

Achievable Rates and Algorithms for Group Testing with Runlength Constraints

Stefano Della Fiore

Department of Information Engineering
University of Brescia
Brescia (BS), Italy
s.dellafiore001@unibs.it

Marco Dalai

Department of Information Engineering
University of Brescia
Brescia (BS), Italy
marco.dalai@unibs.it

Ugo Vaccaro

Department of Informatics
University of Salerno
Fisciano (SA), Italy
uvaccaro@unisa.it

Abstract—In this paper, we study bounds on the minimum length of (k, n, d) -superimposed codes introduced by Agarwal *et al.* [1], in the context of Non-Adaptive Group Testing algorithms with runlength constraints. A (k, n, d) -superimposed code of length t is a $t \times n$ binary matrix such that any two 1's in each column are separated by a run of at least d 0's, and such that for any column c and any other $k - 1$ columns, there exists a row where c has 1 and all the remaining $k - 1$ columns have 0. Agarwal *et al.* proved the existence of such codes with $t = \Theta(dk \log(n/k) + k^2 \log(n/k))$. Here we investigate more in detail the coefficients in front of these two main terms as well as the role of lower order terms. We show that improvements can be obtained over the construction in [1] by using different constructions and by an appropriate exploitation of the Lovász Local Lemma in this context. Our findings also suggest $O(n^k)$ randomized Las Vegas algorithms for the construction of such codes. We also extend our results to Two-Stage Group Testing algorithms with runlength constraints.

Index Terms—Lovász Local Lemma, superimposed codes, runlength-constrained codes.

I. INTRODUCTION

Group Testing refers to the scenario in which one has a population I of individuals, and an unknown subset P of I , commonly referred to as “positives”. The goal is to determine the unknown elements of P by performing tests on arbitrary subsets A of I (called *pools*), and the outcome of the test is assumed to return the value 1 (positive) if A contains at least one element of the unknown set P , the value 0 (negative), otherwise. The problem was first introduced by Dorfman [5] during WWII, in the context of mass blood testing. Since then, Group Testing techniques have found applications in a large variety of areas, ranging from DNA sequencing to quality control, data security to network analysis, and much more. We refer the reader to the excellent monographs [6], [9] for an account of the vast literature on the subject.

Group Testing procedures can be adaptive or non-adaptive. In adaptive Group Testing, the tests are performed sequentially, and the content of the pool tested at the generic step i might depend on the previous $i - 1$ test outcomes. Conversely, in non-adaptive Group Testing all pools are a-priori set, and tests are carried out in parallel. Non-adaptive Group Testing (NAGT) schemes typically require more tests to discover the positives, but they are faster since tests can be performed in parallel. To combine the advantages of both techniques, while mitigating

their limitations, it is sometimes preferable to implement a hybrid approach, where a first screening is performed via a NAGT algorithm, followed by a simple one-by-one testing of the members that are identified in the first stage as potentially positives. This latter approach is usually called Two-Stage Group Testing [4].

In NAGT, the algorithm to determine the positives is usually represented by means of a $t \times n$ binary matrix M , where each row of M represents a test while each column is associated to a distinct member of the population $I = \{1, 2, \dots, n\}$. More precisely, we have $M_{ij} = 1$ if and only if the member $j \in I$ belongs to the i -th test. In general, one assumes a known upper bound k on the cardinality of the unknown set of positives P . Having said that, the property one usually requires for M to represent a correct (and efficiently decodable) NAGT is the following [6]: for any k -tuple of the n columns of M we demand that for any column c of the given k -tuple, there exists a row $i \in \{1, \dots, t\}$ such that c has symbol 1 in row i and all the remaining $k - 1$ columns of the k -tuple have a 0 in the same row i . This condition renders matrices M with such a property equivalent to the well known superimposed codes introduced in the seminal paper by Kautz and Singleton [10] and, independently, by Erdős *et al.* in [7].

Motivated by applications in topological DNA-based data storage, the authors of [1] introduced an interesting new variant of NAGT, in which the associated test matrix M has to satisfy additional constraints, in order to comply with the biological constraints of the problem they want to solve. Informally, one of the main problems studied in [1] is to show the existence of a superimposed code M with a “small” number t of rows and satisfying the following additional property: any two 1's in each column are separated by a run of at least d 0's. We refer the reader to [1] for the rationale behind this run-length constraint. The main achievability result obtained in [1] says that codes with these properties exist for $t = \Theta(dk \log(n/k) + k^2 \log(n/k))$.

