

# Preserving Privacy while Broadcasting: $k$ -Limited-Access Schemes

Mohammed Karmoose, Linqi Song, Martina Cardone, Christina Fragouli  
University of California Los Angeles, Los Angeles, CA 90095 USA  
Email: {mkarmoose, songlinqi, martina.cardone, christina.fragouli}@ucla.edu

**Abstract**—Index coding employs coding across clients within the same broadcast domain. This typically assumes that all clients learn the coding matrix so that they can decode and retrieve their requested data. However, learning the coding matrix can pose privacy concerns: it may enable clients to infer information about the requests and side information of other clients [1]. In this paper, we formalize the intuition that the achieved privacy can increase by decreasing the number of rows of the coding matrix that a client learns. Based on this, we propose the use of  $k$ -limited-access schemes: given an index coding scheme that employs  $T$  transmissions, we create a  $k$ -limited-access scheme with  $T_k \geq T$  transmissions, and with the property that each client learns at most  $k$  rows of the coding matrix to decode its message. We derive upper and lower bounds on  $T_k$  for all values of  $k$ , and develop deterministic designs for these schemes for which  $T_k$  has an order-optimal exponent for some regimes.

## I. INTRODUCTION

Consider a server broadcasting publicly available messages to clients, for instance YouTube videos. It is well recognized that the use of coding and side information can offer significant bandwidth savings when broadcasting [2]. However, it can also pose privacy concerns [1]: although messages are publicly available, clients may wish to preserve the anonymity of their requests from other clients. A curious client, by leveraging the broadcast information, may be able to infer what the requests and side information of other clients are. In this paper, we propose new schemes that seek to balance the bandwidth benefits that coding offers with privacy considerations.

We pose this problem within the index coding framework [2]. In index coding, a server has  $m$  messages and can losslessly broadcast to  $n$  clients. Each client requests a specific message and may have a subset of the messages as side information. To satisfy all clients with the minimum number of transmissions  $T$ , the server can send coded broadcast transmissions; the clients then use the coding matrix<sup>1</sup> to decode their messages. In [1], we showed that, by knowing the coding matrix, a curious client can infer information about the side information and requests of other clients.

This paper builds on a new observation: it may not be necessary to provide clients with the entire coding matrix, but with only the rows required for them to decode their own message. For example, assume we have  $m = 4$  messages and  $n = 4$  clients, where client  $i \in \{1, 3\}$  has message  $b_i$  and would like to receive message  $b_{i+1}$ , and client  $i \in \{2, 4\}$  has

message  $b_i$  and would like to receive message  $b_{i-1}$ . The server can satisfy all clients with two broadcast coded transmissions, namely  $b_1 + b_2$  and  $b_3 + b_4$ , i.e., it uses a  $2 \times 4$  coding matrix. To decode their message, clients 1 and 2 only need to know the first row of this matrix (the fact that the first combination is  $b_1 + b_2$ ), and similarly clients 3 and 4 only need to know the second row of the matrix. By restricting the access to the coding matrix, we limit the privacy leakage: the less rows a client learns, the less it can infer about other clients.

We turn around this observation and ask: what if we restrict each user to access at most  $k$  rows of the coding matrix? In particular, assume we are given a coding matrix that uses  $T$  transmissions to satisfy all clients. Can we “transform” it into an “equivalent” coding matrix that potentially uses  $T_k \geq T$  transmissions to satisfy all clients, but where now each client needs to learn at most  $k$  rows of it to decode its message? We refer to the coding schemes that satisfy this condition as  $k$ -limited-access schemes and we evaluate their benefits, cost and feasibility. Our main contributions are:

- 1) *Benefits*: we formalize the intuition that the achieved level of privacy can increase by decreasing the number of rows of the coding matrix that a client learns.
- 2) *Cost*: we derive upper and lower bounds on  $T_k$  that highlight the maximum and minimum cost to pay in terms of additional broadcast transmissions as a function of  $k$ .
- 3) *Feasibility*: we propose deterministic designs for  $k$ -limited-access schemes, for all values of  $k$ . For some regimes, our designs provide values of  $T_k$  whose exponents are order-optimal.

The paper is organized as follows. Section II defines the problem setup. Section III presents our main results, i.e., it formalizes the intuition that privacy benefits can be achieved by limiting clients’ access to the coding matrix, and it provides upper and lower bounds on the number of transmissions needed to satisfy clients when they know only part of the matrix. Section IV proves the upper bounds presented in Section III by designing  $k$ -limited-access schemes and assessing their performance. Section V positions our work with respect to related literature and Section VI concludes the paper.

