

# Coded Cooperative Data Exchange for a Secret Key

Thomas A. Courtade

EECS Department, University of California, Berkeley  
Email: courtade@eecs.berkeley.edu

Thomas R. Halford

Trellisware Technologies, San Diego, CA  
Email: thalford@trellisware.com

**Abstract**—We consider a cooperative data exchange problem with the goal of generating a secret key. Specifically, we investigate the number of public transmissions required for a set of clients to agree on a secret key with probability one, subject to the constraint that it remains private from an eavesdropper.

Although the problems are closely related, we prove that secret key generation with fewest linear transmissions is NP-hard, while it is known that the analogous problem in traditional cooperative data exchange can be solved in polynomial time. In doing this, we completely characterize the best-possible performance of linear coding schemes, and also prove that linear codes can be strictly suboptimal.

## I. INTRODUCTION

In an asymptotic setting, the relationship between secret key (SK) generation and communication for omniscience was revealed in the pioneering work [1]. In [1], Csiszár and Narayan left characterizing the minimum communication rate required to generate a maximum-rate SK as an open problem.

In [2], El Rouayheb *et al.* introduced a non-asymptotic, combinatorial version of Csiszár’s communication for omniscience problem, which they called (*coded*) *cooperative data exchange* (CCDE). Since its introduction, this problem has received significant attention (cf. [3]–[5] and references therein).

In [5], it was shown that the omniscience-secrecy relations in [1] translate nicely to the combinatorial CCDE setting. Building on this, we investigate the number of public transmissions required for a set of clients to agree on a SK with probability one, subject to the constraint that it remains private from an eavesdropper. In doing so, we address a combinatorial analog of Csiszár’s open problem.

Related to the present work is [6], which studies the minimum communication *rate* required to generate a SK in the two-terminal, asymptotic setting. Specifically, a multi-letter expression characterizing this minimum rate is given in terms of the *r-rounds interactive common information*. Despite the similarity in spirit, the asymptotic setting of [6] and the combinatorial setting of the present paper are considerably different in nature, and the proof techniques used are orthogonal.

Finally, the problem of weakly secure network coding is also related to our setting. The primary distinction is that we only aim to generate a SK, and do not require that the nodes communicate for omniscience. We refer the reader to [7] and references therein for details.

This work is supported by DARPA-DSO under contracts W15P7T-12-C-5013 and W911QX-13-C-0010. T. Courtade is supported in part by the NSF Center for Science of Information under grant agreement CCF-0939370.

Approved for public release; distribution is unlimited. DISCLAIMER: The views expressed are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government.

## Our Contributions

Given the close connection between CCDE and SK generation, we show two surprising results. First, we prove that finding an optimal (linear) coding scheme for SK generation is NP-hard, while it is known that the analogous CCDE problem is in P. In doing this, we completely characterize the attainable performance for linear coding schemes. Second, despite linear codes being optimal for CCDE, we demonstrate that they can be strictly suboptimal for SK generation. Several ancillary results are also proved.

This paper is organized as follows. Section II formally defines our system model and reviews relevant results on CCDE. In Section III, we state and prove our main results. Section IV delivers concluding remarks.

## II. SYSTEM MODEL AND PRELIMINARIES

We first establish basic notation. Throughout, we use calligraphic notation to denote sets. For two sets  $\mathcal{A} \subset \mathcal{B}$ , we write  $\mathcal{B} \setminus \mathcal{A}$  to denote those elements in  $\mathcal{B}$ , but not in  $\mathcal{A}$ . If  $\mathcal{A}$  is a singleton set (i.e.,  $\mathcal{A} = \{a\}$ ), then we often use the notation  $\mathcal{B} - a \triangleq \mathcal{B} \setminus \{a\}$  for convenience. We define  $\mathbb{Z}$  to be the set of integers. For positive  $m \in \mathbb{Z}$ , we use the shorthand notation  $[m] \triangleq \{1, 2, \dots, m\}$ . Finally, for random variables  $X, Y$ , we write  $I(X; Y)$  for the mutual information between  $X$  and  $Y$ .

### A. System Model

Throughout, we consider networks defined by a set of  $n$  clients  $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$ , a positive integer  $m$ , and a family of finite sets  $\{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_n\}$  (each  $\mathcal{I}_j \subseteq [m]$  and  $\cup_{j=1}^n \mathcal{I}_j = [m]$ ) in the following way. Define the random (column) vector  $\underline{X} \triangleq [X_1, X_2, \dots, X_m]^T$ , where each  $X_i$  is a discrete random variable with equiprobable distribution on a finite field  $\mathbb{F}$ , and  $(X_1, X_2, \dots, X_m)$  are mutually independent<sup>1</sup>. The random variables  $\{X_i\}_{i=1}^m$  are called *messages*, and  $\{X_i : i \in \mathcal{I}_j\}$  is the set of messages initially held by client  $c_j \in \mathcal{C}$ . In other words,  $\mathcal{I}_j$  defines the indices of messages initially held by client  $c_j$ , for  $j = 1, \dots, n$ . Throughout,  $n$  will always denote the number of clients; since the sets  $\mathcal{I}_j$  are always indexed by  $j \in [n]$ , we will use the shorthand notation  $\{\mathcal{I}_j\}$  to denote the family  $\{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_n\}$ .

