

# Performance Bounds for Group Testing With Doubly-Regular Designs

Nelvin Tan, Way Tan, and Jonathan Scarlett

**Abstract**—In the group testing problem, the goal is to identify a subset of defective items within a larger set of items based on tests whose outcomes indicate whether any defective item is present. This problem is relevant in areas such as medical testing, DNA sequencing, and communications. In this paper, we study a doubly-regular design in which the number of tests-per-item and the number of items-per-test are fixed. We analyze the performance of this test design alongside the Definite Defectives (DD) decoding algorithm in several settings, namely, (i) the sub-linear regime  $k = o(n)$  with exact recovery, (ii) the linear regime  $k = \Theta(n)$  with approximate recovery, and (iii) the size-constrained setting, where the number of items per test is constrained. Under setting (i), we show that our design together with the DD algorithm, matches an existing achievability result for the DD algorithm with the near-constant tests-per-item design, which is known to be asymptotically optimal in broad scaling regimes. Under setting (ii), we provide novel approximate recovery bounds that complement a hardness result regarding exact recovery. Lastly, under setting (iii), we improve on the best known upper and lower bounds in scaling regimes where the maximum test size grows with the total number of items.

**Index Terms**—Group testing, sparsity, performance bounds, randomized test designs, information-theoretic limits

## I. INTRODUCTION

In the group testing problem, the goal is to identify a small subset of defective items of size  $k$  within a larger set of items of size  $n$ , based on a number  $T$  of tests. This problem is relevant in areas such as medical testing, DNA sequencing, and communication protocols [2, Sec. 1.7], and has recently found utility in COVID-19 testing [3].

In non-adaptive group testing, the placements of items into test can be represented by a binary test matrix of size  $T \times n$ . The strongest known theoretical guarantees are based on the idea of generating this matrix at random and analyzing the average performance. Starting from early studies of group testing, a line of works led to a detailed understanding of the i.i.d. test design [4]–[7], and more recently, further improvements were shown for a near-constant tests-per-item design [8], [9], whose asymptotic optimality in sub-linear sparsity regimes was established in [10]. Recently, there has been increasing evidence that *doubly-regular designs* (i.e., both

constant tests-per-item and items-per-test) also play a crucial role in various settings of interest:

- The work of Mezárd *et al.* [11] uses heuristic arguments from statistical physics to suggest that doubly-regular designs achieve the same optimal threshold as the near-constant tests-per-item design when  $k = \Theta(n^\theta)$ , at least when  $\theta$  is not too small.
- In constrained settings where the number of items per test cannot exceed a pre-specified threshold, doubly-regular designs have been used to obtain performance bounds that appear to be difficult or impossible to obtain using the other designs mentioned above [12], [13].
- In the linear sparsity regime (i.e.,  $k = \Theta(n)$ ), various two-stage adaptive designs were studied in [14], and using a doubly-regular design in the first stage led to strict improvements over the other designs. See also Appendix C for analogous observations with non-adaptive testing and approximate recovery.
- From a more practical viewpoint, doubly-regular designs have found utility in application-driven settings, e.g., see [15], [16] for recent studies relating to medical testing.

In this paper, motivated by these developments, we seek to provide a more detailed understanding of non-adaptive doubly-constant test designs, particularly when paired with the Definite Defectives (DD) algorithm [6]. Briefly, our contributions are as follows: (i) For  $\theta \in [\frac{1}{2}, 1)$ , we rigorously prove the above-mentioned result shown heuristically in [11], albeit with a slightly different version of the doubly-regular design; (ii) We establish new performance bounds for non-adaptive group testing in the linear regime ( $k = \Theta(n)$ ) with some false negatives allowed in the reconstruction, complementing strong impossibility results for exact recovery [17], [18]; (iii) We provide improved upper and lower bounds on the number of tests for the constrained setting with at most  $\rho$  items per test. Our bounds apply to general scaling regimes beyond the regime  $\rho = O(1)$  recently studied in [13], and our consideration of the DD algorithm leads to strict improvements over the COMP algorithm considered in [12].

## A. Problem Setup

Let  $n$  denote the number of items, which we label as  $[n] = \{1, \dots, n\}$ . Let  $\mathcal{K} \subset [n]$  denote the fixed set of defective items, and let  $k = |\mathcal{K}|$  be the number of defective items. We adopt the combinatorial prior [2], where  $\mathcal{K}$  is chosen uniformly at random from  $\binom{[n]}{k}$  sets of size  $k$ . We let  $T = T(n)$  be the

N. Tan is with the Department of Engineering, University of Cambridge. W. Tan and J. Scarlett are with the Department of Computer Science, National University of Singapore (NUS), and also with the Department of Mathematics, NUS. e-mails: [tcnt2@cam.ac.uk](mailto:tcnt2@cam.ac.uk); [e0174826@u.nus.edu](mailto:e0174826@u.nus.edu); [scarlett@comp.nus.edu.sg](mailto:scarlett@comp.nus.edu.sg).

This work was presented in part at the 2021 IEEE International Symposium on Information Theory (ISIT) [1].

This work was supported by an NUS Early Career Research Award. N. Tan and W. Tan contributed equally to this work.

number of tests performed. The  $i$ -th test takes the form

$$Y^{(i)} = \bigvee_{j \in \mathcal{K}} X_j^{(i)}, \quad (1)$$

where the test vector  $X^{(i)} = (X_1^{(i)}, \dots, X_n^{(i)}) \in \{0, 1\}^n$  indicates which items are included in the test, and  $Y^{(i)} \in \{0, 1\}$  is the resulting observation, which indicates whether at least one defective item was included in the test. The goal of group testing is to design a sequence of tests  $X^{(1)}, \dots, X^{(T)}$ , with  $T$  ideally as small as possible, such that the outcomes can be used to reliably recover the defective set  $\mathcal{K}$ . We focus on the non-adaptive setting, in which all tests  $X^{(1)}, \dots, X^{(T)}$  must be designed prior to observing any outcomes.

Next, we introduce the main defining features distinguishing the settings we consider:

- Regarding the scaling of  $k$ , we consider the following:
  - **Sub-linear:** We have  $k = \Theta(n^\theta)$  for some constant  $\theta \in (0, 1)$ . This is the regime where defectivity is rare, and also where group testing typically exhibits the greatest gains.
  - **Linear:** We have  $k = pn$  for some prevalence rate  $p \in (0, 1)$ . This regime may potentially be of greater relevance in certain practical situations, e.g., with  $p$  representing the prevalence of a disease.
- Regarding the constraints (or lack thereof), we consider the following:
  - **Unconstrained:** There is no restriction on the number of tests-per-item or the number of items-per-test in the test design.
  - **Size-constrained:** Tests are size-constrained and thus contain no more than  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  items per test, for some constant  $\beta \in (0, 1)$ . Note that if each test comprises of  $\Theta(n/k)$  items, then  $\Theta(k \log n)$  tests suffice for group testing algorithms with asymptotically vanishing error probability [6]–[8], [19]. One can alternatively use exactly  $n$  tests via one-by-one testing, and it has recently been shown that taking the better of the two (i.e.,  $\Theta(\min\{k \log n, n\})$  tests) gives asymptotically optimal scaling [18]. Hence, to avoid essentially reducing to the unconstrained setting, the parameter regime of interest in the size-constrained setting is  $\rho = o(n/k)$ , which justifies the scaling  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$ .
- Regarding recovery criteria, we consider the following:
  - **Exact recovery:** We seek to develop a testing strategy and decoder that produces an estimate  $\hat{\mathcal{K}}$  such that the error probability  $P_e = \mathbb{P}[\hat{\mathcal{K}} \neq \mathcal{K}]$  is asymptotically vanishing as  $n \rightarrow \infty$ .
  - **Approximate recovery:** We seek to characterize the per-item false-positive rate (FPR) and false-negative rate (FNR), which are defined as the probability that a non-defective item (picked uniformly at random) is declared defective (i.e.,  $\text{FPR} = \frac{\mathbb{E}[|\hat{\mathcal{K}} \setminus \mathcal{K}|]}{[n] \setminus \mathcal{K}}$ ), and the probability that a defective item (picked uniformly at random) is declared non-defective (i.e.,  $\text{FNR} = \frac{\mathbb{E}[|\mathcal{K} \setminus \hat{\mathcal{K}}|]}{|\mathcal{K}|}$ ), respectively.

---

**Algorithm 1** COMP and DD algorithms [6], [19].

---

**Require:**  $T$  tests.

- 1: Initialize two empty sets  $\mathcal{PD}$  (possibly defective set) and  $\mathcal{DD}$  (definitely defective set).
  - 2: Label any item in a negative test as definitely non-defective, and add all remaining items to  $\mathcal{PD}$ .
  - 3: **for** each test **do**
  - 4:   If the test contains exactly one item from  $\mathcal{PD}$ , then add that item to the set  $\mathcal{DD}$ .
  - 5: **return**  $\hat{\mathcal{K}} = \mathcal{PD}$  for COMP, or  $\hat{\mathcal{K}} = \mathcal{DD}$  for DD.
- 

With these definitions in place, the settings that we focus on are (i) the unconstrained sub-linear regime with exact recovery, (ii) the unconstrained linear regime with approximate recovery, and (iii) the size-constrained sub-linear regime with exact recovery. More specifically, in the second of these, we only consider the FNR; the FPR is not considered, since the DD algorithm that we study never declares a non-defective item to be defective.<sup>1</sup> While the above definitions lead to  $2^3 = 8$  possible settings of interest, our focus is on three that we believe to be suitably representative and of the most interest given what is already known in existing works (e.g., due to the hardness of exact recovery in the linear regime [17], [18]).

**Notation.** Throughout the paper, the function  $\log(\cdot)$  has base  $e$ , and we make use of Bachmann-Landau asymptotic notation (i.e.,  $O$ ,  $o$ ,  $\Omega$ ,  $\omega$ ,  $\Theta$ ).

## B. Related Work

We focus on non-adaptive and noiseless group testing with a combinatorial prior. We begin by introducing two common decoding algorithms, Combinatorial Orthogonal Matching Pursuit (COMP) and Definite Defectives (DD), in Algorithm 1. A key difference between the COMP and DD algorithm is that the COMP algorithm produces only false positives (i.e., no false negatives), while the DD algorithm produces only false negatives (i.e., no false positives).

Next, we introduce some test designs [2, Section 1.3], which will be useful for purposes of comparison later:

- **Bernoulli design:** Each item is randomly included in each test independently with some fixed probability.
- **Near-constant tests-per-item:** Each item is included in some fixed number of tests, with the tests for each item chosen uniformly at random *with replacement*, independent from the choices for all other items.
- **Constant tests-per-item:** Each item is included in some fixed number of tests, with the tests for each item chosen uniformly at random *without replacement*, independent from the choices for all other items.
- **Doubly-regular design:** Both the number of tests-per-item and the number of items-per-test are fixed to pre-specified values. Previous works predominantly considered the uniform distribution over all designs satisfying these conditions, though our own results will use a slightly different block-structured variant from [20].

<sup>1</sup>See also Appendix C for a result regarding the COMP algorithm with only false positives and no false negatives.

We proceed to review the related work for each setting.

1) *Unconstrained Sub-Linear Regime*: In the unconstrained setting with sub-linear sparsity, the following number of tests multiplied with  $(1 + \epsilon)$  (where  $\epsilon$  is any positive constant) are sufficient to attain asymptotically vanishing error probability:

- Bernoulli testing & COMP decoding [21], [22]:  $ek \log n \approx 2.72k \log n$ ;
- Bernoulli testing & DD decoding [21], [22]:  $e \max\{\theta, 1 - \theta\}k \log n \approx 2.72 \max\{\theta, 1 - \theta\}k \log n$ ;
- Near-constant tests-per-item & COMP decoding [8]:  $\frac{k \log n}{\log^2 2} \approx 2.08k \log n$ ;
- Near-constant tests-per-item & DD decoding [8]:  $\frac{\max\{\theta, 1 - \theta\}}{\log^2 2} k \log n \approx 2.08 \max\{\theta, 1 - \theta\}k \log n$ .

These results indicate that the near-constant tests-per-item is superior to Bernoulli testing, with the intuition being that the former avoids over-testing or under-testing items. We also observe that DD decoding is superior to COMP decoding, with the intuition being that the information from positive tests is “wasted” in the latter. Further improvements for information-theoretically optimal decoding are discussed below.

Additionally, converse results have been proven for each of these test designs: In the sub-linear regime with unconstrained tests, *any* decoding algorithm that uses the following number of tests multiplied by  $1 - \epsilon$  (where  $\epsilon$  is arbitrarily small) is unable to attain asymptotically vanishing error probability:

- Bernoulli testing [21]:  $(\log 2 \cdot \max_{\nu > 0} \min \{H_2(e^{-\nu}), \frac{\nu e^{-\nu}}{\log 2} \frac{1 - \theta}{\theta}\})^{-1} (1 - \theta)k \log n$ ;
- Near-constant tests-per-item [8], [9]:  $\max \left\{ \frac{\theta}{\log^2 2}, \frac{1 - \theta}{\log 2} \right\} k \log n$ .

For both designs, the DD algorithm’s performance matches the converse for  $\theta \geq \frac{1}{2}$ .

For information-theoretically optimal decoding, exact thresholds on the required number of tests were characterized for Bernoulli testing in [7], [21], and for the near-constant tests-per-item design in [8], [9]. Perhaps most importantly among these, it was shown in [10] that the above converse for the near-constant tests-per-item design extends to *arbitrary* non-adaptive designs, and that a matching achievability threshold holds for a certain spatially-coupled test design with *polynomial-time decoding*. Hence,  $\max \left\{ \frac{\theta}{\log^2 2}, \frac{1 - \theta}{\log 2} \right\} k \log n$  is the optimal threshold for all  $\theta \in (0, 1)$ , and for  $\theta \geq \frac{1}{2}$  the DD algorithm with near-constant tests-per-item is asymptotically optimal.

Mézard *et al.* [11] considered doubly-regular designs, and designs with constant tests-per-item only. Their analysis used heuristics from statistical physics to suggest that such designs can improve on Bernoulli designs, and match the above near-constant tests-per-item bound for the DD algorithm. The analysis in [11] contains some non-rigorous steps; in particular, they make use of a “no short loops” assumption that is only verified for  $\theta > \frac{5}{6}$  and conjectured for  $\theta \geq \frac{2}{3}$ , while experimentally being shown to fail for certain smaller values. One of our contributions in this paper is to establish a rigorous version of their result for  $\theta \geq \frac{1}{2}$ .

A distinct line of works has sought designs that not only require a low number of tests, but also near-optimal decoding complexity (e.g.,  $k$  poly( $\log n$ )) [23]–[28]. However, our focus

in this paper is on the required number of tests, for which the existing guarantees of such algorithms contain loose constants or extra logarithmic factors.

2) *Linear Regime*: Under the exact recovery guarantee, the known methods for deriving achievability bounds on  $T$  in the unconstrained sub-linear regime do not readily extend to the linear regime. In fact, as the following results assert, individual testing is optimal for exact recovery in the linear regime.

- **Weak converse [17]**: In the linear regime with prevalence  $p \in (0, 1)$ , if we use  $T < n - 1$  tests, there exists  $\epsilon = \epsilon(p) > 0$  independent of  $n$  such that  $P_e \geq \epsilon$ .
- **Strong converse [18]**: In the linear regime with any fixed  $p \in (0, 1)$ , if  $T \leq (1 - \epsilon)n$  for some constant  $\epsilon > 0$ , then  $P_e \rightarrow 1$  as  $n \rightarrow \infty$ .

These results imply that individual testing is an asymptotically optimal non-adaptive strategy for exact recovery. Hence, for this regime, we will instead investigate the FNR.

A recent study of two-stage adaptive algorithms in [14] turns out to be relevant to our setup (albeit not directly applicable in our non-adaptive setup). There, the high-level approach was the following, which was also considered in earlier works (e.g., see [29], [30]): (i) Conduct non-adaptive testing and identify a set of definitely non-defective items, leaving only the possibly defective items. (ii) Conduct individual testing on the remaining possibly defective items. It was shown in [14] that using a doubly-regular design from [20] in the first stage gives us the lowest expected number of tests required to attain zero error probability, with strict improvements over the near-constant tests-per-item design. This motivates us to analyze the FNR of the DD algorithm with a doubly-regular design.

3) *Size-Constrained Sub-Linear Regime*: The results most relevant to this setting are described as follows, where  $k = \Theta(n^\theta)$  with  $\theta \in [0, 1)$  throughout:

- **Converse [12]**: For  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  with  $\beta \in [0, 1)$ , and an arbitrarily small  $\epsilon > 0$ , any non-adaptive algorithm with error probability at most  $\epsilon$  requires  $T \geq \frac{1 - 6\epsilon}{1 - \beta} \cdot \frac{n}{\rho}$ , for sufficiently large  $n$ .
- **Improved Converse for  $\beta = 0$  [13]**: For  $\rho = \Theta(1)$  satisfying  $\rho \geq 1 + \lfloor \frac{\theta}{1 - \theta} \rfloor$ , if  $T \leq (1 - \epsilon) \max \left\{ \left(1 + \lfloor \frac{\theta}{1 - \theta} \rfloor\right) \frac{n}{\rho}, \frac{2n}{\rho + 1} \right\}$  for some constant  $\epsilon > 0$ , then any non-adaptive algorithm fails (with  $1 - o(1)$  probability if  $\frac{\theta}{1 - \theta}$  is a non-integer, and with  $\Omega(1)$  probability otherwise).
- **Achievability [12]**: Under a doubly-regular random test design and the COMP algorithm, for  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  with  $\beta \in [0, 1)$ , and an arbitrarily small  $\epsilon > 0$ , the error probability is asymptotically vanishing when  $T \geq \left\lceil \frac{1 + \epsilon}{(1 - \theta)(1 - \beta)} \right\rceil \cdot \left\lceil \frac{n}{\rho} \right\rceil$ .
- **Improved Achievability for  $\beta = 0$  [13]**: In the regime  $\rho = \Theta(1)$ , under a suitably-chosen near-regular random test design (which slightly differs depending on whether or not  $\theta \geq \frac{1}{2}$ ), the error probability is asymptotically vanishing when  $T \geq (1 + \epsilon) \max \left\{ \left(1 + \lfloor \frac{\theta}{1 - \theta} \rfloor\right) \frac{n}{\rho}, \frac{2n}{\rho + 1} \right\}$ . This is achieved using the Sequential COMP (SCOMP) algorithm [6], which starts with the DD solution and then iteratively refines it.