## II. SETUP AND PROBLEM FORMULATION

**Notation.** Calligraphic letters indicate sets; boldface lower case letters denote vectors and boldface upper case letters indicate matrices;  $|\mathcal{X}|$  is the cardinality of  $\mathcal{X}$ ;  $[n]$  is the set

<sup>1</sup>The coding matrix has size  $T \times m$  and collects in each row the coding coefficients used for the corresponding broadcast transmission.

of integers  $\{1, \dots, n\}$ ; for all  $x \in \mathbb{R}$ , the floor and ceiling functions are denoted with  $\lfloor x \rfloor$  and  $\lceil x \rceil$ , respectively;  $\mathbf{0}_j$  is the all-zero row vector of dimension  $j$ ;  $\mathbf{1}_j$  denotes a row vector of dimension  $j$  of all ones and  $\mathbf{I}_j$  is the identity matrix of dimension  $j$ ;  $\mathbf{e}_i^j$  is the all-zero row vector of length  $j$  with a 1 in position  $i$ ; logarithms are in base 2.

**Index Coding.** We consider an index coding instance, where a server has a database  $\mathcal{B}$  of  $m$  messages  $\mathcal{B} = \{\mathbf{b}_{\mathcal{M}}\}$ , where  $\mathcal{M} = [m]$  is the set of message indices, and  $\mathbf{b}_j \in \mathbb{F}_2^F, j \in \mathcal{M}$ , with  $F$  being the message size. The server is connected through a broadcast channel to a set of clients  $\mathcal{C} = \{c_{\mathcal{N}}\}$ , where  $\mathcal{N} = [n]$  is the set of client indices. We assume that  $m \geq n$ . Each client  $c_i, i \in \mathcal{N}$ , has a subset of the messages  $\{\mathbf{b}_{\mathcal{S}_i}\}$ , with  $\mathcal{S}_i \subset \mathcal{M}$ , as side information and requests a new message  $\mathbf{b}_{q_i}$  with  $q_i \in \mathcal{M} \setminus \mathcal{S}_i$  that it does not have. We assume that the server employs a *linear code*, i.e., it designs a set of broadcast transmissions that are linear combinations of the messages in  $\mathcal{B}$ . The linear index code can be represented as  $\mathbf{A}\mathbf{B} = \mathbf{Y}$ , where  $\mathbf{A} \in \mathbb{F}_2^{T \times m}$  is the coding matrix,  $\mathbf{B} \in \mathbb{F}_2^{m \times F}$  is the matrix of all the messages and  $\mathbf{Y} \in \mathbb{F}_2^{T \times F}$  is the resulting matrix of linear combinations. Upon receiving  $\mathbf{Y}$ , client  $c_i, i \in \mathcal{N}$ , employs linear decoding to retrieve  $\mathbf{b}_{q_i}$ .

**Problem Formulation.** In [2], it was shown that the index coding problem is equivalent to the rank minimization of an  $n \times m$  matrix  $\mathbf{G} \in \mathbb{F}_2^{n \times m}$  whose  $i$ -th row  $\mathbf{g}_i, i \in [n]$ , has the following properties: (i) has a 1 in the position  $q_i$ , (ii) has a 0 in the  $j$ -th position for all  $j \in \mathcal{M} \setminus \mathcal{S}_i$ , (iii) can have either 0 or 1 in all the remaining positions. With this representation,  $c_i$  can successfully decode  $\mathbf{b}_{q_i}$  using a linear combination of the messages corresponding to the non-zero entries of  $\mathbf{g}_i$ . Finding an optimal linear coding scheme (i.e., with minimum number of transmissions) is equivalent to completing  $\mathbf{G}$  so that it has the minimum possible rank. Once we have one such  $\mathbf{G}$ , we can use a basis of the row space of  $\mathbf{G}$  (of size  $T = \text{rank}(\mathbf{G})$ ) as coding matrix  $\mathbf{A}$ . In this case, in fact, client  $c_i$  can construct  $\mathbf{g}_i$  as a linear combination of the rows of  $\mathbf{A}$ , i.e.,  $c_i$  performs the decoding operation  $\mathbf{d}_i \mathbf{A} \mathbf{B} = \mathbf{d}_i \mathbf{Y}$ , where  $\mathbf{d}_i \in \mathbb{F}_2^T$  is the decoding row vector of  $c_i$  chosen such that  $\mathbf{d}_i \mathbf{A} = \mathbf{g}_i$ . We remark that any index coding scheme that satisfies all clients with  $T$  transmissions (where  $T$  is not necessarily optimal) – and can be obtained by any index code design algorithm [3]–[5] – corresponds to a completion of  $\mathbf{G}$  (i.e., given  $\mathbf{A} \in \mathbb{F}_2^{T \times m}$ , we can create a corresponding  $\mathbf{G}$  in polynomial time).

In our problem formulation we assume we start with a given matrix  $\mathbf{G}$  of rank  $T$ , i.e., we are given  $n$  *distinct* vectors that belong to a  $T$ -dimensional subspace. Using a basis of the row space of the given  $\mathbf{G}$ , we construct  $\mathbf{A} \in \mathbb{F}_2^{T \times m}$ . Then, we ask: *Given  $n$  distinct vectors  $\mathbf{g}_i, i \in [n]$ , in a  $T$ -dimensional space, can we find a minimum-size set  $\mathcal{A}_k$  with  $T_k \geq T$  vectors, such that each  $\mathbf{g}_i$  can be expressed as a linear combination of at most  $k$  vectors in  $\mathcal{A}_k$  (with  $1 \leq k \leq T$ )?*

The vectors in  $\mathcal{A}_k$  form the rows of the coding matrix  $\mathbf{A}_k$  we will employ. We can equivalently restate this as follows.