We adopt the communication model which is standard in index coding and CCDE problems. That is, we consider transmission schemes consisting of a finite number of communication rounds. In each round, a single client broadcasts an element of  $\mathbb{F}$  (which can be a function of the messages initially held by that client and all previous transmissions) to all other clients over an error-free channel. It is further assumed that

<sup>1</sup>For technical reasons, we assume  $|\mathbb{F}| > 2n$ .

all clients have knowledge of the index sets  $\mathcal{I}_1, \dots, \mathcal{I}_n$ , and thus follow a protocol which is mutually agreed upon. We will elaborate on the definition of a transmission protocol in the next subsection.

### B. Transmission Protocols Defined

For a network defined by  $\{\mathcal{I}_j\}$ , a transmission protocol  $\mathbf{P}$  (or simply, a protocol  $\mathbf{P}$ ) consisting of  $t$  communication rounds is defined by  $n$  encoding functions  $\{f_1, f_2, \dots, f_n\}$ , and a  $t$ -tuple  $(i_1, i_2, \dots, i_t)$ , where  $i_k \in [n]$  indicates which client transmits during communication round  $k$ . More specifically, during communication round  $k$ , client  $c_{i_k}$  transmits

$$f_{i_k}(\{X_j : j \in \mathcal{I}_{i_k}\}, k, \{f_{i_\ell}\}_{\ell=1}^{k-1}) \in \mathbb{F}, \quad (1)$$

where we have abbreviated the transmitted symbols in rounds  $\ell \in [k-1]$  by  $\{f_{i_\ell}\}_{\ell=1}^{k-1}$ . For a given transmission protocol  $\mathbf{P}$  requiring  $t$  communication rounds, we let  $\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P}) \in \mathbb{F}^t$  be the column vector with  $k^{\text{th}}$  entry equal to  $f_{i_k}(\{X_j : j \in \mathcal{I}_{i_k}\}, k, \{f_{i_\ell}\}_{\ell=1}^{k-1})$ . Letting  $\|\cdot\|$  be the length function, we have  $\|\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})\| = t$ . Note that  $\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})$  is a random variable since it is a function of the random vector  $\underline{\mathbf{X}}$ . Generally, the transmission protocol under consideration will be clear from context. Hence, we abbreviate  $\mathbf{T}(\underline{\mathbf{X}}) \triangleq \mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})$  for convenience when there is no ambiguity.

A transmission protocol is said to be *linear* (over  $\mathbb{F}$ ) if the encoding functions  $\{f_1, f_2, \dots, f_n\}$  are of the form

$$f_{i_k}(\{X_j : j \in \mathcal{I}_{i_k}\}, k, \{f_{i_\ell}\}_{\ell=1}^{k-1}) = \sum_j \alpha_j^{(k)} X_j, \quad (2)$$

where  $\alpha_j^{(k)} \in \mathbb{F}$  can be interpreted as the encoding coefficient for message  $j$  during communication round  $k$ . In this case, we can express  $\mathbf{T}(\underline{\mathbf{X}}) = A\underline{\mathbf{X}}$ , where  $A \in \mathbb{F}^{t \times m}$  assuming the definitions  $t \triangleq \|\mathbf{T}(\underline{\mathbf{X}})\|$  and  $m \triangleq |\cup_j \mathcal{I}_j|$ . Hence, the *encoding matrix*  $A$  provides a succinct description of a linear transmission protocol. Note that the order of transmissions corresponding to a linear protocol is inconsequential.

### C. Transmission Protocols for Omniscience

Before proceeding, let  $M^*(\{\mathcal{I}_j\})$  denote the optimal value of the following integer linear program (ILP), which arises in the CCDE problem (see [5] for details).

$$\begin{aligned} \text{minimize : } & \sum_{j \in [n]} a_j & (3) \\ \text{subject to : } & \sum_{j \in \mathcal{S}} a_j \geq \left| \bigcap_{j \in \bar{\mathcal{S}}} \bar{\mathcal{I}}_j \right| \quad \forall \text{ nonempty } \mathcal{S} \subset [n] \\ & a_j \in \mathbb{Z} \quad \text{for all } j \in [n], \end{aligned}$$

where  $\bar{\mathcal{I}}_i \triangleq (\cup_j \mathcal{I}_j) \setminus \mathcal{I}_i$  and  $\bar{\mathcal{S}} \triangleq [n] \setminus \mathcal{S}$ . Though offered here without explanation, the quantity  $M^*(\{\mathcal{I}_j\})$  will play an important role in our treatment and will be needed shortly.

A transmission protocol  $\mathbf{P}$  is said to achieve *omniscience* if there exist decoding functions  $\{g_1, g_2, \dots, g_n\}$  which satisfy

$$g_j(\{X_i : i \in \mathcal{I}_j\}, \mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})) = \underline{\mathbf{X}} \quad \text{for each } j \in [n] \quad (4)$$

with probability 1.

**Theorem 1** (See [5, Theorem 2]). *If a protocol  $\mathbf{P}$  achieves omniscience, then  $\|\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})\| \geq M^*(\{\mathcal{I}_j\})$ . Conversely, there always exists a linear protocol  $\mathbf{P}_L$  that achieves omniscience and has  $\|\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P}_L)\| = M^*(\{\mathcal{I}_j\})$ .*

Theorem 1 addresses the central issue in the CCDE problem, which primarily investigates the number of transmissions required to achieve omniscience.

