table 3).
- **Active querying algorithm is free of model parameters:** For passive querying, the knowledge of  $p - q$  and  $n_{\min}$ , or other side information, is needed to a priori set the number of queries to be made that can guarantee the exact recovery of clusters (unless we make all  $\binom{n}{2}$  queries). On the other hand, our active querying algorithm does not require the knowledge of  $p, q, K$  or the cluster sizes ahead of time to guarantee exact recovery. The only assumption is that the workers are better than random guessers ( $p > 1/2 > q$ ).
- **Active vs. passive querying sample complexity:** The sufficient number of queries to guarantee exact recovery of clusters for our active querying algorithm is  $\mathcal{O}\left(\frac{nK}{\Delta^2} \log n \log \frac{1}{\Delta}\right)$ , where  $\Delta = \min\{p - \frac{1}{2}, \frac{1}{2} - q\}$ . Let us compare this to the state-of-the-art sufficient conditions for exact recovery of clusters via graph clustering for SBM (see (Chen et al. 2014; Vinayak, Oymak, and Hassibi 2014; Jalali et al. 2015) & references there in).
  - When the smallest cluster is  $\Theta(\sqrt{n})$ : Passive graph clustering can guarantee exact recovery using at most  $\mathcal{O}\left(n^2/(p - q)^2\right)$  random queries. In the case of  $\sqrt{n}$

clusters of size  $\Theta(\sqrt{n})$ , passive graph clustering takes at most  $\mathcal{O}(n^2/(p-q)^2)$  random queries to guarantee success while our algorithm takes at most  $\mathcal{O}\left(\frac{n^{1.5}}{\Delta^2} \log n \log \frac{1}{\Delta}\right)$ . So, there is room for  $\sqrt{n}$  gain in sample complexity for active clustering.

- When the smallest cluster is very large, i.e., when all the clusters are of size  $\Theta(n)$ : In this case, since  $K$  is a constant, our algorithm takes at most  $\mathcal{O}\left(\frac{n}{\Delta^2} \log n \log \frac{1}{\Delta}\right)$ . Passive graph clustering algorithms can obtain correct clustering by using at most  $\mathcal{O}(n(\log n)^2/(p-q)^2)$  random queries. So, the relative advantage of active clustering might be limited here and might not kick in until the dataset sizes are very large, depending on the hidden constants in these bounds. From our experiments on real crowdsourcing platform (see Section 6, Table 2), we observe that passive querying followed by graph clustering can provide very good clustering outcomes with much fewer queries compared to active clustering when cluster sizes are large. However, the active clustering seems to pick up more granular differences within each cluster. See Section 6 for more details.

In summary, our active clustering algorithm 1 has advantages in terms of being agnostic to model parameters, and that its success does not depend on a cluster size bottleneck. Active clustering is competitive or better than random queries order-wise, but the advantages can be realized in practice in the regime when the cluster sizes are small, which is a hard scenario for passive algorithms.

**Remark 3.** *Here, we make a few remarks comparing our results with related works.*

- *Near-optimality:* The bounds we obtain match the bounds for the optimal algorithm in (Mazumdar and Saha 2017a) up to a  $\mathcal{O}(\log \log \frac{1}{\Delta})$  factor. A similar comparison can be made with the asymptotic optimality results in (Yun and Proutiere 2014) by observing that the Kullback-Leibler divergence between Bernoulli random variable with parameter  $\frac{1}{2}$  and  $\frac{1}{2} + \Delta$  is approximately  $\Delta^2$ . Thus, the cost for not knowing the error probabilities is, at most, a  $\log \log \frac{1}{\Delta}$  factor.
- *Computational efficiency:* The optimal algorithm in (Mazumdar and Saha 2017a) is not computationally efficient and needs the knowledge of error probability  $p = 1 - q$  to set the number of queries a priori. The hardness in their setting arises from not allowing repeated queries, leading to the challenge of initializing clusters that are of size  $\Omega(\log n/(1-2p)^2)$ . In contrast, our algorithm is computationally efficient while being near-optimal in query complexity.
- *Knowledge of error probabilities:* The setting in many related works (Mazumdar and Saha 2017a; Mukherjee, Peng, and Zhang 2022) assumes the knowledge of error probability  $p$ . This is a very strong assumption as we do not know this in practice and it is difficult to estimate it due to heterogeneity in the errors between different pairs. The key parts of previous algorithms rely on knowing the error probabilities, e.g., the first phase of the algorithm