*Given a coding matrix  $\mathbf{A}$ , can we find  $\mathbf{P} \in \mathbb{F}_2^{T_k \times T}$ , with  $T_k$  as small as possible, such that  $\mathbf{A}_k = \mathbf{P}\mathbf{A}$  and each row in  $\mathbf{G}$*

*can be reconstructed by combining at most  $k$  rows of  $\mathbf{A}_k$ ?*

Note that  $k = T$  corresponds to the conventional transmission scheme of an index coding problem for which  $\mathbf{P} = \mathbf{I}_T$ .

**Transmission Overhead.** We note that the server can privately share the (at most)  $k$  coding vectors that each  $c_i$  needs by using a private secret key or a dedicated channel (e.g., the same channel used by  $c_i$  to convey the request  $q_i$  to the server). Thus, using a  $k$ -limited-access scheme incurs an extra transmission overhead to privately convey the coding vectors. In particular, the total number of transmitted bits  $C_k$  is upper bounded by  $C_k \leq nkm + T_k F$ , while the total number of transmitted bits  $C$  using a conventional scheme is  $C = T(F + m)$ . We observe that the extra overhead incurred is negligible in comparison to the broadcast transmissions that convey the encoded messages when  $n$  and  $m$  are both  $o(F)$ , which is a reasonable assumption for large file sizes (for instance, when sharing YouTube videos).

### III. MAIN RESULTS

Consider the setup in the previous section and suppose that client  $c_1$  is curious, i.e., by leveraging the  $k$  (linearly independent) rows of  $\mathbf{A}_k$  that it receives, it seeks to infer information about  $c_i, i \in [n], i \neq 1$ . We are interested in quantifying the amount of information that  $c_1$  can obtain about  $q_i$  (i.e., the identity of the request of  $c_i$ ) as a function of  $k$ .

As a first step towards this end, we define our privacy metric as follows. We assume that the index coding instance is random and we let  $L$  (respectively,  $L_1$ ) be the random variable associated with the subspace spanned by the  $T$  rows of the coding matrix  $\mathbf{A} \in \mathbb{F}_2^{T \times m}$  (respectively, spanned by the  $k$  vectors given to  $c_1$ ). Assume that  $c_1$  knows  $T$ . Then,

**Definition III.1.** The privacy metric is defined as  $H(L|L_1, T)$ , i.e., it quantifies the amount of uncertainty (entropy) that  $c_1$  has about the subspace spanned by the  $T$  rows of the index coding matrix  $\mathbf{A}$ .

The main motivation behind our choice of the privacy metric is that it offers a yardstick for evaluating the amount of information that  $c_1$  can obtain about  $q_i$ . This is because  $\mathbf{g}_i \in \mathbb{F}_2^m$  (that  $c_i$  needs to recover  $\mathbf{b}_{q_i}$ ) lies in the subspace spanned by the  $T$  rows of  $\mathbf{A}$ . Then, given the specific realizations  $T = t$  and  $L_1 = \ell_1$ , we compute

$$P_k = H(L|L_1 = \ell_1, T = t) \stackrel{(a)}{=} \log(|\mathcal{L}(t, \ell_1)|) \\ \stackrel{(b)}{=} \log\left(\prod_{\ell=0}^{t-k-1} \frac{2^m - 2^{k+\ell}}{2^t - 2^{k+\ell}}\right) \stackrel{m \gg t}{\approx} m(t-k), \quad (1)$$

where: (i) in (a) we let  $\mathcal{L}(t, \ell_1)$  represent the set of subspaces  $L_t \subset \mathbb{F}_2^m$  of dimension  $t$  such that  $\ell_1 \subset L_t$ ; moreover, the equality follows by assuming that the underlying system maintains a uniform distribution across all feasible  $t$ -dimensional subspaces of  $\mathbb{F}_2^m$ ; (ii) the equality in (b) follows by standard counting arguments used to characterize the number of distinct subspaces of a given dimension in a vector space. It is clear that, when  $m \gg t$ , then  $P_k$  in (1) decreases linearly with  $k$ , i.e., the less rows of the coding matrix  $c_1$  learns, the less it

can infer about the subspace spanned by the  $T$  rows of the coding matrix  $\mathbf{A}$ . This suggests that, by increasing  $k$ ,  $c_1$  has more uncertainty about  $q_i$ . It is also clear that  $P_k$  in (1) is zero when  $k = t$ ; this is because, under this condition,  $c_1$  receives the entire index coding matrix and hence it will be able to perfectly reconstruct the subspace spanned by its rows. However, although  $P_k$  in (1) is zero when  $k = t$ ,  $c_1$  might still have uncertainty about  $q_i$  [1]. Quantifying this uncertainty is an interesting open problem that does not appear to be an easy task; this uncertainty, in fact, depends on the underlying system, e.g., on the index code used by the server and on the distribution with which the index coding matrix is selected.