The above results for  $\beta = 0$  (i.e.,  $\rho = \Theta(1)$ ) strictly improve on the results for general  $\beta$ , and enjoy matching achievability

and converse thresholds. Thus, it is natural to ask whether we can attain similar improvements for the case that  $\rho = \Theta((n/k)^\beta)$ , for  $\beta = (0, 1)$  (i.e., large  $\rho$ ). We partially answer this question in the affirmative.

Regarding our use of a doubly-regular design, the column weight restriction is not strictly imposed by the testing constraints, but helps in avoiding “bad” events where some items are not tested enough (or even not tested at all). For example, in the case of the COMP algorithm, the doubly-regular design helps to reduce the number of tests by a factor of  $O(\log n)$  compared to i.i.d. testing.

Additional results are given for the adaptive setting in [13], [31], and for the noisy non-adaptive setting in [12, Section 7.2]. Another notable type of sparsity constraint is that of finitely divisible items (i.e., bounded tests-per-item) constraint, which are studied in [12], [13], [31]. The main reason that we do not consider such constraints here is that it was already studied extensively in [13] without any analog of the above-mentioned restrictive assumption  $\rho = O(1)$ .

### C. Contributions

Our main contributions are as follows:

- **Unconstrained sub-linear regime with exact recovery:** We provide an achievability result for the DD algorithm with a doubly-regular design that matches a result of [8] for the DD algorithm with a near-constant tests-per-item design, which is asymptotically optimal when  $\theta \geq \frac{1}{2}$ . Thus, for this range of  $\theta$ , we provide a rigorous counterpart to the result shown heuristically in [11].
- **Unconstrained linear regime with approximate recovery:** We provide an asymptotic bound on the FNR for the DD algorithm with a doubly-regular design, and further characterize the low-sparsity limit analytically, while evaluating various higher-sparsity regimes numerically.
- **Size-constrained sub-linear regime with exact recovery:** Motivated by recently-shown gains for the regime  $\rho = O(1)$  [13], we show that analogous gains are also possible for more general  $\rho = o(\frac{n}{k})$ . We improve on the best known achievability and converse bounds in such regimes, in particular using the DD algorithm to improve over known results for the COMP algorithm.

Our analysis techniques build on the existing works outlined above, but also come with several new aspects and challenges; specific comparisons are deferred to Remarks 1 and 2.

## II. MAIN RESULTS

We first describe a randomized construction of a doubly-regular  $T \times n$  test matrix  $X$ , with  $T = \frac{rn}{s}$ , where  $r$  and  $s$  are variables to be chosen according to the setting being studied. We select the test matrix in the following manner:

- Sample  $r$  matrices  $X_1, \dots, X_r$  independently, where each matrix  $X_j$  ( $j \in \{1, \dots, r\}$ ) is sampled uniformly from all  $\frac{n}{s} \times n$  binary matrices with exactly  $s$  items per test (i.e., a row weight of  $s$ ) and one item per column (i.e., a column weight of one).<sup>2</sup>

<sup>2</sup>We perform our analysis assuming that  $\frac{n}{s}$  is an integer, since the effect of rounding is asymptotically negligible.

- Form  $X$  by concatenating  $X_1, \dots, X_r$  vertically.

In other words, we perform  $r$  independent rounds of testing, where each round randomly partitions the items into  $\frac{n}{s}$  tests of size  $s$ . This approach was proposed in [20], and was used as the first step of a two-stage procedure in [14]. As we hinted in the previous section, it is distinct from the design that follows the uniform distribution over *all* matrices with row weight  $s$  and column weight  $r$  (e.g., see [11]). The above design is considered primarily to facilitate the analysis; we expect the two designs to behave similarly, but we leave it as an open problem as to whether there exist settings in which one provably outperforms the other.

After testing the items using our test matrix, we run the DD algorithm, shown in Algorithm 1, to attain our estimate  $\widehat{K}$  of the defective set.

### A. Unconstrained Testing in the Sub-Linear Regime

We state our first main result as follows, and prove it in Section III-A.

**Theorem 1.** *Under the doubly-regular design described above with parameters  $s = \frac{n \log 2}{k}$  and  $r = c \log n$  for some constant  $c > 0$ ,<sup>3</sup> when there are  $k = \Theta(n^\theta)$  defective items with constant  $\theta \in (0, 1)$ , the DD algorithm attains vanishing error probability if*

$$T \geq (1 + \epsilon) \frac{\max\{\theta, 1 - \theta\}}{\log^2 2} k \log n, \quad (2)$$

where  $\epsilon$  is an arbitrarily small positive constant (i.e., when the constant  $c$  in  $r = c \log n$  satisfies  $c \geq (1 + \epsilon) \frac{\max\{\theta, 1 - \theta\}}{\log 2}$ , noting that  $T = \frac{rn}{s}$ ).

### B. Approximate Recovery in the Linear Regime

We state our main result for the linear regime as follows, and prove it in Section III-B. In Appendix C, we also provide a counterpart of this result for the case that there are false negatives but no false positives.

**Theorem 2.** *Under the doubly-regular design described above with parameters  $s$  and  $r$ , when there are  $k = pn$  defective items with constant  $p \in (0, 1)$ , the DD algorithm attains  $\text{FNR} \leq \min\{1, \text{FNR}_{\max}(1 + o(1))\}$ , where*

$$\text{FNR}_{\max} = \left( (1 - (1 - p)^{s-1}) + \frac{(1 - p)(1 - (1 - p)^{s-1})^r}{p} \right)^r. \quad (3)$$

The corresponding number of tests is  $T = \frac{rn}{s}$ .

Next, we pause momentarily to introduce the following definition to help us evaluate our result:

**Definition 1 (Rate).** Under the combinatorial prior with  $n$  items,  $k$  defectives and  $T$  tests, the rate is equal to  $\frac{\log_2 \binom{n}{k}}{T}$ , which measures the average number of bits of information that would be gained per test if  $\mathcal{K}$  were recovered perfectly. In the linear regime with  $k = pn$ , this can be simplified to  $\frac{nH_2(p)}{T}$

<sup>3</sup>This result holds regardless of whether we round these values up or down.

(up to a  $1 + o(1)$  equivalence), where  $H_2(p) = -p \log_2 p - (1-p) \log_2 (1-p)$  is the binary entropy function.

Returning to (3), we observe that this expression is minimized when  $r$  is large and  $s$  is small. However, this conflicts with our goal of minimizing the number of tests  $T = \frac{rn}{s}$ , which implies that we should be aiming for small  $r$  and large  $s$  instead. To balance these conflicting goals, we evaluate our results in terms of their achievable rates, subject to a maximum permissible FNR. We first partially do so analytically, and then turn to numerical evaluations.

The following corollary concerns the limit of a small proportion of defectives, i.e., the low-sparsity regime.

**Corollary 1.** *Under the setup of Theorem 2, there exist choices of  $r$  and  $s$  (depending on  $p$ ) such that, in the limit as  $p \rightarrow 0$  (after having taking  $n \rightarrow \infty$ ), we have (i)  $\text{FNR}_{\max}$  approaches 0, (ii) the rate approaches  $\log 2$ , and (iii) it holds that  $s = \frac{\log 2}{p}(1 + o(1))$  and  $r = \frac{\log(\frac{1}{p})}{\log 2}(1 + o(1))$ .*

*Proof.* We start by choosing  $s = \frac{\log 2}{p}$  (the effect of rounding is negligible as  $p \rightarrow 0$ ). By Theorem 2, we have

$$\text{FNR}_{\max} = \left( 1 - (1-p)^{s-1} + \frac{(1-p)(1 - (1-p)^{s-1})^r}{p} \right)^r \quad (4)$$

$$\stackrel{(a)}{=} \left( 1 - e^{-p(s-1)(1+o(1))} + \frac{1}{p} (1 - e^{-p(s-1)(1+o(1))})^r \right)^r \quad (5)$$

$$\stackrel{(b)}{=} \left( \frac{1+o(1)}{2} + \frac{1}{p} \left( \frac{1+o(1)}{2} \right)^r \right)^r, \quad (6)$$

where (a) applies  $(1-p)^{s-1} = e^{-p(s-1)(1+o(1))}$  as  $p \rightarrow 0$ , and (b) substitutes  $s = \frac{\log 2}{p}$ . The above expression approaches zero if  $\frac{1}{p} \left( \frac{1+o(1)}{2} \right)^r \leq \frac{1}{2} - \delta$  for some positive constant  $\delta \in (0, 0.5)$ , because this gives  $\text{FNR}_{\max} = (1 - \delta + o(1))^{\omega(1)} = o(1)$ . Making  $r$  the subject, we obtain

$$\begin{aligned} r &\geq \frac{\log p + \log(\frac{1}{2} - \delta)}{-\log 2 + o(1)} \\ &= \frac{\log(\frac{1}{p}) - \log(\frac{1}{2} - \delta)}{\log 2 + o(1)} = \frac{\log(\frac{1}{p})}{\log 2} (1 + o(1)). \end{aligned} \quad (7)$$

Hence, we have  $r = \frac{\log(\frac{1}{p})}{\log 2} (1 + o(1))$  as desired. Finally,

$$\begin{aligned} \text{rate} &= \frac{\log_2 \binom{n}{k}}{T} \stackrel{(a)}{=} \frac{s \log_2 \binom{n}{pn}}{rn} \\ &\stackrel{(b)}{=} \frac{sp \log(\frac{1}{p})}{r \log 2} (1 + o(1)) \stackrel{(c)}{=} (\log 2)(1 + o(1)), \end{aligned} \quad (8)$$

where (a) substitutes  $T = \frac{rn}{s}$ , (b) applies  $\log_2 \binom{n}{pn} = \frac{1}{\log 2} pn(1 + o(1)) \log(\frac{n}{pn})$  for  $p = o(1)$ , and (c) substitutes  $s = \frac{\log 2}{p}$  and  $r = \frac{\log(\frac{1}{p})}{\log 2} (1 + o(1))$  along with some simplifications.  $\square$

Corollary 1 is consistent with the fact that, in the sub-linear regime  $k = o(n)$ , DD can achieve a rate of  $\log 2$  bits/test for an arbitrarily small target FNR [2, Sec. 5.1]. Essentially, these

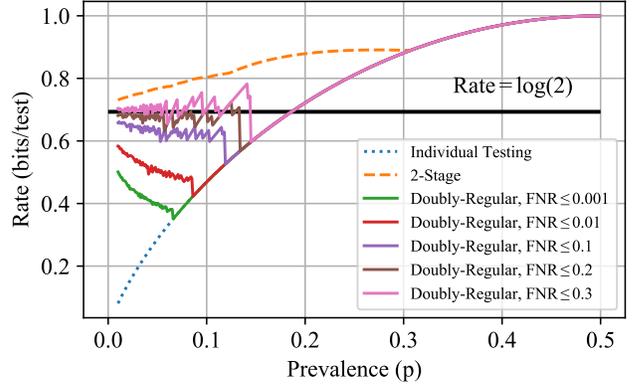


Fig. 1: Achievable rates for DD decoding with the doubly-regular design and approximate recovery (along with individual testing and a two-stage design [14], both of which attain exact recovery).

two results can both be viewed as taking the limits  $\frac{k}{n} \rightarrow 0$  and  $\text{FNR} \rightarrow 0$ , but in the opposite order.

**Numerical evaluation and comparison:** To numerically evaluate the result given in Theorem 2, we perform the following:

- 1) Select a value  $\alpha$  to be the maximum permissible FNR, and select values of  $p$  from the interval  $(0, 0.5]$  to evaluate.
- 2) For each  $p$ , numerically optimize the free parameters  $(s, r)$  to minimize the *aspect ratio*  $\frac{T}{n} = \frac{r}{s}$ , subject to  $\text{FNR} \leq \alpha$ .
- 3) Compute the rate  $\frac{nH_2(p)}{T}$ , and plot this over the chosen values of  $p$ .

The rates attained by the doubly-regular design (from Theorem 2) are shown in Figure 1.

From Figure 1, we observe that doubly-regular testing with the DD algorithm attains strictly higher rates than individual testing for smaller values of  $p$ , while reducing to individual testing for larger  $p$ . The extent of improvement increases as the target FNR increases. However, the rates obtained by conservative 2-stage testing (with exact recovery) remain higher. We are not aware of analogous results on the FNR for other non-adaptive designs such as Bernoulli or near-constant column weight, but in Appendix C we compare to those designs under COMP decoding using the FPR instead of FNR, and see that the doubly-regular design almost always outperforms them.

The discontinuities in the plot can be explained by the fact that  $r$  (tests per item) and  $s$  (items per test) must both be integers. Since  $\text{FNR}_{\max}$  is increasing in  $p$ , the pair  $(r, s)$  will change whenever  $\text{FNR}_{\max}$  exceeds  $\alpha$ . This leads to a downward jump in rate, albeit with  $\text{FNR}_{\max}$  potentially being significantly smaller than  $\alpha$ .

Figure 2 illustrates the optimal values of  $r$  and  $s$ . Generally,  $r$  is small, since  $\text{FNR}_{\max}$  decreases exponentially as  $r$  increases. Although the curve for  $s$  is not smooth, it empirically satisfies  $s = \Theta(p^{-1})$ , which is consistent with Corollary 1.

At this stage, one may wonder why some of the curves in Figure 1 appear to approach a value strictly smaller than

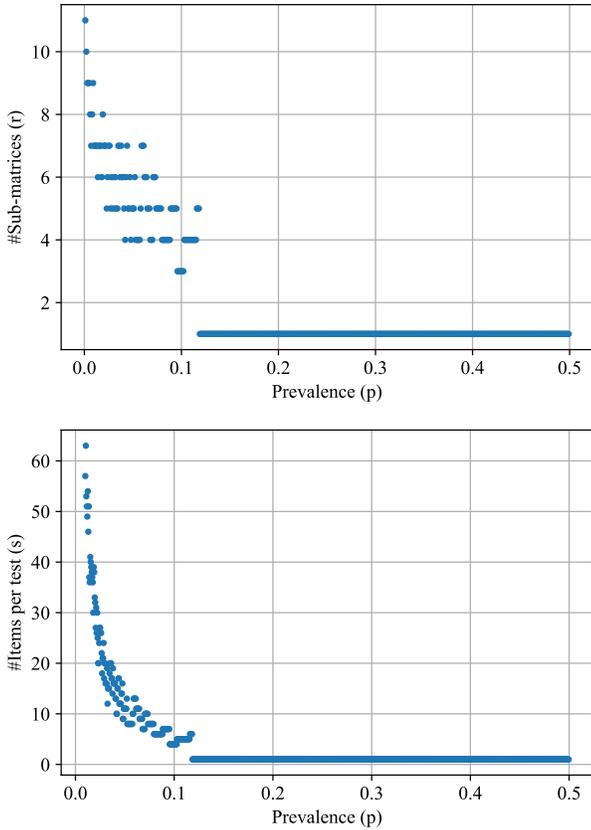


Fig. 2: Optimal  $r$  (Left) and  $s$  (Right) for  $\alpha = 0.1$ .

$\log 2 \approx 0.693$  despite Corollary 1. The reason is that this behavior is only observed for *extremely* small  $p$ , and the rate almost instantaneously drops to a significantly smaller value as  $p$  increases. In fact, we found that if we set  $s = \frac{\log 2}{p}$  and  $r = \frac{\log(\frac{1}{p})}{\log 2}$  as suggested by Corollary 1, our upper bound on the FNR exceeds one even when  $p$  is brought down to the order of  $10^{-14}$ . These findings suggest that asymptotic results for the regime  $k = o(n)$  should be interpreted with caution when it comes to practical problem sizes.

**Discussion on possible converse results.** It is difficult to gauge the tightness of our achievability result, due to the lack of converse results in this setting. While we do not attempt to make any formal statements addressing this, we believe that the constant  $\log 2$  in Corollary 1 is likely to be the best possible. This is supported by the following:

- Under the near-constant tests-per-item design, the DD algorithm is known to fail to attain  $\text{FNR} \rightarrow 0$  at rates exceeding  $\log 2$  bits/test [9]. Furthermore, it appears unlikely that any algorithm could outperform DD for the goal of attaining both  $\text{FPR} = 0$  and  $\text{FNR} \rightarrow 0$ .<sup>4</sup>
- In Appendix C, we study the “opposite” goal of attaining  $\text{FNR} = 0$  and  $\text{FPR} \rightarrow 0$ , and in that case one can rigorously show that  $\log 2$  bits/test is the best possible.

<sup>4</sup>On the other hand, an improved rate of 1 bit/test is possible when we only require  $\text{FPR} \rightarrow 0$  and  $\text{FNR} \rightarrow 0$  [7].

### C. Size-Constrained Sub-Linear Regime

We state our achievability result below, and prove it in Section III-C.

**Theorem 3.** For  $k = \Theta(n^\theta)$  with  $\theta \in [0, 1)$ , and  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  with  $\beta \in (0, 1)$ , for any integer  $r$  satisfying:

- If  $\theta \geq \frac{1}{2}$ :  $r > \frac{\theta}{(1-\theta)(1-\beta)}$  and  $r \geq \frac{2-\beta}{1-\beta}$ ;
- If  $\theta < \frac{1}{2}$ :  $r \geq \frac{1-\theta\beta}{(1-\theta)(1-\beta)}$ ;

the DD algorithm with  $T = \frac{rn}{\rho}$  tests, chosen according to the above randomized doubly-regular design with  $s = \rho$ , recovers the defective set with asymptotically vanishing error probability.

