Community Recovery in Hypergraphs

Kwangjun Ahn¹⁰, Kangwook Lee¹⁰, and Changho Suh¹⁰, Member, IEEE

Abstract-Community recovery is a central problem that arises in a wide variety of applications such as network clustering, motion segmentation, face clustering, and protein complex detection. The objective of the problem is to cluster data points into distinct communities based on a set of measurements, each of which is associated with the values of a certain number of data points. While most prior works focus on a setting in which the number of data points involved in a measurement is two, this paper explores a generalized setting in which the number can be more than two. Motivated by applications particularly in machine learning and channel coding, we consider two types of measurements: 1) homogeneity measurement that indicates whether or not the associated data points belong to the same community and 2) parity measurement that denotes the modulo-2 sum of the values of the data points. Such measurements are possibly corrupted by Bernoulli noise. We characterize the fundamental limits on the number of measurements required to reconstruct the communities for the considered models.

Index Terms— Clustering algorithms, channel coding, hypergraph clustering, information-theoretic limits, generalized censored block model (GCBM).

I. INTRODUCTION

C LUSTERING of data is one of the central problems, and it arises in many fields of science and engineering. Among many related problems, *community recovery in graphs* has received considerable attention with applications in numerous domains such as social networks [3]–[5], computational biology [6], and machine learning [7], [8]. The goal of the problem is to cluster data points into different communities based on *pairwise* information. Among a variety of models for the community recovery problem, the stochastic block model (SBM) [9] and the censored block model (CBM) [10] have received significant attention in recent years. In SBM, two data points in the same communities are more likely to be connected by an edge than the other edges. In the case of CBM, each measurement returns the modulo-2 sum of

Manuscript received September 7, 2017; accepted May 11, 2019. Date of publication June 4, 2019; date of current version September 13, 2019. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea Government (MSIT) under Grant 2018R1A1A1A05022889. This paper was presented in part at the 54th Annual Allerton Conference on Communication, Control, and Computing 2016 [1], and the IEEE International Symposium on Information Theory 2017 [2]. (Kwangjun Ahn and Kangwook Lee contributed equally to this work.)

K. Ahn is with the Department of Mathematical Sciences, Korea Advanced Institute of Science and Technology, Daejeon 34141, South Korea (e-mail: kjahnkorea@kaist.ac.kr).

K. Lee and C. Suh are with the School of Electrical Engineering, Korea Advanced Institute of Science and Technology, Daejeon 34141, South Korea (e-mail: kw1jjang@kaist.ac.kr; chsuh@kaist.ac.kr).

Communicated by C. Caramanis, Associate Editor for Machine Learning. Color versions of one or more of the figures in this article are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TIT.2019.2920637

the values assigned to the two nodes, possibly corrupted by Bernoulli noise.

While these models reflect interactions between a pair of two nodes, there are numerous applications in which interactions occur across more than two nodes [11]-[13]. One such application is a folksonomy, a social network in which users can annotate items with different tags [11]. In this application, the graph consists of nodes corresponding to different users, different items, and different tags. When user *i* annotates item j with tag k, one can view this as a hyperedge connecting node i, node j and node k. Therefore, in order to cluster nodes of such a graph based on such interactions, one needs a model that can capture such three-way interactions. Another application is molecular biology, in which multiway interactions between distinct systems capture complex molecular interactions [12]. There are also a broad range of applications in other domains including computer vision [14], VLSI circuits [15], and categorical databases [16].

These applications naturally motivate us to investigate a hypergraph setting in which measurements are of multiway information type. Specifically, we consider a simple yet practically-relevant model, which we name the generalized censored block model (GCBM). In the GCBM, the n data points are modeled as nodes in a hypergraph, and their interactions are encoded as hyperedges across the nodes. The n nodes are divided into two communities. Each node takes a label between 0 or 1 depending on its affiliation. As our measurement model, we consider a random d-uniform hypergraph in which each subset of d nodes is sampled as a hyperedge with probability p. Each sampled hyperedge is then assigned with a binary label which is a function on the labels of d nodes involved. Inspired by applications in machine learning and channel coding, we study the following two types of measurements:

- *the homogeneity measurement* which reveals whether or not the *d* nodes are from the same community, i.e., whether or not the *d* nodes are having the same label; and
- *the parity measurement* which computes the modulo-2 sum of the labels of the *d* nodes.

We also investigate a noisy measurement setting in which the label of each hyperedge can possibly be flipped with probability, say $\theta \in [0, 1]$.

A. Main Contributions

Specialized to the d = 2 case, the above two measurement models both reduce to the CBM, in which the informationtheoretic limit on the expected number of edges required for

0018-9448 © 2019 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

TABLE I

Summary of Main Results. The Information-Theoretic Limits on Sample Complexity $p\binom{n}{d}$ Are summarized. Here *n* Denotes the Number of Nodes, θ Denotes the Noise Flipping Probability, and d Denotes the Size of Hyperedges. An Observation that the Sample COMPLEXITY DECREASES IN *d* FOR THE PARITY MEASUREMENT CASE MOTIVATES US TO STUDY THE CASE d = f(n) WHERE *f* IS SOME Increasing Function. Here " $\Theta_{\theta,d}$ " Means That the Multiplicative Constant Factor Depends on θ and d

	d = 2	d > 2 (const.)	d = f(n)
Homogeneity	$\frac{1}{2} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$	$\frac{2^{d-2}}{d} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$	
Parity	$\frac{1}{2} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$	$\frac{1}{d} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$	$\Theta_{\theta,d}\left(\max\left\{n, \frac{n\log n}{d}\right\}\right)$

exact recovery is characterized as $p\binom{n}{2} = \frac{1}{2} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$ [17].

On the other hand, the information-theoretic limits for the case of arbitrary d has not been settled. This precisely sets the goal of our paper: We seek to characterize the information-theoretic limits on the sample complexity for exact recovery under the two models. A summary of our findings is as follows. For a fixed constant d, the information-theoretic limits are:

- (the homogeneity measurement case) $p\binom{n}{d} = \frac{2^{d-2}}{d}$. • (the nonsection) $\frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}; \text{ and}$ • (the parity measurement case) $p\binom{n}{d} = \frac{1}{d} \cdot \frac{n \log n}{\left(\sqrt{1-\theta} - \sqrt{\theta}\right)^2}$ if

d is a fixed constant.

One interesting observation is that the sample complexity $p\binom{n}{d}$ decreases in d for the parity measurement case. This naturally motivates us to ask whether one can further reduce the sample complexity with a larger d. We also address this question. Suppose that d can arbitrarily scale with n, i.e., d =f(n) for some increasing function f. Then, the informationtheoretic limits for the parity measurement case¹ read as follows:

• $p\binom{n}{d} = \Omega\left(\frac{n\log n}{d}\right)$ if $d = o(\log n)$; and • $p\binom{n}{d} = \Omega(n)$ if $d = \Omega(\log n)$.

Our findings provide some interesting implications to applications such as subspace clustering and channel coding. For instance, the results offer concrete guidelines as to how to choose d that minimizes sample complexity while ensuring successful clustering. See details in Sec. II-A and Sec. III.

B. Related Work

1) The d = 2 Case: The exact recovery problem in standard graphs (d = 2) has been studied in great generality. In SBM, both the fundamental limits and computationally efficient algorithms are investigated initially for the case of two communities [17]-[19], and recently for the case of an arbitrary number of communities [20]. In CBM, [21] characterizes the sample complexity limit, and [17] develops a computationally efficient algorithm that achieves the limit.

Another important recovery requirement is *detection*, which asks whether one can recover the clusters better than a random guess. The modern study of the detection problem in SBM is initiated by a paper by Decelle et al. [22], which conjectures

phase transition phenomena for the detection problem.² This conjecture is initially tackled for the case of two communities. The impossibility of the detection below the conjectured threshold is established in [23], and it is proved in [24]-[26] that the conjectured threshold can be achieved efficiently. The achievability part of the conjecture for an arbitrary number of communities is recently settled by Abbe and Sandon [27]. For another line of researches, minimax-optimal rates are derived in [28], and algorithms that achieve the rates are developed in [29]. We refer to a recent survey by Abbe [30] for more exhaustive information.

2) The Homogeneity Measurement Case: Recently, [31], [32] develop efficient algorithms for hypergraph community recovery together with theoretical guarantees for the algorithms. They derive an upper bound on sample complexity for almost exact recovery, which allows for a vanishing fraction of misclassified nodes. Although the models considered therein are analogues of SBM rather than CBM in a hypergraph setting, one can apply their algorithm to our model by making small adjustments in our measurement.³ Applying their algorithm to our model, their upper bound reads $p\binom{n}{d}$ = $\Omega(n \log^2 n)$. Whether or not the sufficient condition is also necessary has been unknown. In this work, we show that it is not the case, demonstrating that the minimal sample complexity even for exact recovery is $\Theta(n \log n)$.

We note that the homogeneity measurement case is closely related to subspace clustering, one of the popular problems in computer vision [14], [33], [34]; See Sec. II-A.1 for details.

3) The Parity Measurement Case: This case has been explored by [35] in the context of random constraint satisfaction problems. The case of d = 3 has been well-studied: it is shown that the maximum likelihood decoder succeeds if $p\binom{n}{3} \ge 2 \cdot \frac{n \log n}{(0.5 - \theta)^2}$ [35]. Unlike the prior result which only considers the case of d = 3, we cover an arbitrary constant d, and characterize the sharp threshold on the sample complexity.

Abbe-Montanari [10] relate the parity measurement model to a channel coding problem in which random LDGM codes with a constant right-degree d are employed. By proving the concentration phenomenon of the mutual information between channel input and output, they demonstrate the existence of phase transition for an even d. Our results span any fixed d, and hence fully settle the phase transition (see Sec. III).

¹For the homogeneity measurement setting, we anticipate that the sample complexity would keep growing with an increase in d, which is an undesirable situation. Hence, the case d = f(n) is not studied in depth under the setting.

²In the paper, it is also conjectured that an information-computation gap might exist for the case of more than 3 communities $(k \ge 4)$.

³Specifically, we first remove hyperedges with label 0. Then, the remaining hypergraph will only have hyperedges with label 1. We regard this hypergraph as an unlabeled hypergraph by disregarding labels.

4) Detection in the Stochastic Block Model for Hypergraphs: There are several works which study the community detection under SBM for hypergraphs. In [36], the authors explore the case of two equal-sized communities.⁴ Applying their algorithm to our homogeneity measurement model, their result shows that detection is possible if $\binom{n}{d}p = \Omega(n)$. Moreover, [37] recently conjectures phase transition thresholds for detection. Lastly, [38] derives the minimax-optimal error rates, and generalizes the results in [28] to the hypergraph case.

5) Other Relevant Problems: Community recovery in hypergraphs bears similarities to other inference problems, in which the goal is to reconstruct data from multiple queries. Those problems include crowdsourced clustering [39], [40], group testing [41] and data extraction from histogram-type information [42], [43]. Here one can make a connection to our problem by viewing each query as a hyperedge measurement. However, a distinction lies in the way that queries are collected. For instance, an adaptive measurement model is considered in the crowdsourced setting [39], [40] unlike our non-adaptive setting in which hyperedges are sampled uniformly at random. Histogram-type information acts as a query in [41]–[43].

Lastly, it is also worth mentioning our follow-up work [44]. The key distinction relative to this paper (which focuses on characterizing the information-theoretic limits) is that [44] develops efficient algorithms. While [44] does not provide any sharp-threshold result, the algorithms therein can be applied to more general settings in which there could be more than two communities, and/or hyperedge measurements are in [0, 1], not limited to the binary values.

C. Paper Organization

Sec. II introduces the considered model; in Sec. III, our main results are presented along with some implications; in Sec. IV, V and VI, we provide the proofs of the main theorems; Sec. VII presents experimental results that corroborate our theoretical findings and discuss interesting aspects in view of applications; and in Sec. VIII, we conclude the paper with some future research directions.

D. Notations

For any two sequences f(n) and g(n): $f(n) = \Omega(g(n))$ if there exists a positive constant c such that $f(n) \ge cg(n)$; f(n) = O(g(n)) if there exists a positive constant c such that $f(n) \le cg(n)$; $f(n) = \omega(g(n))$ if $\lim_{n\to\infty} \frac{f(n)}{g(n)} = \infty$; f(n) = o(g(n)) if $\lim_{n\to\infty} \frac{f(n)}{g(n)} = 0$; and $f(n) \le g(n)$ or $f(n) = \Theta(g(n))$ if there exist positive constants c_1 and c_2 such that $c_1g(n) \le f(n) \le c_2g(n)$.

For a set A and an integer $m \leq |A|$, we denote $\binom{A}{m} := \{B \subset A : |B| = m\}$. Let [n] denote $\{1, \ldots, n\}$. Let \mathbf{e}_i be the *i*th standard unit vector. Let **0** be the all-zero-vector and **1** be the all-one-vector. We use $\mathbb{I}\{\cdot\}$ to denote an indicator function. Let $\mathsf{D}_{\mathsf{KL}}(p||q)$ be the Kullback-Leibler (KL)

divergence between Bern(p) and Bern(q), i.e., $D_{\text{KL}}(p||q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$. We shall use $\log(\cdot)$ to indicate the natural logarithm. We use $H(\cdot)$ to denote the binary entropy function.

II. GENERALIZED CENSORED BLOCK MODELS

Consider a collection of *n* nodes $\mathcal{V} = [n]$, each represented by a binary variable $X_i \in \{0, 1\}, 1 \leq i \leq n$. Let $\mathbf{X} := \{X_i\}_{1 \leq i \leq n}$ be the ground-truth vector. Let *d* denote the size of a hyperedge. Samples are obtained as per a *measurement hypergraph* $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ where $\mathcal{E} \subset {\binom{[n]}{d}}$. We assume that each element in ${\binom{[n]}{d}}$ belongs to \mathcal{E} independently with probability $p \in [0, 1]$. Sample complexity is defined as the number of hyperedges in a random measurement hypergraph, which is concentrated around $p {\binom{n}{d}}$ in the limit of *n*. Each sampled edge $E \in \mathcal{E}$ is associated with a noisy binary measurement Y_E :

$$Y_E = f(X_{i_1}, X_{i_2}, \dots, X_{i_d}) \oplus Z_E,$$
 (1)

where $f : \{0, 1\}^d \to \{0, 1\}$ is some binary-valued function, \oplus denotes modulo-2 sum, and $Z_E \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(\theta)$ is a random variable with noise rate $0 \le \theta < \frac{1}{2}$.⁵ For the choice of f, we focus on the two cases:

• the homogeneity measurement:

$$f_h(X_{i_1}, X_{i_2}, \dots, X_{i_d}) = \mathbb{I}\{X_{i_1} = X_{i_2} = \dots = X_{i_d}\};$$

• the parity measurement:

$$f_p(X_{i_1}, X_{i_2}, \ldots, X_{i_d}) = X_{i_1} \oplus X_{i_2} \oplus \cdots \oplus X_{i_d}.$$

Let $\mathbf{Y} := \{Y_E\}_{E \in \mathcal{E}}$. We remark that when d = 2, this reduces to CBM [21].