We now build on the analysis above – that shows the benefits of limiting the access of the clients to the coding matrix – and focus on finding conditions that guarantee that  $\mathbf{P}$  can be constructed while ensuring that each client  $c_i, i \in [n]$ , successfully decodes its request  $\mathbf{b}_{q_i}$  using at most  $k$  transmissions. Towards this end, we derive upper and lower bounds on  $T_k$ . In particular, our main result is stated in the theorem below.

**Theorem III.1.** *Given an index coding matrix  $\mathbf{A} \in \mathbb{F}_2^{T \times m}$  with  $T \geq 2$ , it is possible to transform it into  $\mathbf{A}_k = \mathbf{P}\mathbf{A}$  with  $\mathbf{P} \in \mathbb{F}_2^{T_k \times T}$ , such that each client can recover its request by combining at most  $k$  rows of it, if and only if*

$$T_k \geq T^* = \min \left\{ T_k : \sum_{i=1}^k \binom{T_k}{i} \geq n \right\}. \quad (2)$$

Moreover, there exist constructions of  $\mathbf{P}$  such that:

- When  $\lceil T/2 \rceil \leq k < T$ , then

$$T_k \leq \min \{n, T + 1\}; \quad (3)$$

- When  $1 \leq k < \lceil T/2 \rceil$ , then

$$T_k \leq \min \{n, T_{ub}\}, \quad (4)$$

where

$$T_{ub} = \begin{cases} 2^T & \text{if } k = 1 \\ 2(2k + 1)^{\lceil \frac{T}{2k-1} \rceil - 1} & \text{if } T_{last} = 1 \\ (2k + 1)^{\lceil \frac{T}{2k-1} \rceil - 1} (T_{last} + 2) & \text{otherwise} \end{cases} \quad (5)$$

$$= 2^{O(\frac{T}{k} \log k)} \quad \text{if } k \neq 1,$$

$$\text{where } T_{last} = T - (2k - 1) \left( \left\lceil \frac{T}{2k-1} \right\rceil - 1 \right).$$

We provide the proof of the lower bound in (2) in the Appendix, while in Section IV we give explicit constructions for  $\mathbf{P}$  for the two regimes in Theorem III.1, hence proving the upper bounds on  $T_k$  in (3) and (4). The results in Theorem III.1 also imply the following lemma (see also the Appendix).

**Lemma III.2.** *Consider the regime  $n = 2^T - 1$ . We have*

- When  $\lceil T/2 \rceil \leq k < T$ , the bounds in (2) and (3) coincide, i.e., the provided construction of  $\mathbf{P}$  is optimal;
- When  $1 \leq k < \lceil T/2 \rceil$ , then the bound in (2) becomes

$$T_k \geq \frac{k}{e} \left( \frac{2^T - 1}{k} \right)^{1/k} = 2^{\Omega(\frac{T}{k} + \alpha \log k)}, \quad \alpha = \frac{k-1}{k}. \quad (6)$$

We now conclude this section with some comparisons between the lower and upper bounds on  $T_k$  for the case  $n = 2^T - 1$ . According to Lemma III.2, a construction of  $\mathbf{P}$  with  $T_k = T + 1$  (provided in Section IV) is optimal for  $k \geq \lceil T/2 \rceil$ . In other words, by adding only one more transmission to the original index code, clients need *at most* half of the transmissions to recover their request; this enhances the attained level of privacy. Differently, for  $1 \leq k < \lceil T/2 \rceil$ , the orders of the lower bound in (6) and upper bound in (4) are different. This implies that the construction of  $\mathbf{P}$  for this regime (see Section IV for the details) is not optimal. However, we show next that there exist some regimes of  $k$  where the two bounds are close in order. In particular,

- **$k$  is constant.** In this case, we have  $T_k = 2^{\Theta(T)}$ , i.e., the upper and lower bounds have the same order in the exponent.
- **$k = T/c$  where  $c > 1$  is constant.** In this case, we have  $T_k = 2^{\Theta(\log T)}$ , i.e., this represents another regime where the upper and lower bounds have the same order in the exponent.

#### IV. CONSTRUCTIONS OF $k$ -LIMITED-ACCESS SCHEMES

In this section, we give explicit constructions of the  $\mathbf{P}$  matrix and prove the two upper bounds on  $T_k$  in (3) and (4). Our design of  $\mathbf{P}$  allows to reconstruct any of the  $2^T$  vectors of size  $T$ . Recall that  $\mathbf{A}$  is full rank and that the  $i$ -th row of  $\mathbf{G}$  can be expressed as  $\mathbf{g}_i = \mathbf{d}_i \mathbf{A}$ , where  $\mathbf{d}_i \in \mathbb{F}_2^T$  is the coefficients row vector associated with  $\mathbf{g}_i$ .