We additionally provide a converse result, which is stated as follows and proved in Section IV.

**Theorem 4.** Suppose that  $k = \Theta(n^\theta)$ , for  $\theta \in (0, 1)$ , and let  $X$  be a non-adaptive test matrix such that each test contains at most  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  items, for  $\beta \in (0, 1)$ . Given<sup>5</sup>

$$r = \max \left\{ 2, \frac{1}{1-\beta}, \left\lceil \frac{1 - (1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)} \right\rceil \right\}, \quad (9)$$

for an arbitrary constant  $\epsilon > 0$ , if there are  $(1-\epsilon)\frac{rn}{\rho}$  or fewer tests, then any decoder has error probability  $1 - o(1)$ .

The plots of the constant  $r$  against  $\theta$  are displayed in Figure 3. Our main results do not apply to the case that  $\beta = 0$  exactly (which was already handled in [13]), but we can plot the relevant limits as  $\beta \rightarrow 0$ . Similarly, the results of [13] do not apply when  $\beta > 0$ , but we can plot the relevant limits as  $\rho \rightarrow \infty$ . In Figure 3 (top-left), the same curve is obtained in both limits, and in both cases we have matching achievability and converse bounds.

For the other values of  $\beta$  shown, we observe strict improvements of DD over COMP, and the gap widens as  $\beta$  increases. For the converse, we similarly observe a strict improvement over the previous converse for  $\beta \in (0, 1)$ .

## III. ACHIEVABILITY ANALYSIS

In this section, we prove the three achievability results stated above. In general, doubly-regular designs come with more complicated dependencies that are difficult to handle. The construction that we consider (i.e., concatenating independent doubly-regular sub-matrices with column weight one) allows us to simplify the analysis of the DD algorithm by extending the analysis of one sub-matrix to the entire test matrix.

We proceed by outlining the key steps of the analysis and introducing the relevant notation. This applies to all three settings that we consider, and the differences lie in the specific details (e.g., the choices of  $r$  and  $s$ ) and the subsequent parts that extend from these steps. The key steps are:

- 1) **Determine concentration of #non-defective items in  $\mathcal{PD}$ :** Let  $\mathcal{G}$  be the set of non-defective items in  $\mathcal{PD}$ , and  $G = |\mathcal{G}|$ . Furthermore, for each non-defective  $i$ , let  $\text{PD}_i = \mathbb{1}\{i \in \mathcal{PD}\}$ . Then,  $\sum_{i \in [n] \setminus \mathcal{K}} \text{PD}_i = G$ . This step concerns the concentration of  $G$  around its mean.

<sup>5</sup>The  $\frac{1}{1-\beta}$  term comes from the converse in [12], and need not be integer-valued.

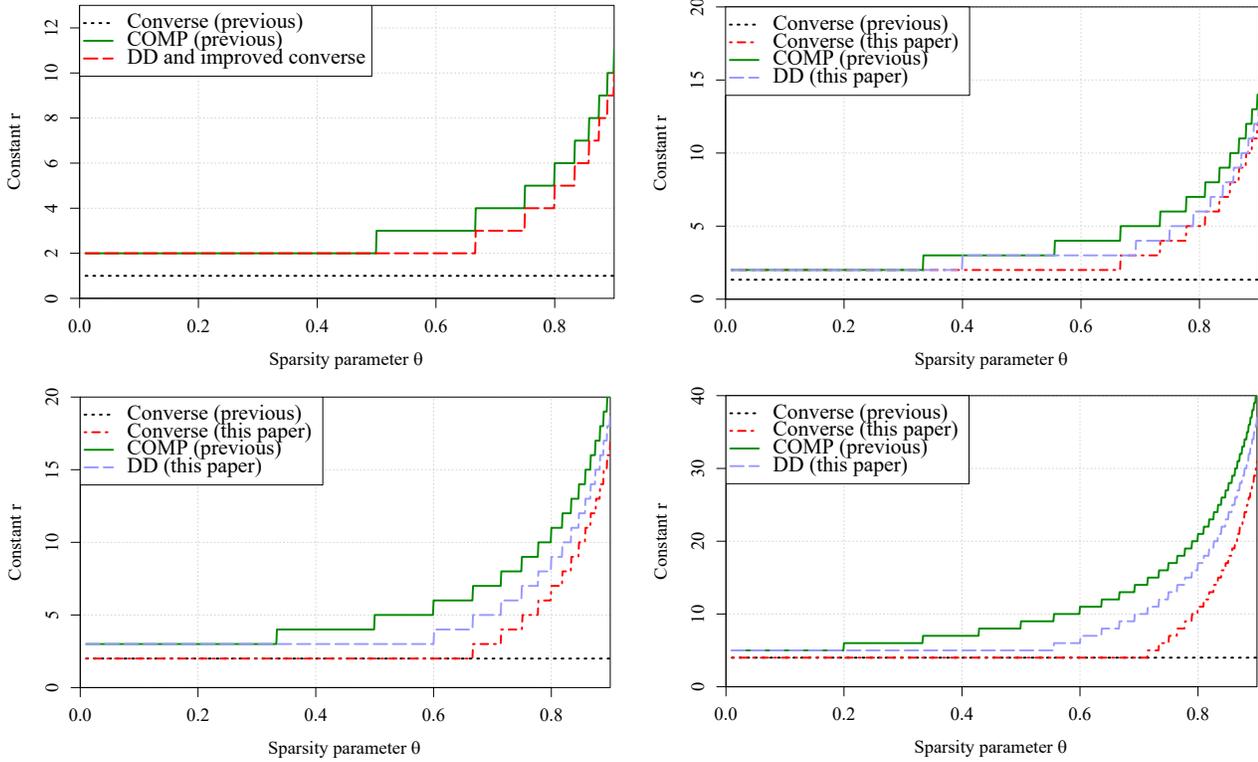


Fig. 3: Plots of  $r$  vs.  $\theta$  for  $\beta \rightarrow 0$  (top-left),  $\beta = 0.25$  (top-right),  $\beta = 0.5$  (bottom-left), and  $\beta = 0.75$  (bottom-right).

We start by considering the probability that a particular non-defective item  $i$  is found in a negative test in a single sub-matrix  $X_j$ . Consider the unique test which item  $i$  is found in. The probability that all the other  $s - 1$  items in the test are non-defective is

$$\begin{aligned} \frac{\binom{n-k-1}{s-1}}{\binom{n-1}{s-1}} &= \frac{(n-k-1)!(n-s)!}{(n-1)!(n-k-s)!} \\ &= \prod_{i=1}^{s-1} \frac{n-i-k}{n-i} = \prod_{i=1}^{s-1} \left(1 - \frac{k}{n-i}\right). \end{aligned} \quad (10)$$

Hence, consider all of the  $r$  sub-matrices, we obtain

$$\begin{aligned} \mathbb{E}[\text{PD}_i] &\stackrel{(a)}{=} \mathbb{P}[i \in \mathcal{PD}] \\ &\stackrel{(b)}{=} \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k}{n-i}\right)\right)^r, \end{aligned} \quad (11)$$

where (a) applies  $\mathbb{E}[\mathbb{1}\{A\}] = \mathbb{P}[A]$ , and (b) applies (10). Hence, we have

$$\begin{aligned} \mathbb{E}[G] &= (n-k)\mathbb{E}[\text{PD}_i] \\ &= (n-k) \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k}{n-i}\right)\right)^r. \end{aligned} \quad (12)$$

Depending on the setting, we use an appropriate concentration inequality (e.g., Chebyshev's inequality) to show that  $G$  is close to this value with probability  $1 - o(1)$ .

- 2) **Determine concentration of #masked defective items:** We introduce two further definitions.

**Definition 2.** Consider a defective item  $i$  and a set  $\mathcal{L}$

that does not include  $i$ . We say that defective item  $i$  is *masked* by  $\mathcal{L}$  if every test that includes  $i$  also includes at least one item from  $\mathcal{L}$ .

**Definition 3.** We call an item  $i$  *masked* if it is masked by  $\mathcal{K} \setminus \{i\}$ .<sup>6</sup> If the matrix is not specified then this property is defined with respect to the full test matrix  $X$ , but we will also use the same terminology with respect to a given sub-matrix  $X_j$ .

For a given defective item  $i$ , let  $M_i^j = \mathbb{1}\{i \text{ is masked in } X_j\}$ , and let  $M^j = \sum_{i \in \mathcal{K}} M_i^j$  be the number of masked defective items in  $X_j$ . Without loss of generality (WLOG), suppose that items  $1, \dots, k$  are defective. We have

$$\mathbb{P}[M_i^j = 0] = \frac{\binom{n-k}{s-1}}{\binom{n-1}{s-1}} = \prod_{i=1}^{s-1} \left(1 - \frac{k-1}{n-i}\right), \quad (13)$$

where the last equality uses the same steps as (10). Summing over all  $k$  defective items, it follows that

$$\begin{aligned} \mathbb{E}[M^j] &= k(1 - \mathbb{P}[M_i^j = 0]) \\ &= k \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k-1}{n-i}\right)\right). \end{aligned} \quad (14)$$

Depending on the setting, we use an appropriate technique to show that  $M^j$  concentrates around this value for all  $X_j$  simultaneously, with probability  $1 - o(1)$ .

- 3) **Establish conditional independence:** In this step, we condition on the preceding high-probability bounds on

<sup>6</sup>Here we are concerned with the case that  $i$  is defective, but later we will use this definition where  $i$  is non-defective, and hence  $\mathcal{K} \setminus \{i\} = \mathcal{K}$ .

$G$  and  $M^j$  holding. To facilitate the analysis, it is useful to not only condition on such events, but to condition on more specific events that ensure such concentration (apart from these, one final conditioning event will also be given below):

- We condition on a fixed set of tests being positive, and the remaining tests being negative. This fixed set has no explicit constraints.
- We condition on a fixed realization of the defective set  $\mathcal{K}$ .
- We condition on  $\mathcal{G}$  (the non-defectives that are marked as possibly defective) being a fixed set with some fixed size  $G$ . This value of  $G$  is assumed to satisfy the above-established concentration behavior.
- For  $j = 1, \dots, r$ , we condition on the *number* of masked defectives in sub-matrix  $X_j$  being  $M^j$ , whose values again satisfy our established concentration results. Note that unlike with  $\mathcal{G}$ , we do not condition on the specific set of masked item indices, but instead only on the total number.

Once we show that the conditional error probability is suitably small, it follows that the same is true after averaging (over all realizations of positive tests, all possible  $\mathcal{K}$ , and so on).

By the symmetry of the test design, without loss of generality, we can consider the following:

- The defective items are indexed by  $1, \dots, k$ .
- The items in  $\mathcal{G}$  indexed by  $k+1, k+2, \dots, k+G$ .

A slight issue here is that if we naively condition on the above sets of size  $G$  and  $M^j$ , the resulting sub-matrices  $X_1, \dots, X_r$  will typically not be conditionally independent, since the sub-matrices are coupled via  $\mathcal{G}$ . Specifically, conditioning on  $\mathcal{G}$  amounts to conditioning on (i) items  $k+1, k+2, \dots, k+G$  only being in positive tests in *all* sub-matrices, and (ii) each of items  $k+G+1, k+G+2, \dots, n$  being in a negative test in *some* sub-matrix. The latter of these means, for example, that if we are told that item  $k+G+1$  is in a negative test in submatrix 1, it reduces the probability of the same being true in submatrix 2 (thus violating independence).

To overcome this difficulty, we additionally condition on the *fixed and specific* placements of items  $k+G+1, k+G+2, \dots, n$  into negative tests (but not positive tests). This must be done subject to each of them being in some negative test, but the placements are otherwise arbitrary and do not impact our analysis.

When this additional conditioning is done, the overall conditioning involving  $G$  and  $M^j$  becomes an “AND” of  $r$  sub-events (one per sub-matrix), with each sub-event requiring that (i)  $M^j$  defectives are masked in sub-matrix  $j$ , (ii) items  $k+1, k+2, \dots, k+G$  are only included in positive tests, and (iii) items  $k+G+1, k+G+2, \dots, n$  are included in the specific negative tests indicated. Due to this “AND” structure, conditional independence is maintained among the  $r$  sub-matrices. For completeness, we provide the mathematical details of this finding in Appendix A.

- 4) **Bound the total error probability:** Recall that  $M_i^j$  is the indicator event that defective  $i$  is masked by  $\mathcal{K} \setminus \{i\}$  in  $X_j$ , and let  $MG_i^j$  be the event that defective item  $i$  is masked by  $\mathcal{G}$  in  $X_j$ . For a given  $X_j$ , we have the following, where  $\mathcal{E}$  is the event described in step 3 above and the subscript  $\mathbb{P}_{\mathcal{E}}[\cdot]$  indicates conditioning:

$$\begin{aligned} & \mathbb{P}_{\mathcal{E}}[M_i^j \cup MG_i^j] \\ &= \mathbb{P}_{\mathcal{E}}[M_i^j] + \mathbb{P}_{\mathcal{E}}[MG_i^j \cap \neg M_i^j] \end{aligned} \quad (15)$$

$$= \mathbb{P}_{\mathcal{E}}[M_i^j] + \mathbb{P}_{\mathcal{E}}[\neg M_i^j] \mathbb{P}_{\mathcal{E}}[MG_i^j | \neg M_i^j] \quad (16)$$

$$\stackrel{(a)}{=} \frac{M^j}{k} + \left(1 - \frac{M^j}{k}\right) \mathbb{P}_{\mathcal{E}}[MG_i^j | \neg M_i^j] \quad (17)$$

$$\stackrel{(b)}{\leq} \frac{M^j}{k} + \left(1 - \frac{M^j}{k}\right) \frac{G}{k - M^j} \quad (18)$$

$$= \frac{M^j}{k} + \frac{G}{k}, \quad (19)$$

where:

- (a) uses the fact that that we conditioned on having  $M^j$  masked defective items, so by the symmetry of the randomized test design, we have  $\mathbb{P}[M_i^j] = \frac{M^j}{k}$ .
- (b) holds because conditioned on  $\neg M_i^j$ , item  $i$  is in one of the  $k - M^j$  tests with a single defective item (recall that  $k - M^j$  is the number of non-masked defective items). Again using symmetry, each non-defective item has at most  $\frac{1}{k - M^j}$  probability of being in that particular test (due to there existing  $k - M^j$  equally likely options). Taking the union bound over all items in  $\mathcal{G}$  gives the desired bound.

By (19) and the conditional independence of the  $X_j$ 's, the probability of defective item  $i$  being masked by  $\mathcal{P}\mathcal{D} \setminus \{i\}$  in all sub-matrices satisfies

$$\prod_{j=1}^r \mathbb{P}_{\mathcal{E}}[M_i^j \cup MG_i^j] \leq \left(\frac{M^j}{k} + \frac{G}{k}\right)^r. \quad (20)$$

**Remark 1.** (*Comparison to existing techniques*) Our approach is distinct from existing analyses of the DD algorithm [6], [8], [13], which use a “globally symmetric” random matrix with no sub-matrix block structure. Like all of these works, we still adopt the the high-level steps of characterizing  $G$  and then characterizing  $M^j$  conditioned on  $G$ , but the details are largely different, including the unique aspect of transferring results regarding simpler sub-matrices to the entire test matrix. Moreover, regarding the size-constrained regime, we note that the analysis of [13] relies strongly on the assumption  $\rho = O(1)$ , which is what led us to adopting a distinct approach. Compared to [13], we believe that our test design and analysis are relatively simpler, though the regimes handled are different and neither subsumes the other ( $\beta = 0$  vs.  $\beta > 0$ ).

#### A. Unconstrained Sub-Linear Regime With Exact Recovery (Theorem 1)

We follow the four steps described above to obtain our result, and begin by selecting the relevant parameters. We choose  $s = \frac{n \log 2}{k}$  and  $r = c \log n$ , where  $c$  is a constant to

be specified shortly<sup>7</sup>. This results in each  $X_j$  being a  $\frac{k}{\log 2} \times n$  sub-matrix, and  $X$  being a  $\frac{ck \log n}{\log 2} \times n$  matrix.

**Step 1:** Setting  $s = \frac{n \log 2}{k}$  and  $r = c \log n$  in (12), we have

$$\mathbb{E}[G] = \Theta \left( n \left( 1 - \prod_{i=1}^{(n/k) \log 2 - 1} \left( 1 - \frac{k}{n-i} \right) \right)^{c \log n} \right) \quad (21)$$

$$\stackrel{(a)}{=} \Theta \left( n \left( 1 - \left( 1 - \frac{k}{n(1+o(1))} \right)^{(n/k) \log 2 - 1} \right)^{c \log n} \right) \quad (22)$$

$$\stackrel{(b)}{=} \Theta \left( n \left( 1 - \exp \left( - \frac{k n \log 2}{n} (1 + o(1)) \right) \right)^{c \log n} \right) \quad (23)$$

$$= \Theta \left( n \left( \frac{1}{2} (1 + o(1)) \right)^{c \log n} \right) \quad (24)$$

$$= \Theta \left( n^{1 - (c \log 2)(1 + o(1))} \right), \quad (25)$$

where (a) uses  $i \leq \frac{n \log 2}{k} - 1 = o(n)$ , and (b) uses the fact  $(1+a)^b = e^{ab(1+o(1))}$  when  $|a| = o(1)$ . Applying Markov's inequality, we have

$$\mathbb{P} \left[ G \geq \frac{k}{\log n} \right] \leq \frac{\mathbb{E}[G]}{k/\log n} \stackrel{(a)}{=} \frac{\Theta(n^{1 - (c \log 2)(1 + o(1))})}{k^{1 - o(1)}} \stackrel{(b)}{=} \Theta(n^{1 - c \log 2 - \theta + o(1)}), \quad (26)$$

where (a) uses (25) and  $\frac{k}{\log n} = k^{1 - o(1)}$ , and (b) substitutes  $k = \Theta(n^\theta)$  and further simplifies the expression. Observe that the power of  $n$  is below zero when  $c \geq (1 + \epsilon) \frac{1 - \theta}{\log 2}$ , where  $\epsilon$  is any positive constant. This implies that  $G = o(k)$  with probability  $1 - o(1)$  when  $c \geq (1 + \epsilon) \frac{1 - \theta}{\log 2}$ .