Method	VI↓ ( $K^* = 3$ )	TQ	VI↓ ( $K^* = 30$ )	TQ
Active 1 (this paper)	<b>0.09 ± 0.06</b> ( $K = 5$ )	41,014	<b>0.52 ± 0.07</b> ( $K = 40$ )	277,370
(Yun and Proutiere 2014)	2.74 ± 0.37 ( $K = 14$ )	52,225	3.63 ± 0.33 ( $K = 24$ )	294,530
Adaptive (Yun and Proutiere 2014)	3.41 ± 0.15 ( $K = 14$ )	52,225	3.84 ± 0.32 ( $K = 13$ )	294,530
Passive				
K-means, Passive (Vinayak, Oymak, and Has- sibi 2014)	0.31 ± 0.35 ( $K = 3$ ) <b>0 ± 0</b> ( $K = 3$ )	41,014 41,014	3.4 ± 0 ( $K = 1$ ) 3.4 ± 0 ( $K =$ 1)	277,370 277,370
+Kmeans, Passive				
Spectral, Passive (Vinayak, Oymak, and Has- sibi 2014)	<b>0.04 ± 0</b> ( $K = 3$ ) <b>0.07 ± 0.22</b> ( $K = 3$ )	41,014 41,014	3.4 ± 0 ( $K = 1$ ) 3.4 ± 0 ( $K =$ 1)	277,370 277,370
+Spectral, Passive				

Table 1: VI for the clustering outcome and the total number of queries, denoted as TQ in the column header, made after running Algorithm 1 and passive clustering on simulated datasets.

in (Mazumdar and Saha 2017a) involves initializing clusters of suitable sizes which relies on these parameters. In the second phase, the membership of each remaining item is decided by querying with  $\mathcal{O}(\log n/(1-2p)^2)$  items from each cluster using Hoeffding’s inequality, which again relies on knowing  $p$ .

We would like to emphasize that the challenge we address in this paper concerns how to decide a cluster membership without knowing  $p$  and  $q$ . Our routine that uses time-varying confidence intervals can make the second phase of the algorithms in (Mazumdar and Saha 2017a) practical. However, the initial phase in their algorithms would still need the knowledge of error probability.

## 5 Simulations: Passive vs. Active Querying

Here, we compare the performance of active and passive clustering algorithms under easy and difficult settings on simulated data. With  $n = 900$  items to be clustered, we consider two scenarios varying the number of clusters  $K^*$  with each cluster of equal size: (1) *Easy case* with  $K^* = 3$  with large cluster sizes, of the order of,  $\Theta(n)$ , (2) *difficult case* with  $K^* = 30$  with small cluster sizes around the threshold of  $\Theta(\sqrt{n})$ . To simulate crowdworkers’ answers, we construct a confusion matrix  $P \in [0, 1]^{n \times n}$  for each scenario. Each entry  $P_{ij}$  denotes the probability of

observing an edge between item  $i$  and  $j$ . We draw  $P_{ij} \sim \text{Uniform}[0.6, 0.85]$  if item  $i$  and  $j$  belong to the same cluster, otherwise  $P_{ij} \sim \text{Uniform}[0.1, 0.35]$ . We use variation of information (VI) (Meila 2007) to measure the difference between the output clustering and ground truth clustering. Note that  $VI \geq 0$  and the smaller the VI, the better with  $VI = 0$ , indicating a perfect match.