Our results. We study the achievable coefficients in front of the $dk \log(n/k)$ and $k^2 \log(n/k)$ terms and on lower order terms in k , whose values are of importance as they determine the achievable rates of such codes for fixed values of d and k . We show that better results than those derived in [1] can be obtained using different random coding constructions which

also admit simpler analyses. Also, we show that improved results can be obtained by an appropriate use of the Lovász Local Lemma (see, e.g., [2]). By exploiting the celebrated result by Moser and Tardos [12], this directly implies a $O(n^k)$ randomized Las Vegas algorithm to construct such codes. In the final part of this paper, we also extend our results to Two Stages Group Testing algorithms.

II. NEW UPPER BOUNDS

Throughout the paper, the logarithms without subscripts are in base two, and we denote with $\ln(\cdot)$ the natural logarithm. For notation convenience we denote with $[a, b]$ the set $\{a, a+1, \dots, b\}$.

Definition II.1. [1] Let k, n, d be positive integers, $k \leq n$. A (k, n, d) -superimposed code is a $t \times n$ binary matrix M such that any two 1's in each column of M are separated by a run of at least d 0's, and for any k -tuple of the columns of M we have that for any column \mathbf{c} of the given k -tuple, there exists a row $i \in [1, t]$ such that \mathbf{c} has symbol 1 in row i and all the remaining $k-1$ columns of the k -tuple are equal to 0. The number of rows t of M is called the length of the (k, n, d) -superimposed code.

Definition II.2. A (k, n, d, w) -superimposed code is a (k, n, d) -superimposed with the additional constraint that each column has weight w (number of ones).

First, we need the following enumerative lemma.

Lemma II.3. Let $S \subseteq \{0, 1\}^t$ be the set of all distinct binary vectors of length t such that each vector has Hamming weight $w \geq 1$ and any two 1's in each vector are separated by a run of at least d 0's. If $t \geq (w-1)d + w$, then

$$|S| = \binom{t - (w-1)d}{w}.$$

Proof. Let A be the set of all distinct binary vectors of length $t - (w-1)d$ and weight w . One can see that $|S| = |A|$ since each element of S can be obtained from an element $a \in A$ by adding between each pair of consecutive ones in a exactly d 0's. Conversely, each element of A can be obtained from an element $s \in S$ by removing between each pair of consecutive ones in s exactly d 0's. \square

We also need the following technical lemma and an easy corollary, which have been proved in [8]. We report here the proofs for the reader's convenience.

Lemma II.4. Let a, b, c be positive integers such that $c \leq a \leq b$. We have that

$$\frac{a}{b} \cdot \frac{a-c}{b-c} \leq \left(\frac{a - \frac{c}{2}}{b - \frac{c}{2}} \right)^2.$$

Proof. Clearly $a(a-c)c^2 \leq b(b-c)c^2$. Then adding the quantity $4ab(a-c)(b-c)$ to both members implies that $a(a-c)(2b-c)^2 \leq b(b-c)(2a-c)^2$. Therefore Lemma II.4 follows. \square

Corollary II.5. Let a, b, c be positive integers such that $c \leq a \leq b$. We have that

$$\frac{\binom{a}{c}}{\binom{b}{c}} \leq \left(\frac{a - \frac{c-1}{2}}{b - \frac{c-1}{2}} \right)^c. \quad (1)$$

Proof. Expanding the LHS of (1) we get

$$\frac{\binom{a}{c}}{\binom{b}{c}} = \frac{a}{b} \cdot \frac{a-1}{b-1} \cdots \frac{a-c+1}{b-c+1}. \quad (2)$$

Let us group the terms in (2) into pairs as follows

$$\frac{a-i}{b-i} \cdot \frac{a-(c-i-1)}{b-(c-i-1)} \text{ for } i = 0, \dots, \left\lceil \frac{c-1}{2} \right\rceil - 1. \quad (3)$$

If c is odd then we leave alone the term $(a - \frac{c-1}{2}) / (b - \frac{c-1}{2})$. By Lemma II.4, each term in (3) can be upper bounded by

$$\frac{a-i}{b-i} \cdot \frac{a-(c-i-1)}{b-(c-i-1)} \leq \left(\frac{a - \frac{c-1}{2}}{b - \frac{c-1}{2}} \right)^2.$$

Hence Corollary II.5 follows. \square

The main tool to prove Theorem II.7 is the Lovász Local Lemma for the symmetric case. We state here the lemma.