**Step 2:** Setting  $s = \frac{n \log 2}{k}$  and  $r = c \log n$  in (14), we obtain

$$\mathbb{E}[M^j] = k \left( 1 - \prod_{i=1}^{(n/k) \log 2 - 1} \left( 1 - \frac{k-1}{n-i} \right) \right) = \frac{k}{2} (1 + o(1)), \quad (27)$$

where the last equality uses the same steps as those in (21)–(24) to simplify the product to  $\frac{1}{2}(1 + o(1))$ . Next we have

$$\begin{aligned} \text{Var}[M^j] &= \text{Var} \left[ \sum_{i=1}^k M_i^j \right] \\ &= k \text{Var}[M_i^j] + k(k-1) \text{Cov}[M_1^j, M_2^j]. \end{aligned} \quad (28)$$

We proceed to evaluate the variance and covariance terms separately. The calculation is the same for any value of  $i$ , so we fix  $i = 1$  and write

$$\begin{aligned} \text{Var}[M_1^j] &= \mathbb{E}[(M_1^j)^2] - \mathbb{E}[M_1^j]^2 \\ &= \mathbb{P}[M_1^j = 1] - \mathbb{P}[M_1^j = 1]^2 \\ &= \frac{1}{4} (1 + o(1)), \end{aligned} \quad (29)$$

since the equality  $\mathbb{P}[M_1^j = 1] = \frac{1}{2}(1 + o(1))$  is implicit in (27). Next, we define  $S_{12}^j$  to be the event that defective items

1 and 2 are in the same test in  $X_j$ , yielding

$$\mathbb{P}[S_{12}^j] = \frac{\binom{n-2}{\frac{n \log 2}{k} - 2}}{\binom{n-1}{\frac{n \log 2}{k} - 1}} = \frac{(n/k) \log 2 - 1}{n-1} \leq \frac{\log 2}{k}. \quad (30)$$

In addition, we have

$$\begin{aligned} \text{Cov}[M_1^j, M_2^j] &= \mathbb{E}[M_1^j, M_2^j] - \mathbb{E}[M_1^j] \mathbb{E}[M_2^j] \\ &= \mathbb{P}[M_1^j = 1, M_2^j = 1] - \mathbb{P}[M_1^j = 1] \mathbb{P}[M_2^j = 1]. \end{aligned} \quad (31)$$

Focusing on the first term, by the law of total probability, we have (32)–(36) at the top of the next page, where (a) substitutes (30) and uses  $\mathbb{P}[M_1^j = 1, M_2^j = 1 | S_{12}^j] \leq 1$ , (b) uses  $\mathbb{P}[\neg S_{12}^j] \leq 1$  and  $\mathbb{P}[A \cap B] = 1 - \mathbb{P}[\neg A \cup \neg B]$  for any two events  $A$  and  $B$ , (c) uses the inclusion-exclusion principle, and (d) calculates  $\mathbb{P}[M_1^j = 0 | \neg S_{12}^j]$  by counting the number of ways to choose the other  $\frac{n \log 2}{k} - 1$  items in the test containing item 1 (this calculation also holds for  $\mathbb{P}[M_2^j = 0 | \neg S_{12}^j]$ ). Likewise,  $\mathbb{P}[M_1^j = 0, M_2^j = 0 | \neg S_{12}^j]$  is calculated by counting the number of ways to choose the other  $\frac{n \log 2}{k} - 2$  items in the tests (rows) of item 1 and item 2. Substituting (36) into (31), we have (37)–(39) at the top of the next page, where (a) computes the probabilities in the same way as (36), and (b) follows by expanding the square and simplifying.

The combinatorial terms in (39) can be bounded using manipulations that are elementary but tedious, so we state the resulting bound here and defer the proof to Appendix B.

**Lemma 1.** *Under the preceding setup, we have  $\text{Cov}[M_1^j, M_2^j] \leq O(\max\{\frac{1}{k}, \frac{k}{n}\})$ .*

Applying (29) and Lemma 1 in (28), we obtain  $\text{Var}[M^j] = O(\max\{k, \frac{k^3}{n}\})$ . By Chebyshev's inequality, it follows that

$$\begin{aligned} \mathbb{P} \left[ \left| M^j - \mathbb{E}[M^j] \right| \geq \frac{\mathbb{E}[M^j]}{\log n} \right] &\leq \frac{\text{Var}[M^j]}{(\mathbb{E}[M^j]/\log n)^2} \\ &= \frac{O(\max\{k, \frac{k^3}{n}\}) \log^2 n}{\Theta(k^2)} = \frac{\log^2 n}{n^{\Omega(1)}}, \end{aligned} \quad (40)$$

where we used the fact that  $k = \Theta(n^\theta)$  with  $\theta \in (0, 1)$ . Hence,  $M^j = \mathbb{E}[M^j](1 + O(\frac{1}{\log n})) = \frac{k}{2}(1 + o(1))$  for all  $X_j$  simultaneously with probability

$$\left( 1 - \frac{\log^2 n}{n^{\Omega(1)}} \right)^{c \log n} \stackrel{(a)}{=} 1 - O\left( \frac{c \log^3 n}{n^{\Omega(1)}} \right) = 1 - o(1), \quad (41)$$

where (a) uses the fact  $(1+a)^b = (1+ab)(1+o(1))$  when  $|a| < 1$  and  $ab = o(1)$ .

**Step 3:** We condition on the events described in the general description of Step 3 following (14). Here the high-probability bounds dictate that  $G = o(k)$  and  $M^j = \frac{k}{2}(1 + o(1))$ .

**Step 4:** Substituting  $r = c \log n$  into (20), the conditional probability of defective item  $i$  being masked by  $\mathcal{PD} \setminus \{i\}$  in all sub-matrices is

$$\prod_{j=1}^r \mathbb{P}_{\mathcal{E}} [M_i^j \cup \text{MG}_i^j] \leq \left( \left( \frac{M^j}{k} + \frac{G}{k} \right) (1 + o(1)) \right)^{c \log n}$$

<sup>7</sup>Rounding  $r$  has a negligible effect on  $c$  (only changing by a factor of  $1 + o(1)$ ) and is thus ignored in our analysis.

$$\mathbb{P}[M_1^j = 1, M_2^j = 1] = \mathbb{P}[S_{12}^j] \mathbb{P}[M_1^j = 1, M_2^j = 1 | S_{12}^j] + \mathbb{P}[\neg S_{12}^j] \mathbb{P}[M_1^j = 1, M_2^j = 1 | \neg S_{12}^j] \quad (32)$$

$$\stackrel{(a)}{\leq} \frac{\log 2}{k} + \mathbb{P}[\neg S_{12}^j] \mathbb{P}[M_1^j = 1, M_2^j = 1 | \neg S_{12}^j] \quad (33)$$

$$\stackrel{(b)}{\leq} \frac{\log 2}{k} + 1 - \mathbb{P}[M_1^j = 0 \cup M_2^j = 0 | \neg S_{12}^j] \quad (34)$$

$$\stackrel{(c)}{\leq} \frac{\log 2}{k} + 1 - \mathbb{P}[M_1^j = 0 | \neg S_{12}^j] - \mathbb{P}[M_2^j = 0 | \neg S_{12}^j] + \mathbb{P}[M_1^j = 0, M_2^j = 0 | \neg S_{12}^j] \quad (35)$$

$$\stackrel{(d)}{=} \frac{\log 2}{k} + 1 - 2 \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1}} + \frac{\binom{n-k}{\frac{n \log 2}{k} - 1} \binom{n-k - \frac{n \log 2}{k} + 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1} \binom{n - \frac{n \log 2}{k} - 1}}, \quad (36)$$

$$\text{Cov}[M_1^j, M_2^j] \leq \frac{\log 2}{k} + 1 - 2 \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1}} + \frac{\binom{n-k}{\frac{n \log 2}{k} - 1} \binom{n-k - \frac{n \log 2}{k} + 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1} \binom{n - \frac{n \log 2}{k} - 1}} - \mathbb{P}[M_1^j = 1] \mathbb{P}[M_2^j = 1] \quad (37)$$

$$\stackrel{(a)}{=} \frac{\log 2}{k} + 1 - 2 \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1}} + \frac{\binom{n-k}{\frac{n \log 2}{k} - 1} \binom{n-k - \frac{n \log 2}{k} + 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1} \binom{n - \frac{n \log 2}{k} - 1}} - \left(1 - \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-1}{\frac{n \log 2}{k} - 1}}\right)^2 \quad (38)$$

$$\stackrel{(b)}{=} \frac{\log 2}{k} + 2 \underbrace{\left( \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-1}{\frac{n \log 2}{k} - 1}} - \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1}} \right)}_{\text{first part}} + \underbrace{\left( \frac{\binom{n-k}{\frac{n \log 2}{k} - 1} \binom{n-k - \frac{n \log 2}{k} + 1}}{\binom{n-2}{\frac{n \log 2}{k} - 1} \binom{n - \frac{n \log 2}{k} - 1}} - \frac{\binom{n-k}{\frac{n \log 2}{k} - 1}}{\binom{n-1}{\frac{n \log 2}{k} - 1}} \right)^2}_{\text{second part}}, \quad (39)$$

$$\begin{aligned} &\stackrel{(a)}{=} \left( \frac{1}{2} (1 + o(1)) + \frac{o(k)}{k} \right)^{c \log n} \\ &= \left( \frac{1}{2} (1 + o(1)) \right)^{c \log n} \\ &= \exp((-c \log 2 + o(1)) \log n) \\ &= n^{-c \log 2 + o(1)}, \end{aligned} \quad (42)$$

where (a) substitutes  $M^j = \frac{k}{2}(1 + o(1))$  and  $G = o(k)$ .

Taking the union bound over all defective items, the probability of any defective item being masked by  $\mathcal{PD} \setminus \{i\}$  in all sub-matrices is at most  $kn^{-c \log 2 + o(1)} = \Theta(n^{\theta - c \log 2 + o(1)})$ . The power of  $n$  is below zero when  $c \geq (1 + \epsilon) \frac{\theta}{\log 2}$  for some positive constant  $\epsilon$ . Together with our previous bound on  $c$  (see (26)), this requires us to choose  $c \geq (1 + \epsilon) \frac{\max\{\theta, 1 - \theta\}}{\log 2}$ , which means it suffices to have

$$T \geq (1 + \epsilon) \frac{\max\{\theta, 1 - \theta\}}{\log^2 2} k \log n. \quad (43)$$

This completes the proof of Theorem 1.

### B. Unconstrained Linear Regime with Approximate Recovery (Theorem 2)

We again follow the four-step procedure to obtain the desired result. We assume that  $s = \Theta(1)$  and  $r = \Theta(1)$ , but their values are otherwise generic. Recall that each  $X_j$  is a  $\frac{n}{s} \times n$  sub-matrix, and  $X$  is a  $\frac{rn}{s} \times n$  matrix. It suffices to prove the theorem for  $s \geq 2$ , since otherwise each  $X_j$  amounts to one-by-one testing and the FNR is trivially zero.

The following technical lemma will be used throughout the analysis.

**Lemma 2.** For any positive constants  $a$  and  $b$  (not depending on  $n$ ) satisfying  $a \leq b$ , we have

$$\prod_{i=a}^b \left(1 - \frac{k}{n-i}\right) = (1-p)^{b-a+1} \left(1 - O\left(\frac{1}{n}\right)\right), \quad (44)$$

where  $p = \frac{k}{n}$ .

*Proof.* We have

$$\begin{aligned} \prod_{i=a}^b \left(1 - \frac{k}{n-i}\right) &= \left(1 - \frac{k}{n(1 + O(\frac{1}{n}))}\right)^{b-a+1} \\ &\stackrel{(a)}{=} \left(1 - p \left(1 + O\left(\frac{1}{n}\right)\right)\right)^{b-a+1} \\ &= (1-p) \left(1 - \frac{p}{1-p} \cdot O\left(\frac{1}{n}\right)\right)^{b-a+1} \\ &\stackrel{(b)}{=} (1-p)^{b-a+1} \left(1 - O\left(\frac{1}{n}\right)\right), \end{aligned} \quad (45)$$

where (a) substitutes  $k = pn$  and uses  $(1 - O(\frac{1}{n}))^{-1} = 1 + O(\frac{1}{n})$ , and (b) applies  $(1 - \frac{p}{1-p} O(\frac{1}{n}))^{b-a+1} = 1 - O(\frac{1}{n})$ , noting that  $\frac{p}{1-p}$  and  $b-a+1$  are both constants.  $\square$

We now proceed with the main steps.

**Step 1:** Continuing from (12), we have

$$\mathbb{E}[G] = (n-k) \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k}{n-i}\right)\right)^r \quad (46)$$

$$\stackrel{(a)}{=} (n-k) \left(1 - (1-p)^{s-1} \left(1 - O\left(\frac{1}{n}\right)\right)\right)^r \quad (47)$$

$$\stackrel{(b)}{=} (n-k)(1-(1-p)^{s-1})^r \left(1 + O\left(\frac{1}{n}\right)\right), \quad (48)$$

where (a) applies Lemma 2, and (b) applies  $(1 - O(\frac{1}{n}))^r = (1 - O(\frac{1}{n}))$  for constant  $r$ . Next, we have

$$\begin{aligned} \text{Var}[G] &= \text{Var}\left[\sum_{i \in [n] \setminus \mathcal{K}} \text{PD}_i\right] \\ &= (n-k)\text{Var}[\text{PD}_1] \\ &\quad + (n-k)(n-k-1)\text{Cov}[\text{PD}_1, \text{PD}_2], \end{aligned} \quad (49)$$

where here and in the rest of this step (step 1), we assume for notational convenience that items 1 and 2 are non-defective.<sup>8</sup> We proceed to evaluate the variance and covariance terms separately. We have

$$\begin{aligned} \text{Var}[\text{PD}_1] &= \mathbb{E}[\text{PD}_1^2] - \mathbb{E}[\text{PD}_1]^2 \\ &= \mathbb{P}[\text{PD}_1 = 1] - \mathbb{P}[\text{PD}_1 = 1]^2 \\ &\stackrel{(a)}{=} ((1 - (1-p)^{s-1})^r - (1 - (1-p)^{s-1})^{2r})(1 + o(1)) \\ &= \Theta(1), \end{aligned} \quad (50)$$

where (a) uses the same steps as (10)–(11) followed by Lemma 2. Next, we have

$$\begin{aligned} \text{Cov}[\text{PD}_1, \text{PD}_2] &= \mathbb{E}[\text{PD}_1 \text{PD}_2] - \mathbb{E}[\text{PD}_1]\mathbb{E}[\text{PD}_2] \\ &= \mathbb{P}[\text{PD}_1 = 1, \text{PD}_2 = 1] \\ &\quad - (1 - (1-p)^{s-1})^{2r} \left(1 + O\left(\frac{1}{n}\right)\right). \end{aligned} \quad (51)$$

Focusing on the first term, we note that the event  $(\text{PD}_1 = 1) \cap (\text{PD}_2 = 1)$  holds if items 1 and 2 are both contained in a positive test in each sub-matrix. Since the sub-matrices are sampled independently, we can consider them separately. Let  $\text{PD}_1^j, \text{PD}_2^j$  respectively be the events that items 1, 2 are contained in a positive test in sub-matrix  $X_j$ . Then  $\mathbb{P}[\text{PD}_1 = 1, \text{PD}_2 = 1] = \mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1]^r$ .

Let  $S_{12}^j$  be the event that items 1, 2 are placed into the same test in  $X_j$ . Similar to before (see (30)), we have  $\mathbb{P}[S_{12}^j] = \frac{s-1}{n-1}$ . Conditioning on event  $S_{12}^j$ , we have:

$$\begin{aligned} \mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1 | S_{12}^j] &= 1 - \frac{\binom{n-k-2}{s-2}}{\binom{n-2}{s-2}} \\ &= (1 - (1-p)^{s-2}) \left(1 + O\left(\frac{1}{n}\right)\right), \end{aligned} \quad (53)$$

where the final equality uses the same steps as (10) followed by Lemma 2. On the other hand, under the complement event, we have (54)–(58) at the top of the next page, where (a) uses the inclusion-exclusion principle, (b) uses the same steps as (10) followed by Lemma 2, and (c) uses the fact that  $s$  and  $p$  are constant. Hence, by the law of total probability, we have

<sup>8</sup>This should not be confused with the convention  $\mathcal{K} = \{1, \dots, k\}$  and  $G = \{k+1, \dots, k+G\}$  used when analyzing the defective items in other steps. The precise indices are inconsequential and merely a matter of notational convenience, so there is no contradiction in using indices 1 and 2 for non-defectives here but for defectives in other parts.

(59)–(62) at the top of the next page, where (a) substitutes  $\mathbb{P}[S_{12}^j] = \frac{s-1}{n-1}$  along with (53) and (58), (b) applies  $\frac{s-1}{n-1} = O(\frac{1}{n})$ , and both (b) and (c) use the fact that  $s$  and  $p$  are constant. Combining (62) with (52) gives

$$\text{Cov}[\text{PD}_1, \text{PD}_2] = O\left(\frac{1}{n}\right), \quad (63)$$

since the leading terms cancel and only leave the  $O(\frac{1}{n})$  remainder. Substituting (50) and (63) into (49), we obtain

$$\text{Var}[G] \leq \Theta(n)\Theta(1) + \Theta(n^2)O(n^{-1}) = \Theta(n). \quad (64)$$

Since  $\mathbb{E}[G] = \Theta(n)$  (see (48)), it follows from Chebyshev's Inequality that

$$\mathbb{P}\left[|G - \mathbb{E}[G]| \geq \frac{\mathbb{E}[G]}{\log n}\right] \leq \frac{\text{Var}[G] \log^2 n}{(\mathbb{E}[G])^2} = \frac{\log^2 n}{n}, \quad (65)$$

implying that  $G = (n-k)(1 - (1-p)^{s-1})^r(1 + o(1))$  with probability  $1 - o(1)$ .