The goal of this problem is to recover **X** from **Y**. In this work, we will focus on the case of even *d*. One can obtain the same results also for the odd *d* case using the same proof techniques. This will be clearer later while proving the results. When *d* is even, the conditional distribution $\mathbf{Y}|\mathbf{X} \oplus \mathbf{1}$. Hence, given a recovery scheme ψ , the probability of error is defined as⁶

$$P_e(\psi) := \max_{\mathbf{X} \in \{0,1\}^n} \Pr\left(\psi(\mathbf{Y}) \notin \{\mathbf{X}, \ \mathbf{X} \oplus \mathbf{1}\}\right).$$

We intend to characterize the minimum sample complexity, above which there exists a recovery algorithm ψ such that $P_e(\psi) \rightarrow 0$ as *n* tends to infinity, and under which $P_e(\psi) \not\rightarrow 0$ for all algorithms.

A. Relevant Applications

⁵Note that the condition $\theta < 1/2$ is to ensure that the maximum likelihood estimator is equivalent to the minimum distance decoder. The other case $\frac{1}{2} < \theta \le 1$ can be readily dealt with by simply flipping all the binary measurements. Moreover, the choice of Bernoulli distribution is only for simplicity. One can choose a different noise distribution as long as the tail of the distribution behaves nicely. This will only change the constant factor in the characterization of information-theoretic limits. See Sec. III.

⁶Notice in the parity measurement case that for an odd *d*, the definition of $P_e(\psi)$ should read instead: $P_e(\psi) := \max_{\mathbf{X} \in \{0,1\}^n} \Pr(\psi(\mathbf{Y}) \neq \mathbf{X})$. This is because the conditional distribution $\mathbf{Y}|\mathbf{X}$ is no longer equal to $\mathbf{Y}|\mathbf{X} \oplus \mathbf{1}$ under the case.

⁴Actually, the main model in the paper is *the bipartite stochastic block model*, which is not a hypergraph model. However, the result for the hypergraph case follows as a corollary (see Theorem 5 therein).



Fig. 1. Connection to subspace clustering. Subspace clustering is illustrated for a simple scenario in which the entire signal space is two-dimensional and data points are approximately lying on a union of two 1-dimensional affine spaces (lines). A common procedure in the existing algorithms includes construction of a *d*-th order affinity tensor ($d \ge 2$) each entry of which represents a quantity that captures a level of similarity across *d* data points, so taking either 0 or 1 depending on the similarity level. For instance, the four points involved in E_1 in the figure lie near the same affine space, so the similarity measure is decided as 1; on the other hand, the four points in E_2 span different affine spaces, so the similarity measure is decided as 0. Since each data point does not exactly lie in a subspace, an error can occur in the decision—the similarity measurement can be noisy. Hence one can view this problem as the GCBM under the homogeneity measurement model.

1) Subspace Clustering Homogeneity and the Measurement: Subspace clustering is a popular problem of which the task is to cluster *n* data points that approximately lie in a union of lower-dimensional affine spaces. The problem arises in a variety of applications such as motion segmentation [45] and face clustering [46], where data points corresponding to the same class (tracked points on a moving object or faces of a person) lie on a single lower-dimensional subspace; for details, see [47] and references therein. A common procedure of the existing algorithms for subspace clustering [34], [48]–[50] begins construction of a d-th order affinity tensor $(d \ge 2)$ whose entries represent *similarities* between every d data points. Since this construction incurs a complexity that scales like n^d , sampling-based approaches are proposed in [14], [33], [34].

A similarity between d data points in prior works [14], [33], [34] is defined such that it tends to 1 if all of the d points can be well-fitted by a single low-dimensional affine space and 0 otherwise. Hence, restricted to the two-subspace case, one can view a similarity over a d-tuple E as a homogeneity measurement.⁷ By setting the probability of each entry being sampled as p, one can relate this to our homogeneity measurement model; see Fig. 1 for visual illustration.

2) Channel Coding and the Parity Measurement: The community recovery problem has an inherent connection with channel coding problems [18], [21]. To see this, consider a communication setting which employs random LDGM codes with a constant right-degree *d*. To make a connection, we begin by constructing a random *d*-uniform hypergraph with *n* nodes,



Fig. 2. **Connection to channel coding.** GCBM with the parity information can be seen as a channel coding problem which employs random LDGM codes with a constant right-degree d. To see this, we first draw a random d-uniform hypergraph with n nodes, where each edge of size d appears with probability p. Given the input sequence of n information bits, the parity bits corresponding to all the sampled hyperedges are concatenated, forming a codeword. The noisy measurement can be mapped to the output of a binary symmetric channel (BSC) with crossover probability θ , when fed by the codeword. A recovery algorithm ψ corresponds to the decoder which wishes to infer the n information bits from the received signals. One can then see that recovering communities in hypergraphs is equivalent to the above channel coding problem.

where each edge of size *d* appears with probability *p*. Given the input sequence of *n* information bits, we then concatenate the parity bits with respect to the sampled hyperedges to form a codeword of average length $p\binom{n}{d}$. Note that the expected code rate is $\frac{n}{p\binom{n}{d}}$. The noisy measurement can be mapped to the output of a binary symmetric channel (BSC) with crossover probability θ , when fed by the codeword. A recovery algorithm ψ corresponds to the decoder which wishes to infer the *n* information bits from the received signals. One can then see that recovering communities in hypergraphs is equivalent to the above channel coding problem; see Fig. 2 for visual illustration.

III. MAIN RESULTS

A. The Homogeneity Measurement Case

Theorem 1. Let $d \ge 2$ be a constant. Under the homogeneity measurement case, the following holds for any $\epsilon > 0$:

$$\inf_{\psi} P_e(\psi) \to 0 \quad if \binom{n}{d} p \ge (1+\epsilon) \frac{2^{d-2}}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2}; \\ \inf_{\psi} P_e(\psi) \not\to 0 \quad if \binom{n}{d} p \le (1-\epsilon) \frac{2^{d-2}}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2}.$$

Proof: See Sec. IV.

We first make a comparison to the result in [31]. While [31] models a fairly general similarity measurement, it considers a more relaxed performance metric, so called almost exact recovery, which allows for a vanishing fraction of misclassified nodes; and provides a sufficient condition on sample complexity under the setting [52]. On the other hand, we identify the sufficient and necessary condition for *exact* recovery, thereby characterizing the fundamental limit. Specializing their result to the model of our interest, the sufficient condition in [31]

⁷In subspace clustering, similarities can sometimes be noisy in that even though the *d* data points are from the same (different) subspace, similarity can be 0 (1). Note that Z_E in (1) captures this noise. Nonetheless, our noise model cannot fully respect the noise effect that arises in subspace clustering in which the noises are not necessarily i.i.d. In fact, more realistic noise models were taken into consideration in [51], which however takes non-sampling approaches.

reads $\Omega(n \log^2 n)$, which comes with an extra log *n* factor gap to the optimality.

One interesting observation in Theorem 1 is that the sample complexity limit is proportional to $\frac{2^{d-2}}{d}$. This suggests that the amount of information that one hyperedge reveals on average is approximately $\frac{d}{2^{d-2}}$ bits. To understand why this is the case, consider a setting in which $\theta = 0$ and an hyperedge $E = \{i_1, i_2, \ldots, i_d\}$ is observed. The case of $Y_E = 1$ implies $X_{i_1} = X_{i_2} = \cdots = X_{i_d}$, in which there are only two uncertain cases (all zeros and all ones), i.e., the d-1 bits of information are revealed. On the other hand, the case of $Y_E = 0$ provides much less information as it rules out only two possible cases $(X_{i_1} = X_{i_2} = \cdots = X_{i_d} = 0$ and $X_{i_1} = X_{i_2} = \cdots = X_{i_d} = 1)$ out of 2^d possible candidates. This amounts to roughly $d \cdot \frac{2}{2^d}$ bits. Since $Y_E = 1$ occurs with probability $\frac{1}{2^{d-1}}$, the amount of information that one hyperedge can carry on average should read about $\frac{1}{2^{d-1}}(d-1) + (1 - \frac{1}{2^{d-1}}) \frac{d}{2^{d-2}} \approx \frac{d}{2^{d-2}}$. Relying on the connection to subspace clustering elaborated

Relying on the connection to subspace clustering elaborated in Sec. II-A, one can make an interesting implication from Theorem 1. The result offers a detailed guideline as to how to choose d for sample-efficient subspace clustering. In the case where the measurement quality reflected in θ is irrelevant of the number d of data points involved in a measurement, the limit increases in d. In practical applications, however, θ may depend on d. Actually, the quality of similarity measure can improve as more data points get involved, making θ decrease as d increases. In this case, choosing d as small as possible minimizes $\frac{2^{d-2}}{d}$ but may make θ too large. Hence, there might be a *sweet spot* on d that minimizes the sample complexity. It turns out this is indeed the case in practice. Actually we identify such optimal d* for motion segmentation application; see Sec. VII-A for details.

B. The Parity Measurement Case

Theorem 2. Let $d \ge 2$ be a constant. Under the parity measurement case, the following holds for any $\epsilon > 0$:

$$\begin{cases} \inf_{\psi} P_e(\psi) \to 0 \quad if \binom{n}{d} p \ge (1+\epsilon) \frac{1}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2}; \\ \inf_{\psi} P_e(\psi) \not\to 0 \quad if \binom{n}{d} p \le (1-\epsilon) \frac{1}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2}. \end{cases}$$

Proof: See Sec.V.

Notice that for a fixed θ and *n*, the minimum sample complexity is proportional to $\frac{1}{d}$, hence decreases in *d* unlike the homogeneity measurement case.

In view of the connection made in Sec. II-A, a natural question that arises in the context of channel coding is to ask how far the rate of the random LDGM code is from the capacity of the BSC channel. The connection can help immediately answer the question. We see from Theorem 2 that the rate of the LDGM code is

$$\frac{n}{p\binom{n}{d}} = \frac{d(\sqrt{1-\theta} - \sqrt{\theta})^2}{\log n}$$

This suggests that the code rate increases in d. Note that as long as d is constant, the rate vanishes, being far from the capacity of BSC channel $1 - H(\theta)$. On the other hand, it is not clear as to whether or not the random LDGM code can

achieve a non-vanishing code rate possibly by increasing the value of d. To check this, we explore the case where d can scale with n. By symmetry, it suffices to consider the case $2 \le d \le n/2$. Moreover, to avoid pathological cases where d fluctuates as n increases, we assume that d is a monotone function.

Theorem 3. Fix d, a monotone function of n such that $2 \le d \le n/2$, and $\epsilon > 0$. Under the parity measurement case,

• (upper bound) $\inf_{\psi} P_e(\psi) \to 0$ if

$$\binom{n}{d}p \ge (1+\epsilon)\frac{5/2}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2} \quad and \qquad (2)$$
$$\binom{n}{d}p \ge (1+\epsilon)5\log 2\frac{n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}; \qquad (3)$$

- $\binom{d}{p} \ge (1+\epsilon)5\log 2\frac{1}{(\sqrt{1-\theta}-\sqrt{\theta})^2};$
- (lower bound) $\inf_{\psi} P_e(\psi) \not\rightarrow 0$ if

$$\binom{n}{d}p \le (1-\epsilon)\frac{1}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2} \text{ or }$$
(4)

$$\binom{n}{d}p \le \frac{n}{1 - H(\theta)}.$$
(5)

Proof: See Sec. VI.

To see what these results mean, consider the two cases: $d = \Omega(\log n)$ and $d = o(\log n)$. In the case $d = \Omega(\log n)$, the theorem says that for a fixed θ ,

$$\inf_{\psi} P_e(\psi) \to 0 \text{ if } \binom{n}{d} p > \beta_1 n \text{ and}$$
$$\inf_{\psi} P_e(\psi) \not\to 0 \text{ if } \binom{n}{d} p < \beta_2 n,$$

where $\beta_1 = \max\left\{\frac{5/2\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2 d}, \frac{5\log 2}{(\sqrt{1-\theta}-\sqrt{\theta})^2}\right\} \approx 1$ and $\beta_2 = \max\left\{\frac{\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2 d}, \frac{1}{1-H(\theta)}\right\} \approx 1$. This suggests that as long as *d* grows asymptotically larger than $\log n$, we can achieve an order-wise tight sample complexity that is linear in *n*. On the other hand, in the case $d = o(\log n)$, the theorem asserts that

$$\inf_{\psi} P_e(\psi) \to 0 \text{ if } \binom{n}{d} p > \frac{5/2}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2} \text{ and}$$
$$\inf_{\psi} P_e(\psi) \not\to 0 \text{ if } \binom{n}{d} p < \frac{1}{d} \frac{n \log n}{(\sqrt{1-\theta} - \sqrt{\theta})^2}.$$

This implies that one cannot achieve the linear-order sample complexity if d grows slower than $\log n$. The implication of the above two can be formally stated as follows.

Corollary 1. For $d = o(\log n)$, reliable recovery is impossible with linear-order sample complexity, while it is possible for $d = \Omega(\log n)$.

From this, we see that the random LDGM code can achieve a constant rate as soon as $d = \Omega(\log n)$.

IV. PROOF OF THEOREM 1

The achievability and converse proofs are streamlined with the help of Lemmas 1 and 2, of which the proofs are left in Appendix A. For illustrative purpose, we focus on the noisy case $(\theta > 0)$ and assume that *n* is even. For a vector $\mathbf{V} := \{V_i\}_{1 \le i \le n} \in \{0, 1\}^n$, we define

$$\begin{cases} f_{\{i_1,i_2,\dots,i_d\}}(\mathbf{V}) & := f(V_{i_1}, V_{i_2},\dots, V_{i_d}); \\ \mathbf{F}(\mathbf{V}) & := \{f_E(\mathbf{V})\}_{E \in \mathcal{E}}; \\ d_{\mathsf{H}}(\mathbf{V}) & := \|\mathbf{Y} - \mathbf{F}(\mathbf{V})\|_1. \end{cases}$$
(6)

Let ψ_{ML} be the maximum likelihood (ML) decoder. One can easily verify that

$$\psi_{\mathrm{ML}}(\mathbf{Y}) = \arg\min_{\mathbf{V} \in \{0,1\}^n} \mathsf{d}_{\mathsf{H}}(\mathbf{V}),\tag{7}$$

where ties are randomly broken.

A. Achievability Proof

We intend to prove that

$$\max_{\mathbf{X}\in\{0,1\}^n} \Pr(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{1}\}) \to 0$$

under the claimed condition. Let $\mathbf{A} \in \{0, 1\}^n$ be the groundtruth vector. Without loss of generality, assume that the first k coordinates are 0's and the next n - k coordinates are 1's, where $0 \le k \le n/2$.

Let $\mathcal{A}_{i,j}$ denote the collection of all vectors whose coordinates are different from that of **A** in *i* many positions among the first *k* coordinates and in *j* many positions among the next n-k coordinates. Note that $\mathcal{A}_{0,0} = \{\mathbf{A}\}$ and $\mathcal{A}_{k,n-k} = \{\mathbf{A} \oplus \mathbf{1}\}$. Thus, a decoding algorithm ψ is successful if and only if the output $\psi(\mathbf{Y}) \in \mathcal{A}_{0,0} \cup \mathcal{A}_{k,n-k}$. Let $\mathcal{I} := \{(i, j) : (i, j) \notin \{(0, 0), (k, n-k)\}, 0 \le i \le k, \text{ and } 0 \le j \le n-k\}$. We also define

$$\mathbf{V}_{i,j} := (\underbrace{1,\ldots,1}_{i}, 0, \ldots, 0, \underbrace{0,\ldots,0}_{j}, 1, \ldots, 1), \underbrace{1}_{k}, \underbrace{1}_{n-k}$$

which is a representative vector of $A_{i,j}$.