Lemma II.6. [2] Let E_1, E_2, \dots, E_m be events in an arbitrary probability space. Suppose that each event E_i is mutually independent of a set of all other events E_j but at most d , and that $\Pr(E_i) \leq p$ for all $1 \leq i \leq m$. If

$$edp \leq 1$$

then $\Pr(\cap_{i=1}^m \overline{E}_i) > 0$.

Now, we are ready to state our main result.

Theorem II.7. There exists a (k, n, d, w) -superimposed code of length t , where t is the minimum integer such that the following inequality holds

$$ek \left[\binom{n}{k} - \binom{n-k+1}{k} \right] \left(\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \right)^w \leq 1. \quad (4)$$

Proof. Let M be a $t \times n$ binary matrix, where each column \mathbf{c} is picked uniformly at random between the set of all distinct binary vectors of length t such that each column has weight w and any two 1's in each column of M are separated by a run of at least d 0's. Therefore by Lemma II.3 we have that

$$\Pr(\mathbf{c}) = \binom{t - (w-1)d}{w}^{-1}.$$

For a given index $i \in [1, n]$ and a set of column-indices B , $|B| = k-1$, $i \notin B$, let $E_{i,B}$ be the event such that for every row in which \mathbf{c}_i (the i -th column) has 1, there exists an index $j \in B$ such that \mathbf{c}_j has 1 in that same row. We can write this event in terms of supports as $\text{Supp}(\mathbf{c}_i) \subseteq \text{Supp}(\mathbf{c}_B)$. There

are $n \binom{n-1}{k-1}$ such events. We can express the probability of such an event as follows

$$\Pr(E_{i,B}) = \sum_{c'=(c'_1, \dots, c'_{k-1})} \Pr(\mathbf{c}_B = c') \cdot \Pr(\text{Supp}(\mathbf{c}_i) \subseteq \text{Supp}(\mathbf{c}_B) | \mathbf{c}_B = c'), \quad (5)$$

where we have denoted with \mathbf{c}_B the vector $(\mathbf{c}_{j_1}, \dots, \mathbf{c}_{j_{k-1}})$ in which j_1, \dots, j_{k-1} are the elements of B . The sum in (5) is over all the possible configurations of $k-1$ vectors of length t , weight w and the distance between ones in each column is at least d . Then, we can upper bound (5) by the maximum of $\Pr(\text{Supp}(\mathbf{c}_i) \subseteq \text{Supp}(\mathbf{c}_B) | \mathbf{c}_B = c')$ over all $k-1$ vectors $c' = (c'_1, \dots, c'_{k-1})$. Therefore, we can consider the worst-case scenario where the $k-1$ columns of M with indices in B maximize this probability. It can be seen that the maximum is achieved when the $w(k-1)$ ones of the $k-1$ columns indexed by B are placed in $w(k-1)$ different rows. Hence,

$$\Pr(E_{i,B}) \leq \frac{\binom{w(k-1)}{w}}{\binom{t-(w-1)d}{w}}. \quad (6)$$

Using Corollary II.5 we upper bound (6) as follows

$$\Pr(E_{i,B}) \leq \left(\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \right)^w. \quad (7)$$

Proceeding as in [8], it can be proved that an arbitrary event $E_{i,A}$ is mutually independent from all the events $E_{j,C}$, where $C \subseteq [1, n] \setminus (A \cup \{i\})$ and $j \notin C$. Since the number of events $E_{j,C}$ is equal to

$$\binom{n-k}{k-1} (n-k+1) = k \binom{n-k+1}{k},$$

each event $E_{i,A}$ is dependent of at most

$$f = k \left[\binom{n}{k} - \binom{n-k+1}{k} \right] \quad (8)$$

other events. If the probability that none of the events $E_{i,A}$ occurs is strictly positive then there exists a matrix M that is a (k, n, d, w) -superimposed code of length t . Therefore, using Lemma II.6 and taking p equal to the RHS of (7) and f as defined in equation (8), Theorem II.7 follows. \square