**Step 2:** Continuing from (14), we have

$$\begin{aligned} \mathbb{E}[M^j] &= k \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k-1}{n-i}\right)\right) \\ &= k(1 - (1-p)^{s-1})(1 + o(1)), \end{aligned} \quad (66)$$

where the last equality uses the same steps as those in (46)–(48). We can use a near-identical approach as Step 1 to establish that  $\text{Var}[M^j] = \Theta(n)$ , and again applying Chebyshev's inequality, we have  $M^j = k(1 - (1-p)^{s-1})(1 + o(1))$  for all  $X_j$  simultaneously with probability  $(1 - o(1))^r = 1 - o(1)$  (since  $r = \Theta(1)$ ).

**Step 3:** We again condition on the events described in the general description of Step 3 following (14). Here the high-probability bounds dictate that  $G = k(1 - (1-p)^{s-1})^r(1 + o(1))$  and  $M^j = k(1 - (1-p)^{s-1})(1 + o(1))$ .

**Step 4:** Continuing from (20), the conditional probability of defective item  $i$  being masked by  $\mathcal{PD} \setminus \{i\}$  in all sub-matrices (which is precisely the FNR) satisfies

$$\begin{aligned} &\prod_{j=1}^r \mathbb{P}_{\mathcal{E}}[M_i^j \cup \text{MG}_i^j] \\ &\leq \left(\left(\frac{M^j}{k} + \frac{G}{k}\right)(1 + o(1))\right)^r \end{aligned} \quad (67)$$

$$\stackrel{(a)}{=} \left((1 - (1-p)^{s-1}) + \frac{(n-k)(1 - (1-p)^{s-1})^r}{k}\right)^r \times (1 + o(1)) \quad (68)$$

$$\stackrel{(b)}{=} \left((1 - (1-p)^{s-1}) + \frac{(1-p)(1 - (1-p)^{s-1})^r}{p}\right)^r \times (1 + o(1)), \quad (69)$$

where (a) substitutes  $M^j = k(1 - (1-p)^{s-1})(1 + o(1))$  and  $G = (n-k)(1 - (1-p)^{s-1})^r(1 + o(1))$ , and (b) substitutes  $k = pn$ . This completes the proof of Theorem 2.

### C. Size-Constrained Sub-Linear Regime with Exact Recovery (Theorem 3)

In this setting, to reduce the number of constants throughout, we consider  $k = n^\theta$  and  $\rho = \left(\frac{n}{k}\right)^\beta$  (i.e., implied constants

$$\mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1 | \neg S_{12}^j] = 1 - \mathbb{P}[\text{PD}_1^1 = 0 \cup \text{PD}_2^1 = 0 | \neg S_{12}^j] \quad (54)$$

$$\stackrel{(a)}{=} 1 - \left( \mathbb{P}[\text{PD}_1^1 = 0 | \neg S_{12}^j] + \mathbb{P}[\text{PD}_2^1 = 0 | \neg S_{12}^j] - \mathbb{P}[\text{PD}_1^1 = 0, \text{PD}_2^1 = 0 | \neg S_{12}^j] \right) \quad (55)$$

$$= 1 - \left( 2 \frac{\binom{n-k-2}{s-1}}{\binom{n-2}{s-1}} - \frac{\binom{n-k-2}{s-1}}{\binom{n-2}{s-1}} \cdot \frac{\binom{n-k-s-1}{s-1}}{\binom{n-s-1}{s-1}} \right) \quad (56)$$

$$\stackrel{(b)}{=} 1 - \left( 2(1-p)^{s-1} - (1-p)^{2(s-1)} \right) \left( 1 - O\left(\frac{1}{n}\right) \right) \quad (57)$$

$$\stackrel{(b)}{=} \left( 1 - (1-p)^{s-1} \right)^2 \left( 1 + O\left(\frac{1}{n}\right) \right), \quad (58)$$

$$\mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1] = \mathbb{P}[S_{12}^j] \mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1 | S_{12}^j] + \mathbb{P}[\neg S_{12}^j] \mathbb{P}[\text{PD}_1^1 = 1, \text{PD}_2^1 = 1 | \neg S_{12}^j] \quad (59)$$

$$\stackrel{(a)}{=} \left( \frac{s-1}{n-1} \cdot (1 - (1-p)^{s-2}) + \left( 1 - \frac{s-1}{n-1} \right) \cdot (1 - (1-p)^{s-1})^2 \right) \left( 1 + O\left(\frac{1}{n}\right) \right) \quad (60)$$

$$\stackrel{(b)}{=} \left( O\left(\frac{1}{n}\right) + \left( 1 - O\left(\frac{1}{n}\right) \right) (1 - (1-p)^{s-1})^2 \right) \left( 1 + O\left(\frac{1}{n}\right) \right) \quad (61)$$

$$\stackrel{(c)}{=} \left( 1 - (1-p)^{s-1} \right)^2 \left( 1 + O\left(\frac{1}{n}\right) \right), \quad (62)$$

of one in their scaling laws), but the general case follows with only minor changes. We again follow the above four-step procedure to obtain our result, and begin by selecting the relevant parameters. We choose  $s = \rho$ , and let  $r$  be a constant (not depending on  $n$ ) to be chosen later. This results in each  $X_j$  having size  $\frac{n}{\rho} \times n$ , and  $X$  having size  $\frac{rn}{\rho} \times n$ .

**Step 1:** Substituting  $s = \rho$  into (12), we obtain

$$\mathbb{E}[G] = \Theta \left( n \left( 1 - \prod_{i=1}^{\rho-1} \left( 1 - \frac{k}{n-i} \right) \right)^r \right) \quad (70)$$

$$\stackrel{(a)}{=} \Theta \left( n \left( 1 - \left( 1 - \frac{k}{n(1-o(1))} \right)^{\rho-1} \right)^r \right) \quad (71)$$

$$\stackrel{(b)}{=} \Theta \left( n \left( \frac{k\rho}{n} \right)^r \right), \quad (72)$$

where (a) uses  $i \leq \rho - 1 = o(n)$ , and (b) uses the fact  $(1+a)^b = (1+ab)(1+o(1))$  when  $|a| < 1$  and  $ab = o(1)$  (recall that  $\rho \rightarrow \infty$  with  $\rho = o(\frac{n}{k})$ ). We now consider two cases:

- For  $\theta \geq \frac{1}{2}$ , we use Markov's inequality to obtain

$$\mathbb{P}\left[G \geq \frac{k^2\rho}{n} \log n\right] \leq \frac{\mathbb{E}[G]}{\frac{k^2\rho}{n} \log n} \leq \frac{n \left(\frac{k\rho}{n}\right)^r}{\frac{k^2\rho}{n} \log n} \stackrel{(a)}{=} \frac{n^{(1-\theta)(2-\beta)-r(1-\theta)(1-\beta)}}{\log n}, \quad (73)$$

where (a) substitutes  $\rho = (n/k)^\beta$  and  $k = n^\theta$ . Note that when  $r \geq \frac{2-\beta}{1-\beta}$ , the expression in (73) scales as  $O\left(\frac{1}{\log n}\right)$ .

- For  $\theta < \frac{1}{2}$ , we similarly have

$$\mathbb{P}\left[G \geq k^\beta \log n\right] \leq \frac{\mathbb{E}[G]}{k^\beta \log n} \leq \frac{n \left(\frac{k\rho}{n}\right)^r}{k^\beta \log n} \stackrel{(a)}{=} \frac{n^{1-r(1-\theta)(1-\beta)-\theta\beta}}{\log n}, \quad (74)$$

where (a) substitutes  $\rho = (n/k)^\beta$  and  $k = n^\theta$ . When  $r \geq \frac{1-\theta\beta}{(1-\theta)(1-\beta)}$ , the expression in (74) scales as  $O\left(\frac{1}{\log n}\right)$ . As a side-note, to see why we split  $\theta$  into two cases here, we note that  $\frac{2-\beta}{1-\beta} \leq \frac{1-\theta\beta}{(1-\theta)(1-\beta)}$  if and only if  $\theta \geq \frac{1}{2}$  (for any  $\beta \in [0, 1)$ ).

**Step 2:** Substituting  $s = \rho$  into (14), we have

$$\mathbb{E}[M^j] = k \left( 1 - \prod_{i=1}^{\rho-1} \left( 1 - \frac{k-1}{n-i} \right) \right) = \frac{k^2\rho}{n} (1 + o(1)), \quad (75)$$

where the last equality uses the same steps as those in (70)–(72). By Markov's inequality, it follows that

$$\mathbb{P}\left[M^j \geq \frac{k^2\rho}{n} \log n\right] \leq \frac{\mathbb{E}[M^j]}{\frac{k^2\rho}{n} \log n} \leq \frac{1 + o(1)}{\log n}. \quad (76)$$

Turning to the desired event for all  $r$  sub-matrices simultaneously, we have  $M^j < \frac{k^2\rho}{n} \log n$  with probability at least  $(1 - \frac{1+o(1)}{\log n})^r = 1 - O\left(\frac{1}{\log n}\right)$ , using the fact that  $r = \Theta(1)$ .

**Step 3:** We again condition on the events described in the general description of Step 3 following (14). Here the high-probability bounds dictate that

$$G \leq \begin{cases} \frac{k^2\rho}{n} \log n, & \text{if } \theta \geq \frac{1}{2} \\ k^\beta \log n, & \text{if } \theta < \frac{1}{2}, \end{cases} \quad (77)$$

and  $M^j \leq \frac{k^2\rho}{n} \log n$  for all  $X_j$  simultaneously.

**Step 4:** Continuing from (20), the probability of defective item  $i$  being masked by  $\mathcal{PD} \setminus \{i\}$  in all sub-matrices is

$$\prod_{j=1}^r \mathbb{P}_{\mathcal{E}}[M_i^j \cup MG_i^j]$$

$$\leq \left( \left( \frac{M^j}{k} + \frac{G}{k} \right) (1 + o(1)) \right)^r \quad (78)$$

$$\leq \begin{cases} \left( O\left( \frac{k\rho}{n} \log n \right) \right)^r & \text{if } \theta \geq \frac{1}{2} \\ \left( O\left( \frac{k\rho}{n} \log n + k^{-(1-\beta)} \log n \right) \right)^r & \text{if } \theta < \frac{1}{2}, \end{cases} \quad (79)$$

where we substituted  $M^j \leq \frac{k^2\rho}{n} \log n$  and  $G$  as shown in (77) and applies some simplifications.

Taking the union bound over all  $k$  defective items and equating the resulting bound with a target value of  $\frac{1}{\log n}$ , we obtain the following conditions for the desired events to hold with probability at least  $1 - \frac{1}{\log n}$ :

$$k \left( O\left( \frac{k\rho}{n} \log n \right) \right)^r \leq \frac{1}{\log n} \text{ if } \theta \geq \frac{1}{2} \quad (80)$$

$$k \left( O\left( \frac{k\rho}{n} \log n + k^{-(1-\beta)} \log n \right) \right)^r \leq \frac{1}{\log n} \text{ if } \theta < \frac{1}{2}. \quad (81)$$

In other words, there exist constants  $C_1$  and  $C_2$  such that it suffices to have

$$\begin{aligned} k \left( C_1 \frac{k\rho}{n} \log n \right)^r &\leq \frac{1}{\log n} \text{ if } \theta \geq \frac{1}{2} \\ k \left( C_2 \max \left\{ \frac{k\rho}{n} \log n, k^{-(1-\beta)} \log n \right\} \right)^r &\leq \frac{1}{\log n} \text{ if } \theta < \frac{1}{2}. \end{aligned} \quad (82)$$

Substituting  $\rho = \left(\frac{n}{k}\right)^\beta$  and  $k = n^\theta$ , and performing some simplifications, we obtain the following conditions:

- For  $\theta \geq \frac{1}{2}$ :

$$n^\theta \left( C_1 n^{-(1-\theta)(1-\beta)+o(1)} \right)^r \leq \frac{1}{\log n}. \quad (83)$$

- For  $\theta < \frac{1}{2}$ :

$$\begin{aligned} n^\theta \left( C_2 \max \left\{ n^{-(1-\theta)(1-\beta)+o(1)}, n^{-\theta(1-\beta)+o(1)} \right\} \right)^r \\ \leq \frac{1}{\log n}. \end{aligned} \quad (84)$$

Observe that for any integer  $r$  satisfying

$$r > \frac{\theta}{(1-\theta)(1-\beta)}, \quad (85)$$

the exponent of  $n$  on the left-hand side of (83) becomes negative, so that (83) is satisfied. Similarly, observe that for any integer  $r$  satisfying

$$r > \max \left\{ \frac{\theta}{(1-\theta)(1-\beta)}, \frac{1}{1-\beta} \right\}, \quad (86)$$

the exponent of  $n$  on the left-hand side of (84) becomes negative, so that (84) is indeed satisfied. Hence, together with the bounds on  $r$  from step 1, we require  $r$  to be an integer such that the following conditions hold:

- For  $\theta \geq \frac{1}{2}$ :

$$r > \frac{\theta}{(1-\theta)(1-\beta)}, \quad r \geq \frac{2-\beta}{1-\beta}. \quad (87)$$

- For  $\theta < \frac{1}{2}$ :

$$r > \max \left\{ \frac{\theta}{(1-\theta)(1-\beta)}, \frac{1}{1-\beta} \right\}, \quad r \geq \frac{1-\theta\beta}{(1-\theta)(1-\beta)}. \quad (88)$$

Note that (88) simplifies to  $r \geq \frac{1-\theta\beta}{(1-\theta)(1-\beta)}$  alone, because for any  $\theta \in [0, \frac{1}{2}]$  and  $\beta \in [0, 1)$ , we have (i)  $1-\theta\beta > \theta$ , which implies that  $\frac{1-\theta\beta}{(1-\theta)(1-\beta)} > \frac{\theta}{(1-\theta)(1-\beta)}$ , and (ii)  $\frac{1-\theta\beta}{1-\theta} > 1$ , which implies that  $\frac{1-\theta\beta}{(1-\theta)(1-\beta)} > \frac{1}{1-\beta}$ .

Finally, we used a total of  $rn/\rho$  tests and incurred a total error probability of at most  $O\left(\frac{1}{\log n}\right)$ . This completes the proof of Theorem 3.

#### IV. CONVERSE ANALYSIS FOR THE SIZE-CONSTRAINED SETTING

In this section, we prove Theorem 4. We prove the results corresponding to  $r = \lceil \frac{1-(1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)} \rceil$  and  $r = 2$  separately in Section IV-A and Section IV-B respectively; the remaining term  $\frac{1}{1-\beta}$  is already known from [12].

**Remark 2.** Compared to the achievability part, the analysis in this section builds more closely on that of the prior work [13] handling the regime  $\rho = O(1)$ . In particular, for the first part (Section IV-A), we mostly follow [13] but with different choices of parameters to carefully account for the scaling of  $\rho$ . On the other hand, for the second part (Section IV-B), more substantial changes are needed: In [13], the scaling  $\rho = O(1)$  makes it relatively easier to identify positive tests with multiple items of degree one (leading to failure), whereas in our setting with  $\rho = \omega(1)$  the argument is somewhat more lengthy and technical, though still adopts a similar high-level approach.

##### A. Part I

To reduce the number of constants throughout, we consider  $k = n^\theta$  and test sizes  $\rho = \left(\frac{n}{k}\right)^\beta = n^{\beta(1-\theta)}$  (i.e., implied constants of one in their scaling laws), but the general case follows with only minor changes. Without loss of generality, we may assume that  $\theta > \frac{1+\beta}{2+\beta}$ ; this is because if  $\theta \leq \frac{1+\beta}{2+\beta}$ , then  $\lceil \frac{1-(1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)} \rceil \leq 1$ , which implies that it will not be the maximum among the three constant terms in (9). We make use of the notion of masked items in Definition 3, and we are now interested in characterizing both masked defectives and masked non-defectives.

Let the true defectivity vector be  $\mathbf{u} \in \{0, 1\}^n$ , where  $u_i = 1$  indicates that the  $i$ -th item is defective, and  $u_i = 0$  otherwise. Furthermore, let the test outcomes be represented by the vector  $\mathbf{y} \in \{0, 1\}^T$ , where  $y_i = 1$  denotes that the  $i$ -th test is positive and  $y_i = 0$  otherwise. We also define

$$\begin{aligned} d^- &= \left\lceil \frac{1-(1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)} \right\rceil - 1, \text{ and} \\ d^+ &= \left\lceil \frac{1-(1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)} \right\rceil, \end{aligned} \quad (89)$$

and we note that our assumption  $\theta > \frac{1+\beta}{2+\beta}$  implies that  $d^- > 0$ .

We can visualize any non-adaptive group testing design matrix  $\mathbf{X}$  as a bipartite graph  $\mathcal{G} = (V \cup F, E)$  with  $|F| = T$

factor nodes  $\{a_1, \dots, a_T\}$  (tests) and  $|V| = n$  variable nodes  $\{u_1, \dots, u_n\}$  (items). An edge between an item  $u_i$  and test  $a_j$  indicates that  $u_i$  takes part in  $a_j$ . We let  $\{\partial_{\mathcal{G}} a_1, \dots, \partial_{\mathcal{G}} a_T\}$  and  $\{\partial_{\mathcal{G}} u_1, \dots, \partial_{\mathcal{G}} u_n\}$  denote the neighborhoods in  $\mathcal{G}$ ; the sparsity constraint on the test implies  $|\partial_{\mathcal{G}} a_j| \leq \rho$ . Finally, let  $V_{1+}(\mathcal{G})$  be the set of defectives that are masked,  $V_{0+}(\mathcal{G})$  be the set of non-defectives that are masked, and  $V_+(\mathcal{G}) = V_{1+}(\mathcal{G}) \cup V_{0+}(\mathcal{G})$ . We have the following auxiliary results.