We run our active clustering algorithm 1 and adaptive algorithm in (Yun and Proutiere 2014). We also tried running the active clustering algorithm in (Mazumdar and Saha 2017a), but it failed to run as the initial stage did not yield any clusters even after searching over hyper-parameters. For passive algorithms, we ran the random query algorithm from (Yun and Proutiere 2014), k-means, spectral clustering (McSherry 2001), and convex algorithms (Vinayak and Hassibi 2016). Each algorithm is run 10 times and the results are shown in Table 1, and discussed below:

- In the *easy setting* with  $K = 3$  clusters (large clusters), the passive clustering is on par with our active clustering algorithm 1. In particular, the spectral clustering and convex clustering algorithms obtain nearly perfect clustering.
- In the *difficult setting* (small clusters), our active clustering algorithm 1 outperforms all the other algorithms. This setting is around the threshold where the passive algorithms struggle to recover the clusters. We also note that the adaptive algorithm in (Yun and Proutiere 2014) is designed to work for the easy setting where the cluster sizes are at least  $\Theta(n)$ . So, it is not surprising that it does not perform well in the difficult setting. However, both the active and passive versions of the algorithms in (Yun and Proutiere 2014) do not perform well in the easy setting either. This is because they rely on estimating a homogeneous error parameter for each block to decide a priori the number of repeated queries to make with each cluster, which affects the accuracy of decisions of cluster memberships. This highlights the problem of relying on problem parameters for deciding cluster memberships.

## 6 Experiments Using Real Data

In this section, we present experimental results using real datasets and real crowdworkers. We have also conducted experiments using synthetic datasets and simulated workers.

### Experiments on a Real Crowdsourcing Platform

For experiments with real crowdworkers, we use Amazon Mechanical Turk (AMT 2005) platform where crowdworkers answered pairwise queries (Figure 1a). The instructions we provided are shown in Figure 2. We note that we did not enforce the gold standard questions. We used all the data we obtained and paid all the workers who participated regardless of accuracy. Histograms of worker error rates for all the experiments with real crowdworkers are shown in Figure 1b.

The backend and frontend for this active querying system are implemented using Node.js, embedded Javascript, CSS, and Bootstrap. We implemented a batched version of Algorithm 1 for efficiency in terms of time to run the experiments on AMT. Instead of querying one item at a time and waiting

Method	Pair Err.%	VI ↓	mean T	TQ
active (from scratch)	12.5%	1.85	21.98	43, 572
passive (Vinayak and Hassibi 2016)	20%	<b>0.23</b>	N/A	<b>17, 626</b>
active, initialized (non-random repeat)	14.27%	1.42	23.20	29, 189
active, initialized (random repeat)	14.14%	1.12	22.14	28, 824

Table 2: The percentage of node pairs in error, variation of information (VI) for the clustering outcome, the average number of repetitions per query, and the total number of queries, denoted as TQ, made after running Algorithm 1 on Dogs3 dataset run with the help of real crowdworkers on AMT. The 2nd row shows the best clustering result from (Vinayak and Hassibi 2016) (in Table 4) for the same dataset for passive clustering.

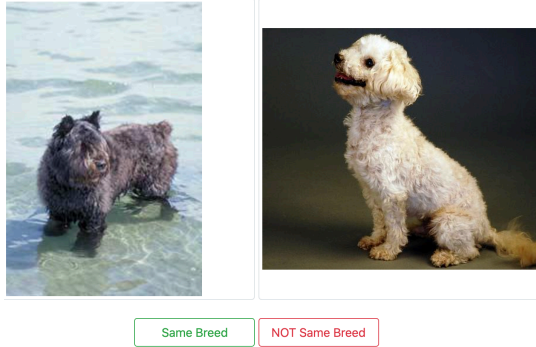
for its cluster membership to be decided, we maintain an active querying batch with 30 images to be queried (until the end, where only a few items remain to be clustered). We also maintain a yet-to-be-queried set and a clustered set. When a decision is arrived at for an image in the active querying batch as to which cluster it belongs to or to form a new cluster, it is moved from the batch to the clustered set, and a randomly chosen image from the yet-to-be-queried set is added to the batch. Each crowdworker is shown 30 pairs of images to cluster (Figure 1a shows an example of a query).