Remark II.8. We note that in Theorem II.7 we could use the union bound instead of the Local Lemma. Since the total number of events is $n \binom{n-1}{k-1}$, we have that there exists a (k, n, d, w) -superimposed code of length t , provided that

$$n \binom{n-1}{k-1} \left(\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \right)^w < 1. \quad (9)$$

In [1] the authors proved that a (k, n, d, w) -superimposed code of length t exists, provided that

$$n \binom{n-1}{k-1} \left(\frac{w(k-1)}{t - (2d+1)(w-1)} \right)^w < 1. \quad (10)$$

It is clear that our bound given in Remark II.8 is better than the bound given in (10) since

$$\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \leq \frac{w(k-1)}{t - (2d+1)(w-1)}$$

for all positive integers w, k, d .

If we compare the bounds of Theorem II.7 and Remark II.8 then it has been proved in [8] that

$$ek \left[\binom{n}{k} - \binom{n-k+1}{k} \right] \leq n \binom{n-1}{k-1} \quad (11)$$

for all $k \leq 0.667\sqrt{n}$. Therefore when k is much smaller than n (which is indeed the case in circumstances of interest), the Local Lemma performs better than the union bound. It is important to note that a conjecture of Erdős, Frankl and Füredi [7] says that for $k \geq \sqrt{n}$ optimal superimposed codes have length equal to n . The current best known result has been proved in [13] which states that if $k \geq 1.157\sqrt{n}$ then the minimum length of superimposed codes is equal to n .

Corollary II.9. *There exists a (k, n, d, w) -superimposed code of length t , where*

$$t \leq \left[(w-1)d + \frac{w-1}{2} + \left(w(k-1) - \frac{w-1}{2} \right) \cdot \left(\min \left\{ n \binom{n-1}{k-1}, ek \left[\binom{n}{k} - \binom{n-k+1}{k} \right] \right\} \right)^{\frac{1}{w}} \right]. \quad (12)$$

Proof. It easily follows rearranging the terms in equation (4) and in equation (9). \square

Corollary II.10. *There exists a (k, n, d) -superimposed code of length t with $k \leq n/e$, where*

$$t \leq \ln 2 \cdot dk \log(n/k) + e^2 \cdot k^2 \log(n/k) - \frac{(3e^2 - \ln 2)}{2} k \log(n/k) - d + O(1).$$

Proof. Substitute $w = k \ln(n/k)$ in (12) and upper bound

$$\min \left\{ n \binom{n-1}{k-1}, ek \left[\binom{n}{k} - \binom{n-k+1}{k} \right] \right\} < k \left(\frac{en}{k} \right)^k.$$

Therefore we obtain

$$t \leq d(k \ln(n/k) - 1) + \frac{k}{2} \ln(n/k) + e \cdot (ke^k)^{\frac{1}{k \ln(n/k)}} k \left(k - \frac{3}{2} \right) \ln(n/k) + O(1). \quad (13)$$

Hence Corollary II.10 follows since $n \geq ek$ and $k^{1/k} \leq \frac{1}{\ln 2}$ for every integer $k \geq 1$. \square

We note that in the explicit bound given in Corollary II.10 the leading coefficient of the term $k \log(n/k)$ can be improved, for $k \leq 0.667\sqrt{n}$, by using a better estimation of the minimum in equation (12) that comes from the use of the Local Lemma.

By exploiting the celebrated result by Moser and Tardos [12], this directly implies a $O(n^k)$ randomized Las Vegas algorithm to construct the codes of Corollary II.10

From the inequality (10) we can derive an explicit upper bound on the length of the codes whose existence was showed in [1] when $w = k \ln(n/k)$ by upper bounding $n \binom{n-1}{k-1}$ with $k \left(\frac{en}{k}\right)^k$. We report here the obtained result.

Theorem II.11. [1] *There exists a (k, n, d) -superimposed code of length t , where*

$$t \leq 2d(k \ln(n/k) - 1) + k \ln(n/k) + e \cdot (ke^k)^{\frac{1}{k \ln(n/k)}} k(k-1) \ln(n/k) + O(1).$$

It is clear that our result given in equation (13) improves the one of Theorem II.11.

Remark II.12. *We note that it was proved in [1] that every (k, n, d) -superimposed code of length t must satisfy*

$$t \geq \min \{n, 1 + (k-1)(d+1)\}.$$

This implies that if $k \geq \frac{n-1}{d+1} + 1$ then $t = n$, so we cannot construct a (k, n, d) -superimposed code of length t that is better than the identity matrix of size $n \times n$.