**Case I:**  $\lceil T/2 \rceil \leq k < T$ . When  $n \geq T + 1$ , let

$$\mathbf{P} = \begin{bmatrix} \mathbf{I}_T \\ \mathbf{1}_T \end{bmatrix}, \quad (7)$$

which results in a matrix  $\mathbf{A}_k$  with  $T_k = T + 1$ , matching the bound in (3). We now show that each  $\mathbf{g}_i = \mathbf{d}_i \mathbf{A}$ ,  $i \in [n]$ , can be reconstructed by combining up to  $k$  vectors of  $\mathbf{A}_k$ . Let  $w(\mathbf{d}_i)$  be the Hamming weight of  $\mathbf{d}_i$ . If  $w(\mathbf{d}_i) \leq \lceil T/2 \rceil$ , then we can reconstruct  $\mathbf{g}_i$  as  $\mathbf{g}_i = [\mathbf{d}_i \mathbf{0}] \mathbf{A}_k$ , which involves adding  $w(\mathbf{d}_i) \leq \lceil T/2 \rceil \leq k$  rows of  $\mathbf{A}_k$ . Differently, if  $w(\mathbf{d}_i) \geq \lceil T/2 \rceil + 1$ , then we can reconstruct  $\mathbf{g}_i$  as  $\mathbf{g}_i = [\mathbf{d}_i \mathbf{1}] \mathbf{A}_k$ , where  $\bar{\mathbf{d}}_i$  is the bitwise complement of  $\mathbf{d}_i$ . In this case, reconstructing  $\mathbf{g}_i$  involves adding  $T - w(\mathbf{d}_i) + 1 \leq \lfloor T/2 \rfloor \leq k$  rows of  $\mathbf{A}_k$ .

When  $n < T + 1$ , then it is sufficient to send  $n$  uncoded transmissions, where the  $i$ -th transmission satisfies  $c_i, i \in [n]$ . In this case  $c_i$  has access only to the  $i$ -th transmission, i.e.,  $k = 1$ . This completes the proof of the upper bound in (3).

**Case II:**  $1 \leq k < \lceil T/2 \rceil$ . First, we consider  $n \geq T_{ub}$ , where  $T_{ub}$  is defined in (5). For this, we provide a construction for  $\mathbf{P}$  that is based on multiple uses of the construction in Case I. In what follows, we let  $T_c = \lceil \frac{T}{2k-1} \rceil$ . Consider the following sets of *distinct* vectors (i.e., by omitting replicated vectors)

$$\mathcal{P}_j = \{ \mathbf{0}_{2k-1}, \mathbf{1}_{2k-1}, \mathbf{e}_i^{2k-1}, \forall i \in [2k-1] \}, j \in [T_c - 1], \quad (8a)$$

$$\mathcal{P}_{T_c} = \{ \mathbf{0}_{T_{last}}, \mathbf{1}_{T_{last}}, \mathbf{e}_i^{T_{last}}, \forall i \in [T_{last}] \}, \quad (8b)$$

where  $T_{last} = T - (2k - 1)(T_c - 1)$ . Then, our construction of  $\mathbf{P}$  is based on different concatenations of various elements of the above sets as we explain in what follows. Let

$$\mathcal{P} = \mathcal{P}_1 \times \mathcal{P}_2 \times \dots \times \mathcal{P}_{T_c}$$

$$\mathbf{P} = \begin{bmatrix} 0 & 0 & 0 & | & 0 & 0 & 0 & | & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 & | & 1 & 1 \\ 0 & 0 & 0 & | & 0 & 0 & 0 & | & 1 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 & | & 0 & 1 \\ 0 & 0 & 0 & | & 1 & 1 & 1 & | & 0 & 0 \\ \vdots & & & & \vdots & & & & \vdots & \\ 0 & 0 & 1 & | & 0 & 0 & 1 & | & 0 & 0 \\ 0 & 0 & 1 & | & 0 & 0 & 1 & | & 1 & 1 \\ 0 & 0 & 1 & | & 0 & 0 & 1 & | & 1 & 0 \\ 0 & 0 & 1 & | & 0 & 0 & 1 & | & 0 & 1 \end{bmatrix}$$

Figure 1. Construction of  $\mathbf{P}$  for  $T = 8$  and  $k = 2$ .

be the Cartesian product of the sets defined in (8) and  $\mathcal{P}(i)$ ,  $i \in [T_k]$ , be the  $i$ -th tuple of  $\mathcal{P}$ . Then, the  $i$ -th row vector  $\mathbf{p}_i$  of  $\mathbf{P}$  is constructed by concatenating the elements of the tuple  $\mathcal{P}(i)$  in their respective order (i.e., the first element of  $\mathcal{P}(i)$  is the left-most part of  $\mathbf{p}_i$ , the second element of  $\mathcal{P}(i)$  is the second left-most part of  $\mathbf{p}_i$  and so on). It is not difficult to see that with this construction  $\mathbf{p}_i$  has length  $T$ . Since from (8) we have that, for  $j \in [T_c - 1]$ ,

$$|\mathcal{P}_j| = \begin{cases} 2 & k = 1 \\ 2k + 1 & k > 1 \end{cases} \quad |\mathcal{P}_{T_c}| = \begin{cases} 2 & T_{\text{last}} = 1 \\ T_{\text{last}} + 2 & T_{\text{last}} > 1 \end{cases},$$

then we have  $\prod_{j=1}^{T_c} |\mathcal{P}_j|$  possible different ways of concatenating vectors from these sets. This gives the bound in (4). To illustrate this process consider the following example.