**Lemma 3.** [13, Claim 2.3] *Conditioned on any non-empty realizations of the sets  $V_{1+}(\mathcal{G})$  and  $V_{0+}(\mathcal{G})$ , any decoding procedure fails with probability at least  $1 - \frac{1}{|V_{1+}(\mathcal{G})| + |V_{0+}(\mathcal{G})|}$ .*

**Lemma 4.** (Multiplicative Chernoff bound [32, Thm. 4.2]) *Suppose  $X_1, \dots, X_m$  are independent random variables taking values in  $\{0, 1\}$ . Let  $X = \sum_{i=1}^m X_i$ . Then for any  $\epsilon \in (0, 1)$ ,*

$$\mathbb{P}[X \leq (1 - \epsilon)\mathbb{E}[X]] \leq \exp\left(-\frac{\epsilon^2 \mathbb{E}[X]}{2}\right). \quad (90)$$

Our goal is to show that with  $T = (1 - \epsilon)\frac{d^+ n}{\rho}$  for some constant  $\epsilon > 0$ , we have  $|V_{1+}(G)|, |V_{0+}(G)| = \omega(1)$  with high probability (i.e., with probability  $1 - o(1)$ ), which implies (using Lemma 3) that  $P_e = 1 - o(1)$ . As a stepping stone to our result for the combinatorial prior, we first study the i.i.d. prior, whose defectivity vector we denote by  $\mathbf{u}^*$ , such that each entry is one independently with probability  $p = \frac{k - \sqrt{k \log n}}{n}$ . The following existing result allows us to transfer from the latter prior to the former.

**Lemma 5.** [13, Corollary 3.6] *Given non-negative integers  $C_1, C_2$  and fixed  $\epsilon_1, \epsilon_2 \in (0, 1)$ , if the modified model (i.i.d. prior) satisfies*

$$\begin{aligned} \mathbb{P}[|V_{1+}(\mathcal{G}, \mathbf{u}^*)| > 2C_1] &\geq 1 - \epsilon_1 \text{ and} \\ \mathbb{P}[|V_{0+}(\mathcal{G}, \mathbf{u}^*)| > 2C_2] &\geq 1 - \epsilon_2, \end{aligned} \quad (91)$$

*then the original model (combinatorial prior) satisfies*

$$\begin{aligned} \mathbb{P}[|V_{1+}(\mathcal{G}, \mathbf{u})| > C_1] &\geq 1 - \epsilon_1 - o(1) \text{ and} \\ \mathbb{P}[|V_{0+}(\mathcal{G}, \mathbf{u})| > C_2] &\geq 1 - \epsilon_2 - o(1). \end{aligned} \quad (92)$$

In view of this result, we proceed by working with  $\mathbf{u}^*$  instead of  $\mathbf{u}$ . We will also use the following key fact [9]: Whenever items have distance at least 6 in the underlying graph, the events of being masked are independent. Leveraging on this fact, we introduce a procedure for obtaining a set  $\mathcal{B}$  of size  $N$  (to be specified shortly), such that each item has a degree of at most  $d^-$ , and the pairwise distance between items is at least 6. The procedure is as follows:

- 1) Create  $\mathcal{G}_0$  from  $\mathcal{G}$  by executing the following:
  - a) Remove all items whose degree in  $\mathcal{G}$  is greater than  $\log n$ .
  - b) Identify all tests whose degree in  $\mathcal{G}$  is less than  $\frac{\rho}{\log n}$ , and simultaneously remove all of those tests and their respective items.
- 2) For  $i = 1, \dots, N$ , where  $N = \frac{n^{1-2\beta(1-\theta)}}{\log^3 n}$ :
  - a) Arbitrarily select an item in  $\mathcal{G}_{i-1}$  whose degree in  $\mathcal{G}$  is at most  $d^-$  (assuming one exists).
  - b) Extract the item and add it to the set  $\mathcal{B}$ .

- c) Remove the extracted item's tests in the first and third neighborhood, its items in the second and fourth neighborhood, and the extracted item itself. Let  $\mathcal{G}_i$  be the resulting graph.

We now proceed to analyze each of the steps.

**Analysis of Step 1:** In Step 1(a), we remove all items with degree greater than  $\log n$ . Then, in Step 1(b), we further remove tests with degree less than  $\frac{\rho}{\log n}$  and their items. We are left with a sub-graph  $\mathcal{G}_0$  of the original graph, whose nodes' degree properties turn out to be convenient for the analysis.

Before proceeding with  $\mathcal{G}_0$ , we need to understand its number of items. To do so, we investigate how many items of various kinds could have been removed from  $\mathcal{G}$ . We first note that the number of edges in  $\mathcal{G}$  contributed by items with degree greater than  $\log n$  is upper bounded by the total number of edges, which equals

$$\sum_{u \in V(\mathcal{G})} |\partial_{\mathcal{G}} u| = \sum_{a \in F(\mathcal{G})} |\partial_{\mathcal{G}} a| \leq T\rho = (1 - \epsilon)d^+ n. \quad (93)$$

Hence, the number of items in  $\mathcal{G}$  with degree greater than  $\log n$  (and hence removed in Step 1(a)) is at most  $\frac{(1-\epsilon)d^+ n}{\log n} = o(n)$ .

Next, we note that the number of tests with degree less than  $\frac{\rho}{\log n}$  removed in Step 1(b) is trivially no higher than the total number of tests  $\frac{(1-\epsilon)d^+ n}{\rho}$ . Hence, the number of edges contributed by those tests is at most  $\frac{\rho}{\log n} \cdot \frac{(1-\epsilon)d^+ n}{\rho} = o(n)$ . This implies that the number of items that are removed in Step 1(b) is at most  $o(n)$ . Combining these observations, we conclude that  $\mathcal{G}_0$  has  $n(1 - o(1))$  items, all with degree at most  $\log n$  in  $\mathcal{G}_0$ . Moreover, every test in  $\mathcal{G}_0$  has degree between  $\frac{\rho}{\log n}$  and  $\rho$  in  $\mathcal{G}$  (though their degree in  $\mathcal{G}_0$  itself could be below  $\frac{\rho}{\log n}$ ).

**Lemma 6.** *Under the preceding setup with  $T \leq (1 - \epsilon)\frac{d^+ n}{\rho}$ , the number of items in  $\mathcal{G}_0$  with degree at most  $d^-$  in  $\mathcal{G}$  is  $\Theta(n)$ .*

*Proof.* Suppose for contradiction that the number of items in  $\mathcal{G}_0$  with degree at most  $d^-$  in  $\mathcal{G}$  is  $o(n)$ . Then, the number of items in  $\mathcal{G}_0$  with degree at least  $d^+$  in  $\mathcal{G}$  is  $n(1 - o(1))$ . It follows that the number of edges (in  $\mathcal{G}$ ) contributed by these items is at least  $d^+ n(1 - o(1))$ , which implies that the number of tests required for those edges is at least

$$\frac{d^+ n(1 - o(1))}{\rho} > (1 - \epsilon)\frac{d^+ n}{\rho}, \quad (94)$$

for sufficiently large  $n$ , contradicting the assumption on  $T$ .  $\square$

**Analysis of Step 2:** Due to the degree upper bounds (in  $\mathcal{G}_0$ ) established above, in each iteration, we remove at most

$$\rho \log n + \rho^2 \log^2 n = \Theta(n^{2\beta(1-\theta)} \log^2 n) \quad (95)$$

items. This scales as  $o(n)$  when  $\theta > 1 - \frac{1}{2\beta}$  with  $\beta \in (0, 1)$ , which is satisfied due to our assumption  $\theta > \frac{1+\beta}{2+\beta} = 1 - \frac{1}{2+\beta}$  (note that  $2 + \beta \geq 2\beta$ ). In total, after  $N = \frac{n^{1-2\beta(1-\theta)}}{\log^3 n}$  iterations, we removed a number of items scaling as

$$N \cdot \Theta(n^{2\beta(1-\theta)} \log^2 n) = \Theta\left(\frac{n}{\log n}\right) = o(n). \quad (96)$$

Since we have  $\Theta(n)$  items with degree at most  $d^-$  in  $\mathcal{G}_0$  (see Lemma 6), and we remove only  $o(n)$  items (from (96)), it follows that we will always be able to extract  $N$  items, i.e., the item described in Step 2(a) is always guaranteed to exist under our assumptions.

**Analysis of the masking probability.** We now focus our attention on the set  $\mathcal{B}$ . For each item  $u \in \mathcal{B}$ , we have

$$\mathbb{P}[u \in V_+(\mathcal{G}, \mathbf{u}^*)] \stackrel{(a)}{\geq} \prod_{a \in \partial_{\mathcal{G}} u} \left(1 - (1-p)^{|\partial_{\mathcal{G}} a| - 1}\right) \quad (97)$$

$$\stackrel{(b)}{=} \prod_{a \in \partial_{\mathcal{G}_0} u} \left(1 - (1-p)^{|\partial_{\mathcal{G}} a| - 1}\right) \quad (98)$$

$$\stackrel{(c)}{=} \prod_{a \in \partial_{\mathcal{G}_0} u} \left(p(|\partial_{\mathcal{G}} a| - 1)(1 + o(1))\right) \quad (99)$$

$$\stackrel{(d)}{\geq} \left(\frac{p\rho}{\log n}(1 + o(1))\right)^{|\partial_{\mathcal{G}_0} u|} \quad (100)$$

$$\stackrel{(e)}{\geq} \left(\frac{p\rho}{\log n}(1 + o(1))\right)^{d^-}, \quad (101)$$

where:

- (a) follows from the fact that the events of  $u$  being masked in each of its tests satisfy an increasing property<sup>9</sup> [17, Lemma 4]; this is a sufficient condition for applying the FKG inequality [33], which lower bounds the probability by the expression that would hold under independence across  $a \in \partial_{\mathcal{G}} u$ . Note also that  $|\partial_{\mathcal{G}} a| \geq \frac{\rho}{\log n} \geq 2$  because our procedure (Step 1(b)) ensures that item  $u$  is not an item that is individually tested in  $\mathcal{G}$ .
- (b) applies  $\partial_{\mathcal{G}} u = \partial_{\mathcal{G}_0} u$ , which holds because transforming  $\mathcal{G}$  into  $\mathcal{G}_0$  causes  $u$  to stay in the same tests. This is because whenever we remove a test in Step 1 of the procedure, we also remove all the items in the test.
- (c) follows by first noting that  $p|\partial_{\mathcal{G}} a| \leq \frac{k\rho}{n}(1 - o(1)) = o(1)$ , where the inequality holds by substituting  $p = \frac{k}{n}(1 - o(1))$  and  $|\partial_{\mathcal{G}} a| \leq \rho$ , and the equality uses  $\rho = o(\frac{n}{k})$ . We then use a Taylor expansion in (98), followed by some simplifications.
- (d) uses the fact that  $|\partial_{\mathcal{G}} a| \geq \frac{\rho}{\log n} = \omega(1)$  for each test  $a$  in  $\mathcal{G}_0$ .
- (e) uses the fact that  $\frac{p\rho}{\log n}(1 + o(1)) = O(\frac{k\rho}{n \log n}) = \frac{n^{-\Omega(1)}}{\log n} < 1$  (implying that a higher power in the exponent makes the overall expression smaller), and  $|\partial_{\mathcal{G}_0} u| \leq |\partial_{\mathcal{G}} u| \leq d^-$  by construction in our extraction procedure.

We now turn to the number of masked defective items in  $\mathcal{B}$ . Recall that for any two items  $u, u' \in \mathcal{B}$ , the events of being masked are independent due to the pairwise distances being at least 6. Furthermore, each item is defective independently with probability  $p$  under  $\mathbf{u}^*$ , and the property of any given item being masked is independent of the item's defectivity status (as noted in [10]). Hence,  $|V_{1+}(\mathcal{G}, \mathbf{u}^*)|$  stochastically dominates Binomial( $N, p \cdot (\frac{p\rho}{\log n}(1 + o(1)))^{d^-}$ ), which has an

<sup>9</sup>Namely, marking any additional item(s) as defective can only increase (or keep unchanged) the probability of the masking event associated with each test  $a \in \partial_{\mathcal{G}} u$ .

expectation of

$$N \cdot p \cdot \left(\frac{p\rho}{\log n}(1 + o(1))\right)^{d^-} \\ = \frac{n^{1-2\beta(1-\theta)}}{\log^3 n} \cdot \frac{k}{n} \cdot \left(\frac{k\rho}{n \log n}\right)^{d^-} (1 + o(1)) \quad (102)$$

$$= \frac{n^{1-(1-\theta)(2\beta+1)-d^-(1-\theta)(1-\beta)}}{(\log n)^{3+d^-}} (1 + o(1)). \quad (103)$$

$$= \frac{n^{\Omega(1)}}{(\log n)^{3+d^-}} (1 + o(1)), \quad (104)$$

where (103) substitutes  $k = n^\theta$ ,  $\rho = (\frac{n}{k})^\beta$ , and  $p = \frac{k}{n}(1 + o(1))$ , and (104) holds since  $d^- < \frac{1-(1-\theta)(2\beta+1)}{(1-\theta)(1-\beta)}$  by definition. By Lemma 4 with  $\epsilon = 1 - \frac{2(\log n)^{4+d^-}}{n^{\Omega(1)}}(1 - o(1))$ , we find that  $|V_{1+}(\mathcal{G}, \mathbf{u}^*)| > 2 \log n$  with probability  $1 - o(1)$ . The same analysis can be repeated to show that  $|V_{0+}(\mathcal{G}, \mathbf{u}^*)| > 2 \log n$  with probability  $1 - o(1)$  (intuitively, the same readily follows for non-defectives because  $k = o(n)$ , i.e., the number of non-defectives is much higher).

Finally, we transfer our findings from the i.i.d. prior  $\mathbf{u}^*$  to the combinatorial prior  $\mathbf{u}$ . Specifically, using Lemma 5, we immediately obtain  $|V_{1+}(\mathcal{G}, \mathbf{u})|, |V_{0+}(\mathcal{G}, \mathbf{u})| > \log n = \omega(1)$  with probability  $1 - o(1)$ . By conditioning on  $|V_{1+}(\mathcal{G}, \mathbf{u})| = \omega(1)$  and  $|V_{0+}(\mathcal{G}, \mathbf{u})| = \omega(1)$ , and applying Lemma 3, we get a conditional error probability of  $1 - o(1)$ , which completes the first part of the proof of Theorem 4.

## B. Part II

We again work in the regime where  $k = n^\theta$  and  $\rho = (\frac{n}{k})^\beta = n^{\beta(1-\theta)}$ , with  $\beta \in (0, 1)$ , i.e., the constant factors are set to one. We start by stating an auxiliary result that will be used later.

**Lemma 7.** [34, Section 5] *For a random variable  $Z_1 \sim \text{Hypergeometric}(n, z, k)$  (i.e.,  $\mathbb{P}[Z_1 = z_1] = \binom{z}{z_1} \binom{n-z_1}{k-z_1} \binom{n}{k}^{-1}$ ),  $p = \frac{z}{n}$ , and  $0 < t < p$ , we have*

$$\mathbb{P}[Z_1 \leq (p-t)k] \leq e^{-kD(p-t||t)} \leq e^{-2t^2k}, \quad (105)$$

where  $D(a||b) = a \log \frac{a}{b} + (1-a) \log \frac{1-a}{1-b}$  is the binary KL divergence.

We wish to show that with  $T = (1 - \epsilon) \frac{2n}{\rho}$ ,<sup>10</sup> where  $\epsilon$  is a positive constant, no algorithm successfully recovers  $\mathbf{u}$  from  $(\mathbf{y}, \mathcal{G})$  with probability  $\Omega(1)$ . We proceed with a proof by contradiction, assuming that the opposite is true, i.e., that with  $T = (1 - \epsilon) \frac{2n}{\rho}$ , there exists an algorithm that successfully recovers  $\mathbf{u}$  from  $(\mathbf{y}, \mathcal{G})$  with probability  $\Omega(1)$ . Based on this assumption, we will show that  $\epsilon = o(1)$ , which gives us a contradiction.

**Lemma 8.** *If  $T = (1 - \epsilon) \frac{2n}{\rho}$  and the success probability is  $\Omega(1)$ , then there are at least  $2\epsilon n(1 - o(1))$  degree-one items.*

*Proof.* Let  $\alpha_0 n$  and  $\alpha_1 n$  be the number of items with degree zero and one respectively. We start by establishing that

<sup>10</sup>We may assume equality for the purpose of proving a converse, since reducing the number of tests only makes recovery even more difficult.

$\alpha_0 = o(1)$ . This is because if there were  $\Omega(n)$  degree-zero items, then with high probability there would be  $\omega(1)$  degree-zero defectives and  $\omega(1)$  degree-zero non-defectives (e.g., using the variant of Hoeffding's bound for the hypergeometric distribution [35]). Since these are impossible to distinguish from each other, this implies  $o(1)$  success probability, in contradiction with the lemma assumption.