To avoid excessive cost by repetition of *difficult* to cluster images, we set the maximum number of repetitions to 80. If an image takes more than 80 queries to decide whether it belongs to a cluster and happens for all the clusters, it is considered a *difficult* to cluster image and put in a separate bucket of such hard images. We set  $\zeta = 0.0001, \delta = 0.3$  for all the experiments unless specified otherwise.

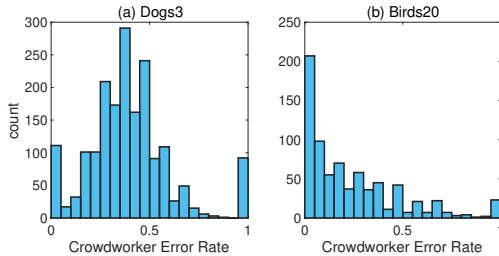
**Easy setting: Large cluster sizes Dogs3 dataset** (Khosla et al. 2011; Vinayak and Hassibi 2016) has 473 dogs of 3 different breeds (see Figure 3): Norfolk Terrier (172 images), Toy Poodle (151 images) and Bouvier des Flanders (150 images). This dataset has larger cluster sizes.

We ran the following three experiments on AMT for the Dogs3 dataset:

1. Starting from no images being clustered (referred to as *from scratch*).
2. Starting from initialized clusters where we start with the three clusters initialized with 50 images randomly chosen from respective breeds. When querying for an image  $i$ 's membership with a cluster( $j$ ),
  - (a) Choosing a random image from cluster( $j$ ) as a representative and repeating the same query to different crowdworkers (referred to as *initialized, non-random repeat*).
  - (b) Picking a randomly chosen image from cluster( $j$ ) for



(a) Sample of the pairwise queries displayed to the crowdworkers on Amazon Mechanical Turk (AMT).



(b) Histograms of crowdworkers error rate on AMT for Dogs3 and Birds20 datasets.

#### Instructions :

- You will be shown 30 pairs of images or less with dogs in them.
- For each pair of images, you will see two buttons.
- If you think the dogs in the two images are of the same breed, then click on **Same Breed** button. If not, click on **NOT Same Breed** button.
- You need to answer all the questions.
- Of the 30, there are 5 **GOLD STANDARD** questions. You need to get at least 4 of them correct to get the answers accepted.

Start

Figure 2: Sample of instructions shown for pair queries. Note that we did not enforce the gold standard questions. We used all the data we obtained and paid all the workers who participated in the tasks.

each repetition to different crowdworkers (referred to as *initialized, random repeat*).

The initialization and the order in which the images were picked to be added to the active query batch were the same for both the *non-random* and *random repeat* versions.

The results are summarized in Table 2. For comparison with passive clustering, we will refer to the results in (Vinayak and Hassibi 2016) for the same dataset. A total of 134 images for *from scratch*, 104 images for *initialized non-random repeat* and 38 images for *initialized random repeat* experiments respectively remained as *difficult* to cluster in these experiments. The output clusters for *from scratch* setting are shown in Figure 4.

Comparing the results for *active from scratch* and the best passive clustering result from (Vinayak and Hassibi 2016), we make the following observations. The clustering outcome for passive querying followed by graph cluster-



Figure 3: Sample images from the three clusters in the Dogs3 Dataset

ing seems to significantly outperform active querying with just 40% of the number of queries. Recall from discussions in Section 4 that the theoretical bounds on total query complexity for active algorithm in large cluster regime is  $\mathcal{O}\left(\frac{n}{\Delta^2} \log n \log \frac{1}{\Delta}\right)$  which when compared to the bound of  $\mathcal{O}\left(n(\log n)^2/(p-q)^2\right)$  for passive clustering is only order-wise better marginally. The data sizes we are working with here might be too small for such a slight advantage to get reflected depending on the hidden constants in these bounds. We further note that the clustering outcome from *active from scratch* has overall 3 large clusters (corresponding to the three breeds), 5 very small clusters that pick up two groups of poodles that look very different from the rest, two groups of terriers that are slightly darker and those with ears pointed when imaged, and a group of Bouvier des Flanders and 6 outliers (with only one image per cluster). So, while the clustering outcome of Algorithm 1 overall does not match the ground truth of three clusters very well compared to passive querying, it does seem to capture more granular nuances in the images.