Using these notations and the union bound, we get:

$$\begin{aligned} &\Pr(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{1}\} \mid \mathbf{X} = \mathbf{A}) \\ &\stackrel{(a)}{\leq} \Pr\left(\bigcup_{(i,j)\in\mathcal{I}} \bigcup_{\mathbf{V}\in\mathcal{A}_{i,j}} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{V}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right]\right) \\ &\leq \sum_{(i,j)\in\mathcal{I}} \sum_{\mathbf{V}\in\mathcal{A}_{i,j}} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{V}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right) \\ &= \sum_{(i,j)\in\mathcal{I}} \binom{k}{i} \binom{n-k}{j} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{V}_{i,j}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right), \quad (8) \end{aligned}$$

where the step (*a*) follows from the fact that the ML decoder outputs $V \notin \{A, A \oplus 1\}$ if $d_H(V) \le d_H(A)$.

To compare $d_{H}(V_{i,j})$ with $d_{H}(A)$, we define the set of *distinctive* hyperedges, i.e., the set of hyperedges such that $f_{E}(A) \neq f_{E}(V_{i,j})$:

$$\mathcal{F}_{i,j} := \left\{ E \in \binom{[n]}{d} : f_E(\mathbf{A}) \neq f_E(\mathbf{V}_{i,j}) \right\}$$
(9)

and $\mathcal{E}_{i,j} := \mathcal{E} \cap \mathcal{F}_{i,j}$. By definition, for $E \in \mathcal{E}_{i,j}$, $Y_E = f_E(\mathbf{A})$ if $Z_E = 0$; $Y_E = f_E(\mathbf{V}_{i,j})$ otherwise. Hence,

 $d_{\mathsf{H}}(\mathbf{V}_{i,j}) \leq d_{\mathsf{H}}(\mathbf{A})$ if and only if $\sum_{E \in \mathcal{E}_{i,j}} Z_E \geq \frac{|\mathcal{E}_{i,j}|}{2}$. This leads to:

where (a) is due to Hoeffding [53]. By letting $p' := (1 - e^{-D(0.5 ||\theta)})p$ and applying this to (8), we get:

$$\Pr(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{I}\} \mid \mathbf{X} = \mathbf{A})$$

$$\leq \sum_{(i,j)\in\mathcal{I}} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|}.$$
(12)

To give a tight upper bound on (12), one needs a tight lower bound on the size of the set of distinctive hyperedges, i.e., $|\mathcal{F}_{i,j}|$. It turns out that bounding $|\mathcal{F}_{i,j}|$ when d > 2requires non-trivial combinatorial counting. Note that this was not the case when d = 2, since $|\mathcal{F}_{i,j}|$ can be exactly computed via simple counting. Indeed, one of our main technical contributions lies in the derivation of tight bounds on $|\mathcal{F}_{i,j}|$, which we detail below.

Fact 1. The number of distinctive hyperedges can be calculated as follows:

$$|\mathcal{F}_{i,j}| = \sum_{\ell=1}^{d-1} {i \choose \ell} {k-i \choose d-\ell} + \sum_{\ell=1}^{d-1} {j \choose \ell} {n-k-j \choose d-\ell} + \sum_{\ell=1}^{d-1} {i \choose \ell} {n-k-j \choose d-\ell} + \sum_{\ell=1}^{d-1} {k-i \choose \ell} {j \choose d-\ell}.$$
 (13)

Proof: Consider a hyperedge $E = \{i_1, i_2, \ldots, i_d\}$ such that $f_E(\mathbf{A}) = 1$. That is, the hyperedge is connected only to a subset of the first k nodes or only to a subset of the last n - k nodes. That is, $\{i_1, i_2, \ldots, i_d\} \subset \{1, 2, \ldots, k\}$ or $\{i_1, i_2, \ldots, i_d\} \subset \{k + 1, k + 2, \ldots, n\}$. Consider the first case, i.e., $\{i_1, i_2, \ldots, i_d\} \subset \{1, 2, \ldots, k\}$. In order for this hyperedge to be distinctive, i.e., $f_E(\mathbf{V}_{i,j}) = 0$, at least one element of *E* must be in $\{1, 2, \ldots, k\}$. Thus, the total number of such distinctive hyperedges is $\sum_{\ell=1}^{d-1} {i \choose \ell} {k-i \choose d-\ell}$. Similarly, one can count the number of distinctive hyperedges for the case $\{i_1, i_2, \ldots, i_d\} \subset \{k + 1, k + 2, \ldots, n\}$: $\sum_{\ell=1}^{d-1} {i \choose \ell} {n-k-j \choose d-\ell}$. By considering the opposite case where $f_E(\mathbf{A}) = 0$ and $f_E(\mathbf{V}_{i,j}) = 1$, one can also obtain the remaining two terms, proving the statement.

By symmetry, we see that $|\mathcal{F}_{i,j}| = |\mathcal{F}_{k-i,n-k-j}|$. Hence,

$$\begin{split} &\sum_{(i,j)\in\mathcal{I}} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|} \\ &\leq \sum_{(i,j)\in\mathcal{I}, \ j\leq\lfloor\frac{n-k}{2}\rfloor} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|} \\ &+ \sum_{(i,j)\in\mathcal{I}, \ j\geq\lfloor\frac{n-k}{2}\rfloor} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|} \\ &= \sum_{(i,j)\in\mathcal{I}, \ j\leq\lfloor\frac{n-k}{2}\rfloor} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|} \\ &+ \sum_{(i,j)\in\mathcal{I}, \ j\leq\lfloor\frac{n-k}{2}\rfloor} \binom{k}{k-i} \binom{n-k}{n-k-j} (1-p')^{|\mathcal{F}_{k-i,n-k-j}|} \\ &= 2 \sum_{(i,j)\in\mathcal{I}, \ j\leq\lfloor\frac{n-k}{2}\rfloor} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|} =: 2\Phi. \end{split}$$

In order to show $\Phi \to 0$, for a fixed constant $\delta > 0$, we define the following index sets: $\mathcal{I}_{\text{big}} := \{(i, j) \in \mathcal{I} : [j \leq \frac{n-k}{2}] \cap ([i \geq \delta n] \cup [j \geq \delta n])\}$ and $\mathcal{I}_{\text{small}} := \{(i, j) \in \mathcal{I} : [j \leq \frac{n-k}{2}] \cap ([i < \delta n] \cap [j < \delta n])\}$. Then, using the index sets, one can express Φ as $\Phi_{\text{big}} + \Phi_{\text{small}}$, where

$$\Phi_{\text{big}} := \sum_{(i,j)\in\mathcal{I}_{\text{big}}} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|}$$

and

$$\Phi_{\text{small}} := \sum_{(i,j)\in\mathcal{I}_{\text{small}}} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|}.$$

Let us first consider Φ_{big} . Without loss of generality, assume $i \ge \delta n$. Then it follows from Fact 1 that

$$\begin{aligned} |\mathcal{F}_{i,j}| &\geq \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{n-k-j}{d-\ell} \stackrel{(a)}{\geq} \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{n/4}{d-\ell} \\ &\geq \binom{i}{1} \binom{n/4}{d-1} \geq \delta n \binom{n/4}{d-1} = \Omega(n^d), \end{aligned}$$

where (a) follows from the hypothesis that $j \leq \frac{n-k}{2}$ and $k \leq \frac{n}{2}$. Then it is easy to show that $\Phi_{\text{big}} \rightarrow 0$:

$$\Phi_{\text{big}} \leq \sum_{(i,j)\in\mathcal{I}} \binom{k}{i} \binom{n-k}{j} e^{-p'\Omega(n^d)}$$

$$\stackrel{(a)}{=} e^{-\Omega(n\log n)} \sum_{(i,j)\in\mathcal{I}} \binom{k}{i} \binom{n-k}{j} \leq e^{-\Omega(n\log n)} 2^n \to 0,$$

where (a) follows from the fact that $p'\Omega(n^d) \asymp p\binom{n}{d} = \Omega(n \log n)$.

Now we consider Φ_{small} . The following lemma gives a tight lower bound on $|\mathcal{F}_{i,j}|$ for this case:

Lemma 1. For $i < \delta n$ and $j < \delta n$,

$$|\mathcal{F}_{i,j}| \ge (i+j) \cdot \frac{(1-2\delta)^{d-1}}{2^{d-2}} \binom{n-1}{d-1}.$$

Proof: See Sec. A-A.

$$\sum_{\substack{(i,j)\in\mathcal{I}_{\text{small}}\\\leq}} \binom{k}{i} \binom{n-k}{j} (1-p')^{|\mathcal{F}_{i,j}|}$$

$$\stackrel{(a)}{\leq} \sum_{\substack{(i,j)\in\mathcal{I}_{\text{small}}\\\leq}} n^{i} n^{j} e^{-p'(i+j)\cdot\frac{(1-2\delta)^{d-1}}{2^{d-2}}\binom{n-1}{d-1}}}{\sum_{\substack{(i,j)\in\mathcal{I}_{\text{small}}\\\leq}}} \exp\left((i+j)\left\{\log n - \frac{p'(1-2\delta)^{d-1}\binom{n-1}{d-1}}{2^{d-2}}\right\}\right),$$
(14)

where (a) follows due to $\binom{k}{i} \leq n^i$, $\binom{n-k}{j} \leq n^j$ and Lemma 1. A straightforward computation yields $(1 - e^{-\mathsf{D}_{\mathsf{KL}}(0.5 \parallel \theta)}) = (\sqrt{1 - \theta} - \sqrt{\theta})^2$, so the claimed condition

$$\binom{n}{d}p \ge (1+\epsilon)\frac{2^{d-2}}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}$$

becomes

$$\binom{n}{d}p' \ge (1+\epsilon)\frac{2^{d-2}}{d}n\log n .$$
(15)

Under this claimed condition, we get:

$$\frac{p'(1-2\delta)^{d-1}\binom{n-1}{d-1}}{2^{d-2}} = \frac{p'(1-2\delta)^{d-1}\binom{n}{d}\frac{d}{n}}{2^{d-2}}$$
$$\stackrel{(a)}{\geq} (1+\epsilon)(1-2\delta)^{d-1}\log n$$
$$\stackrel{(b)}{\geq} (1+\epsilon/2)\log n,$$

where (a) follows from (15); (b) follows by choosing δ sufficiently small $((1-2\delta)^{d-1} \rightarrow 0 \text{ as } \delta \rightarrow 0)$. Thus, the RHS of (14) converges to 0 as *n* tends to infinity. This completes the proof.

B. Converse Proof

Let $V_{1/2}$ be the collection of *n*-dimensional vectors, each consisting of n/2 number of 0's and n/2 number of 1's. Moreover, let $\mathbf{X}_{1/2}$ be the random vector sampled uniformly at random over $V_{1/2}$. For any scheme ψ , by definition of $P_e(\psi)$, we see that

$$\Pr\left(\psi(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{1}\} \mid \mathbf{X} = \mathbf{X}_{1/2}\right) \le P_e(\psi)$$

and hence

$$\inf_{\psi} \Pr\left(\psi(\mathbf{Y}) \notin \{\mathbf{X}, \ \mathbf{X} \oplus \mathbf{1}\} \mid \mathbf{X} = \mathbf{X}_{1/2}\right) \leq \inf_{\psi} P_e(\psi).$$

Relying on this inequality, our proof strategy is to show that the left hand side is strictly bounded away from 0. Note that the infimum in the left hand side is achieved by $\psi_{ML,1/2}$:

$$\psi_{\mathrm{ML},1/2}(\mathbf{Y}) = \arg\min_{\mathbf{V}\in\mathcal{V}_{1/2}} \mathsf{d}_{\mathsf{H}}(\mathbf{V}) \,.$$

By letting $\mathbf{A} = (\underbrace{0, \dots, 0}_{n/2}, \underbrace{1, \dots, 1}_{n/2})$, we obtain:

$$\begin{aligned} &\Pr\left(\psi_{\mathrm{ML},1/2}(\mathbf{Y})\notin\{\mathbf{X},\ \mathbf{X}\oplus\mathbf{1}\}\mid\mathbf{X}=\mathbf{X}_{1/2}\right)\\ &=\Pr\left(\psi_{\mathrm{ML},1/2}(\mathbf{Y})\notin\{\mathbf{A},\ \mathbf{A}\oplus\mathbf{1}\}\mid\mathbf{X}=\mathbf{A}\right).\end{aligned}$$

Let *S* be the success event:

$$S := \bigcap_{\mathbf{V} \in \mathcal{V}_{1/2} \setminus \{\mathbf{A}, \mathbf{A} \oplus \mathbf{I}\}} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{V}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \right]$$

One can show that $\Pr(\psi_{ML,1/2}(\mathbf{Y}) \notin {\mathbf{A}, \mathbf{A} \oplus \mathbf{I}} | \mathbf{X} = \mathbf{A}) \ge \frac{1}{3} \Pr(S^c)$. This is due to the fact that given S^c , there are more than two candidates for $\arg\min_{\mathbf{V} \in \mathcal{V}_{1/2}} \mathbf{d}_{\mathsf{H}}(\mathbf{V})$, so

$$\Pr\left(\psi_{\mathrm{ML},1/2}(\mathbf{Y})\notin\{\mathbf{A},\ \mathbf{A}\oplus\mathbf{1}\}\mid\mathbf{X}=\mathbf{A},\ S^{c}\right)\geq\frac{1}{3}.$$

Hence, it suffices to show $Pr(S) \rightarrow 0$. To give a tight upper bound on Pr(S), we construct a subset of nodes such that any two nodes in the subset do not share the same hyperedge. To this end, we use the deletion technique (alteration technique) [54]. We first choose a big subset

$$\mathcal{R}_{\text{big}} = \{1, 2, \dots, r\} \bigcup \left\{\frac{n}{2} + 1, \frac{n}{2} + 2, \dots, \frac{n}{2} + r\right\},\$$

where $r = \lceil \frac{n}{\log^7 n} \rceil$; then erase every node in \mathcal{R}_{big} which shares hyperedges with other nodes in \mathcal{R}_{big} to obtain \mathcal{R}_{res} . The following lemma guarantees that \mathcal{R}_{res} has a comparable size as that of \mathcal{R}_{big} with high probability. For the later usage, we allow *d* to scale with *n*.

Lemma 2. Suppose $\binom{n}{d}p = O(n \log n)$ and $d = O(\log n)$. Let \mathcal{R}_{big} be a subset of [n] and \mathcal{R}_{res} be a subset obtained from \mathcal{R}_{big} by deleting every node which shares hyperedges with other nodes in \mathcal{R}_{big} . If $|\mathcal{R}_{big}| = O(n/\log^7 n)$, then with probability approaching 1,

$$|\mathcal{R}_{res}| = (1 - o(1))|\mathcal{R}_{big}|.$$

Proof: See Sec. A-B.