By Remark II.12, it is clear that the constraint $k \leq n/e$ in Corollary II.10 is reasonable since $1 + (k-1)(d+1) \geq ek$ for every $k, d \geq 2$.

We also note that a simple generalization of the method given by Cheng et al. in [3] provide the following result.

Theorem II.13. *There exists a (k, n, d) -superimposed code of length t , $t \leq \frac{1}{B_k} (k \log(n/k) + \log(ke^k))$, where*

$$B_k = \max_{q \geq 2} B_{k,q}, \quad (14)$$

$$B_{k,q} = \frac{-\log \left[1 - \left(1 - \frac{1}{q}\right)^{k-1} \right]}{q + d}.$$

For $k \rightarrow \infty$, the point q that maximize (14) is linear in k .

The proof of Theorem II.13 is similar to the one in [3], we only need to ensure that when we construct a binary matrix M starting from a random q -ary matrix each column \mathbf{c} of M has a run of at least d 0's between any two 1's. This can be done by mapping each q -ary symbol into a binary vector of length $q + d$ where the last d elements are fixed to 0.

If we lower bound B_k with $B_{k,q}$ for the choice $q = \frac{1}{\ln 2}(k-1)$ then, for k sufficiently large, Theorem II.13 gives the following explicit bound on the minimum length t of (k, n, d) -superimposed codes

$$t \leq dk \log(n/k) + \frac{1}{\ln 2} \cdot k(k-1) \log(n/k) + \left(\frac{1}{\ln 2}(k-1) + d \right) \log(ke^k) + O(1). \quad (15)$$

One can see that this bound already improves, for k sufficiently large, the one given in Theorem II.11 but not the one obtained in Corollary II.10 for $k < d$.

III. SELECTORS

Selectors were introduced in [4] and they can be seen as a generalization of superimposed codes. Like superimposed codes, selectors find applications in many circumstances, like Group Testing [4], efficient conflict resolution in the transmission model of [11], etc.. In this section, we introduce selectors in which the weight of each column is equal to some fixed value w and where any two 1's in each column of M are separated by a run of at least d 0's, so that they can be applied to the scenario of [1]. Successively, we will show that selectors can be used to construct efficient two-stage procedure for Group Testing with runlength constraints, that require a much smaller number of tests, with respect to the NAGT considered in [1] and in the previous section of the present paper. Let us start by giving some definitions.

Definition III.1. Let k, n, d, p be positive integers, $1 \leq p \leq k \leq n$. A (k, n, d, p) -selector is a $t \times n$ binary matrix M such that any two 1's in each column of M are separated by a run of at least d 0's, and for any k -tuple of the columns of M we have that at least p rows of the identity matrix of size $k \times k$ are contained in that k -tuple of columns. The number of rows t of M is called the length of the (k, n, d, p) -selector.

One can see that for $p = k$ we get the definition of (k, n, d) -superimposed codes studied in Section II.

Definition III.2. A (k, n, d, p, w) -selector is a (k, n, d, p) -selector with the additional constraint that each column has weight w .

It can be seen (see [8, Lemma 2]) that Definition III.1 is equivalent to requiring that for any k -tuple of columns of a (k, n, d, p) -selector and any $k-p+1$ columns among the selected k -tuple, there exists a row of the identity matrix of size $k \times k$ where the 1 is contained in one of the $k-p+1$ columns. Therefore, thanks to this equivalence, we can generalize the proof of Theorem II.7 to obtain the following.

Theorem III.3. *There exists a (k, n, d, p, w) -selector of length t , where t is the minimum integer such that the following inequality holds*

$$e \binom{k}{p-1} \left[\binom{n}{k} - \binom{n-k}{k} \right] \cdot \left(\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \right)^{w(k-p+1)} \leq 1. \quad (16)$$