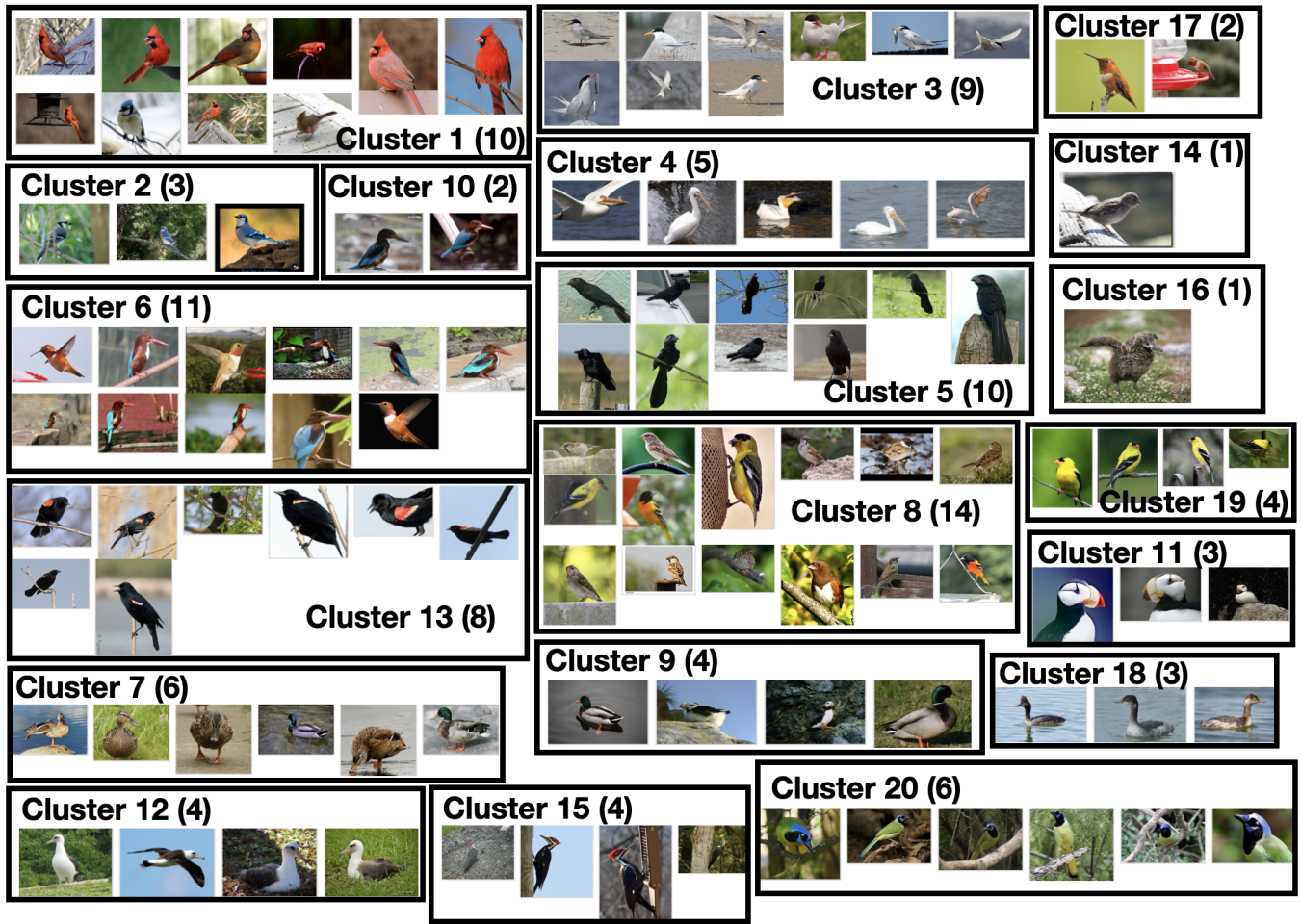
By comparing the results for the setup of *initialized non-random repeat* and *random repeat*, we note that there is no large difference in the percentage of pairs that were in error and the average number of repetitions made per query. An issue that could arise when the representative of a cluster is fixed, as is the case in the *non-random repeat* setting, is that if we are unlucky to pick a bad/atypical example from the cluster as the representative, it can lead to either error or exceeding the difficult query repeat limit. Whereas, in the *random repeat* set up, this is usually ameliorated as a random representative is chosen for each repetition. This is also reflected in the clustering outcome, where the random repeat setting performs slightly better than that of the non-random repeat. We note that in both cases, if the image being queried itself is a difficult image, then it is hard to avoid a large number of repetitions.