Let Δ be the event that $|\mathcal{R}_{\text{res}}| \geq (1 - o(1))|\mathcal{R}_{\text{big}}|$. Given the event Δ , both $\{1, 2, ..., n/2\} \cap \mathcal{R}_{\text{res}}$ and $\{\frac{n}{2} + 1, \frac{n}{2} + 2, ..., n\} \cap \mathcal{R}_{\text{res}}$ contain more than r/2 elements. We collect r/2 elements from each of these sets and denote by $\{b_1, b_2, ..., b_{r/2}\}$ and $\{c_1, c_2, ..., c_{r/2}\}$, respectively. Suppose that there exist (k, ℓ) such that $\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_k}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})$ and $\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{c_\ell}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})$. Conditioning on Δ , there are no hyperedges that contain both b_k and c_ℓ , so $\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{c_\ell}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})$. Hence conditioning on Δ ,

$$S \subset \bigcap_{k=1}^{r/2} [\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_k}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A})] \bigcup \bigcap_{k=1}^{r/2} [\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{c_k}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A})]$$

=: S'.

Since the event Δ occurs with probability approaching 1 and $S \subset S'$, $\Pr(S) \simeq \Pr(S \mid \Delta) \le \Pr(S' \mid \Delta)$. Hence,

$$\Pr(S) \lesssim \Pr\left(S' \mid \Delta\right)$$

$$\leq 2 \Pr\left(\bigcap_{k=1}^{r/2} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_k}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right] \mid \Delta\right)$$

$$\stackrel{(a)}{=} 2 \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta\right)^{r/2},$$

where (a) follows from the fact that the events $\{[\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_k}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A})]\}_{1 \le k \le r/2}$ are mutually independent conditioned on Δ . Let $p' = (1 - e^{-\mathsf{D}_{\mathsf{KL}}(0.5 \parallel \theta)})p$ as in the achievability proof. We intend to give an upper bound

on $\Pr(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta)$, i.e., a lower bound on $\Pr(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) \le \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta)$. Recall from the proof of achievability (see (11)) that

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{V}_{i,j}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right) \leq (1 - (1 - e^{-\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)})p)^{|\mathcal{F}_{i,j}|}.$$

For the case of $\mathbf{V}_{i,j} = \mathbf{A} \oplus \mathbf{e}_{b_1}$, $|\mathcal{F}_{i,j}| = \binom{n/2-1}{d-1} + \binom{n/2}{d-1}$ (note that k = n/2, i = 1, j = 0). So we get:

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) \le \mathsf{d}_{\mathsf{H}}(\mathbf{A})\right) \le e^{-p'\left\binom{n/2-1}{d-1} + \binom{n/2}{d-1}\right)}.$$
 (16)

On the other hand, what we need for the converse proof is a lower bound. In what follows, we will show that (16) is tight enough, more precisely,

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) \le \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta\right) \ge (1 - o(1))e^{-2p'\binom{n/2-1}{d-1}}.$$
(17)

What this means at a high level is that Chernoff-Hoeffding is tight enough. Let us condition on the event Δ for the time being. As in (9), we define the following sets:

$$\mathcal{F}_{b_1} := \left\{ E \in \binom{[n]}{d} : f_E(\mathbf{A}) \neq f_E(\mathbf{A} \oplus \mathbf{e}_{b_1}) \right\}$$

and $\mathcal{E}_{b_1} := \mathcal{E} \cap \mathcal{F}_{b_1}$. By definition, for $E \in \mathcal{E}_{b_1}$, $Y_E = f_E(\mathbf{A})$ if $Z_E = 0$; $Y_E = f_E(\mathbf{A} \oplus \mathbf{e}_{b_1})$ otherwise. We see that

$$\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \Leftrightarrow \sum_{E \in \mathcal{E}_{b_1}} Z_E \geq \frac{|\mathcal{E}_{b_1}|}{2}$$

Now we want to manipulate $\Pr(\mathbf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_1}) \leq \mathbf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta)$ as we did in (10). However, here we need to give a careful attention to the range of summation as \mathcal{E}_{b_1} cannot be equal to \mathcal{F}_{b_1} . This is because every hyperedge in \mathcal{E}_{b_1} should intersect \mathcal{R}_{big} at exactly one node, which is b_1 . Indeed, for any hyperedge E in \mathcal{E}_{b_1} , b_1 should be in E to satisfy $f_E(\mathbf{A}) \neq f_E(\mathbf{A} \oplus \mathbf{e}_{b_1})$, while if another node from \mathcal{R}_{big} is contained in E, it contradicts the fact that $b_1 \in \mathcal{R}_{\text{res}}$. This implies that \mathcal{E}_{b_1} is always contained in a proper subset \mathcal{G}_{b_1} of \mathcal{F}_{b_1} defined as:

$$\mathcal{G}_{b_1} := \mathcal{F}_{b_1} \setminus \left\{ E \in \binom{[n]}{d} : |E \cap \mathcal{R}_{\text{big}}| \ge 2 \right\}.$$
(18)

Now a manipulation similar to (10) yields:

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_{1}}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta\right)$$

=
$$\sum_{\ell=1}^{|\mathcal{G}_{b_{1}}|} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_{1}}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid |\mathcal{E}_{b_{1}}| = \ell, \Delta\right)$$
$$\cdot \Pr(|\mathcal{E}_{b_{1}}| = \ell \mid \Delta).$$

Since the event Δ is related to the occurrence of edges in

$$\left\{ E \in \binom{[n]}{d} : |E \cap \mathcal{R}_{\text{big}}| \ge 2 \right\}$$

and \mathcal{E}_{b_1} is subject to (18), Δ and $[|\mathcal{E}_{b_1}| = \ell]$ are independent. and thus $(21) = o(1) \cdot (20)$. Thus, we get:

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_{1}}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta\right)$$

= $\sum_{\ell=1}^{|\mathcal{G}_{b_{1}}|} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_{1}}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid |\mathcal{E}_{b_{1}}| = \ell, \ \Delta\right) \Pr(|\mathcal{E}_{b_{1}}| = \ell)$
= $\sum_{\ell=1}^{|\mathcal{G}_{b_{1}}|} \Pr\left(\sum_{E \in \mathcal{E}_{b_{1}}} Z_{E} \geq \frac{\ell}{2} \mid |\mathcal{E}_{b_{1}}| = \ell\right) \binom{|\mathcal{G}_{b_{1}}|}{\ell} \frac{p^{\ell}}{(1-p)^{\ell-|\mathcal{G}_{b_{1}}|}}.$
(19)

By the reverse Chernoff-Hoeffding bound [53], for a fixed $\delta > 0$, there exists $n_{\delta} > 0$ such that

$$\Pr\left(\sum_{E \in \mathcal{E}_{b_1}} Z_E \ge \frac{\ell}{2} \middle| |\mathcal{E}_{b_1}| = \ell\right) \ge e^{-(1+\delta)\ell \mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}$$

for all $\ell \ge n_{\delta}$. Let g_n be a sequence (to be determined) such that $g_n \to \infty$ as $n \to \infty$. For a sufficiently large n,

$$(19) \geq \sum_{\ell=1}^{|\mathcal{G}_{b_1}|} {|\mathcal{G}_{b_1}| \choose \ell} \frac{(e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)} p)^{\ell}}{(1-p)^{\ell-|\mathcal{G}_{b_1}|}} \qquad (20)$$
$$-\sum_{\ell=1}^{g_n-1} {|\mathcal{G}_{b_1}| \choose \ell} \frac{(e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)} p)^{\ell}}{(1-p)^{\ell-|\mathcal{G}_{b_1}|}}. \qquad (21)$$

Actually one can choose g_n so that (21) is negligible compared to (20). To see this, we consider:

$$\frac{(21)}{(20)} \leq \frac{(1-p)^{|\mathcal{G}_{b_1}|} \sum_{\ell=1}^{g_n-1} \left(|\mathcal{G}_{b_1}| \frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \right)^{\ell}}{(1-p)^{|\mathcal{G}_{b_1}|} \sum_{\ell=1}^{|\mathcal{G}_{b_1}|} \binom{|\mathcal{G}_{b_1}|}{\ell} \left(\frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \right)^{\ell}}{\left(1+\frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \right)^{|\mathcal{G}_{b_1}|}} \\ \stackrel{(a)}{=} \frac{\sum_{\ell=1}^{g_n-1} \left(|\mathcal{G}_{b_1}| \frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \right)^{|\mathcal{G}_{b_1}|}}{(1+o(1))\exp\left(|\mathcal{G}_{b_1}| \frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \right)}{1-p}} \\ =: \frac{\sum_{\ell=1}^{g_n-1} q^{\ell}}{(1+o(1))e^{q}},$$
(22)

where (a) follows from the fact that $\lim_{x\to 0+} \frac{1+x}{e^x} = 1$, and the last equation is due to the following definition: $q := |\mathcal{G}_{b_1}| \frac{pe^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p}.$ One can easily verify that $|\mathcal{F}_{b_1}| = \binom{n/2-1}{d-1} + \binom{n/2}{d-1}$ and $|\mathcal{G}_{b_1}| = \binom{n/2-1-r}{d-1} + \binom{n/2-r}{d-1}.$ Since r = o(n), $\lim_{n\to\infty} |\mathcal{G}_{b_1}|/|\mathcal{F}_{b_1}| \to 1$. Thus,

$$q = |\mathcal{G}_{b_1}| \frac{p e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p}$$
(23)

$$\asymp |\mathcal{F}_{b_1}| \frac{p e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}}{1-p} \asymp n^{d-1} p = \Omega(\log n) \,. \tag{24}$$

Therefore, if one chooses $g_n = \lfloor \log q \rfloor$,

$$\frac{(21)}{(20)} = \frac{\sum_{\ell=1}^{g_n-1} q^{\ell}}{e^q} \le \frac{g_n q^{g_n}}{e^q} \le \frac{\log q \cdot q^{\log q}}{e^q} \le \frac{\log q \cdot q^{\log q}}{e^q} \to 0,$$

Hence, we get:

$$(19) = (20) - (21)$$

$$\geq (1 - o(1)) \sum_{\ell=1}^{|\mathcal{G}_{b_1}|} {|\mathcal{G}_{b_1}| \choose \ell} \frac{(e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}p)^{\ell}}{(1 - p)^{\ell - |\mathcal{G}_{b_1}|}}$$

$$= (1 - o(1)) \left(1 - (1 - e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)})p\right)^{|\mathcal{G}_{b_1}|}$$

$$\stackrel{(a)}{\geq} (1 - o(1)) \left(1 - (1 - e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)})p\right)^{2\binom{n/2}{d-1}}$$

$$\stackrel{(b)}{\equiv} (1 - o(1)) \exp\left(-2\binom{n/2}{d-1}(1 - e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)})p\right),$$

where (a) follows since $|\mathcal{G}_{b_1}| \leq |\mathcal{F}_{b_1}| \leq 2\binom{n/2}{d-1}$; (b) follows from the fact that $\lim_{x\to 0+} \frac{1+x}{e^x} = 1$. As $\delta > 0$ can be chosen arbitrarily small, the term $e^{-(1+\delta)\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}$ can be made arbitrarily close to $e^{-\mathsf{D}_{\mathsf{KL}}(0.5\|\theta)}$, which in turn ensures that the last term is essentially equal to

$$(1-o(1))e^{-2p'\binom{n/2}{d-1}}.$$

Applying this to the previous upper bound on Pr(S), we get:

$$\begin{aligned} \Pr(S) &\leq \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{e}_{b_{1}}) > \mathsf{d}_{\mathsf{H}}(\mathbf{A}) \mid \Delta\right)^{r/2} \\ &\leq \left(1 - (1 - o(1))e^{-2p'\binom{n/2}{d-1}}\right)^{r/2} \\ &\leq \exp\left(-(1 - o(1))\frac{r}{2}e^{-2p'\binom{n/2}{d-1}}\right) \\ &= \exp\left(-(1 - o(1))\frac{n}{2\log^{7}n}e^{-(1 + o(1))\cdot\frac{p'd\binom{n}{d}}{2d-2n}}\right), \end{aligned}$$

where the last equality follows from the fact that

$$\lim_{n \to \infty} \frac{2p'\binom{n/2}{d-1}}{p'd\binom{n}{d}/2^{d-2}n} \to 1 \text{ and } r = \left\lceil \frac{n}{\log^7 n} \right\rceil.$$

The last term converges to 0 as $p' \leq (1-\epsilon)\frac{2^{d-2}}{d}\frac{n\log n}{\binom{n}{2}}$.

V. PROOF OF THEOREM 2

In this section, we prove a similar statement for the parity measurement case.

A. Achievability Proof

Note that the parity measurement is symmetric in a sense that for any vector A, we have

$$\Pr(\psi_{ML}(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{1}\} \mid \mathbf{X} = \mathbf{A})$$
$$= \Pr(\psi_{MI}(\mathbf{Y}) \notin \{\mathbf{0}, \mathbf{1}\} \mid \mathbf{X} = \mathbf{0}).$$

Indeed, this follows from a simple observation that for any vector **B**, there is a trivial coupling between the conditional distribution $\{Y|X = B\}$ and $\{Y^{\oplus A}|X = B \oplus A\}$, where $Y^{\oplus A} =$ $\{(Y^{\oplus \mathbf{A}})_E\}_{E \in \mathcal{E}}$ is defined as

$$(Y^{\oplus \mathbf{A}})_E := Y_E \oplus \bigoplus_{i \in E} \mathbf{A}_i$$

for each $E \in \mathcal{E}$.

By this symmetry, we have

$$\max_{\mathbf{X} \in \{0,1\}^n} \Pr(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{X}, \mathbf{X} \oplus \mathbf{1}\})$$

= $\Pr(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{0}, \mathbf{1}\} \mid \mathbf{X} = \mathbf{0}) ,$

and hence, it suffices to prove that

Pr (
$$\psi_{ML}$$
(**Y**) ∉ {**0**, **1**} | **X** = **0**) → 0.