Proof. Let M be a $t \times n$ binary matrix, where each column \mathbf{c} is picked uniformly at random between the set of all distinct binary vectors of length t such that each column has weight w and with distance between ones in each column at least d . As in Theorem II.7, by Lemma II.3 we have that

$$\Pr(\mathbf{c}) = \binom{t - (w-1)d}{w}^{-1}.$$

For a given pair of sets $B_1, B_2 \subseteq [1, n]$ where $|B_1| = k-p+1$, $|B_2| = p-1$ and $B_1 \cap B_2 = \emptyset$, let E_{B_1, B_2} be the event such

that for each column \mathbf{c}_i with $i \in B_1$ and every row r where $\mathbf{c}_i(r) = 1$ there exists an index $j \in (B_1 \cup B_2 \setminus \{i\})$ such that \mathbf{c}_j has 1 in that same row r . There are $\binom{k}{p-1} \binom{n}{k}$ such events. Then, by the same argument used in the proof of Theorem II.7 we can easily upper bound the probability of such events as follows

$$\Pr(E_{B_1, B_2}) \leq \left(\frac{\binom{w(k-1)}{w}}{\binom{t-(w-1)d}{w}} \right)^{k-p+1}. \quad (17)$$

Using Corollary II.5 we upper bound (17) as follows

$$\Pr(E_{B_1, B_2}) \leq \left(\frac{w(k-1) - \frac{w-1}{2}}{t - (w-1)d - \frac{w-1}{2}} \right)^{w(k-p+1)}. \quad (18)$$

Let us fix an arbitrary event E_{A_1, A_2} then it is easy to see that it is mutually independent from all the events $E_{A'_1, A'_2}$ such that $A'_1 \subseteq [1, n] \setminus (A_1 \cup A_2)$, $A'_2 \subseteq [1, n] \setminus (A_1 \cup A_2 \cup A'_1)$. The number of events $E_{A'_1, A'_2}$ is equal to $\binom{k}{p-1} \binom{n-k}{k}$. Therefore each event E_{A_1, A_2} is dependent of at most

$$f = \binom{k}{p-1} \left[\binom{n}{k} - \binom{n-k}{k} \right] \quad (19)$$

other events. If the probability that none of the events E_{A_1, A_2} occurs is strictly positive then there exists a matrix M that is a (k, n, d, p, w) -selector of length t . Using Lemma II.6 and taking p equal to the RHS of (18) and f as defined in equation (19), Theorem III.3 follows. \square

Corollary III.4. *There exists a (k, n, d, p, w) -selector of length t , where*

$$t \leq \left\lceil \left((w-1)d + \frac{w-1}{2} + \left(w(k-1) - \frac{w-1}{2} \right) \cdot \left(e \binom{k}{p-1} \left[\binom{n}{k} - \binom{n-k}{k} \right] \right)^{\frac{1}{w(k-p+1)}} \right) \right\rceil. \quad (20)$$

Proof. It follows rearranging the terms in equation (16). \square

Again, by exploiting the result by Moser and Tardos [12], we get a $O(n^k)$ randomized Las Vegas algorithm to construct the codes of Corollary III.4

Thanks to Corollary III.4 we obtain the following upper bound on the minimum length of (k, n, d, p) -selectors.

Corollary III.5. *There exists a (k, n, d, p) -selector of length t with $k \leq n/e$, where*

$$t \leq \ln 2 \cdot \frac{dk}{k-p+1} \log(n/k) + \ln 2 \cdot e^{3+\frac{1}{e}} \frac{k^2}{k-p+1} \log(n/k) + O(k \log(n/k)).$$

Proof. Substituting $w = \frac{k}{k-p+1} \ln(n/k)$ in (20) and using the well-known inequality $\binom{m}{s} \leq \left(\frac{em}{s}\right)^s$, we get

$$t \leq \frac{dk}{k-p+1} \ln(n/k) + e \left[e^{1+\frac{p}{k}} \left(\frac{k}{p-1} \right)^{\frac{p-1}{k}} \right]^{\frac{1}{\ln(n/k)}} \frac{k^2}{k-p+1} \ln(n/k) + O(k \ln(n/k)).$$

Hence Corollary III.5 follows since $p \leq k$, $n \geq ek$ and since the function $x^{1/x}$ takes its maximum at $x = e$. \square

A. Application of (k, n, d, p) -selectors to two-stage Group Testing with runlength constraints

We need the following result, whose proof for "classical" selectors (that is, for selectors without the runlength constraint studied in this paper) is implicit in the discussion before Theorem 3 of [4]. It is trivial to see that the proof carries out also in the present scenario.