**Hard setting: Small cluster sizes** **Birds20** is a dataset we created using a subset of Caltech-UCSD Birds dataset (Wah et al. 2011). It has 125 images of birds from 20 different species: American Goldfinch (6), Arctic Tern (5), Baltimore Oriole (7), Blue Jay (4), Cardinal (10), Eared Grebe (3), Eastern Towhee (5), Fish Crow (4), Green Jay (6), Groove Billed Ani (6), Horned Puffin (5), House Sparrow (10), Laysan Albatross (5), Least Tern (5), Mallard (10), Pileated Woodpecker (4), Red Winged Blackbird (10), Rufous Hummingbird (5), White Breasted Kingfisher (10), and White Pelican (5). The number of images in each species cluster is shown in the bracket. This dataset has very small cluster sizes and allows us to investigate the performance of Algo-

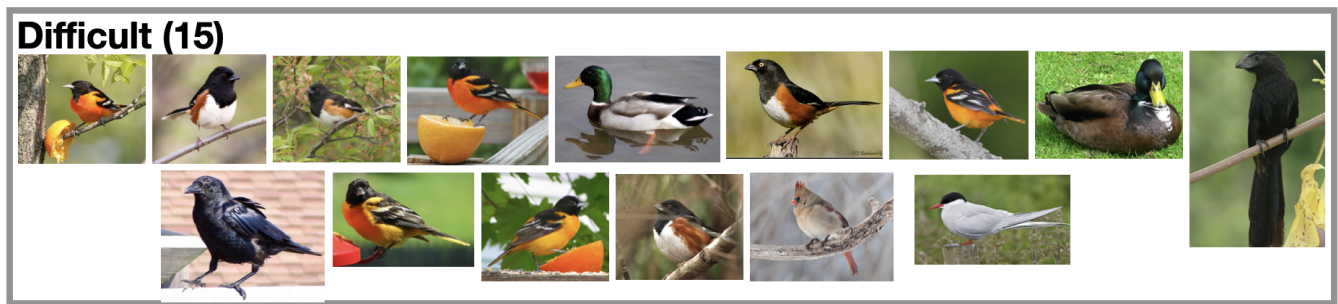


Figure 4: Clusterings obtained by running Algorithm 1 on the Dogs3 dataset from scratch.





(a)



(b)

Figure 5: Clustering obtained by running Algorithm 1 on Birds20 dataset from scratch.

Method	Pair err.%	VI ↓	mean T	TQ
active, from scratch	1.69%	<b>0.88</b> (K = 20)	12.34	15, 160
passive full, (7750 edges × 1)	15.6%	1.64 ± 0.11 (K = 6)	N/A	7, 750
passive subset repeat, (5054 edges × 3)	18.4%	1.64 ± 0.13 (K = 11)	N/A	15, 162

Table 3: The percentage of node pairs in error, VI for the clustering outcome, the average number of repetitions per query, and the total number of queries, denoted as TQ, made after running Algorithm 1 and passive clustering on Birds20 dataset run with the help of real crowdworkers on AMT.

rithm 1 in small-cluster-regime in practice.