*Example.* Let  $T = 8$  and  $k = 2$  for which  $T_c = 3$  and  $T_{\text{last}} = 2$ . Then, we have

$$\begin{aligned} \mathcal{P}_1 = \mathcal{P}_2 &= \{[0 \ 0 \ 0], [1 \ 1 \ 1], [1 \ 0 \ 0], \\ &\quad [0 \ 1 \ 0], [0 \ 0 \ 1]\}, \\ \mathcal{P}_3 &= \{[0 \ 0], [1 \ 1], [1 \ 0], [0 \ 1]\}. \end{aligned}$$

Figure 1 shows how  $\mathbf{P}$  is then constructed.

We now need to prove that any  $\mathbf{g}_i$ ,  $i \in [n]$ , can be reconstructed using at most  $k$  rows of  $\mathbf{A}_k$ . Notice that this is equivalent to showing that we need at most  $k$  rows of  $\mathbf{P}$  to reconstruct  $\mathbf{d}_i$ ,  $i \in [n]$ . This is because, if this holds, then  $\mathbf{d}_i = \mathbf{d}_i^* \mathbf{P}$  where the row vector  $\mathbf{d}_i^* \in \mathbb{F}_2^{T_k}$  has at most  $k$  non-zero elements. Then, this would imply  $\mathbf{g}_i = \mathbf{d}_i \mathbf{A} = \mathbf{d}_i^* \mathbf{P} \mathbf{A} = \mathbf{d}_i^* \mathbf{A}_k$ , i.e.,  $\mathbf{g}_i$  is reconstructed by using at most  $k$  rows of  $\mathbf{A}_k$ . In what follows, we therefore prove that any  $\mathbf{d}_i$ ,  $i \in [n]$ , can be reconstructed by using at most  $k$  rows of  $\mathbf{P}$ . As a running example to illustrate the different steps of our proof we use the case in Figure 1 with  $\mathbf{d}_i = [1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]$ .

**Step 1.** Starting from the left-most bit, we split  $\mathbf{d}_i$ ,  $i \in [n]$ , into  $(T_c - 1)$  parts of length  $(2k - 1)$  and one last part of length  $T_{\text{last}}$ . We denote by  $\mathbf{d}_i(j)$  the  $j$ -th part with  $j \in [T_c]$ .

*Running example.* We have

$$\mathbf{d}_i(1) = [1 \ 1 \ 0], \quad \mathbf{d}_i(2) = [0 \ 1 \ 0], \quad \mathbf{d}_i(3) = [0 \ 0].$$

**Step 2.** We leverage our proof of Case I, where we showed that  $\mathbf{P}$  in (7) can be used to reconstruct any vector of length  $T$  using  $\lceil T/2 \rceil \leq k < T$  rows. This, in fact, implies that: (i) any  $\mathbf{d}_i(j)$ ,  $j \in [T_c - 1]$ , can be reconstructed by adding at most

$k$  elements of  $\mathcal{P}_j$  (excluding the first element), and (ii) any  $\mathbf{d}_i(T_c)$  can be reconstructed by adding at most  $k$  elements of  $\mathcal{P}_{T_c}$  (excluding the first element). We let  $\mathcal{R}_j$ ,  $j \in [T_c]$ , be the set of elements of  $\mathcal{P}_j$  needed to reconstruct  $\mathbf{d}_i(j)$ . Clearly,  $|\mathcal{R}_j| \leq k$ ,  $j \in [T_c]$ . Let  $R^* = \max_{j \in [T_c]} |\mathcal{R}_j|$ . Then, we further populate  $\mathcal{R}_j$ ,  $j \in [T_c]$  with  $R^* - |\mathcal{R}_j|$  zero vectors, so that all  $\mathcal{R}_j$  have the same cardinality.