Now, by lower bounding the number of edges contributed by the items, and upper bounding the number of edges contributed by the tests, we obtain

$$\begin{aligned} \alpha_1 n + 2(1 - \alpha_0 - \alpha_1)n &\leq \sum_{u \in V(G)} |\partial_{\mathcal{G}} u| = \sum_{a \in F(G)} |\partial_{\mathcal{G}} a| \\ &\leq T\rho = (1 - \epsilon)2n. \end{aligned} \quad (106)$$

Combining the left-most and right-most expressions and making  $\alpha_1$  the subject and applying  $\alpha_0 = o(1)$ , we obtain  $\alpha_1 \geq 2\epsilon - 2\alpha_0 = 2\epsilon(1 - o(1))$ .  $\square$

Before proceeding, we introduce some further notation:

- $\mathcal{T}'$  denotes the set of tests with at least  $\frac{\rho}{\log n}$  items of degree one, and  $T' = |\mathcal{T}'|$ . We observe that the lower bound  $T' \geq \frac{n}{\rho \log n}$  must hold, since otherwise the total number of items of degree one would be at most  $(1 - \epsilon)\frac{2n}{\rho} \cdot \frac{\rho}{\log n} + \frac{n}{\rho \log n} \cdot \rho \leq \frac{3n}{\log n}$ , contradicting Lemma 8.
- $T'_+$  denotes the number of tests in  $\mathcal{T}'$  with exactly one defective item of degree one.
- $z$  denotes the number of items with degree one that appear in the tests of  $\mathcal{T}'$ . Note that  $z$  satisfies  $z \geq \frac{T'_+ \rho}{\log n} \geq \frac{n}{\log^2 n}$  because each test in  $\mathcal{T}'$  has at least  $\frac{\rho}{\log n}$  items of degree one, and  $z \leq T'_+ \rho$  because of the  $\rho$ -sized test constraint. These inequalities will be used in our analysis later.
- $Z_1$  denotes the number of defective items with degree one that appear in the tests of  $\mathcal{T}'$ .

Note that among the four quantities introduced above, only  $T'_+$  and  $Z_1$  are random (i.e., affected by the randomness of the defective set).

**Lemma 9.** *When the number  $T'$  of tests containing at least  $\frac{\rho}{\log n}$  items of degree one is at least  $\frac{n}{\rho \log n}$ , any algorithm recovering  $\mathbf{u}$  from  $(\mathbf{y}, \mathcal{G})$  has a success probability of  $o(1)$ .*

Before proving this lemma, we state some additional auxiliary results.

**Lemma 10.** *Under the combinatorial prior with our assumed scaling laws, the size of the set  $V_{0+}$  of masked non-defective items scales as  $\omega(1)$  with probability  $1 - o(1)$ .*

*Proof.* We momentarily return to the i.i.d. prior  $\mathbf{u}^*$ , where  $p = \frac{k - \sqrt{k \log n}}{n} = \frac{k}{n}(1 - o(1))$ , and study how  $T'_+$  scales. Under  $\mathbf{u}^*$ , the number of defective items of degree one in a test containing  $\rho' \geq \frac{\rho}{\log n}$  items of degree one is distributed as  $\text{Binomial}(\rho', \frac{k}{n}(1 - o(1)))$ , which implies

$$\begin{aligned} &\mathbb{P}[1 \text{ def. item of degree one}] \\ &= \binom{\rho'}{1} \left(\frac{k}{n}(1 - o(1))\right) \left(1 - \frac{k}{n}(1 - o(1))\right)^{\rho'-1} \end{aligned} \quad (107)$$

$$\stackrel{(a)}{\geq} \frac{\rho}{\log n} \left(\frac{k}{n}(1 - o(1))\right) \left(1 - \frac{k}{n}(1 - o(1))\right)^{\rho-1} \quad (108)$$

$$\stackrel{(b)}{=} \frac{k\rho}{n \log n} \left(1 - \frac{k\rho}{n}(1 - o(1))\right) (1 - o(1)) \quad (109)$$

$$\stackrel{(c)}{=} \frac{k\rho}{n \log n} (1 - o(1)), \quad (110)$$

where (a) uses  $\rho' \geq \frac{\rho}{\log n}$  and  $\rho' \leq \rho$ , (b) uses  $\frac{k\rho}{n} = o(1)$  and a Taylor expansion (followed by some simplifications), and (c) uses  $\frac{k\rho}{n} = o(1)$ . Observe that the event of a test in  $\mathcal{T}'$  having at exactly one defective item of degree one is independent of the same event for the other tests in  $\mathcal{T}'$ , since the items are defective independently and cannot be shared across the tests (due to the degree of one). Hence,  $T'_+$  stochastically dominates  $\text{Binomial}(T', \frac{k\rho}{n \log n}(1 - o(1)))$ , which implies

$$\mathbb{E}[T'_+] \geq T' \cdot \frac{k\rho}{n \log n} \cdot (1 - o(1)) = \frac{k}{\log^2 n} (1 - o(1)), \quad (111)$$

by applying  $T' \geq \frac{n}{\rho \log n}$ . By Lemma 4 with  $\epsilon = 1 - \frac{1+o(1)}{\log n}$ , it follows that

$$\mathbb{P}\left[T'_+ \leq \frac{k}{\log^3 n}\right] \leq \exp\left(-\frac{k}{2 \log^2 n} (1 + o(1))\right) = o(1), \quad (112)$$

which implies that we have  $T'_+ > \frac{k}{\log^3 n} = \omega(1)$  with probability  $1 - o(1)$ . This further implies that  $|V_{0+}| > \frac{k}{\log^3 n} \left(\frac{\rho}{\log n} - 1\right) = \omega(1)$  (by counting the masked non-defective items with degree one in the  $T'_+$  tests) with probability  $1 - o(1)$  under the i.i.d. prior. By Lemma 5, it then follows that  $|V_{0+}| = \omega(1)$  with probability  $1 - o(1)$  under the combinatorial prior.  $\square$

**Lemma 11.** *Under the combinatorial prior, the number  $Z_1$  of defective items with degree one that appear in the tests in  $\mathcal{T}'$  scales as  $\omega(1)$  with probability  $1 - o(1)$ .*

*Proof.* Due to fact that the defective set  $\mathcal{K}$  is uniformly distributed, we have  $\mathbb{P}[Z_1 = z_1] = \binom{z}{z_1} \binom{n-z}{k-z_1} \binom{n}{k}^{-1}$ , i.e.,  $Z_1 \sim \text{Hypergeometric}(n, z, k)$ . Using this fact, we have

$$\mathbb{P}\left[Z_1 \leq \frac{\mathbb{E}[Z_1]}{\log n}\right] = \mathbb{P}\left[Z_1 \leq \frac{zk}{n \log n}\right] \quad (113)$$

$$\stackrel{(a)}{\leq} \exp\left(-2\left(\frac{z}{n}\left(1 - \frac{1}{\log n}\right)\right)^2 k\right) \quad (114)$$

$$\stackrel{(b)}{\leq} \exp\left(-2\left(\frac{1}{\log^2 n}\left(1 - \frac{1}{\log n}\right)\right)^2 k\right) \quad (115)$$

$$= \exp\left(-\frac{2k}{\log^4 n} (1 - o(1))\right), \quad (116)$$

where (a) uses Lemma 7 with  $t = \frac{z}{n}\left(1 - \frac{1}{\log n}\right)$ , and (b) uses  $z \geq \frac{n}{\log^2 n}$ . This proves that  $Z_1 > \frac{\mathbb{E}[Z_1]}{\log n}$  with probability  $1 - o(1)$ . Since  $\frac{\mathbb{E}[Z_1]}{\log n} = \frac{zk}{n \log n} \geq \frac{k}{\log^3 n}$  (by applying  $z \geq \frac{n}{\log^2 n}$ ), we have  $Z_1 > \frac{k}{\log^4 n} = \omega(1)$  with probability  $1 - o(1)$ .  $\square$

With the auxiliary results in place, we turn to the proof of Lemma 9.

*Proof of Lemma 9.* Our aim is to show that with  $T' \geq \frac{n}{\rho \log n}$ , the success probability is  $o(1)$ . From this point onwards, we condition on  $|V_{0+}| = \omega(1)$  and  $Z_1 = \omega(1)$ , both occurring

with probability  $1 - o(1)$  (see Lemma 10 and Lemma 11). Under these conditioned events, we have the following possible cases: (1)  $|V_{1+}| = 0$  and  $|V_{0+}| = \omega(1)$ ; (2)  $|V_{1+}| > 0$  and  $|V_{0+}| = \omega(1)$ . We proceed to analyze each case separately:

- **Case 1:** By assumption, we have  $Z_1 = \omega(1)$ , i.e., an unbounded number of defective items with degree one in the tests of  $\mathcal{T}'$ . The condition  $|V_{1+}| = 0$  implies that the tests those  $Z_1$  items are in cannot include any other defective items, meaning that there are  $Z_1$  corresponding tests with exactly one degree-one defective, so that  $T'_+ \geq Z_1 = \omega(1)$ . For each of these tests, the decoder can do no better than to guess the defective item uniformly at random from all degree-one items in the test, of which there are at least  $\frac{\rho}{\log n}$ . Hence, the success probability in each test is at most  $\frac{\log n}{\rho} = o(1)$ .
- **Case 2:** In this case, a direct application of Lemma 3 gives us an error probability of  $1 - o(1)$  (i.e., a success probability of  $o(1)$ ).

Since the success probability is  $o(1)$  in both cases regardless of the specific values of  $|V_{+}| = \omega(1)$  and  $Z_1 = \omega(1)$  being conditioned on, it follows that the unconditional success probability is also  $o(1)$ . This proves Lemma 9.  $\square$

In view of Lemma 9, for the assumed condition stated following Lemma 7 to hold (with an algorithm attaining success probability  $\Omega(1)$ ), it must be the case that  $T' < \frac{n}{\rho \log n}$ . We proceed by assuming that this is true, and showing that we arrive at a contradiction.

The condition  $T' < \frac{n}{\rho \log n}$  implies that the number of edges contributed by the  $T'$  tests is below  $\frac{n}{\rho \log n} \cdot \rho = \frac{n}{\log n} = o(n)$ , which further implies that  $o(n)$  items of degree one participate in the tests of  $\mathcal{T}'$ . Recalling Lemma 8, this leaves us with greater than  $2\epsilon n(1 - o(1))$  items of degree one that participate in tests containing fewer than  $\frac{\rho}{\log n}$  items of degree one. Comparing the overall degrees, we find that the number of such tests is greater than  $\frac{2\epsilon n(1 - o(1))}{\rho / \log n}$ , and combining this with the assumption  $T = (1 - \epsilon) \frac{2n}{\rho}$  gives

$$\begin{aligned} \frac{2\epsilon n(1 - o(1))}{\rho / \log n} &\leq (1 - \epsilon) \frac{2n}{\rho} \\ \implies \epsilon &\leq \frac{1}{(\log n)(1 - o(1)) + 1} = o(1), \end{aligned} \quad (117)$$

which gives the desired contradiction (we assumed  $\epsilon = \Theta(1)$ ). This completes the second part of the proof of Theorem 4.

## V. CONCLUSION

In this paper, we have analyzed the performance of doubly-regular group testing designs in several settings of interest, with our main results summarized as follows:

- In the unconstrained setting with sub-linear sparsity, the block-structured doubly-regular design with the DD algorithm matches the achievability result of the DD algorithm with near-constant tests-per-item, which is known to be optimal for  $\theta \geq \frac{1}{2}$ .
- In the unconstrained setting with linear sparsity, we complemented hardness results for exact recovery by

deriving achievability results with only false negatives in the reconstruction.

- In the setting of  $\rho$ -sized tests with sub-linear sparsity and exact recovery, we improved on previously-known upper and lower bounds in regimes of the form  $\rho = \Theta\left(\left(\frac{n}{k}\right)^\beta\right)$  with  $\beta \in (0, 1)$ , complementing recent improvements that only hold for  $\beta = 0$ .

An immediate direction for future research is to establish converse bounds for approximate recovery in the linear regime with no false positives, and also to establish a more detailed understanding of the entire FPR vs. FNR trade-off. In addition, in the size-constrained setting, the optimal thresholds still remain unknown in general, despite the gap now being narrower.

## APPENDIX A

### DETAILS OF CONDITIONAL INDEPENDENCE ARGUMENT

Here we provide a more mathematical description of the conditioning argument used in Step 3 in Section III. Before defining the conditioning events for sub-matrices  $X_1, \dots, X_r$ , we first consider fixing several quantities:

- The defective set  $\mathcal{K}$  (e.g.,  $\mathcal{K} = \{1, \dots, k\}$ );
- The masked non-defective set  $\mathcal{G}$  (e.g.,  $\mathcal{G} = k+1, \dots, k+G$ , where  $G$  is the number of masked non-defectives);
- The set of positive tests  $\mathcal{P}_1, \dots, \mathcal{P}_r$  for the  $r$  sub-matrices (e.g., with  $X_1$  corresponding to tests  $1, \dots, \frac{n}{s}$ , we could have  $\mathcal{P}_1$  being a fixed size- $\frac{n}{2s}$  subset of those tests);
- The counts of masked defectives  $M^1, \dots, M^r$  (e.g.,  $M^1 = 0, M^2 = 3$ , etc.);
- The placements of non-masked non-defectives into negative tests (e.g., item  $k+G+1$  is in the first negative test for  $X_1$ , item  $k+G+2$  is in the fifth negative test for  $X_7$ , etc.). We represent this by a set FixedNeg containing triplets  $(i, j, t)$  with  $i \in \{1, \dots, n\} \setminus (\mathcal{K} \cup \mathcal{G})$  indexing the (non-defective) item,  $j \in \{1, \dots, r\}$  indexing the sub-matrix, and  $t \in \{1, \dots, T\}$  indexing the (negative) test.

With the sets of positive tests  $\{\mathcal{P}_j\}_{j=1}^r$  fixed, the corresponding sets of negative tests are also fixed, and are denoted by  $\{\mathcal{N}_j\}_{j=1}^r$ . Then, with  $\mathcal{K}, \mathcal{G}, \{\mathcal{P}_j\}_{j=1}^r, \{\mathcal{N}_j\}_{j=1}^r, \{M^j\}_{j=1}^r$  and FixedNeg fixed, we can now list the conditioning events for  $j = 1, \dots, r$ :

- Let  $\mathcal{A}_j$  be the event that each test in  $\mathcal{P}_j$  contains at least one item from  $\mathcal{K}$ , and that each test in  $\mathcal{N}_j$  contains no items from  $\mathcal{K}$ .
- Let  $\mathcal{B}_j$  be the event that each test in  $\mathcal{N}_j$  contains no items from  $\mathcal{G}$ .
- Let  $\mathcal{C}_j$  be the event that there are exactly  $M^j$  items in  $\mathcal{K}$  whose (only) test in  $X_j$  contains at least one other defective (i.e., it contains two or more in total);
- Let  $\mathcal{D}_j$  be the event that for all  $(i, t)$  satisfying  $(i, j, t) \in \text{FixedNeg}$ , it holds that entry  $(i, t)$  in  $X_j$  equals one.

The final conditioning event is  $\bigcap_{j=1}^r (\mathcal{A}_j \cap \mathcal{B}_j \cap \mathcal{C}_j \cap \mathcal{D}_j)$ . Then, since  $\mathcal{A}_j \cap \mathcal{B}_j \cap \mathcal{C}_j \cap \mathcal{D}_j$  is an event depending only on  $X_j$  (for  $j = 1, \dots, r$ ), we observe that the independence of the matrices  $\{X_j\}_{j=1}^r$  before conditioning is still maintained after conditioning, as desired.

$$\frac{(n-k)!(n-\frac{n \log 2}{k})!}{(n-\frac{n \log 2}{k}-k+1)!(n-1)!} - \frac{(n-k)!(n-\frac{n \log 2}{k}-1)!}{(n-\frac{n \log 2}{k}-k+1)!(n-2)!} \quad (118)$$

$$= \frac{(n-k)!((n-\frac{n \log 2}{k})!(n-2)! - (n-\frac{n \log 2}{k}-1)!(n-1)!)}{(n-\frac{n \log 2}{k}-k+1)!(n-1)!(n-2)!} \quad (119)$$

$$= \frac{(n-k)!(n-\frac{n \log 2}{k}-1)!(n-2)!(n-\frac{n \log 2}{k}-n+1)}{(n-\frac{n \log 2}{k}-k+1)!(n-1)!(n-2)!} \quad (120)$$

$$\stackrel{(a)}{=} \frac{(n-\frac{n \log 2}{k}-1) \dots (n-\frac{n \log 2}{k}-k+2)}{(n-1) \dots (n-k+1)} \left( -\frac{n \log 2}{k} + 1 \right) \quad (121)$$

$$= \left( -\frac{n \log 2}{k} (1 + o(1)) \cdot \frac{1}{n-k+1} \right) \prod_{i=1}^{k-2} \frac{n-\frac{n \log 2}{k}-i}{n-i} \quad (122)$$

$$\stackrel{(b)}{=} -\frac{\log 2}{k} (1 + o(1)) \left( 1 - \frac{\log 2}{k} (1 + o(1)) \right)^{k-2} \quad (123)$$

$$\stackrel{(c)}{=} O\left(\frac{1}{k}\right), \quad (124)$$

$$\frac{(n-k)!(n-k-\frac{n \log 2}{k}+1)!(n-\frac{n \log 2}{k}-1)!(n-\frac{2n \log 2}{k})!}{(n-2)!(n-\frac{n \log 2}{k}-1)!(n-k-\frac{n \log 2}{k}+1)!(n-k-\frac{2n \log 2}{k}+2)!} - \left( \frac{(n-k)!(n-\frac{n \log 2}{k})!}{(n-1)!(n-k-\frac{n \log 2}{k}+1)!} \right)^2 \quad (125)$$

$$= \frac{(n-k)!(n-\frac{2n \log 2}{k})!}{(n-2)!(n-k-\frac{2n \log 2}{k}+2)!} - \left( \frac{(n-k)!(n-\frac{n \log 2}{k})!}{(n-1)!(n-k-\frac{n \log 2}{k}+1)!} \right)^2 \quad (126)$$

$$\stackrel{(a)}{=} \frac{(n-\frac{2n \log 2}{k}) \dots (n-k-\frac{2n \log 2}{k}+3)}{(n-2) \dots (n-k+1)} - \frac{(n-\frac{n \log 2}{k})^2 \dots (n-\frac{n \log 2}{k}-k+2)^2}{(n-1)^2 \dots (n-k+1)^2} \quad (127)$$

$$= \prod_{i=0}^{k-3} \frac{n-\frac{2n \log 2}{k}-i}{n-i-2} - \prod_{i=0}^{k-2} \frac{(n-\frac{n \log 2}{k}-i)^2}{(n-i-1)^2} \quad (128)$$

$$= \prod_{i=0}^{k-3} \left( 1 - \frac{2n \log 2 - 2}{n-i-2} \right) - \prod_{i=0}^{k-2} \left( 1 - \frac{n \log 2 - 1}{n-i-1} \right)^2 \quad (129)$$

$$= \left( 1 - \frac{2 \log 2}{k} \left( 1 + O\left(\frac{k}{n}\right) \right) \right)^{k-2} - \left( 1 - \frac{\log 2}{k} \left( 1 + O\left(\frac{k}{n}\right) \right) \right)^{2(k-1)} \quad (130)$$

$$\stackrel{(b)}{=} \exp\left( - (2 \log 2) \left( 1 + O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right) \right) \right) - \exp\left( - (2 \log 2) \left( 1 + O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right) \right) \right) \quad (131)$$

$$\stackrel{(c)}{=} O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right), \quad (132)$$

## APPENDIX B

### PROOF OF LEMMA 1 (COVARIANCE CALCULATION)

It suffices to show that the first and second parts of (39) simplify to  $O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right)$ .