### D. Transmission Protocols for Secret Keys

A transmission protocol (with corresponding transmission sequence  $\mathbf{T}(\underline{\mathbf{X}})$ ) generates a secret key (SK) if there exist decoding functions  $\{k_1, k_2, \dots, k_n\}$  which satisfy the following three properties:

(i) For all  $j \in [n]$ , and with probability 1,

$$k_j(\{X_i : i \in \mathcal{I}_j\}, \mathbf{T}(\underline{\mathbf{X}})) = k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}})).$$

(ii)  $k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}}))$  is equiprobable on  $\mathbb{F}$ .

(iii)  $I(k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}})); \mathbf{T}(\underline{\mathbf{X}})) = 0$ .

In words, requirement (iii) guarantees that the public transmissions  $\mathbf{T}(\underline{\mathbf{X}})$  reveal no information about  $k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}}))$ . Requirement (i) asserts that all clients  $c_j \in \mathcal{C}$  can compute  $k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}}))$ . For these reasons,  $k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\underline{\mathbf{X}}))$  is called a *secret key*. Naturally, a secret key should be equiprobable on its domain to make guessing difficult, thus motivating requirement (ii).

It is not immediately clear whether any protocol  $\mathbf{P}$  generates a SK. However, it turns out that such protocols exist in great abundance. In particular, the existence of protocols that generate a SK depends solely on the family  $\{\mathcal{I}_j\}$ .

**Theorem 2** (See [5, Theorem 6]). *For a network defined by  $\{\mathcal{I}_j\}$ , there exists a protocol  $\mathbf{P}$  which generates a SK if and only if*

$$|\cup_j \mathcal{I}_j| \geq M^*(\{\mathcal{I}_j\}) + 1. \quad (5)$$

Despite the fact that  $M^*(\{\mathcal{I}_j\})$  corresponds to the optimal value of an ILP, it can be computed in time polynomial in  $|\cup_j \mathcal{I}_j|$  (see [3], [5]). Therefore, for any family  $\{\mathcal{I}_j\}$ , we can efficiently test whether (5) holds. Hence, the essential remaining question is: “How many transmissions are needed to generate a SK?”

To this end, let  $\mathcal{P}(\{\mathcal{I}_j\})$  denote the set of protocols for  $\{\mathcal{I}_j\}$  that generate a SK, and define

$$S(\{\mathcal{I}_j\}) \triangleq \min \left\{ \|\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})\| : \mathbf{P} \in \mathcal{P}(\{\mathcal{I}_j\}) \right\}. \quad (6)$$

That is,  $S(\{\mathcal{I}_j\})$  is the minimum number of transmissions needed to generate a SK. Similarly, let  $\mathcal{P}_L(\{\mathcal{I}_j\})$  denote the set of linear protocols for  $\{\mathcal{I}_j\}$  that generate a SK, and define

$$S_L(\{\mathcal{I}_j\}) \triangleq \min \left\{ \|\mathbf{T}(\underline{\mathbf{X}}, \mathbf{P})\| : \mathbf{P} \in \mathcal{P}_L(\{\mathcal{I}_j\}) \right\}. \quad (7)$$

In words,  $S_L(\{\mathcal{I}_j\})$  is the minimum number of transmissions required to generate a SK when we restrict our attention to linear protocols. If  $\{\mathcal{I}_j\}$  does not satisfy (5), then we set  $S(\{\mathcal{I}_j\}) = S_L(\{\mathcal{I}_j\}) = \infty$ .

**Remark 1.** *We will often write “ $\{\mathcal{I}_j\}$  generates a SK” instead of the more accurate, but cumbersome, “For the network defined by  $\{\mathcal{I}_j\}$ , there exists a protocol  $\mathbf{P}$  which generates a SK” whenever (5) holds.*

## III. MAIN RESULTS

In this section, we investigate the number of transmissions required to generate a SK. In particular, we completely characterize  $S_L(\{\mathcal{I}_j\})$ , and make progress toward characterizing  $S(\{\mathcal{I}_j\})$ .

As demonstrated in the previous section, the CCDE and SK-generation problems are closely connected through the quantity  $M^*(\{\mathcal{I}_j\})$ . Since Theorem 1 and the tractability of ILP (3) essentially resolve the CCDE problem, it is natural to conjecture that a similar result should hold for  $S(\{\mathcal{I}_j\})$  and  $S_L(\{\mathcal{I}_j\})$ . Unfortunately, there is a fundamental difference between the problems, which is revealed by the following two negative results<sup>2</sup>.

**Theorem 3.** *For some families  $\{\mathcal{I}_j\}$ ,  $S_L(\{\mathcal{I}_j\}) > S(\{\mathcal{I}_j\})$ .*

**Theorem 4.** *Computing  $S_L(\{\mathcal{I}_j\})$  is NP-hard.*

For the CCDE problem, Theorem 1 asserts that linear protocols achieve optimal performance. Furthermore, the number of transmissions required by linear protocols is easily computed. For the problem of SK generation, the opposite is true. That is, linear protocols can be suboptimal, and the number of transmissions required by linear protocols is generally difficult to compute. This situation is parallel to that of multicast network coding and index coding. The two problems are closely related (cf. [8]), but exhibit the same dichotomy.

Despite the negative results offered by Theorems 3 and 4, we can characterize several properties of  $S_L(\{\mathcal{I}_j\})$ ,  $S(\{\mathcal{I}_j\})$ , and  $M^*(\{\mathcal{I}_j\})$ . Some of these are demonstrated in the following results, which are needed as we progress toward proving Theorem 4. A complete characterization of  $S_L(\{\mathcal{I}_j\})$  will be given in Theorem 5.