Conditioning on $\mathbf{X} = \mathbf{0}$,

$$\begin{aligned} &\Pr\left(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{0}, \mathbf{1}\}\right) \\ &\leq \Pr\left(\bigcup_{A \neq \mathbf{0}, \mathbf{1}} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{A}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right]\right) \\ &= \Pr\left(\bigcup_{k=1}^{n-1} \bigcup_{\|\mathbf{A}\|_{1}=k} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{A}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right]\right) \\ &\leq \sum_{k=1}^{n-1} \sum_{\|\mathbf{A}\|_{1}=k} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right) \\ &\stackrel{(a)}{=} 2 \cdot \sum_{k=1}^{n/2} \sum_{\|\mathbf{A}\|_{1}=k} \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{A}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right) \\ &\stackrel{(b)}{=} 2 \cdot \sum_{k=1}^{n/2} \binom{n}{k} \Pr\left(\mathsf{d}_{\mathsf{H}}\left(\sum_{i=1}^{k} \mathbf{e}_{i}\right) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right), \end{aligned} \tag{25}$$

where (a) follows form the fact that $\Pr(\mathsf{d}_{\mathsf{H}}(\mathbf{A}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})) = \Pr(\mathsf{d}_{\mathsf{H}}(\mathbf{A} \oplus \mathbf{1}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0})); (b)$ follows due to symmetry. To compare $\mathsf{d}_{\mathsf{H}}\left(\sum_{i=1}^{k} \mathbf{e}_{i}\right)$ and $\mathsf{d}_{\mathsf{H}}(\mathbf{0})$, we define

$$\mathcal{F}_k := \left\{ E \in \binom{[n]}{d} : f_E(\mathbf{0}) \neq f_E\left(\sum_{i=1}^k \mathbf{e}_i\right) \right\}$$

and $\mathcal{E}_k := \mathcal{E} \cap \mathcal{F}_k$. As in (11), we obtain

$$\Pr\left(\mathsf{d}_{\mathsf{H}}\left(\sum_{i=1}^{k} \mathbf{e}_{i}\right) \le \mathsf{d}_{\mathsf{H}}(\mathbf{0})\right)$$
$$\le (1 - (1 - e^{-\mathsf{D}_{\mathsf{KL}}(0.5||\theta)})p)^{|\mathcal{F}_{k}|} = (1 - p')^{|\mathcal{F}_{k}|}$$

yielding

$$\frac{1}{2} \cdot (25) \le \sum_{k=1}^{n/2} \binom{n}{k} (1-p')^{|\mathcal{F}_k|}.$$
(26)

We again count $|\mathcal{F}_k|$ in an effort to obtain a tight upper bound on (26). Notice that $E \in \mathcal{F}_k$ if $|E \cap [k]|$ is odd, and hence

$$|\mathcal{F}_k| = \sum_{\substack{i \le d\\ i \text{ is odd}}} \binom{k}{i} \cdot \binom{n-k}{d-i}.$$
(27)

Let $\delta > 0$ be a small constant that will be determined later. For the case $k \ge \delta n$, it follows that

$$|\mathcal{F}_k| \ge \binom{k}{1}\binom{n-k}{d-1} \ge \delta n\binom{n/2}{d-1} = \Omega(n^d).$$

Then it is easy to show $(26) \rightarrow 0$ for this case:

$$\sum_{k=\delta n}^{n/2} \binom{n}{k} (1-p')^{|\mathcal{F}_k|} \le \sum_{k=\delta n}^{n/2} \binom{n}{k} e^{-p'\Omega(n^d)}$$
$$\stackrel{(a)}{=} e^{-\Omega(n\log n)} \sum_{k=\delta n}^{n/2} \binom{n}{k} \le e^{-\Omega(n\log n)} 2^n \to 0$$

where (a) follows from the fact that $p'\Omega(n^d) \simeq p\binom{n}{d} = \Omega(n \log n)$. For the case $k < \delta n$, we see that

$$|\mathcal{F}_k| \ge \binom{k}{1} \binom{n-k}{d-1} \ge k \binom{(1-\delta)n}{d-1}$$
$$\stackrel{(a)}{\underset{n\to\infty}{\equiv}} (1+o(1))k(1-\delta)^{d-1} \binom{n-1}{d-1}, \qquad (28)$$

where (a) follows since

$$\lim_{n \to \infty} \frac{\alpha^{d-1} \binom{n-1}{d-1}}{\binom{an}{d-1}} = 1$$
(29)

holds for a fixed d and $\alpha \in (0, 1)$. Hence, we get

$$\sum_{k=1}^{on} \binom{n}{k} (1-p')^{|\mathcal{F}_k|} \leq \sum_{k=1}^{on} n^k e^{-(1+o(1))p'k(1-\delta)^{d-1}\binom{n}{d-1}} = \sum_{k=1}^{\delta n} e^{k \cdot \left\{ \log n - (1+o(1))p'(1-\delta)^{d-1}\binom{n}{d-1} \right\}}.$$
(30)

By choosing δ arbitrarily small, under the claimed condition, one can make

$$p'(1-\delta)^{d-1} \binom{n}{d-1} = (1+o(1))(1-\delta)^{d-1} \binom{n}{d} p' \frac{d}{n}$$

$$\geq (1+\epsilon/2) \log n ,$$

which implies that (30) converges to 0 as n tends to infinity.

B. Converse Proof

As the parity measurement is symmetric,

$$\inf_{\psi} P_e(\psi) = \Pr\left(\psi_{\mathrm{ML}}(\mathbf{Y}) \notin \{\mathbf{X}, \ \mathbf{X} \oplus \mathbf{1}\} \mid \mathbf{X} = \mathbf{0}\right) \,.$$

As before, we define the success event as:

$$S := \bigcap_{\mathbf{V} \neq \mathbf{0}, \mathbf{1}} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{V}) > \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \right]. \tag{31}$$

Again, it suffices to show that $Pr(S) \rightarrow 0$, and to this end, we construct a subset of nodes such that any two nodes in the subset do not share the same hyperedge. Unlike the previous case, the subset is now defined as:

$$\mathcal{R}_{\text{big}} := \{1, 2, \dots, r\}$$
 (32)

where $r = \lceil \frac{n}{\log^7 n} \rceil$, and we erase every node in \mathcal{R}_{big} which shares hyperedges with other nodes in \mathcal{R}_{big} to obtain \mathcal{R}_{res} . In view of Lemma 2, we have $|\mathcal{R}_{res}| \ge (1 - o(1))r$ almost surely; let Δ be such event. Conditioning on Δ , we enumerate r/2 many elements of \mathcal{R}_{res} by $b_1, \ldots, b_{r/2}$. As there are no hyperedges that connect two nodes in \mathcal{R}_{res} , the events $\{[d_H(e_{b_k}) > d_H(0)]\}_{1 \le k \le r/2}$ are mutually independent conditioned on Δ . Hence, we get:

$$\Pr(S) \lesssim \Pr(S \mid \Delta)$$

$$\leq \Pr\left(\bigcap_{k=1}^{r/2} \left[\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_{k}}) > \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \right] \mid \Delta \right)$$

$$= \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_{1}}) > \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \mid \Delta\right)^{r/2}.$$
(33)

Let $p' = (1 - e^{-D_{\mathsf{KL}}(0.5 \| \theta)})p$ as before. Using similar arguments used in the previous section, we have

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_1}) \le \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \mid \Delta\right) \ge (1 - o(1))e^{-p'\binom{n-1}{d-1}}.$$
 (34)

This gives:

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_{1}}) > \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \mid \Delta\right)^{r/2} \le \left(1 - (1 - o(1))e^{-p'\binom{n-1}{d-1}}\right)^{r/2} \le \exp\left(-(1 - o(1))\frac{r}{2}\exp\left\{-p'\binom{n-1}{d-1}\right\}\right) \le \exp\left(-(1 - o(1))\frac{n}{2\log^{7} n}\exp\left\{-(1 + o(1)) \cdot \frac{p'\binom{n}{d}d}{n}\right\}\right)$$

Notice that the last term converges to 0 as $\binom{n}{d}p' \leq (1 - \epsilon)\frac{n \log n}{d}$, which completes the proof.

VI. PROOF OF THEOREM 3

When d scales with n, a technical challenge arises, and we will focus on such technical difficulties, skipping most of the redundant parts.

A. Proof of the Upper Bound

From (26) and (27), we get

$$P_e(\psi_{\rm ML}) \le \sum_{k=1}^{n/2} {n \choose k} (1-p')^{N_k},$$
 (35)

where

$$N_k := \sum_{\substack{1 \le i \le d \\ i \text{ is odd}}} \binom{k}{i} \cdot \binom{n-k}{d-i}$$
(36)

and $p' := (\sqrt{1-\theta} - \sqrt{\theta})^2 p$. Let us focus on counting N_k . When $d \approx 1$, $\binom{n}{d} \approx \frac{n^d}{d!}$ suffices to obtain a proper bound on N_k . However, in the general case where *d* scales with *n*, one needs a more delicate bounding technique to obtain sharp results. The following lemma presents our new bound.

Lemma 3. Let
$$\beta := \lceil \frac{n-d+1}{2d+1} \rceil < n/2$$
 and $\alpha := \frac{n-d+1}{d}$. Then

$$\sum_{\substack{1 \le i \le d \\ i \text{ is odd}}} \binom{k}{i} \binom{n-k}{d-i} \ge \begin{cases} \frac{2k}{5\alpha} \binom{n}{d}, & k < \beta; \\ \frac{1}{5} \binom{n}{d}, & \beta \le k \le n/2. \end{cases}$$

Proof: See Sec. VI-C. The proof requires an involved combinatorial counting, which is one of our main technical contributions.

Employing Lemma 3, we get:

$$(35) \leq \sum_{k=1}^{\beta-1} \binom{n}{k} (1-p')^{N_k} + \sum_{k=\beta}^{n/2} \binom{n}{k} (1-p')^{N_k}$$

$$\leq \sum_{k=1}^{\beta-1} \binom{n}{k} (1-p')^{\frac{2k}{5\alpha}} \binom{n}{\alpha} + \sum_{k=\beta}^{n/2} \binom{n}{k} (1-p')^{\frac{1}{5}} \binom{n}{\alpha}$$

$$\stackrel{(a)}{\leq} \sum_{k=1}^{\beta-1} \left[n^k e^{-p' \frac{2k}{5\alpha}} \binom{n}{\alpha} \right] + 2^n e^{-\frac{1}{5}p'} \binom{n}{\alpha}$$

$$\leq \sum_{k=1}^{\beta-1} \left[\exp\left\{ k \left(\log n - \frac{2p'\binom{n}{d}}{5\alpha} \right) \right\} \right] \qquad (37)$$

$$+ \exp\left\{ n \log 2 - \frac{1}{5}p'\binom{n}{d} \right\}, \qquad (38)$$

where (a) follows from $\sum_{k=\beta}^{n/2} {n \choose k} \le \sum_{k=0}^{n} {n \choose k} \le 2^{n}$. Note that (38) vanishes due to (3). In order to show that (37) vanishes as well, we consider two cases: d = o(n) and $d \simeq n$. When d = o(n),

$$\sum_{k=1}^{\beta-1} \exp\left\{k\left(\log n - \frac{2p'\binom{n}{d}}{5\alpha}\right)\right\}$$
$$\leq \sum_{k=1}^{\beta-1} \exp\left\{k\left(\log n - \frac{2dp'\binom{n}{d}}{5n}\right)\right\}$$
$$\leq \frac{\exp\left(\log n - \frac{2dp'\binom{n}{d}}{5n}\right)}{1 - \exp\left(\log n - \frac{2dp'\binom{n}{d}}{5n}\right)} \to 0,$$

since
$$\log n - \frac{2dp'\binom{n}{d}}{5n} \to -\infty$$
.
If $d \simeq n$,

$$\sum_{k=1}^{\beta-1} \exp\left\{k\left(\log n - \frac{2p'\binom{n}{d}}{5\alpha}\right)\right\}$$
$$\leq \beta \max_{1 \leq k \leq \beta-1} \exp\left\{k\left(\log n - \frac{2p'\binom{n}{d}}{5\alpha}\right)\right\}$$
$$= \beta \exp\left(\log n - \frac{2p'\binom{n}{d}}{5\alpha}\right),$$

where the last equality holds since $\log n - \frac{2p'\binom{n}{d}}{5\alpha} < 0$, and hence k = 1 achieves the maximum value. Note that this vanishes since β is asymptotically bounded by a constant. Therefore, (37) always vanishes, completing the proof.

B. Proof of the Lower Bound

The lower bound statement can be rewritten as follows: $\inf_{\psi} P_e(\psi) \not\rightarrow 0$ if $\binom{n}{d}p \leq \max((1-\epsilon))^{\frac{1}{d}} \frac{n \log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}, \frac{n}{1-H(\theta)}$. Note that when $d = \omega(\log n)$, the condition reduces to $\binom{n}{d}p \leq \frac{n}{1-H(\theta)}$. Hence, it is sufficient to show the following two statements.

• If $d = O(\log n)$: $\inf_{\psi} P_e(\psi) \not\rightarrow 0$ if $\binom{n}{d}p \leq \max\left((1-\epsilon)\frac{1}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}, \frac{n}{1-H(\theta)}\right)$.

• If $d = \omega(\log n)$: $\inf_{\psi} P_e(\psi) \neq 0$ if $\binom{n}{d} p \leq \frac{n}{1 - H(\theta)}$.

We first show that $\binom{n}{d}p \leq \frac{n}{1-H(\theta)}$ implies $\inf_{\psi} P_e(\psi) \neq 0$ for all *d*. By rearranging terms, we have $\binom{n}{d}p \leq \frac{n}{1-H(\theta)} \Leftrightarrow \frac{n}{\binom{n}{d}p} \geq 1 - H(\theta)$. One can immediately observe that this implies $\inf_{\psi} P_e(\psi) \neq 0$ since $\frac{n}{\binom{n}{d}p}$ (which can be viewed as the rate of a code) cannot exceed the Shannon capacity of the channel $1 - H(\theta)$.

We now prove that $\binom{n}{d}p \leq (1-\epsilon)\frac{1}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}$ implies $\inf_{\psi} P_e(\psi) \not\rightarrow 0$ if $d = O(\log n)$. Further, we will focus on the case of $\binom{n}{d}p \approx \frac{n\log n}{d}$ since this is the regime where the largest amount of information is available. Again, it is enough to show that $\Pr(S) \rightarrow 0$, where S is defined as (31). By defining $\mathcal{R}_{\text{big}}, \mathcal{R}_{\text{res}}, \Delta$ and $b_1, \ldots, b_{r/2}$ as before, we again obtain (33):

$$\Pr(S) \le \Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_1}) > \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \mid \Delta\right)^{r/2} . \tag{39}$$

We finish the proof by showing the following for the considered case:

$$\Pr\left(\mathsf{d}_{\mathsf{H}}(\mathbf{e}_{b_1}) \leq \mathsf{d}_{\mathsf{H}}(\mathbf{0}) \mid \Delta\right) \geq (1 - o(1))e^{-2p'\binom{n-1}{d-1}}.$$

While following the proof of (17), the key technical difficulty arises when checking $q = \Omega(\log n)$ (see (24)): a simple calculation yields $|\mathcal{F}_{b_1}| = \binom{n-1}{d-1}$ and $|\mathcal{G}_{b_1}| = \binom{n-|\mathcal{R}_{\text{big}}|}{d-1}$, but here it is not clear whether $\binom{n-|\mathcal{R}_{\text{big}}|}{d-1} \approx \binom{n-1}{d-1}$ when *d* is not a constant. We resolve this using a careful estimation as follows. As $|\mathcal{R}_{\text{big}}| = \Theta(\frac{n}{\log^7 n})$ and $d = O(\log n)$, it is straightforward to verify

$$1 - \frac{1}{\log^2 n} \le \frac{n - |\mathcal{R}_{\text{big}}| - j}{n - 1 - j}$$

for $0 \le j \le d - 2$. This simple yet crucial inequality concludes:

$$\frac{\binom{n-|\mathcal{R}_{\text{big}}|}{d-1}}{\binom{n-1}{d-1}} = \prod_{j=0}^{d-2} \frac{n-|\mathcal{R}_{\text{big}}|-j}{n-1-j}$$
$$\geq \left(1-\frac{1}{\log^2 n}\right)^{d-1} \approx \exp\left\{-\frac{d-1}{\log^2 n}\right\} \to 1.$$

C. Proof of Lemma 3

Without loss of generality, we prove the lemma assuming that $d \le k$. The proofs for the other cases are similar. We wish to obtain lower bounds on

 $N_{k} = \sum_{\substack{1 \le i \le d \\ i \text{ is odd}}} \binom{k}{i} \binom{n-k}{d-i}$ $= \underbrace{\binom{k}{1}\binom{n-k}{d-1}}_{\text{boundary odd term}} + \underbrace{\sum_{i=3,5,\dots,d-3}\binom{k}{i}\binom{n-k}{d-i}}_{\text{intermediate odd terms}}$ $+ \underbrace{\binom{k}{d-1}\binom{n-k}{1}}_{\text{boundary odd term}}$

in terms of $\binom{n}{d}$. First, observe that

$$\binom{n}{d} = \sum_{0 \le i \le d} \binom{k}{i} \binom{n-k}{d-i}$$
$$= \underbrace{\binom{k}{0}\binom{n-k}{d}}_{\text{boundary term}} + \underbrace{\sum_{1 \le i \le d-1} \binom{k}{i} \binom{n-k}{d-i}}_{\text{intermediate terms}}$$
$$+ \underbrace{\binom{k}{d}\binom{n-k}{0}}_{\text{boundary term}}.$$

Suppose we have the following bounds for some quantities $A_1, A_2, A_3 > 0$.