*Running example.* We have  $R^* = 2$  and

$$\begin{aligned} \mathcal{R}_1 &= \{[1 \ 0 \ 0], [0 \ 1 \ 0]\}, \\ \mathcal{R}_2 &= \{[0 \ 1 \ 0], [0 \ 0 \ 0]\}, \\ \mathcal{R}_3 &= \{[0 \ 0], [0 \ 0]\}. \end{aligned}$$

**Step 3.** We concatenate the different elements of  $\mathcal{R}_j$ ,  $j \in [T_c]$ . In particular, for each  $\ell \in [R^*]$  we create a vector of length  $T$  by concatenating the elements in the  $\ell$ -th position of all  $\mathcal{R}_j$ ,  $j \in [T_c]$ , as follows: we put the  $\ell$ -th element of  $\mathcal{R}_1$  as the left-most part, then we concatenate to it the  $\ell$ -th element of  $\mathcal{R}_2$  and so on until  $\mathcal{R}_{T_c}$ . Thus, we obtain a set  $\mathcal{R}^*$  of  $R^*$  vectors of length  $T$ . Clearly, from our construction of  $\mathbf{P}$ , each element of  $\mathcal{R}^*$  is a row of  $\mathbf{P}$ . Moreover, from our construction in the previous step of  $\mathcal{R}_j$ ,  $j \in [T_c]$ , we have that the sum of the  $R^*$  vectors in  $\mathcal{R}_j$  reconstructs  $\mathbf{d}_i(j)$ . Hence, it is not difficult to see that the sum of the  $R^*$  elements of  $\mathcal{R}^*$  reconstructs  $\mathbf{d}_i$ .

*Running example.* We have

$$\mathcal{R}^* = \left\{ \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \right\}.$$

By adding the two elements of  $\mathcal{R}^*$  we obtain  $[1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]$ , which is precisely the  $\mathbf{d}_i$  we wanted to reconstruct.

When  $n < T_{\text{ub}}$ , then it is sufficient to send  $n$  uncoded transmissions, where the  $i$ -th transmission satisfies  $c_i$ ,  $i \in [n]$ . In this case  $c_i$  has access only to the  $i$ -th transmission, i.e.,  $k = 1$ . This completes the proof of the upper bound in (4).

## V. RELATED WORK

The problem of protecting privacy was initially proposed to enable the disclosure of databases for public access, while maintaining the anonymity of the users [6]. Similar concerns have been raised in the context of *Private Information Retrieval* (PIR), which was introduced in [7] and has received a fair amount of attention [8], [9]. In particular, in PIR the goal is to ensure that no information about clients' requests is revealed to a set of malicious databases when clients are trying to retrieve information from them. Similarly, the problem of *Oblivious Transfer* (OT) was studied [10], [11] to establish, by means of cryptographic techniques, two-way private connections between the clients and the server.

We were here interested in addressing privacy concerns in broadcast domains. In particular, we analyzed this problem within the index coding framework, as we recently proposed in [1]. This problem differs from secure index coding [12], where the goal is to guarantee that each client does not learn any information about the *content* of the messages other than its request. Differently, our goal was to limit the information

that a client can learn about the *identities* of the requests of other clients. Moreover, our approach here has a significant difference with respect to [1]. In fact, while in [1] our goal was to design the encoding matrix to guarantee a high-level of privacy, here we assumed that an index coding matrix (that satisfies all clients) is given to us and we developed methods to increase its achieved levels of privacy.

The solution that we here proposed to limit the privacy leakage is based on finding overcomplete bases. This approach is closely related to compressed sensing and dictionary learning [13], where the goal is to learn a dictionary of signals such that other signals can be *sparsely* and *accurately* represented using atoms from this dictionary. These problems seek lossy solutions, i.e., signal reconstruction is not necessarily perfect. This allows a convex optimization formulation of the problem, which can be solved efficiently [14]. In contrast, our problem was concerned with lossless reconstructions, in which case the optimization problem is no longer convex.

## VI. CONCLUSION

We studied an index coding problem, where clients are eager to learn the identity of the request of other clients. We proposed the use of  $k$ -limited-access schemes to mitigate the privacy risks, which provide clients with only part of the coding matrix and still ensure that they can all recover their requested message. We showed that such approach can achieve higher levels of privacy than conventional schemes (where the entire matrix is broadcast to all users) at the cost of additional number of transmissions. This analysis sheds light on an inherent tradeoff between bandwidth savings and privacy protection in broadcast domains. Future work would include the derivation of tighter upper and lower bounds on the number of transmissions required by a  $k$ -limited-access scheme.

## APPENDIX

In order to prove the lower bound in (2), we establish a connection between our problem and a linear-algebraic one, namely the problem of representing vectors in finite vector spaces. Given a matrix  $\mathbf{A}$ , denote by  $\mathbb{V}_{\mathbf{A}} \subseteq \mathbb{F}_2^T$  the subspace formed by the span of the rows of  $\mathbf{A}$ . It is clear that the dimension of  $\mathbb{V}_{\mathbf{A}}$  is at most  $T$  (exactly  $T$  if  $\mathbf{A}$  is full rank) and that the  $n$  distinct rows of  $\mathbf{G}$  lie in  $\mathbb{V}_{\mathbf{A}}$ . Let  $\mathbf{a}_i \in \mathbb{F}_2^m, i \in [T_k]$ , be the  $i$ -th row of  $\mathbf{A}_k$ . Then this problem is equivalent to the following: *what is a minimum-size set of vectors  $\mathcal{A}_k = \{\mathbf{a}_{[T_k]}\}$  such that any row vector of  $\mathbf{G}$  can be represented by a linear combination of at most  $k$  vectors of  $\mathcal{A}_k$ ?*