**Lemma 1.** *If  $\{\mathcal{I}_j\}$  generates a SK, then*

$$S(\{\mathcal{I}_j\}) \leq S_L(\{\mathcal{I}_j\}) \leq M^*(\{\mathcal{I}_j\}). \quad (8)$$

*Proof.* By definition,  $S(\{\mathcal{I}_j\}) \leq S_L(\{\mathcal{I}_j\})$  since  $\mathcal{P}_L(\{\mathcal{I}_j\}) \subseteq \mathcal{P}(\{\mathcal{I}_j\})$ . The second inequality follows from the proof of [5, Theorem 6], in which a linear transmission protocol  $\mathbf{P}$  is constructed that generates a SK with  $\|\mathbf{T}(\mathbf{X}), \mathbf{P}\| = M^*(\{\mathcal{I}_j\})$  communication rounds.  $\square$

We say that  $\{\mathcal{J}_j\}$  is a *subfamily* of  $\{\mathcal{I}_j\}$  if there is a set  $\mathcal{S} \subset \cup_j \mathcal{I}_j$  such that  $\mathcal{J}_j = \mathcal{I}_j \setminus \mathcal{S}$  for all  $j \in [n]$ .

**Lemma 2.** *If  $\{\mathcal{J}_j\}$  is a subfamily of  $\{\mathcal{I}_j\}$ , then*

$$M^*(\{\mathcal{J}_j\}) \leq M^*(\{\mathcal{I}_j\}), \quad (9)$$

$$S(\{\mathcal{J}_j\}) \geq S(\{\mathcal{I}_j\}), \text{ and} \quad (10)$$

$$S_L(\{\mathcal{J}_j\}) \geq S_L(\{\mathcal{I}_j\}). \quad (11)$$

*Proof.* By De Morgan's law, it is easy to verify that

$$\left| \bigcap_{j \in \mathcal{S}} \bar{\mathcal{I}}_j \right| \geq \left| \bigcap_{j \in \mathcal{S}} \bar{\mathcal{J}}_j \right| \text{ for all nonempty } \mathcal{S} \subset [n], \quad (12)$$

where  $\bar{\mathcal{I}}_i \triangleq (\cup_j \mathcal{I}_j) \setminus \mathcal{I}_i$  and  $\bar{\mathcal{J}}_i \triangleq (\cup_j \mathcal{J}_j) \setminus \mathcal{J}_i$ . Therefore, the constraints in ILP (3) are relaxed, and  $M^*(\{\mathcal{J}_j\}) \leq M^*(\{\mathcal{I}_j\})$  by definition.

To show (10), observe that any transmission protocol which generates a SK for the subfamily  $\{\mathcal{J}_j\}$  also generates a SK for the family  $\{\mathcal{I}_j\}$  by ignoring the set of messages  $\{X_i : i \notin \cup_j \mathcal{J}_j\}$ . Hence, it follows that  $S(\{\mathcal{J}_j\}) \geq S(\{\mathcal{I}_j\})$ . If  $\{\mathcal{J}_j\}$  can not generate a SK, the inequality trivially holds. This argument also proves (11).  $\square$

Lemma 2 demonstrates monotonicity, but offers no insight into whether inequalities (9)-(11) are tight. The following lemma identifies settings under which (11) holds with equality, and will prove useful later on.

**Lemma 3.** *If  $S_L(\{\mathcal{I}_j\}) < M^*(\{\mathcal{I}_j\})$ , then there exists some  $\ell \in \cup_j \mathcal{I}_j$  for which  $S_L(\{\mathcal{I}_j - \ell\}) = S_L(\{\mathcal{I}_j\})$ .*

*Proof.* Define  $m \triangleq |\cup_j \mathcal{I}_j|$ . By definition, there is a linear transmission protocol  $\mathbf{P}_L$  which generates a SK in  $S_L(\{\mathcal{I}_j\})$  communication rounds. Let  $\mathbf{T}(\mathbf{X}, \mathbf{P}_L) = \mathbf{A}\mathbf{X}$  be the sequence of transmissions made by  $\mathbf{P}_L$ , and let  $\{k_1, \dots, k_n\}$  be valid decoding functions.

Since  $\|\mathbf{T}(\mathbf{X}, \mathbf{P}_L)\| = S_L(\{\mathcal{I}_j\}) < M^*(\{\mathcal{I}_j\})$ , Theorem 1 asserts that the protocol  $\mathbf{P}_L$  can not achieve omniscience. Therefore, by a possible permutation of clients, we can assume without loss of generality that there is no function  $g_1$  for which

$$g_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{A}\mathbf{X}) = \mathbf{X} \text{ with probability 1.} \quad (13)$$

As a consequence, there must exist a nonzero vector  $\mathbf{v}$  such that  $A\mathbf{v} = 0$ , and  $v_i = 0$  for all  $i \in \mathcal{I}_1$ . Indeed, if there is no such  $\mathbf{v}$ , then client  $c_1$  can solve a full-rank system of equations to recover  $\mathbf{X}$ , yielding a contradiction. Since  $\mathbf{v}$  is not identically zero, there is some  $\ell$  for which  $v_\ell \neq 0$ .