We ran the following three experiments on AMT for the Birds20 dataset:

1. Active clustering with the batched implementation of Algorithm 1 starting from no images being clustered (referred to as *from scratch*).
2. Passive querying followed by graph clustering with
  - (a) All  $\binom{125}{2} = 7750$  edges queried once (referred to as *passive full*).
  - (b) Randomly chosen subset of edges with each edge queried thrice (referred to as *passive subset repeat*). We chose 5054 edges randomly so that with 3 repetitions, it matches the total queries made in *active from scratch* setting and used majority voting to get the adjacency matrix.

We ran k-means, spectral clustering (McSherry 2001) and improved convex algorithm (Vinayak, Oymak, and Hassibi 2014) (followed by k-means and spectral clustering) for graph clustering on the passively queried adjacency matrices. The results for these experiments with Birds20 dataset are summarized in Table 3. See Figure 4 for the clusters output by our active clustering algorithm. For passive clustering, we present the best results with the number of clusters that are resolved by the respective adjacency matrices.

Comparing the outcome of Algorithm 1 (*active from scratch*) with passive clustering, we make the following observations. In this small cluster regime, active Algorithm 1 provides much better clustering outcomes than passive clustering. We note that Algorithm 1 recovered 20 clusters overall. In contrast, the adjacency matrices filled by *passive full* and *passive repeat* could only resolve 6 and 11 clusters, respectively. This is due to the limitations of efficient clustering algorithms with respect to recovering small clusters (see Section 4).

## 7 Conclusion

In this work, we considered the problem of clustering a set of items into disjoint clusters with the help of noisy crowdworkers who can answer pairwise comparison queries of

type “Are items  $i$  and  $j$  in the same cluster?”. We proposed a practical active clustering algorithm towards this goal and, under mild assumptions, provided bounds on query complexity that guarantee the exact recovery of the clusters. The proposed active algorithm does not need the knowledge of any problem parameters, in particular the error probabilities, number of clusters, or size of the clusters. We implemented this algorithm on a real crowdsourcing platform to demonstrate its efficacy and study its performance in large and small cluster regimes. While the theoretical bound on the query complexity is order-wise better for the active clustering algorithm when the clusters are large, passive algorithms can, in fact, provide very good clustering outcomes with much fewer queries in practice. The main advantage of the active clustering algorithm seems to be in the case when there could be clusters of very small sizes that passive clustering algorithms will fail to recover. A hybrid approach that gets the best of both worlds would be useful to develop, and we leave it to future work.

## Acknowledgements

We thank Nihar B. Shah for insightful discussions and comments on the paper. This work was partially supported by NSF grants NCS-FO 2219903 and NSF CAREER Award CCF 2238876.

## References

- AMT. 2005. Amazon Mechanical Turk. <https://www.mturk.com/>. Accessed: 2022-10-12.
- Chen, Y.; Jalali, A.; Sanghavi, S.; and Xu, H. 2014. Clustering partially observed graphs via convex optimization. *Journal of Machine Learning Research*, 15(1): 2213–2238.
- Condon, A.; and Karp, R. M. 2001. Algorithms for graph partitioning on the planted partition model. *Random Struct. Algorithms*, 18(2): 116–140.
- Gomes, R. G.; Welinder, P.; Krause, A.; and Perona, P. 2011. Crowdclustering. In *Advances in Neural Information Processing Systems 24*, 558–566.
- Heckel, R.; Shah, N. B.; Ramchandran, K.; and Wainwright, M. J. 2019. Active ranking from pairwise comparisons and when parametric assumptions do not help. *The Annals of Statistics*, 47(6): 3099–3126.
- Holland, P. W.; Laskey, K. B.; and Leinhardt, S. 1983. Stochastic blockmodels: First steps. *Social Networks*, 5(2): 109 – 137.
- Ibrahim, S.; and Fu, X. 2021. Mixed Membership Graph Clustering via Systematic Edge Query. *IEEE Transactions on Signal Processing*.
- Jalali, A.; Han, Q.; Dumitriu, I.; and Fazel, M. 2015. Relative Density and Exact Recovery in Heterogeneous Stochastic Block Models. arXiv:1512.04937.
- Jamieson, K. G.; Malloy, M.; Nowak, R.; and Bubeck, S. 2014. lil’ UCB : An Optimal Exploration Algorithm for Multi-Armed Bandits. In *COLT*, 423–439.