Bound 1:
$$\underbrace{\binom{k}{0}\binom{n-k}{d} + \binom{k}{d}\binom{n-k}{0}}_{\text{sum of boundary terms}} \leq A_1\underbrace{\left[\binom{k}{1}\binom{n-k}{d-1} + \binom{k}{d-1}\binom{n-k}{1}\right]}_{\text{sum of boundary odd terms}}$$
Bound 2:
$$\underbrace{\sum_{i=1,2,\dots,d-2,d-1}\binom{k}{i}\binom{n-k}{d-i}}_{\text{intermediate terms}} \leq A_2 \cdot \underbrace{\sum_{i=3,5,\dots,d-3}\binom{k}{i}\binom{n-k}{d-i}}_{i} + A_3N_k.$$

intermediate odd terms

Then, by summing up the two inequalities, one can obtain a lower bound on N_k :

$$\binom{n}{d} \leq A_1 \cdot \left[\binom{k}{1}\binom{n-k}{d-1} + \binom{k}{d-1}\binom{n-k}{1}\right]$$
$$+ A_2 \cdot \sum_{i=3,5,\dots,d-3}\binom{k}{i}\binom{n-k}{d-i} + A_3 \cdot N_k$$
$$\leq \max(A_1, A_2) \cdot N_k + A_3 \cdot N_k$$
$$= (\max(A_1, A_2) + A_3) \cdot N_k .$$

Thus, the proof is completed as long as one can find appropriate quantities A_1 , A_2 and A_3 .

The following lemma asserts that $A_2 = 2$ and $A_3 = 3$ satisfy Bound 2.

Lemma 4. For $1 \le k \le n/2$,

$$\sum_{\substack{i=1,2,\dots,d-2,d-1\\i \le dd}} \binom{k}{i} \binom{n-k}{d-i}$$
$$\leq 2 \cdot \sum_{\substack{3 \le i \le d-3\\i \ge odd}} \binom{k}{i} \binom{n-k}{d-i} + 3N_k$$

Proof: See Sec. A-C. For Bound 1, the following lemma characterizes A_1 .



Fig. 3. Algorithm 1 achieves the optimal sample complexity. We run Monte Carlo simulations to estimate the probability of success when: (a) n = 1000, d = 4, and for various choices of θ ; (b) n = 1000, $\theta = 0.05$, and for various choices of d. For each curve, we normalize the number of samples by the respective information theoretic limits, characterized in Theorem 1. Observe that the probability of success quickly approaches 1 as the normalized sample complexity crosses 1.

Lemma 5. Let
$$\alpha := \frac{n-d+1}{d}$$
 and $\beta := \left\lceil \frac{n-d+1}{2d+1} \right\rceil$. For $\beta \le k \le n/2$,
 $\binom{k}{0}\binom{n-k}{d} + \binom{k}{d}\binom{n-k}{0}$
 $\le 2\left[\binom{k}{1}\binom{n-k}{d-1} + \binom{k}{d-1}\binom{n-k}{1}\right]$.

For $k < \beta$,

$$\binom{k}{0}\binom{n-k}{d} + \binom{k}{d}\binom{n-k}{0}$$
$$\leq \frac{a}{k}\left[\binom{k}{1}\binom{n-k}{d-1} + \binom{k}{d-1}\binom{n-k}{1}\right]$$

and $\frac{\alpha}{k} \geq 2$.

Proof: See Sec. A-D.

That is, $A_1 = 2$ if $\beta \le k \le n/2$, and $A_1 = \frac{\alpha}{k}$ if $k < \beta$.

We are now' ready to prove Lemma 3 with the help of Lemma 4, Lemma, 5 and (40). When $\beta \le k < n/2$,

$$\binom{n}{d} \le 5 \ N_k$$

When $k < \beta$,

$$\binom{n}{d} \leq \left(\max\left(2, \frac{\alpha}{k}\right) + 3\right) N_k \leq \frac{5\alpha}{2k} N_k,$$

where the last inequality holds since $\frac{a}{k} \ge 2$. This completes the proof.

VII. EXPERIMENTAL RESULTS

A. The Homogeneity Measurement Case

1) Efficient Algorithms: We also develop a computationally-efficient algorithm that achieves the information-theoretic limit characterized in Theorem $1.^{8}$

Algorithm 1 An Efficient Algorithm for the Homogeneity Case

1: For $E \in {\binom{[n]}{d}}$, define

$$W_E := \begin{cases} Y_E & \text{if } E \in \mathcal{E}; \\ 0, & \text{otherwise.} \end{cases}$$

2: Apply Hypergraph Spectral Clustering (HSC) [44] to a weighted hypergraph $([n], \{W_E\}_{E \in \binom{[n]}{d}})$ to obtain

$$\mathbf{X}^{(0)} = \{X_i^{(0)}\}_{1 \le i \le n} \in \{0, 1\}^n.$$

3: Compute $\mathbf{X} = \{X_i\}_{1 \le i \le n}$ as follows:

$$X_i = \begin{cases} X_i^{(0)} & \text{if } \mathsf{d}_{\mathsf{H}}(\mathbf{X}^{(0)}) < \mathsf{d}_{\mathsf{H}}(\mathbf{X}^{(0)} \oplus \mathbf{e}_i); \\ X_i^{(0)} \oplus 1 & \text{if } \mathsf{d}_{\mathsf{H}}(\mathbf{X}^{(0)}) \ge \mathsf{d}_{\mathsf{H}}(\mathbf{X}^{(0)} \oplus \mathbf{e}_i), \end{cases}$$

where $d_{H}(\cdot)$ is defined in (6).

4: Output **X**.

Here we only present the algorithm while deferring a detailed analysis to our companion paper [44]. The algorithm operates in two stages, beginning with a decent initial estimate from Hypergraph Spectral Clustering [44] followed by iterative refinement. Detailed procedures are presented in Algorithm 1. Our algorithm is inspired by two-stage approaches that have been applied to a wide variety of problems including matrix completion [55], [56], phase retrieval [57], [58], robust PCA [59], community recovery [18], [20], [29], [60], [61], EM-algorithm [62], and rank aggregation [63].

2) *Performance of Algorithm 1:* We demonstrate the performance of Algorithm 1 by running Monte Carlo simulations.

Each point plotted in Fig. 3a and Fig. 3b indicates an empirical success rate. We take 100 Monte Carlo trials. Fig. 3a shows the probability of success when n = 1000, d = 4, and for various choices of θ . Shown in Fig. 3b is the performance of our algorithm with n = 1000, $\theta = 0.05$, and for various choices of d. For both figures, the x-axis denotes the number of samples normalized by the respective information-theoretic limits, characterized in Theorem 1. One can observe that the success probability due to Algorithm 1 quickly approaches 1

⁸Notice that our theoretical guarantee for the efficient algorithm is only for the balanced communities case, i.e., $\frac{n_1}{n_2} = \Theta(1)$, where n_1 and n_2 denote the sizes of the two communities, respectively. This condition ensures the weak consistency guarantee of HSC (Stage 2).



Fig. 4. Existence of d^* in motion segmentation. (a) We estimate the empirical noise rate $\hat{\theta}$ as a function of d in motion segmentation. (b) We plug $\hat{\theta}$ to the limit characterized in Theorem 1 and verify that $d^* = 6$.



Fig. 5. **Optimal choice of** d when θ decays with d. We run Monte Carlo simulations to estimate the probability of success with the data set shown in (a). We observe that the effective noise rate decreases as d increases. For varying d from 3 to 6, the success probability of Algorithm 1 is shown in (b): the best performance of the algorithm is observed when d = 4.

as the normalized sample complexity crosses 1, which corroborates our theoretical findings.

3) Optimal d for Subspace Clustering: We observe how the fundamental limit varies as a function of d. As we briefly discussed in Sec. III, if the noise rate θ is irrelevant to d, the optimal choice of d would be the minimum possible value of d. However, if the noise quality θ depends on d, there may be a sweet spot for d.

We demonstrate the existence of a sweet spot in one of subspace clustering applications: motion segmentation. We use the benchmark Hopkins 155 [64] dataset to compute an empirical noise rate θ as a function *d* as follows. For each sampled hyperedge $E = \{i_1, \ldots, i_d\}$, we adopt the method proposed in [34] to evaluate similarity between the corresponding *d* data points that we denote by *D*. Then, we set $Y_E = 1$ if and only if *D* is less than a fixed threshold, which is appropriately chosen so that $\Pr(Y_E = 0 \mid i_1, i_2, \ldots, i_d \text{ are from the same line}) \approx \Pr(Y_E = 1 \mid i_1, i_2, \ldots, i_d \text{ are not from the same line}). We estimate the effective noise rate <math>\hat{\theta} := \Pr(Y_E = 0 \mid i_1)$

 i_1, i_2, \ldots, i_d are from the same line) for various d, and observe that $\hat{\theta}$ quickly decreases as d increases; see Fig. 4a. We then plug these $\hat{\theta}$'s to the limit characterized in Theorem 1; see Fig. 4b. Note that d = 5 is not the optimal choice, but d = 6 is the sweet spot.

We also corroborate the existence of a sweet spot in a synthetic data set for subspace clustering, shown in Fig. 5a. Here the goal is to cluster $n \ (= 200)$ 2-dimensional data points approximately lying on a union of two lines (1-dimensional subspaces). We compute Y_E as above and evaluate the performance of Algorithm 1, shown in Fig. 5b. As a result, we observe that the optimal choice of d here is 4 rather than 3. Interestingly, the optimal d^* , as per the experiment, matches with the one computed with the estimated noise rates $\hat{\theta}$'s for each d as per the above procedure, the results read $\hat{\theta} = 0.248, 0.112, 0.051, 0.022, 0.011$ for d = 3, 4, 5, 6, 7, respectively. Plugging these $\hat{\theta}$'s into the limit characterized in Theorem 1, we obtain the estimated



Fig. 6. We run the Monte Carlo simulations to estimate the probability of success for n = 1000, varying d, and $\theta = 0$. For each d, we normalize the number of samples by max $(n, n \log n/d)$. Observe that the probability of success quickly approaches 1 as the normalized sample complexity crosses 1. Here we see much sharper transitions, relative to the ones in Fig. 3. This may be due to the fact that this experiment uses a deterministic algorithm (Gaussian elimination) and is subject to the noiseless case.

sample complexities 5163, 2870, 3033, 3984, 6103 for d = 3, 4, 5, 6, 7, respectively.

B. The Parity Measurement Case

1) Efficient Algorithms: For the parity measurement case, there are two efficient algorithms in the literature [35], [65]. In [35], it is shown that for d = 3, a variant of message passing algorithm successfully recovers the ground-truth vector provided that $\binom{n}{3}p = \Omega(n^2/\log n)$. Another efficient algorithm is based on a low-rank tensor factorization algorithm proposed in [65], and it is proved that reliable community recovery is feasible if $\binom{n}{3}p = \Omega(n^{1.5}\log^4 n)$. In either of the two cases, the sufficient condition comes with a polynomial term (*n* or $n^{1/2}$) to the fundamental limit characterized in Theorem 1. In fact, it is conjectured in [36] (see Conjecture 1 therein) that at least $n^{1.5}$ many samples are required for exact recovery.

On the other hand, focusing on the $\theta = 0$ case, recovering the ground-truth vector from the measurement vector **Y** is essentially the same as solving linear equations over the Galois field of two elements \mathbb{F}_2 . Hence, one can use efficient algorithms for solving linear equations such as Gaussian elimination to recover ground truth in the noiseless case. For subsequent experiments, we will only consider the noiseless case and use Gaussian elimination as an efficient recovery algorithm.

2) Information-Theoretic Limit: We first provide Monte Carlo simulation results which corroborate our theoretical findings in Theorem 2. Each point plotted in Fig. 6 and Fig. 7 is an empirical success rate. All results are obtained with 50 Monte Carlo trials. In Fig. 6, we plot the probability of successful recovery for n = 1000, varying d, and $\theta = 0$. For each d, we normalize the number of samples by $\max(n, n \log n/d)$. One can observe that the probability of success quickly approaches 1 as the normalized sample complexity crosses 1.

3) Minimum d for Linear Sample Complexity: According to Corollary 1, d should be $\Omega(\log n)$ for linear



Fig. 7. We run the Monte Carlo simulations to estimate the probability of success for varying n, varying d, $\theta = 0$, and $p = 1.1n/\binom{n}{d}$. Note that when n increases by a multiplicative factor of 4, the success probability curve shifts rightward by about the same amount. In particular, the minimum d required for a reasonable success probability, say 0.9, increases by about the same amount. These results are consistent with Corollary 1.

sample complexity. We corroborate this through the following experiment. In the experiment, we set n = 50, 200, 800, 3200, i.e., we increase n by a multiplication factor of 4 at each time. For each n, we vary d, while maintaining sample size to be $p\binom{n}{d} = 1.1n$, i.e., a linear sample size. Plotted in Fig. 7 are the experimental results for this setting. Note that when n increases by a multiplicative factor of 4, the success probability curve shifts rightward by about the same amount. In particular, the minimum d required for a reasonable success probability, say 0.9, increases by about the same amount. These results are consistent with Corollary 1 that the minimum d for linear sample complexity is $\Omega(\log n)$.

VIII. CONCLUSION

In this paper, we investigate the problem of community recovery in hypergraphs under the two generalized censored block models (GCBM), one based on the homogeneity measurement and the other based on the parity measurement. For these two models, we fully characterize the information-theoretic limits on sample complexity as a function of the number of nodes n, the size of edges d, the noise rate θ , and the edge observation probability p. We also corroborate our theoretical findings via experiments.