Define  $\hat{X}_\ell \equiv 0$ , and  $\hat{X}_i \triangleq X_i$  for  $i \in \cup_j (\mathcal{I}_j - \ell)$ . Also, define vectors  $\hat{\mathbf{X}} \triangleq [\hat{X}_1, \hat{X}_2, \dots, \hat{X}_m]^T$  and  $\mathbf{X}' \triangleq \hat{\mathbf{X}} + X_\ell \cdot \mathbf{v}$ . First, we note that

$$k_j(\{\hat{X}_i : i \in \mathcal{I}_j\}, \mathbf{A}\hat{\mathbf{X}}) = k_1(\{\hat{X}_i : i \in \mathcal{I}_1\}, \mathbf{A}\hat{\mathbf{X}})$$

for all  $j \in [n]$  since

$$k_j(\{X_i : i \in \mathcal{I}_j\}, \mathbf{A}\mathbf{X}) = k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{A}\mathbf{X})$$

with probability 1 by definition.

Next, observe that:

$$\begin{aligned} I\left(k_1(\{\hat{X}_i : i \in \mathcal{I}_1\}, \mathbf{A}\hat{\mathbf{X}}); \mathbf{A}\hat{\mathbf{X}}\right) \\ = I\left(k_1(\{\hat{X}_i + X_\ell \cdot v_i : i \in \mathcal{I}_1\}, \mathbf{A}\mathbf{X}'); \mathbf{A}\mathbf{X}'\right) \end{aligned} \quad (14)$$

$$= I(k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{A}\mathbf{X}); \mathbf{A}\mathbf{X}) = 0 \quad (15)$$

In the above,

- (14) follows since  $\mathbf{A}\mathbf{X}' = A(\hat{\mathbf{X}} + X_\ell \cdot \mathbf{v}) = \mathbf{A}\hat{\mathbf{X}}$ , and  $v_i = 0$  for all  $i \in \mathcal{I}_1$ .
- (15) follows from the (crucial) observation that  $\mathbf{X}'$  and  $\mathbf{X}$  are equal in distribution, and by definition of  $A$  and  $k_1$ .

<sup>2</sup>The proofs of Theorems 3 and 4 are deferred until the end of this section.

Finally, by similar reasoning, we note that the random variable  $k_1(\{\hat{X}_i : i \in \mathcal{I}_1\}, A\hat{\mathbf{X}})$  is equiprobable on  $\mathbb{F}$  since

$$\begin{aligned} & k_1(\{\hat{X}_i : i \in \mathcal{I}_1\}, A\hat{\mathbf{X}}) \\ &= k_1(\{\hat{X}_i + X_\ell \cdot v_i : i \in \mathcal{I}_1\}, A\mathbf{X}') \end{aligned} \quad (16)$$

$$\stackrel{d}{=} k_1(\{X_i : i \in \mathcal{I}_1\}, A\mathbf{X}), \quad (17)$$

and  $k_1(\{X_i : i \in \mathcal{I}_1\}, A\mathbf{X})$  is equiprobable on  $\mathbb{F}$  by definition. In (17),  $\stackrel{d}{=}$  indicates equality in distribution.

Therefore, we can conclude that a SK can be generated by the subfamily  $\{\mathcal{I}_j - \ell\}_j$  by applying the protocol  $\mathbf{P}_L$  and fixing  $X_\ell \equiv 0$ . This proves that  $S_L(\{\mathcal{I}_j - \ell\}) \leq S_L(\{\mathcal{I}_j\})$ . By Lemma 2, the reverse inequality also holds.  $\square$

In order to proceed, we will need to introduce *critical families*. To this end, let  $\tau \geq 1$  be an integer. A family  $\{\mathcal{I}_j\}$  is  $\tau$ -critical if the following hold:

- (i)  $|\cup_j \mathcal{I}_j| - M^*(\{\mathcal{I}_j\}) = \tau$ , and
- (ii)  $M^*(\{\mathcal{I}_j - i\}) = M^*(\{\mathcal{I}_j\})$  for all  $i \in \cup_j \mathcal{I}_j$ .

It is interesting to note that  $\tau$ -criticality of  $\{\mathcal{I}_j\}$  can be efficiently tested since  $M^*(\{\mathcal{I}_j\})$  is computable in polynomial time. 1-critical families have a threshold property: families  $\{\mathcal{I}_j\}$  that are 1-critical generate a secret key, and no proper subfamilies of  $\{\mathcal{I}_j\}$  generate SKs. This is a consequence of Theorem 2 and the definition of 1-criticality.

A *minimum  $\tau$ -critical subfamily*  $\{\mathcal{J}_j^*\}$  of  $\{\mathcal{I}_j\}$  satisfies

$$|\cup_j \mathcal{J}_j^*| \leq |\cup_j \mathcal{I}_j| \quad (18)$$

for all other  $\tau$ -critical subfamilies  $\{\mathcal{J}_j\}$  of  $\{\mathcal{I}_j\}$ . Note that if  $\{\mathcal{I}_j\}$  is 1-critical, then  $\{\mathcal{I}_j\}$  is its own unique minimum 1-critical subfamily.

The following Theorem demonstrates that minimum 1-critical subfamilies completely characterize  $S_L(\{\mathcal{I}_j\})$ .

**Theorem 5.** *If  $\{\mathcal{I}_j\}$  generates a SK, then*

$$S_L(\{\mathcal{I}_j\}) = M^*(\{\mathcal{J}_j^*\}) = |\cup_j \mathcal{J}_j^*| - 1, \quad (19)$$

where  $\{\mathcal{J}_j^*\}$  is a minimum 1-critical subfamily of  $\{\mathcal{I}_j\}$ .