- Khosla, A.; Jayadevaprakash, N.; Yao, B.; and Fei-Fei, L. 2011. Novel Dataset for Fine-Grained Image Categorization. In *First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition*.
- Mazumdar, A.; and Saha, B. 2017a. Clustering with noisy queries. *Advances in Neural Information Processing Systems*, 30.
- Mazumdar, A.; and Saha, B. 2017b. Query complexity of clustering with side information. *Advances in Neural Information Processing Systems*, 30.
- Mazumdar, A.; and Saha, B. 2017c. A Theoretical Analysis of First Heuristics of Crowdsourced Entity Resolution. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- McSherry, F. 2001. Spectral Partitioning of Random Graphs. In *FOCS*, 529–537. IEEE Computer Society. ISBN 0-7695-1390-5.
- Meila, M. 2007. Comparing Clusterings—an Information Based Distance. *J. Multivar. Anal.*, 98(5): 873–895.
- Mukherjee, C. S.; Peng, P.; and Zhang, J. 2022. Recovering unbalanced communities in the stochastic block model with application to clustering with a faulty oracle. *arXiv preprint arXiv:2202.08522*.
- Raykar, V. C.; Yu, S.; Zhao, L. H.; Valadez, G. H.; Florin, C.; Bogoni, L.; and Moy, L. 2010. Learning From Crowds. *J. Mach. Learn. Res.*, 11: 1297–1322.
- Rohe, K.; Chatterjee, S.; and Yu, B. 2011. Spectral clustering and the high-dimensional stochastic blockmodel. *Annals of Statistics*, 39(4): 1878–1915.
- Sorokin, A.; and Forsyth, D. 2008. Utility data annotation with Amazon Mechanical Turk. In *Computer Vision and Pattern Recognition Workshops, 2008. CVPRW '08. IEEE Computer Society Conference on*, 1–8. IEEE. ISBN 978-1-4244-2339-2.
- Verroios, V.; and Garcia-Molina, H. 2015. Entity Resolution with crowd errors. In *ICDE*, 219–230. IEEE Computer Society.
- Vesdapunt, N.; Bellare, K.; and Dalvi, N. 2014. Crowdsourcing algorithms for entity resolution. *Proceedings of the VLDB Endowment*, 7(12): 1071–1082.
- Vinayak, R. K.; and Hassibi, B. 2016. Crowdsourced Clustering: Querying Edges vs Triangles. In *Advances in Neural Information Processing Systems*, 1316–1324.
- Vinayak, R. K.; Oymak, S.; and Hassibi, B. 2014. Graph Clustering With Missing Data: Convex Algorithms and Analysis. In *Neural Information Processing Systems Conference (NIPS)*.
- Wah, C.; Branson, S.; Welinder, P.; Perona, P.; and Belongie, S. 2011. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology.
- Wang, J.; Kraska, T.; Franklin, M. J.; and Feng, J. 2012. Crowder: Crowdsourcing entity resolution. *Proceedings of the VLDB Endowment*, 5(11): 1483–1494.
- Welinder, P.; Branson, S.; Belongie, S.; and Perona, P. 2010. The Multidimensional Wisdom of Crowds. In *Neural Information Processing Systems Conference (NIPS)*.
- Yun, S.-Y.; and Proutiere, A. 2014. Community detection via random and adaptive sampling. In *Conference on learning theory*, 138–175. PMLR.