**First part:** Expanding the binomial coefficient in terms of factorials, the first term simplifies as in (118)–(124) at the top of the next page, where (a) expands all the factorials and simplifies, (b) uses  $k = o(n)$ , and (c) applies  $(1 - \frac{\log 2}{k}(1 + o(1)))^{k-2} = \exp(-(1 + o(1)) \log 2) = \Theta(1)$ .

**Second part:** Expanding the binomial coefficient in terms of factorials, we have (125)–(132) at the top of the next page, where (a) expands all the factorials and simplifies, (b) uses the fact that  $1 - a = e^{-a+O(a^2)} = e^{-a(1+O(a))}$  and  $(1 + O(a))(1 + O(b)) = 1 + O(\max\{a, b\})$  whenever  $a, b = o(1)$ , and (c) uses the fact that upon applying a Taylor expansion

to each term, the leading terms cancel and only the first-order remainder remains.

**Combining the terms:** Substituting  $O\left(\frac{1}{k}\right)$  and  $O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right)$  into the first and second parts of (39) respectively, we obtain  $\text{Cov}[M_1^j, M_2^j] \leq O\left(\max\left\{\frac{1}{k}, \frac{k}{n}\right\}\right)$ , as desired.

## APPENDIX C

### FALSE POSITIVE RATE OF COMP IN THE LINEAR REGIME

Here we provide an analog of Theorem 2 for the COMP algorithm (see Algorithm 1), which comes essentially “for free” from our analysis of the DD algorithm. However, we note that unlike our DD analysis, the result for COMP would also follow easily from prior work, particularly [14]. Recall that the COMP algorithm only produces false positives, which implies that we only need to look at the FPR (i.e., FNR = 0).

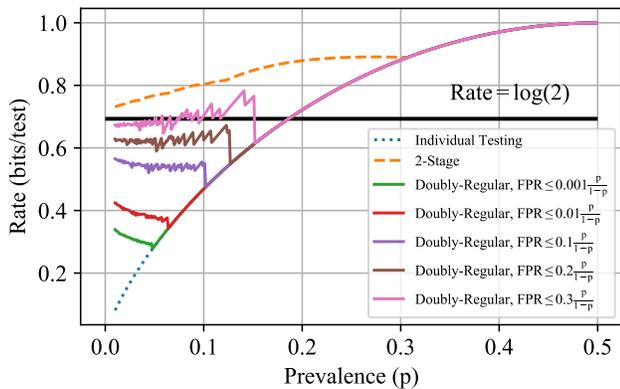


Fig. 4: Achievable rates for COMP decoding with the doubly-regular design and approximate recovery (along with individual testing and a two-stage design [14], both of which attain exact recovery).

**Theorem 5.** Using the COMP algorithm with the block-structured doubly-regular design with fixed parameters  $s$  and  $r$ , when there are  $k = pn$  defective items with constant  $p \in (0, 1)$ , we have  $\text{FPR} \leq \text{FPR}_{\max}(1 + o(1))$  with probability  $1 - o(1)$ , where

$$\text{FPR}_{\max} = (1 - (1 - p)^{s-1})^r. \quad (133)$$

*Proof.* From (11), we have

$$\begin{aligned} \text{FPR} &= \left(1 - \prod_{i=1}^{s-1} \left(1 - \frac{k}{n-i}\right)\right)^r \\ &= \left(1 - \left(1 - \frac{k}{n(1-o(1))}\right)^{s-1}\right)^r \\ &= (1 - (1 - p)^{s-1})^r (1 + o(1)), \end{aligned} \quad (134)$$

where the second and third qualities use the fact that  $s$ ,  $r$ , and  $p$  are all constant with respect to  $n$ .  $\square$

A given FPR value corresponds to an average of  $(n-k)\text{FPR}$  false positives, which may potentially be much larger than  $k$ . To place the number of false positives and the actual number of defectives on the “same scale”, we find it more convenient to work with the normalized quantity  $\text{FPR} \frac{n-k}{k} = \text{FPR} \frac{1-p}{p}$ . Then, this quantity equaling a given value  $\alpha > 0$  corresponds to an average of  $\alpha k$  false positives.

Recalling the notion of rate in Definition 1, we have the following analog of Corollary 1; the proof is similar but simpler, so is omitted.

**Corollary 2.** Under the setup of Theorem 5, there exist choices of  $r$  and  $s$  (depending on  $p$ ) such that, in the limit as  $p \rightarrow 0$  (after having taking  $n \rightarrow \infty$ ), we have (i)  $\text{FPR}_{\max} \frac{1-p}{p} \rightarrow 0$ , (ii) the rate approaches  $\log 2$ , and (iii) it holds that  $s = \frac{\log 2}{p}(1 + o(1))$  and  $r = \frac{\log(\frac{1}{p})}{\log 2}(1 + o(1))$ .

Similarly to the discussion following Corollary 1, this result is consistent with the rate of  $\log 2$  attained for COMP with approximate recovery in the sub-linear regime  $k = o(n)$  [2, Sec. 5.1].

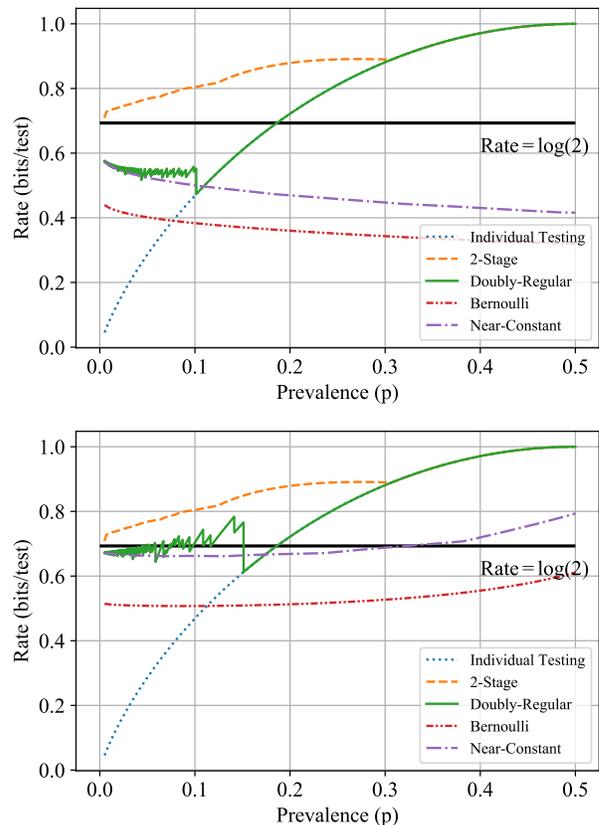


Fig. 5: Comparison of rates for COMP decoding with various tests designs, for approximate recovery with  $\alpha = 0.1$  (Left) and  $\alpha = 0.3$  (Right); the rates for Bernoulli designs and near-constant column weight designs are given in [36]. Individual testing and the two-stage design [14] attain exact recovery.

In this case, we can easily argue that the constant of  $\log 2$  is optimal. To see this, note that if we can could attain at most  $\alpha k$  false positives on average for arbitrarily small  $\alpha > 0$ , then we could use this to construct a two-stage adaptive group testing algorithm where the second stage tests the items outputted in the first stage individually, thus using an average of  $(1 + \alpha)k$  tests or fewer. When  $p$  approaches zero, these additional tests contribute a arbitrarily small fraction compared to the leading  $\Theta(k \log \frac{n}{k}) = \Theta(k \log \frac{1}{p})$  term, so the two-stage design attains zero error probability with an arbitrarily small increase in the rate. However, the converse result of [14] shows that rates above  $\log 2$  are impossible in this two-stage setting (see also [29] for the sublinear regime). This establishes the optimality of the constant  $\log 2$  above.

To visualize the constant- $p$  regime, we plot the rates attained by the COMP algorithm with the doubly-regular design in Figure 4 (see the text following Corollary 1 for discussion on why the behavior of the curves reaching  $\log 2$  as  $p \rightarrow 0$  is not visible; a similar discussion applies here). In Figure 5, we additionally compare against other random non-adaptive test designs, for which bounds on the FPR were given in [36]. We observe that the doubly-regular design consistently outperforms the Bernoulli design, and also outperforms the near-constant column weight design except near certain  $p$

values where the doubly-regular curve is discontinuous.

## REFERENCES

- [1] N. Tan and J. Scarlett, "An analysis of the DD algorithm for group testing with size-constrained tests," in *IEEE Int. Symp. Inf. Theory*, 2021.
- [2] M. Aldridge, O. Johnson, and J. Scarlett, "Group testing: An information theory perspective," *Found. Trend. Comms. Inf. Theory*, vol. 15, no. 3–4, pp. 196–392, 2019.
- [3] M. Aldridge and D. Ellis, *Pooled testing and its applications in the COVID-19 pandemic*. Pandemics: Insurance and Social Protection, 2021.
- [4] M. Malyutov, "The separating property of random matrices," *Math. Notes Acad. Sci. USSR*, vol. 23, no. 1, pp. 84–91, 1978.
- [5] G. Atia and V. Saligrama, "Boolean compressed sensing and noisy group testing," *IEEE Trans. Inf. Theory*, vol. 58, no. 3, pp. 1880–1901, March 2012.
- [6] M. Aldridge, L. Baldassini, and O. Johnson, "Group testing algorithms: Bounds and simulations," *IEEE Trans. Inf. Theory*, vol. 60, no. 6, pp. 3671–3687, June 2014.
- [7] J. Scarlett and V. Cevher, "Phase transitions in group testing," in *Proc. ACM-SIAM Symp. Disc. Alg. (SODA)*, 2016.
- [8] O. Johnson, M. Aldridge, and J. Scarlett, "Performance of group testing algorithms with near-constant tests-per-item," *IEEE Trans. Inf. Theory*, vol. 65, no. 2, pp. 707–723, Feb. 2019.
- [9] A. Coja-Oghlan, O. Gebhard, M. Hahn-Klimroth, and P. Loick, "Information-theoretic and algorithmic thresholds for group testing," in *Int. Colloq. Aut., Lang. and Prog. (ICALP)*, 2019.
- [10] —, "Optimal group testing," in *Conf. Learn. Theory (COLT)*, vol. 125, 2020, pp. 1374–1388.
- [11] M. Mézard, M. Tarzia, and C. Toninelli, "Group testing with random pools: Phase transitions and optimal strategy," *J. Stat. Phys.*, vol. 131, no. 5, pp. 783–801, 2008.
- [12] V. Gandikota, E. Grigorescu, S. Jaggi, and S. Zhou, "Nearly optimal sparse group testing," *IEEE Trans. Inf. Theory*, vol. 65, no. 5, pp. 2760–2773, 2019.
- [13] O. Gebhard, M. Hahn-Klimroth, O. Parczyk, M. Penschuck, M. Rolvien, J. Scarlett, and N. Tan, "Near-optimal sparsity-constrained group testing: Improved bounds and algorithms," *IEEE Trans. Inf. Theory*, vol. 68, no. 5, pp. 3253–3280, 2022.
- [14] M. Aldridge, "Conservative two-stage group testing in the linear regime," 2020, <https://arxiv.org/abs/2005.06617>.
- [15] R. Goenka, S.-J. Cao, C.-W. Wong, A. Rajwade, and D. Baron, "Contact tracing information improves the performance of group testing algorithms," 2021, <https://arxiv.org/abs/2106.02699>.
- [16] S. Ghosh, R. Agarwal, M. A. Rehan, S. Pathak, P. Agarwal, Y. Gupta, S. Consul, N. Gupta, R. Ritika, R. Goenka *et al.*, "A compressed sensing approach to pooled RT-PCR testing for COVID-19 detection," *IEEE Open J. Sig. Proc.*, 2021.
- [17] M. Aldridge, "Individual testing is optimal for nonadaptive group testing in the linear regime," *IEEE Trans. Inf. Theory*, vol. 65, no. 4, pp. 2058–2061, April 2019.
- [18] W. H. Bay, J. Scarlett, and E. Price, "Optimal non-adaptive probabilistic group testing in general sparsity regimes," 2022, article iaab020.
- [19] C. L. Chan, S. Jaggi, V. Saligrama, and S. Agnihotri, "Non-adaptive group testing: Explicit bounds and novel algorithms," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 3019–3035, May 2014.
- [20] A. Z. Broder and R. Kumar, "A note on double pooling tests," 2020, <https://arxiv.org/abs/2004.01684>.
- [21] M. Aldridge, "The capacity of Bernoulli nonadaptive group testing," *IEEE Trans. Inf. Theory*, vol. 63, no. 11, pp. 7142–7148, 2017.
- [22] C. L. Chan, P. H. Che, S. Jaggi, and V. Saligrama, "Non-adaptive probabilistic group testing with noisy measurements: Near-optimal bounds with efficient algorithms," in *Allerton Conf. Comm., Ctrl., Comp.*, Sep. 2011, pp. 1832–1839.
- [23] S. Cai, M. Jahangoshahi, M. Bakshi, and S. Jaggi, "Efficient algorithms for noisy group testing," *IEEE Trans. Inf. Theory*, vol. 63, no. 4, pp. 2113–2136, 2017.
- [24] K. Lee, R. Pedarsani, and K. Ramchandran, "Saffron: A fast, efficient, and robust framework for group testing based on sparse-graph codes," in *IEEE Int. Symp. Inf. Theory*, 2016.
- [25] H. Q. Ngo, E. Porat, and A. Rudra, "Efficiently decodable error-correcting list disjunct matrices and applications," in *Int. Colloq. Automata, Lang., and Prog.*, 2011.
- [26] P. Indyk, H. Q. Ngo, and A. Rudra, "Efficiently decodable non-adaptive group testing," in *ACM-SIAM Symp. Disc. Alg. (SODA)*, 2010.
- [27] E. Price and J. Scarlett, "A fast binary splitting approach to non-adaptive group testing," in *RANDOM*, 2020.
- [28] E. Price, J. Scarlett, and N. Tan, "Fast splitting algorithms for sparsity-constrained and noisy group testing," 2021, <https://arxiv.org/abs/2106.00308>.
- [29] M. Mézard and C. Toninelli, "Group testing with random pools: Optimal two-stage algorithms," *IEEE Trans. Inf. Theory*, vol. 57, no. 3, pp. 1736–1745, 2011.
- [30] A. G. D'yachkov, "Lectures on designing screening experiments," 2014, <https://arxiv.org/abs/1401.7505>.
- [31] N. Tan and J. Scarlett, "Near-optimal sparse adaptive group testing," in *IEEE Int. Symp. Inf. Theory*, 2020.
- [32] R. Motwani and P. Raghavan, *Randomized Algorithms*. Chapman & Hall/CRC, 2010.
- [33] C. M. Fortuin, J. Ginibre, and P. W. Kasteleyn, "Correlation inequalities on some partially ordered sets," *Communications in Mathematical Physics*, vol. 22, no. 2, pp. 89–103, 1971.
- [34] M. Skala, "Hypergeometric tail inequalities: ending the insanity," 2013, <https://arxiv.org/abs/1311.5939>.
- [35] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *J. Amer. Stat. Assoc.*, vol. 58, no. 301, pp. 13–30, 1963.
- [36] W. H. Bay and J. Scarlett, "Optimal non-adaptive probabilistic group testing in general sparsity regimes," 2020, <https://arxiv.org/abs/2006.01325v1>.

**Nelvin Tan** received the B.Comp. degree in computer science and statistics from the National University of Singapore in 2021. He is currently pursuing the Ph.D. degree from the Signal Processing and Communications Group in the Department of Engineering, University of Cambridge. His research interests include information theory and high-dimensional statistics.

**Way Tan** received a double degree in mathematics and computer science from the National University of Singapore in 2021. His research interests include information theory and applied mathematics.

**Jonathan Scarlett** (S'14 – M'15) received the B.Eng. degree in electrical engineering and the B.Sci. degree in computer science from the University of Melbourne, Australia. From October 2011 to August 2014, he was a Ph.D. student in the Signal Processing and Communications Group at the University of Cambridge, United Kingdom. From September 2014 to September 2017, he was post-doctoral researcher with the Laboratory for Information and Inference Systems at the École Polytechnique Fédérale de Lausanne, Switzerland. Since January 2018, he has been an assistant professor in the Department of Computer Science and Department of Mathematics, National University of Singapore. His research interests are in the areas of information theory, machine learning, signal processing, and high-dimensional statistics. He received the Singapore National Research Foundation (NRF) fellowship, and the NUS Early Career Research Award.