*Proof.* By inductively applying Lemma 3, we can find a subfamily  $\{\mathcal{T}_j\}$  of  $\{\mathcal{I}_j\}$  for which

$$S_L(\{\mathcal{I}_j\}) = M^*(\{\mathcal{T}_j\}). \quad (20)$$

Let  $\{\mathcal{J}_j\}$  be any 1-critical subfamily of  $\{\mathcal{T}_j\}$ . We have the following chain of inequalities

$$S_L(\{\mathcal{I}_j\}) \leq S_L(\{\mathcal{J}_j^*\}) \leq M^*(\{\mathcal{J}_j^*\}) \quad (21)$$

$$\leq M^*(\{\mathcal{J}_j\}) \quad (22)$$

$$\leq M^*(\{\mathcal{T}_j\}) \quad (23)$$

$$= S_L(\{\mathcal{I}_j\}). \quad (24)$$

The above steps can be justified as follows:

- (21) follows from Lemmas 1 and 2.
- By definition of  $\tau$ -criticality, (18) is equivalent to  $M^*(\{\mathcal{J}_j^*\}) \leq M^*(\{\mathcal{J}_j\})$ . Thus, (22) follows since  $\{\mathcal{J}_j^*\}$

is a minimum 1-critical subfamily of  $\{\mathcal{I}_j\}$ , and  $\{\mathcal{J}_j\}$  is a 1-critical subfamily of  $\{\mathcal{I}_j\}$ .

- (23) follows from Lemma 2.
- (24) is a repeat of (20).

This proves that  $S_L(\{\mathcal{I}_j\}) = M^*(\{\mathcal{J}_j^*\})$ . Recalling the definition of 1-criticality completes the proof.  $\square$

Theorem 5 implies that  $S_L(\{\mathcal{I}_j\})$  is easily computed if we can identify a minimum 1-critical subfamily of  $\{\mathcal{I}_j\}$ . By Theorem 4, we know this must be NP-hard. In order to prove this, we require the following lemma which lends a hypergraph interpretation to 1-criticality. For a hypergraph  $H = (\mathcal{V}, \mathcal{E})$ , recall that an edge set  $\mathcal{E}' \subseteq \mathcal{E}$  is a *minimal connected dominating edge set* if the subhypergraph  $H' = (\mathcal{V}, \mathcal{E}')$  is connected, and the removal of any edge from  $\mathcal{E}'$  disconnects  $H'$ .

**Lemma 4.** *Consider a hypergraph  $H = (\mathcal{V}, \mathcal{E})$  with vertex set  $\mathcal{V} = \mathcal{C}$ , and edge set  $\mathcal{E} = \cup_j \mathcal{I}_j$ . Define  $H$  as follows: a vertex  $c_j \in \mathcal{V}$  is contained in the edge  $e \in \mathcal{E}$  if and only if  $e \in \mathcal{I}_j$ .  $H$  is connected if and only if*

$$M^*(\{\mathcal{I}_j\}) < |\cup_j \mathcal{I}_j|. \quad (25)$$

*In particular,  $\{\mathcal{I}_j\}$  is 1-critical if and only if  $\mathcal{E}$  is a minimal connected dominating edge set.*

*Proof.* First, suppose  $H$  is not connected. By definition, there must exist a nontrivial partition  $\mathcal{V} = (\mathcal{S}, \bar{\mathcal{S}})$  such that there is no edge  $e \in \mathcal{E}$  which contains vertices from both  $\mathcal{S}$  and  $\bar{\mathcal{S}}$ . Stated another way,  $(\cup_{j \in \mathcal{S}} \mathcal{I}_j) \cap (\cup_{j \in \bar{\mathcal{S}}} \mathcal{I}_j) = \emptyset$ . Hence, ILP (3) includes the two constraints

$$\sum_{j \in \mathcal{S}} a_j \geq \left| \bigcap_{j \in \bar{\mathcal{S}}} \bar{\mathcal{I}}_j \right| = \left| \bigcup_{j \in \mathcal{S}} \mathcal{I}_j \right| \quad (26)$$

$$\sum_{j \in \bar{\mathcal{S}}} a_j \geq \left| \bigcap_{j \in \mathcal{S}} \bar{\mathcal{I}}_j \right| = \left| \bigcup_{j \in \bar{\mathcal{S}}} \mathcal{I}_j \right|, \quad (27)$$

the sum of which imply  $M^*(\{\mathcal{I}_j\}) \geq |\cup_j \mathcal{I}_j|$ . By taking the contrapositive, we have proven

$$M^*(\{\mathcal{I}_j\}) < |\cup_j \mathcal{I}_j| \implies H \text{ is connected.} \quad (28)$$

Next, suppose  $H$  is connected, and assume without loss of generality that  $\mathcal{E} = \cup_j \mathcal{I}_j \triangleq \{1, 2, \dots, m\}$ . Since  $H$  is connected, there is a transmission protocol for which the entries of  $\mathbf{T}(\mathbf{X})$  are precisely  $\{X_1 + X_j\}_{j=2}^m$ . Indeed, by connectivity of  $H$ , there must be some client  $c$  initially holding  $X_1$  and some  $X_e$  (say,  $X_2$  without loss of generality), and can therefore transmit  $X_1 + X_2$  during the first communication round. By induction, assume that  $\{X_1 + X_j\}_{j=2}^{m-1}$  are transmitted during the first  $m-2$  communication rounds (permuting indices of the  $X_i$ 's if necessary). Again, by connectivity of  $H$ , there must be a client  $c'$  which initially holds  $X_m$  and  $X_k$ , where  $k < m$ . Hence, in communication round  $m-1$ , client  $c'$  can transmit  $(X_1 + X_k) - (X_k - X_m) = X_1 + X_m$ . Noting that