We conclude our paper by highlighting a few interesting open problems. One interesting question is whether or not one can sharpen Theorem 3 to characterize exact informationtheoretic limits for the scaling *d* case. From the simulation results in Sec. VII-B, we propose the following conjecture: Under the setting of Theorem 3, the information-theoretic limits is max $\left\{\frac{n}{1-H(\theta)}, \frac{1}{d}\frac{n\log n}{(\sqrt{1-\theta}-\sqrt{\theta})^2}\right\}$. Next open problem is about the computational gap for the parity measurement case: Investigating efficient algorithms for this case would shed some light on the study of information-computation gaps. Lastly, parallel to numerous efforts in the graph case [20], [66]–[68], generalizing our main results to (i) more than two communities or (ii) non-binary labels would be of great interest.

APPENDIX A PROOFS OF LEMMAS

A. Proof of Lemma 1

Recall that

$$\begin{aligned} |\mathcal{F}_{i,j}| &= \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{k-i}{d-\ell} + \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{n-k-j}{d-\ell} \\ &+ \sum_{\ell=1}^{d-1} \binom{j}{\ell} \binom{n-k-j}{d-\ell} + \sum_{\ell=1}^{d-1} \binom{k-i}{\ell} \binom{j}{d-\ell}. \end{aligned}$$

In order to prove the lemma, it is sufficient to prove the following two bounds:

Bound 1:
$$\sum_{\ell=1}^{d-1} {\binom{i}{\ell}} {\binom{k-i}{d-\ell}} + \sum_{\ell=1}^{d-1} {\binom{i}{\ell}} {\binom{n-k-j}{d-\ell}}$$
$$\geq i \cdot \frac{(1-2\delta)^{d-1}}{2^{d-2}} {\binom{n-1}{d-1}}$$

and

Bound 2:
$$\sum_{\ell=1}^{d-1} {j \choose \ell} {n-k-j \choose d-\ell} + \sum_{\ell=1}^{d-1} {k-i \choose \ell} {j \choose d-\ell}$$
$$\geq j \cdot \frac{(1-2\delta)^{d-1}}{2^{d-2}} {n-1 \choose d-1}.$$

Here we will focus on proving Bound 1. We remark that the proof of Bound 2 is essentially identical.

For simplicity, let

$$Z := \sum_{\ell=1}^{d-1} {i \choose \ell} {k-i \choose d-\ell} + \sum_{\ell=1}^{d-1} {i \choose \ell} {n-k-j \choose d-\ell}.$$

Then, since $i < \delta n$ and $j < \delta n$,

$$Z \ge \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{k-\delta n}{d-\ell} + \sum_{\ell=1}^{d-1} \binom{i}{\ell} \binom{n-k-\delta n}{d-\ell}.$$

We further bound Z by considering two cases separately: $k \ge 1$ δn and $k < \delta n$. When $k \ge \delta n$,

$$\sum_{\ell=1}^{d-1} {i \choose \ell} {k-\delta n \choose d-\ell} + \sum_{\ell=1}^{d-1} {i \choose \ell} {n-k-\delta n \choose d-\ell}$$

$$\geq i {k-\delta n \choose d-1} + i {n-k-\delta n \choose d-1}$$

$$\approx i \cdot \left[\left(\frac{k}{n} - \delta\right)^{d-1} + \left(1 - \frac{k}{n} - \delta\right)^{d-1} \right] {n-1 \choose d-1},$$

where the last inequality holds since $\binom{an}{b} \approx a^b \binom{n}{b} \approx a^b \binom{n-1}{b}$ for constants a and b. We then apply Hölder's inequality: Given p, q such that 1/p + 1/q = 1, we have $\sum_{z} |x_{z}y_{z}| \le (\sum_{z} |x_{z}|^{p})^{1/p} (\sum_{z} |y_{z}|^{q})^{1/q}$ for all sequences $\{x_{z}\}$ and $\{y_{z}\}$. By setting $(x_{1}, x_{2}) = (\alpha, \beta), (y_{1}, y_{2}) = (1, 1), p = d - 1, q = d - 1$ $\frac{d-1}{d-2}$, we have

$$\alpha + \beta \le (\alpha^{d-1} + \beta^{d-1})^{\frac{1}{d-1}} 2^{\frac{d-2}{d-1}}$$

Applying this version of Hölder's inequality to the last lower bound, we have

$$i \cdot \left[\left(\frac{k}{n} - \delta\right)^{d-1} + \left(1 - \frac{k}{n} - \delta\right)^{d-1} \right] \binom{n-1}{d-1}$$
$$\geq i \cdot \frac{(1 - 2\delta)^{d-1}}{2^{d-2}} \binom{n-1}{d-1}.$$

When $k < \delta n$, $\sum_{\ell=1}^{d-1} {i \choose \ell} {n-k-j \choose d-\ell}$ becomes the dominant term. Hence,

$$\sum_{\ell=1}^{d-1} {j \choose \ell} {n-k-\delta n \choose d-\ell} + \sum_{\ell=1}^{d-1} {k-i \choose \ell} {\delta n \choose d-\ell}$$

> $i {n-k-\delta n \choose d-1} > i {n-2\delta n \choose d-1}$
 $\approx i \cdot (1-2\delta)^{d-1} {n-1 \choose d-1} > i \cdot \frac{(1-2\delta)^{d-1}}{2^{d-2}} {n-1 \choose d-1}.$
is completes the proof.

This completes the proof.

B. Proof of Lemma 2

Denote by $\mathcal{R}_{\text{big}} \subset [n]$ the set of nodes of size $n/\log^7 n$. One can easily show that with high probability, some nodes of this set are connected by the same hyperedge(s). Denote by \mathcal{R}_{res} the largest subset of \mathcal{R}_{big} , whose elements do not share the same hyperedges. The lemma states that with high probability, $|\mathcal{R}_{res}| \simeq |\mathcal{R}_{big}|$.

We now formally prove this statement. Note that for a hyperedge $E = (i_1, i_2, \dots, i_d), |E \cap \mathcal{R}_{big}|$ is the number of nodes in \mathcal{R}_{big} that are connected by the hyperedge. Hence, if $2 \le |E \cap R_{\text{big}}| \le d$, this hyperedge connects more than one nodes in \mathcal{R}_{big} , and $E \cap \mathcal{R}_{big}$ is the set of the nodes that share the same hyperedge E.

Let us denote by \mathcal{R}_{share} the subset of nodes that are connected by the same hyperedge(s). Then,

$$\mathcal{R}_{\text{share}} := \bigcup_{k=2}^{d} \mathcal{R}_{\text{share}}^{(k)} := \bigcup_{k=2}^{d} \bigcup_{\substack{E \in \mathcal{E}:\\ |E \cap \mathcal{R}_{\text{big}}| = k}} E \cap \mathcal{R}_{\text{big}}.$$
(40)

Our proof strategy is as follows. Since

$$\mathcal{R}_{\text{res}} = \mathcal{R}_{\text{big}} - \mathcal{R}_{\text{share}} = \mathcal{R}_{\text{big}} - \bigcup_{k=2}^{d} \mathcal{R}_{\text{share}}^{(k)},$$
 (41)

it is sufficient to show that

$$\left| \bigcup_{k=2}^{d} \mathcal{R}_{\text{share}}^{(k)} \right| = o(|\mathcal{R}_{\text{big}}|).$$
(42)

More specifically, we will show

$$\Pr\left(\exists k \in \{2, 3, \dots, d\} \text{ s.t. } |\mathcal{R}_{\text{share}}^{(k)}| > \frac{n}{\log^9 n}\right) \to 0. \quad (43)$$

That is, with probability approaching 1, $|\mathcal{R}_{\text{share}}^{(k)}|$ $o(n/\log^8 n)$ for all $k, 2 \leq k \leq d$. Note that this implies (42) since

$$\left| \bigcup_{k=2}^{d} \mathcal{R}_{\text{share}}^{(k)} \right| \le \sum_{k=2}^{d} |\mathcal{R}_{\text{share}}^{(k)}| = O(d) \times o\left(\frac{n}{\log^8 n}\right)$$
$$= o\left(\frac{n}{\log^7 n}\right) = o(|\mathcal{R}_{\text{big}}|).$$

In order to bound (43), we first derive an upper bound on the expected value of $|\mathcal{R}_{\text{share}}^{(k)}|$. By definition,

$$\begin{aligned} |\mathcal{R}_{\text{share}}^{(k)}| &\leq \sum_{\substack{E \in \mathcal{E}:\\|E \cap \mathcal{R}_{\text{big}}|=k}} |E \cap \mathcal{R}_{\text{big}}| \\ &\leq \sum_{\substack{E \in \mathcal{E}:\\|E \cap \mathcal{R}_{\text{big}}|=k}} k = |\{E \in \mathcal{E} : |E \cap \mathcal{R}_{\text{big}}| = k\}| \cdot k. \end{aligned}$$

Observe that $|\{E \in \mathcal{E} : |E \cap \mathcal{R}_{\text{big}}|\}|$ is the sum of $\binom{|\mathcal{R}_{\text{big}}|}{k}\binom{n-|\mathcal{R}_{\text{big}}|}{d-k}$ i.i.d. Bernoulli random variables with probability p. Hence,

$$\mathbb{E}\left\{|\{E \in \mathcal{E} : |E \cap \mathcal{R}_{\text{big}}| = k\}| \cdot k\right\}$$
$$= k \binom{|\mathcal{R}_{\text{big}}|}{k} \binom{n - |\mathcal{R}_{\text{big}}|}{d - k} p$$
$$= |\mathcal{R}_{\text{big}}| \binom{|\mathcal{R}_{\text{big}}| - 1}{k - 1} \binom{n - |\mathcal{R}_{\text{big}}|}{d - k} p.$$

As $|\mathcal{R}_{\text{big}}| = o(n)$, we have $\binom{|\mathcal{R}_{\text{big}}|-1}{k-1} \binom{n-|\mathcal{R}_{\text{big}}|}{d-k} \leq \binom{|\mathcal{R}_{\text{big}}|-1}{1} \binom{n-|\mathcal{R}_{\text{big}}|}{d-2}$, which in turn gives the following upper bound on the last term:

$$\begin{aligned} &|\mathcal{R}_{\text{big}}| \binom{|\mathcal{R}_{\text{big}}|-1}{1} \binom{n-|\mathcal{R}_{\text{big}}|}{d-2} p \\ &\leq 2|\mathcal{R}_{\text{big}}|^2 \binom{n-2}{d-2} p \\ &= 2|\mathcal{R}_{\text{big}}|^2 \binom{n-2}{d-2} \frac{n\log n}{\binom{n}{d}} \\ &= O\left(\frac{|\mathcal{R}_{\text{big}}|^2 d^2 \log n}{n}\right) = O\left(\frac{n}{\log^{11} n}\right), \end{aligned}$$

where the last equality holds since $\binom{n-2}{d-2} \approx \binom{n}{d} \frac{d^2}{n^2}$. Note that this inequality holds for any $2 \leq k \leq d$. Using Markov's inequality,

$$\Pr\left(|\{E \in \mathcal{E} : |E \cap \mathcal{R}_{\text{big}}| = k\}| \cdot k > \frac{n}{\log^9 n}\right)$$
$$\leq \frac{\log^9 n}{n} \cdot O\left(\frac{n}{\log^{11} n}\right) = O\left(\frac{1}{\log^2 n}\right).$$

Applying the union bound over all $2 \le k \le d$, the probability that there exists k = 2, 3, ..., d such that

$$|\{E \in \mathcal{E} : |E \cap \mathcal{R}_{\text{big}}| = k\}| \cdot k > \frac{n}{\log^9 n}$$

is upped bounded by

$$d \cdot O\left(\frac{1}{\log^2 n}\right) = O\left(\frac{1}{\log n}\right)$$

This completes the proof.

C. Proof of Lemma 4

Recall that we assumed $d \leq k$ at the beginning of Sec. VI-C.

Fact 2. For
$$1 \le i \le d-1$$
,
 $\binom{k}{i}\binom{n-k}{d-i} \le 2\binom{k}{i+1}\binom{n-k}{d-(i+1)}$
 $+ 2\binom{k}{i-1}\binom{n-k}{d-(i-1)}$.

Proof: The conclusion follows from the following inequalities:

$$\begin{split} & \frac{\binom{k}{i+1}\binom{n-k}{d-(i+1)} + \binom{k}{i-1}\binom{n-k}{d-(i-1)}}{\binom{k}{i}\binom{n-k}{d-i}} \\ &= \frac{(k-i)(d-i)}{(i+1)(n-k-d+i+1)} + \frac{i(n-k-d+i)}{(k-i+1)(d-i+1)} \\ &\geq 2\sqrt{\frac{(k-i)(d-i)}{(i+1)(n-k-d+i+1)}} \cdot \frac{i(n-k-d+i)}{(k-i+1)(d-i+1)} \\ &= 2\sqrt{\frac{(k-i)}{(k-i+1)}} \cdot \frac{(d-i)}{(d-i+1)} \cdot \frac{i}{i+1} \cdot \frac{(n-k-d+i)}{(n-k-d+i+1)} \\ &\geq 2\sqrt{\left(\frac{1}{2}\right)^4} = \frac{1}{2}, \end{split}$$

where the last inequality follows since k - i, d - i, i, n - k - d + i are all greater than or equal to 1 as long as $1 \le i \le d - 1$ due to the assumption $d \le k$. We conclude the proof using Fact 2:

$$\sum_{\substack{1 \le i \le d-1 \\ i \ge i \le d-2}} \binom{k}{i} \binom{n-k}{d-i} \\ \le \sum_{\substack{2 \le i \le d-2 \\ i \ge \text{ even}}} \binom{k}{i} \binom{n-k}{d-i} + \sum_{\substack{1 \le i \le d-1 \\ i \ge \text{ odd}}} \binom{k}{i} \binom{n-k}{d-i} \\ \stackrel{(a)}{\le} 2 \cdot \binom{k}{1} \binom{n-k}{d-1} + 4 \cdot \sum_{\substack{3 \le i \le d-3 \\ i \ge \text{ odd}}} \binom{k}{i} \binom{n-k}{d-i}$$

$$+2\cdot \binom{k}{d-1}\binom{n-k}{1} + \sum_{\substack{1\leq i\leq d-1\\i: \text{ odd}}}\binom{k}{i}\binom{n-k}{d-i} + 3 N_k,$$
$$=2\cdot \sum_{\substack{3\leq i\leq d-3\\i: \text{ odd}}}\binom{k}{i}\binom{n-k}{d-i} + 3 N_k,$$

where (a) follows from Fact 2.