Lemma III.6. *Let M be a (k, n, d, p) -selector with t rows, and let \mathbf{f} be the $t \times 1$ columns vector obtained by the bitwise OR of at most q , $q \leq p-1$, columns $\mathbf{c}_{i_1}, \dots, \mathbf{c}_{i_q}$ of M . Then, apart from $\mathbf{c}_{i_1}, \dots, \mathbf{c}_{i_q}$, there are at most other $k - q - 1$ columns of M whose 1's are in a subset of the positions in which the vector \mathbf{f} also has 1's.*

Now we proceed as follows. Let k be an upper bound on the number of possible positives in the Group Testing problem. We perform all the tests corresponding to the rows of a $(2k, n, d, k+1)$ -selector M , as explained in the introduction. More precisely, the generic i -th pool T_i , for $i = 1, \dots, t$, contains all elements $j \in [1, n]$ for which $M_{ij} = 1$. After having performed (in parallel) all tests on pools T_1, \dots, T_t , we get a "syndrome" vector \mathbf{f} (of dimension $t \times 1$) equal to the bitwise OR of the (at most) k columns that correspond to the unknown positive elements. The number of columns of M that are "covered" by \mathbf{f} (that is, that have their 1's in a subset of the positions in which the vector \mathbf{f} also has 1's) is upper bounded by $2k$ (by Lemma III.6). In other words, there are at most $s \leq 2k$ potentially positive elements, and the true positive are among them. Hence, one can test individually those s elements to discover the true positives. Altogether, we have used $t + 2k$ tests.

By using Corollary III.5 to estimate t , we get that we can discover all the positive elements by performing a number of tests upper bounded by a quantity that is

$$2d \ln(n/k) + O(k \ln(n/k)). \quad (21)$$

The bound (21) shows that our two-stage Group Testing algorithm outperforms both the NAGT algorithm presented in [1] and also our improved one given in the previous section of the present paper. It is interesting to notice that the bound (21) is information-theoretic optimal, for $d = O(k)$, and that this optimality can be achieved by introducing the least amount of adaptivity in the testing algorithm.

REFERENCES

- [1] A. Agarwal, O. Milenkovic and S. Pattabiraman and J. Ribeiro, Group Testing with Runlength Constraints for Topological Molecular Storage, 2020 IEEE International Symposium on Information Theory, 132-137.
- [2] N. Alon and J. Spencer, *The Probabilistic Method*, Wiley, 2008.
- [3] Y. Cheng, D.-Z. Du, G. Lin, On the upper bounds of the minimum number of rows of disjunct matrices, *Optim. Lett.* 3 (2009), 297–302.
- [4] A. De Bonis, L. Gasieniec and U. Vaccaro, Optimal Two-Stage Algorithms for Group Testing Problems, *SIAM Journal on Computing*, vol. 34, no.5, 1253–1270, 2005.
- [5] R. Dorfman, The Detection of Defective Members of Large Populations, *Ann. Math. Statist.* 14(4): 436-440 (December, 1943).
- [6] D.-Z. Du and F.K. Hwang, *Pooling Designs And Nonadaptive Group Testing: Important Tools For Dna Sequencing*, World Scientific, 2006.
- [7] P. Erdős, P. Frankl, Z. Füredi, Families of finite sets in which no set is covered by the union of r others, *Israel J. Math.* 51, 79–89 (1985).
- [8] L. Gargano, A. A. Rescigno and U. Vaccaro, Low-weight superimposed codes and related combinatorial structures: Bounds and applications, *Theoretical Computer Science* 806 (2020), 655–672.
- [9] O. Johnson, J. Scarlett and M. Aldridge, *Group Testing: An Information Theory Perspective*, Now Publishers 2019.
- [10] W. Kautz and R. Singleton, Nonrandom binary superimposed codes, *IEEE Transactions on Information Theory*, 10, (1964), 363–377.
- [11] J. Komlos and A. Greenberg, An asymptotically fast nonadaptive algorithm for conflict resolution in multiple-access channels, *IEEE Trans. on Information Theory*, 31, Issue: 2, March 1985, 302–306.
- [12] R.A. Moser and G. Tardos, A constructive proof of the general Lovász local lemma, *Journal of the ACM*, 57 (2010), 1–15.
- [13] C. Shangguan and G. Ge, New Bounds on the Number of Tests for Disjunct Matrices, *Transactions on Information Theory*, vol. 62, no. 12, pp. 7518-7521, Dec. 2016.