$$(X_1, X_1 + X_2, \dots, X_1 + X_m) \stackrel{d}{=} (X_1, X_2, \dots, X_m),$$

we have  $I(X_1; \mathbf{T}(\mathbf{X})) = 0$ . If client  $c \in e \in \mathcal{E}$ , then it can recover  $X_1$  from the transmission  $X_1 + X_e$  by simply subtracting  $X_e$ . Since  $H$  is connected, each  $c \in \mathcal{V}$  belongs to some edge in  $\mathcal{E}$ , and therefore all clients can recover  $X_1$  losslessly. Since  $X_1$  is equiprobable on  $\mathbb{F}$  by definition, we can conclude that  $\{\mathcal{I}_j\}$  generates a SK. Theorem 2 asserts that we must have  $M^*(\{\mathcal{I}_j\}) < |\cup_j \mathcal{I}_j|$ , and we have proven

$$M^*(\{\mathcal{I}_j\}) < |\cup_j \mathcal{I}_j| \iff H \text{ is connected.} \quad (29)$$

We now prove the second claim. To this end, suppose  $\{\mathcal{I}_j\}$  is 1-critical. Then  $M^*(\{\mathcal{I}_j\}) = |\cup_j \mathcal{I}_j| - 1$ , which implies  $H$  is connected (and thus  $\mathcal{E}$  is dominating) by (29). Consider the subhypergraph  $H' = (\mathcal{V}, \mathcal{E} \setminus \{e\})$ , which corresponds to the subfamily  $\{\mathcal{I}_j - e\}$  of  $\{\mathcal{I}_j\}$ . Since  $\{\mathcal{I}_j\}$  is 1-critical, we must have  $M^*(\{\mathcal{I}_j - e\}) = M^*(\{\mathcal{I}_j\}) = |\cup_j \mathcal{I}_j| - 1 = |\cup_j (\mathcal{I}_j - e)|$ . By (29),  $H'$  must be disconnected, and therefore  $\mathcal{E}$  is a minimal connected dominating edge set.

On the other hand, suppose  $\mathcal{E}$  is a minimal connected dominating edge set. Since  $H$  is connected, (29) implies

$$M^*(\{\mathcal{I}_j\}) \leq |\cup_j \mathcal{I}_j| - 1. \quad (30)$$

Since  $\mathcal{E}$  is minimal, for any  $e \in \mathcal{E}$ ,  $H' = (\mathcal{V}, \mathcal{E} \setminus \{e\})$  is disconnected, and (29) implies

$$M^*(\{\mathcal{I}_j - e\}) \geq |\cup_j (\mathcal{I}_j - e)| = |\cup_j \mathcal{I}_j| - 1. \quad (31)$$

Applying Lemma 2, we must have  $M^*(\{\mathcal{I}_j\}) = M^*(\{\mathcal{I}_j - e\})$ , and  $|\cup_j \mathcal{I}_j| - M^*(\{\mathcal{I}_j\}) = 1$ , which implies  $\{\mathcal{I}_j\}$  is 1-critical.  $\square$

We are finally in a position to prove Theorem 4.

*Proof of Theorem 4.* Consider a hypergraph  $H = (\mathcal{V}, \mathcal{E})$  with vertex set  $\mathcal{V} = \mathcal{C}$ , and edge set  $\mathcal{E} = \cup_j \mathcal{I}_j$ . Define  $H$  as follows: a vertex  $c_j \in \mathcal{V}$  is contained in the edge  $e \in \mathcal{E}$  if and only if  $e \in \mathcal{I}_j$ .

We can assume  $\{\mathcal{I}_j\}$  generates a SK. By Theorem 5 and Lemma 4, computing  $S_L(\{\mathcal{I}_j\})$  is equivalent to computing the number of edges in a minimum connected dominating edge set (i.e., a minimal connected dominating edge set with fewest possible edges). It is easy to see that the NP-complete SET COVER DECISION PROBLEM is a special case<sup>3</sup>.  $\square$

Theorem 3 is proved by the following example: Let  $n = 7$ , and consider the family  $\{\mathcal{I}_j\}$  defined by  $\mathcal{I}_1 = \{1, 2, 3, 4\}$ , and  $\mathcal{I}_2, \dots, \mathcal{I}_7$  are all  $\binom{4}{2}$  distinct 2-element subsets of  $\{1, 2, 3, 4\}$ . By direct computation, we find that  $\{\mathcal{I}_j - \{1\}\}$  is a minimum 1-critical subfamily, and hence  $S_L(\{\mathcal{I}_j\}) = 2$  by Theorem 5. Suppose  $\mathbb{F} = \{0, 1, \alpha, \beta\}^2 = \text{GF}(4) \times \text{GF}(4)$ . Thus, we can express  $X_j = (X_j^{(1)}, X_j^{(2)})$  for each  $j = 1, \dots, 4$ , where  $X_j^{(1)}, X_j^{(2)}$  are mutually independent, each equiprobable on  $\text{GF}(4)$ . It is readily verified that the single transmission

$$\left( X_1^{(1)} + \alpha X_2^{(1)} + X_3^{(1)}, X_1^{(1)} + \beta X_2^{(1)} + X_4^{(1)} \right) \in \mathbb{F} \quad (32)$$