 β ≤ k ≤ n/2 Since d ≤ n/2 and β < n/2, one can verify the inequality using the following facts:

$$\binom{k}{0}\binom{n-k}{d} \le 2\binom{k}{1}\binom{n-k}{d-1} \Leftrightarrow k \ge \frac{n-d+1}{2d+1}$$

and

$$\binom{k}{d}\binom{n-k}{0} \le 2\binom{k}{d-1}\binom{n-k}{1} \Leftrightarrow k \le n - \frac{n-d+1}{2d+1}$$

We first show that $\alpha/k \ge 2$. Since

$$k \le \left\lceil \frac{n-d+1}{2d+1} \right\rceil - 1 \le \frac{n-d+1}{2d+1},$$

we have

$$\frac{\alpha}{k} = \frac{\left(\frac{n-d+1}{d}\right)}{k} > \frac{\left(\frac{n-d+1}{d}\right)}{\left(\frac{n-d+1}{2d+1}\right)} = \frac{2d+1}{d} \ge 2.$$

Next, the inequality can be checked using the following facts:

$$\binom{k}{d}\binom{n-k}{0} \leq 2\binom{k}{d-1}\binom{n-k}{1}$$
$$\Leftrightarrow k \leq n - \frac{n-d+1}{2d+1}$$

and

$$\frac{\binom{k}{0}\binom{n-k}{d}}{\binom{k}{1}\binom{n-k}{d-1}} = \frac{n-k-d+1}{kd} \le \frac{n-d+1}{kd} = \frac{\alpha}{k} \,.$$

REFERENCES

- K. Ahn, K. Lee, and C. Suh, "Community recovery in hypergraphs," in Proc. Allerton Conf. Commun., Control Comput., 2016, pp. 657–663.
- [2] K. Ahn, K. Lee, and C. Suh, "Information-theoretic limits of subspace clustering," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Aachen, Germany, Jun. 2017, pp. 2473–2477.
- [3] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proc. Nat. Acad. Sci. USA*, vol. 99, no. 12, pp. 7821–7826, Apr. 2002.
- [4] S. Fortunato, "Community detection in graphs," *Phys. Rep.*, vol. 486, nos. 3–5, pp. 75–174, 2010.
- [5] M. A. Porter, J.-P. Onnela, and P. J. Mucha, "Communities in networks," *Notices AMS*, vol. 56, no. 9, pp. 1082–1097, 2009.
- [6] J. Chen and B. Yuan, "Detecting functional modules in the yeast proteinprotein interaction network," *Bioinformatics*, vol. 22, pp. 2283–2290, Sep. 2006.
- [7] Q.-X. Huang and L. Guibas, "Consistent shape maps via semidefinite programming," *Comput. Graph. Forum*, vol. 32, no. 5, pp. 177–186, 2013.
- [8] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 8, pp. 888–905, Aug. 2000.
- [9] P. W. Holland, K. B. Laskey, and S. Leinhardt, "Stochastic blockmodels: First steps," *Social Netw.*, vol. 5, no. 2, pp. 109–137, 1983.
- [10] E. Abbe and A. Montanari, "Conditional random fields, planted constraint satisfaction and entropy concentration," in *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques.* Berlin, Germany: Springer, 2013, pp. 332–346.
- [11] G. Ghoshal, V. Zlatić, G. Caldarelli, and M. E. J. Newman, "Random hypergraphs and their applications," *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 79, no. 6, 2009, Art. no. 066118.
- [12] T. Michoel and B. Nachtergaele, "Alignment and integration of complex networks by hypergraph-based spectral clustering," *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 86, no. 5, 2012, Art. no. 056111.
- [13] X. Dong, P. Frossard, P. Vandergheynst, and N. Nefedov, "Clustering on multi-layer graphs via subspace analysis on Grassmann manifolds," *IEEE Trans. Signal Process.*, vol. 62, no. 4, pp. 905–918, Feb. 2014.
- [14] S. Agarwal, K. Branson, and S. Belongie, "Higher order learning with graphs," in *Proc. ACM ICML*, 2006, pp. 17–24.
- [15] G. Karypis and V. Kumar, "Multilevel k-way hypergraph partitioning," VLSI Des., vol. 11, no. 3, pp. 285–300, 2000.
- [16] D. Gibson, J. Kleinberg, and P. Raghavan, "Clustering categorical data: An approach based on dynamical systems," *Databases*, vol. 8, nos. 3–4, pp. 222–236, Feb. 2000.

- [17] B. Hajek, Y. Wu, and J. Xu, "Achieving exact cluster recovery threshold via semidefinite programming: Extensions," *IEEE Trans. Inf. Theory*, vol. 62, no. 10, pp. 5918–5937, Oct. 2016.
- [18] E. Abbe, A. S. Bandeira, and G. Hall, "Exact recovery in the stochastic block model," *IEEE Trans. Inf. Theory*, vol. 62, no. 1, pp. 471–487, Jan. 2016.
- [19] E. Mossel, J. Neeman, and A. Sly, "Consistency thresholds for the planted bisection model," 2014, arXiv:1407.1591. [Online]. Available: https://arxiv.org/abs/1407.1591
- [20] E. Abbe and C. Sandon, "Community detection in general stochastic block models: Fundamental limits and efficient algorithms for recovery," in *Proc. IEEE FOCS*, Oct. 2015, pp. 670–688.
- [21] E. Abbe, A. S. Bandeira, A. Bracher, and A. Singer, "Decoding binary node labels from censored edge measurements: Phase transition and efficient recovery," *IEEE Trans. Netw. Sci. Eng.*, vol. 1, no. 1, pp. 10–22, Jan. 2014.
- [22] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová, "Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications," *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 84, no. 6, 2011, Art. no. 066106.
- [23] E. Mossel, J. Neeman, and A. Sly, "Reconstruction and estimation in the planted partition model," *Probab. Theory Rel. Fields*, vol. 162, nos. 3–4, pp. 431–461, 2015.
- [24] E. Mossel, J. Neeman, and A. Sly, "A proof of the block model threshold conjecture," 2013, arXiv:1311.4115. [Online]. Available: https://arxiv. org/abs/1311.4115
- [25] L. Massoulié, "Community detection thresholds and the weak Ramanujan property," in *Proc. ACM 46th Annu. ACM Symp. Theory Comput.*, 2014, pp. 694–703.
- [26] C. Bordenave, M. Lelarge, and L. Massoulié, "Non-backtracking spectrum of random graphs: Community detection and non-regular ramanujan graphs," in *Proc. IEEE FOCS*, Oct. 2015, pp. 1347–1357.
- [27] E. Abbe and C. Sandon, "Proof of the achievability conjectures for the general stochastic block model," *Commun. Pure Appl. Math.*, vol. 71, no. 7, pp. 1334–1406, 2017.
- [28] A. Y. Zhang and H. H. Zhou, "Minimax rates of community detection in stochastic block models," *Ann. Statist.*, vol. 44, no. 5, pp. 2252–2280, 2016.
- [29] C. Gao, Z. Ma, A. Y. Zhang, and H. H. Zhou, "Achieving optimal misclassification proportion in stochastic block model," *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 1980–2024, 2017.
- [30] E. Abbe, "Community detection and stochastic block models: Recent developments," J. Mach. Learn. Res., vol. 18, no. 1, pp. 6446–6531, 2017.
- [31] D. Ghoshdastidar and A. Dukkipati, "Uniform hypergraph partitioning: Provable tensor methods and sampling techniques," *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 1638–1678, 2017. [Online]. Available: http://jmlr.org/papers/v18/16-100.html
- [32] D. Ghoshdastidar and A. Dukkipati, "Consistency of spectral hypergraph partitioning under planted partition model," *Ann. Statist.*, vol. 45, no. 1, pp. 289–315, 2017.
- [33] V. M. Govindu, "A tensor decomposition for geometric grouping and segmentation," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2005, pp. 1150–1157.
- [34] G. Chen and G. Lerman, "Spectral curvature clustering (SCC)," Int. J. Comput. Vis., vol. 81, no. 3, pp. 317–330, Mar. 2009.
- [35] O. Watanabe, "Message passing algorithms for MLS-3LIN problem," Algorithmica, vol. 66, no. 4, pp. 848–868, 2013.
- [36] L. Florescu and W. Perkins, "Spectral thresholds in the bipartite stochastic block model," in *Proc. Conf. Learn. Theory (COLT)*, 2016, pp. 943–959.
- [37] M. C. Angelini, F. Caltagirone, F. Krzakala, and L. Zdeborová, "Spectral detection on sparse hypergraphs," in *Proc. Allerton Conf. Commun.*, *Control, Comput.*, Sep./Oct. 2015, pp. 66–73.
- [38] C.-Y. Lin, I. Chien, and I.-H. Wang, "On the fundamental statistical limit of community detection in random hypergraphs," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2017, pp. 2183–2187.
- [39] N. Vesdapunt, K. Bellare, and N. Dalvi, "Crowdsourcing algorithms for entity resolution," *Proc. VLDB Endowment*, vol. 7, no. 12, pp. 1071–1082, 2014.
- [40] H. Ashtiani, S. Kushagra, and S. Ben-David, "Clustering with samecluster queries," in *Proc. Adv. Neural Inf. Process. Syst.*, 2016, pp. 3216–3224.
- [41] R. Dorfman, "The detection of defective members of large populations," Ann. Math. Statist., vol. 14, no. 4, pp. 436–440, Dec. 1943.

- [42] I.-H. Wang, S.-L. Huang, and K.-Y. Lee, "Extracting sparse data via histogram queries," in *Proc. Allerton Conf. Commun., Control, Comput.*, Sep. 2016, pp. 39–45.
- [43] I.-H. Wang, S.-L. Huang, K.-Y. Lee, and K.-C. Chen, "Data extraction via histogram and arithmetic mean queries: Fundamental limits and algorithms," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jun. 2016, pp. 1386–1390.
- [44] K. Ahn, K. Lee, and C. Suh, "Hypergraph spectral clustering in the weighted stochastic block model," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 5, pp. 959–974, Oct. 2018.
- [45] R. Vidal, R. Tron, and R. Hartley, "Multiframe motion segmentation with missing data using powerfactorization and GPCA," *Int. J. Comput. Vis.*, vol. 79, no. 1, pp. 85–105, 2008.
- [46] J. Ho, M.-H. Yang, J. Lim, K.-C. Lee, and D. Kriegman, "Clustering appearances of objects under varying illumination conditions," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 1, Jun. 2003, pp. 1–11.
- [47] R. Vidal, "Subspace clustering," *IEEE Signal Process. Mag.*, vol. 28, no. 2, pp. 52–68, Mar. 2011.
- [48] E. Elhamifar and R. Vidal, "Sparse subspace clustering: Algorithm, theory, and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 11, pp. 2765–2781, Nov. 2013.
- [49] E. L. Dyer, A. C. Sankaranarayanan, and R. G. Baraniuk, "Greedy feature selection for subspace clustering," *J. Mach. Learn. Res.*, vol. 14, no. 1, pp. 2487–2517, 2013.
- [50] R. Heckel and H. Bolcskei, "Robust subspace clustering via thresholding," *IEEE Trans. Inf. Theory*, vol. 61, no. 11, pp. 6320–6342, Nov. 2015.
- [51] E. Arias-Castro, G. Lerman, and T. Zhang, "Spectral clustering based on local PCA," J. Mach. Learn. Res., vol. 18, no. 1, pp. 253–309, 2017.
- [52] B. Hajek, Y. Wu, and J. Xu, "Information limits for recovering a hidden community," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jul. 2016, pp. 1894–1898.
- [53] W. Hoeffding, "Probability inequalities for sums of bounded random variables," J. Amer. Stat. Assoc., vol. 58, no. 301, pp. 13–30, 1963.
- [54] N. Alon and J. H. Spencer, *The Probabilistic Method*. Hoboken, NJ, USA: Wiley, 2004.
- [55] R. H. Keshavan, A. Montanari, and S. Oh, "Matrix completion from a few entries," *IEEE Trans. Inf. Theory*, vol. 56, no. 6, pp. 2980–2998, Jun. 2010.
- [56] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Proc. 45th Annu. ACM Symp. Theory Comput.*, 2013, pp. 665–674.
- [57] P. Netrapalli, P. Jain, and S. Sanghavi, "Phase retrieval using alternating minimization," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 2796–2804.
- [58] E. J. Candès, X. Li, and M. Soltanolkotabi, "Phase retrieval via Wirtinger flow: Theory and algorithms," *IEEE Trans. Inf. Theory*, vol. 61, no. 4, pp. 1985–2007, Apr. 2015.
- [59] X. Yi, D. Park, Y. Chen, and C. Caramanis, "Fast algorithms for robust PCA via gradient descent," in *Proc. Adv. Neural Inf. Process. Syst.*, 2016, pp. 4152–4160.
- [60] Y. Chen, G. Kamath, C. Suh, and D. Tse, "Community recovery in graphs with locality," in *Proc. ICML*, 2016, pp. 689–698.

- [61] P. Chin, A. Rao, and V. Vu, "Stochastic block model and community detection in sparse graphs: A spectral algorithm with optimal rate of recovery," in *Proc. Conf. Learn. Theory (COLT)*, 2015, pp. 391–423.
- [62] S. Balakrishnan, M. J. Wainwright, and B. Yu, "Statistical guarantees for the em algorithm: From population to sample-based analysis," *Ann. Statist.*, vol. 45, no. 1, pp. 77–120, 2017. 10.1214/16-AOS1435
- [63] Y. Chen and C. Suh, "Spectral MLE: Top-k rank aggregation from pairwise comparisons," in *Proc. ICML*, 2015, pp. 371–380.
 [64] R. Tron and R. Vidal, "A benchmark for the comparison of 3-D motion
- [64] R. Tron and R. Vidal, "A benchmark for the comparison of 3-D motion segmentation algorithms," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2007, pp. 1–8.
- [65] P. Jain and S. Oh, "Provable tensor factorization with missing data," in Proc. Adv. Neural Inf. Process. Syst., 2014, pp. 1431–1439.
- [66] Y. Chen, C. Suh, and A. J. Goldsmith, "Information recovery from pairwise measurements," *IEEE Trans. Inf. Theory*, vol. 62, no. 10, pp. 5881–5905, Oct. 2016.
- [67] V. Jog and P.-L. Loh, "Information-theoretic bounds for exact recovery in weighted stochastic block models using the Renyi divergence," 2015, arXiv:1509.06418. [Online]. Available: https://arxiv.org/abs/1509.06418
- [68] S. Heimlicher, M. Lelarge, and L. Massoulié, "Community detection in the labelled stochastic block model," 2012, arXiv:1209.2910. [Online]. Available: https://arxiv.org/abs/1209.2910

Kwangjun Ahn received his B.S. degree in the Department of Mathematical Sciences from Korea Advanced Institute of Science and Technology (KAIST) in 2017. He is currently a military police desk clerk in the U.S. Army as a part of Korean Augmentation to the U.S. Army (KATUSA).

Kangwook Lee is a Research Assistant Professor at Information and Electronics Research Institute at Korea Advanced Institute of Science and Technology (KAIST). He earned his Ph.D. in EECS from UC Berkeley in 2016. He is a recipient of the KFAS Fellowship from 2010 to 2015.

Changho Suh (S'10–M'12) is an Associate Professor in the School of Electrical Engineering at Korea Advanced Institute of Science and Technology (KAIST). He received the B.S. and M.S. degrees in Electrical Engineering from KAIST in 2000 and 2002 respectively, and the Ph.D. degree in Electrical Engineering and Computer Sciences from UC-Berkeley in 2011. From 2011 to 2012, he was a postdoctoral associate at the Research Laboratory of Electronics in MIT. From 2002 to 2006, he had been with the Telecommunication R&D Center, Samsung Electronics.

Dr. Suh received the 2018 IEIE/IEEE Joint Award, the 2015 IEIE Hadong Young Engineer Award, a 2015 Bell Labs Prize finalist, the 2013 IEEE Communications Society Stephen O. Rice Prize, the 2011 David J. Sakrison Memorial Prize, and the 2009 IEEE ISIT Best Student Paper Award.