A lower bound on  $T_k$  can be obtained as follows. Given  $\mathcal{A}_k$ , there must exist a linear combination of at most  $k$  vectors of  $\mathcal{A}_k$  that is equal to each of the  $n$  distinct row vectors of  $\mathbf{G}$ . The number of *distinct* non-zero linear combinations of up to  $k$  vectors is at most equal to  $\sum_{j=1}^k \binom{T_k}{j}$ . Thus, we have

$$\sum_{i=1}^k \binom{T_k}{i} \geq n, \quad (9)$$

which gives precisely the bound in (2).

We now derive the lower bounds in Lemma III.2, i.e., we evaluate (2) for  $n = 2^T - 1$ . From (9), we obtain

$$\sum_{i=1}^k \binom{T_k}{i} \geq 2^T - 1. \quad (10)$$

Since in general  $T_k \geq T$ , to prove that  $T_k \geq T+1$  for  $k < T$ , it is sufficient to show that we have a contradiction for  $T_k = T$ . Indeed, by setting  $T_k = T$ , the bound in (10) becomes

$$\sum_{i=1}^k \binom{T}{i} \geq 2^T - 1 = \sum_{i=1}^T \binom{T}{i},$$

which clearly is not possible since  $k < T$ . Hence,  $T_k \geq T+1$  for all  $k < T$ . However, for  $1 \leq k < \lceil T/2 \rceil$ , we can refine this lower bound as follows

$$\begin{aligned} k \binom{T_k e}{k} &\geq k \binom{T_k}{k} \geq \sum_{i=1}^k \binom{T_k}{i} \geq 2^T - 1 \\ \implies T_k &\geq \frac{k}{e} \left( \frac{2^T - 1}{k} \right)^{1/k}. \end{aligned}$$

This concludes the proof of the lower bounds in Lemma III.2.

## REFERENCES

- [1] M. Karmoose, L. Song, M. Cardone, and C. Fragouli, "Private broadcasting: an index coding approach," to appear in *IEEE International Symposium on Information Theory (ISIT)*, June 2017.
- [2] Z. Bar-Yossef, Y. Birk, T. Jayram, and T. Kol, "Index coding with side information," *IEEE Transactions on Information Theory*, vol. 57, no. 3, pp. 1479–1494, February 2011.
- [3] H. Esfahanizadeh, F. Lahouti, and B. Hassibi, "A matrix completion approach to linear index coding problem," in *IEEE Information Theory Workshop (ITW)*, November 2014, pp. 531–535.
- [4] X. Huang and S. El Rouayheb, "Index coding and network coding via rank minimization," in *IEEE Information Theory Workshop-Fall (ITW)*, October 2015, pp. 14–18.
- [5] M. A. R. Chaudhry and A. Sprintson, "Efficient algorithms for index coding," in *INFOCOM Workshops 2008, IEEE*, April 2008, pp. 1–4.
- [6] C. C. Aggarwal and S. Y. Philip, "A general survey of privacy-preserving data mining models and algorithms," in *Privacy-preserving data mining*. Springer, 2008, pp. 11–52.
- [7] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan, "Private information retrieval," *Journal of the ACM (JACM)*, vol. 45, no. 6, pp. 965–981, November 1998.
- [8] R. Freij-Hollanti, O. Gnilke, C. Hollanti, and D. Karpuk, "Private information retrieval from coded databases with colluding servers," *arXiv:1611.02062*, November 2016.
- [9] K. Banawan and S. Ulukus, "The capacity of private information retrieval from coded databases," *arXiv:1609.08138*, September 2016.
- [10] G. Brassard, C. Crepeau, and J.-M. Robert, "All-or-nothing disclosure of secrets," *Advances in Cryptology: Proceedings of Crypto '86*, Springer-Verlag, pp. 234–238, 1987.
- [11] M. Mishra, B. K. Dey, V. M. Prabhakaran, and S. Diggavi, "The oblivious transfer capacity of the wiretapped binary erasure channel," in *IEEE International Symposium on Information Theory*, June 2014, pp. 1539–1543.
- [12] S. H. Dau, V. Skachek, and Y. M. Chee, "On the security of index coding with side information," *IEEE Transactions on Information Theory*, vol. 58, no. 6, pp. 3975–3988, June 2012.
- [13] G. Chen and D. Needell, "Compressed sensing and dictionary learning," *Finite Frame Theory: A Complete Introduction to Overcompleteness*, vol. 73, p. 201, 2016.
- [14] R. Rubinfeld, A. M. Bruckstein, and M. Elad, "Dictionaries for sparse representation modeling," *Proceedings of the IEEE*, vol. 98, no. 6, pp. 1045–1057, June 2010.