<sup>3</sup>Consider any subsets  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_k$  whose union covers  $\mathcal{U}$ . For  $u' \notin \mathcal{U}$ , define  $\mathcal{U}' = \mathcal{U} \cup \{u'\}$ , and  $\mathcal{A}'_j = \mathcal{A}_j \cup \{u'\}$  for  $j \in [k]$ . Clearly,  $\{\mathcal{A}_{j_i}\}_{i=1}^m$  is a minimum cover of  $\mathcal{U}$  if and only if  $\{\mathcal{A}'_{j_i}\}_{i=1}^m$  is a minimum connected cover of  $\mathcal{U}'$ .

by client  $c_1$  permits reconstruction of the SK  $k_1(\{X_i : i \in \mathcal{I}_1\}, \mathbf{T}(\mathbf{X})) \triangleq (X_3^{(1)}, X_4^{(1)}) \in \mathbb{F}$  at all clients. Hence, we can conclude  $1 = S(\{\mathcal{I}_j\}) < S_L(\{\mathcal{I}_j\}) = M^*(\{\mathcal{I}_j\}) = 2$ .

The above construction is a *vector-linear* transmission protocol, and cannot be realized by a protocol which is linear over  $\mathbb{F}$ . We remark that this construction can be generalized to increase the gap between  $S(\{\mathcal{I}_j\})$  and  $S_L(\{\mathcal{I}_j\})$ , but details are omitted due to space constraints.

#### IV. CONCLUDING REMARKS

Finding a minimum 1-critical subfamily  $\{\mathcal{J}_j^*\}$  of  $\{\mathcal{I}_j\}$  is NP-hard, but we can easily find a 1-critical subfamily  $\{\mathcal{J}_j\}$  satisfying  $|\cup_j \mathcal{J}_j| \leq (2 + \ln(n-1)) |\cup_j \mathcal{J}_j^*|$  by using the greedy algorithm in [9]. Given the connection to the SET COVER PROBLEM pointed out in the proof of Theorem 4, this approximation is best-possible unless NP has  $O(n^{O(\ln n)})$  time algorithms. With  $\{\mathcal{J}_j\}$  in hand, we can implement the linear protocol described in the proof of Lemma 4 requiring  $S_L(\{\mathcal{J}_j\})$  transmissions, which is guaranteed to be roughly within a  $\ln n$  factor of optimal.

Although we have focused exclusively on protocols that generate a single SK, it is natural to consider protocols that generate  $\tau$  independent secret keys. To this end, let  $S_L^{(\tau)}(\{\mathcal{I}_j\})$  denote the minimum number of transmissions required by a linear protocol to generate  $\tau$  independent secret keys. A slight modification of our argument yields:

**Theorem 6.** *Let  $\tau \geq 1$  be an integer. If there is a protocol  $\mathbf{P}$  for  $\{\mathcal{I}_j\}$  which generates  $\tau$  independent secret keys, then*

$$S_L^{(\tau)}(\{\mathcal{I}_j\}) = M^*(\{\mathcal{J}_j^*\}) = |\cup_j \mathcal{J}_j^*| - \tau, \quad (33)$$

where  $\{\mathcal{J}_j^*\}$  is a minimum  $\tau$ -critical subfamily of  $\{\mathcal{I}_j\}$ .

Presumably, it is NP-hard to find a minimum  $\tau$ -critical subfamily of  $\{\mathcal{I}_j\}$  for  $\tau \geq 2$ , which would lead to a natural generalization of Theorem 4. This requires a modification of Lemma 4, and appears challenging. Indeed, the structure of 1-critical families is succinctly captured by minimal connected dominating edge sets, but the structural class of  $\tau$ -critical families is much richer for  $\tau \geq 2$ . We leave this general hardness problem for future work.

#### REFERENCES

- [1] I. Csiszár and P. Narayan, "Secrecy capacities for multiple terminals," *IEEE Trans. on Inf. Theory*, vol. 50, no. 12, pp. 3047–3061, Dec. 2004.
- [2] S. El Rouayheb, A. Sprintson, and P. Sadeghi, "On coding for cooperative data exchange," in *IEEE Information Theory Workshop (ITW)*, Jan. 2010.
- [3] N. Milosavljevic, S. Pawar, S. El Rouayheb, M. Gastpar, and K. Ramchandran, "Optimal deterministic polynomial-time data exchange for omniscience," *arXiv preprint:1108.6046 [cs.IT]*, Aug. 2011.
- [4] M. Gonen and M. Langberg, "Coded cooperative data exchange problem for general topologies," in *Proc. of IEEE International Symposium on Information Theory (ISIT)*, July 2012, pp. 2606–2610.
- [5] T. Courtade and R. Wesel, "Coded cooperative data exchange in multihop networks," *IEEE Trans. Inf. Theory*, vol. 60, no. 2, pp. 1136–1158, 2014.
- [6] H. Tyagi, "Common information and secret key capacity," *IEEE Trans. on Inf. Theory*, vol. 59, no. 9, pp. 5627–5640, 2013.
- [7] M. Yan and A. Sprintson, "Algorithms for weakly secure data exchange," in *Network Coding (NetCod), 2013 Intl. Symp. on.* IEEE, 2013, pp. 1–6.
- [8] M. Effros, S. El Rouayheb, and M. Langberg, "An equivalence between network coding and index coding," *arXiv:1211.6660*, Nov. 2012.
- [9] W. Ren and Q. Zhao, "A note on 'Algorithms for connected set cover problem and fault-tolerant connected set cover problem'," *Theoretical Computer Science*, vol. 412, no. 45, pp. 6451–6454, 2011.