

# Three-Layer Composite Coding for Index Coding

Yucheng Liu<sup>†</sup>, Parastoo Sadeghi<sup>†</sup>, and Young-Han Kim<sup>\*</sup>

<sup>†</sup>Research School of Engineering, Australian National University, Canberra, ACT, 2600, Australia

<sup>\*</sup>Department of Electrical and Computer Engineering, University of California, San Diego, CA 92093, USA

Emails: {yucheng.liu, parastoo.sadeghi}@anu.edu.au, yhk@ucsd.edu

**Abstract**—We extend the composite coding (CC) scheme for the index coding problem from two layers to more layers of random binning. We explicitly introduce the three-layer composite coding (TLCC) scheme and provide the achievable rate region and the error analysis for it. We present a concrete non-trivial example with  $n = 7$  messages where the TLCC strictly outperforms the CC scheme. We also present a number of simplification methods for the TLCC scheme towards better understanding of the scheme, as well as significantly reducing its computational complexity. We further prove that even a simplified version of the TLCC, which can be possibly weaker than the TLCC, still subsumes the CC scheme.

## I. INTRODUCTION

Introduced by Birk and Kol in [1], the index coding problem investigates the optimal broadcast rate of  $n$  messages from a server to multiple receivers via a noiseless finite-capacity channel. Each receiver wants to decode a unique message and has prior knowledge of some other messages. A fundamental goal in index coding is characterizing its full capacity region, which is still an open problem. Various index coding schemes have been developed, which establish inner bounds on the capacity region (or lower bounds on the symmetric capacity) of the index coding problem. The author in [2] showed that linear index coding is optimal for all problems with  $n \leq 5$  messages. Nevertheless, none of the existing schemes is generally optimal. See [2], [3] and the references therein for structured index coding schemes.

The classic random coding method by Shannon still plays a central role in many network information theory problems. *Composite coding* or CC, first introduced in [4] and later enhanced in [5], is an index coding scheme built upon random flat coding. Flat coding is a single-layer random coding scheme where for each message tuple a codeword is generated via random binning. Composite coding is a two-layer random coding scheme, which first generates an index of certain rate for each message subset via random binning and then maps all the indices to a codeword, again via random binning. Composite coding strictly outperforms flat coding and gives tight inner bounds on the capacity region for all index coding problems with  $n \leq 5$  messages, as well as many larger problems with more messages. However, composite coding is in general suboptimal.

A method is proposed in [6] to address this issue by combining structured coding and random coding. However, the structured part of the encoder seems problem-dependent and may involve custom design. A question that arises is whether one can design a general-purpose purely random coding scheme by adding more layers into the CC scheme, such that it strictly outperforms the CC scheme.

In this paper we show that adding more layers to composite coding is possible and indeed beneficial. For ease of exposition and computations, we focus on a *three-layer composite coding* or the TLCC scheme. Section IV describes the codebook generation and the encoding and decoding operations for the TLCC scheme. We characterize the achievable rate region for the TLCC scheme and present the corresponding error analysis. Next in Section V, we present a series of simplification techniques for reducing the computational complexity of the TLCC. We prove that even with some simplifications the TLCC gives a rate region that is no smaller than that given by the CC. A concrete example is provided to show that the TLCC scheme can be strictly superior to the CC scheme.

## II. SYSTEM MODEL

Consider the index coding problem with  $n$  messages,  $x_i \in \{0, 1\}^{t_i}, i \in [n] \doteq \{1, 2, \dots, n\}$ . For brevity, when we say message  $i$ , we mean message  $x_i$ . Let  $X_i$  be the random variable corresponding to  $x_i$ . We assume that  $X_1, \dots, X_n$  are uniformly distributed and independent of each other. For any  $K \subseteq [n]$ , set  $K^c \doteq [n] \setminus K$ , and we use the shorthand notation  $\mathbf{x}_K$  to denote the collection of messages whose index is in  $K$ , and  $2^K$  to denote the collection of all non-empty subsets of  $K$ . In particular, set  $N \doteq 2^{[n]}$ . By convention  $\mathbf{x}_\emptyset = 2^\emptyset = \emptyset$ .

There is a single server that contains all messages  $\mathbf{x}_{[n]}$  and is connected to all receivers via a noiseless broadcast link of normalized capacity  $C = 1$ . Let  $y$  be the output of the server, which is a function of  $\mathbf{x}_{[n]}$ . There are  $n$  receivers, where receiver  $i \in [n]$  wishes to obtain  $x_i$  and knows  $\mathbf{x}_{A_i}$  as side information for some  $A_i \subseteq [n] \setminus \{i\}$ .

We define a  $(\mathbf{t}, r) = ((t_i, i \in [n]), r)$  *index code* by

- An encoder  $\phi : \prod_{i \in [n]} \{0, 1\}^{t_i} \rightarrow \{0, 1\}^r$ , which maps the messages  $\mathbf{x}_{[n]}$  to an  $r$ -bit sequence  $y$ , and
- $n$  decoders, one for each receiver  $i \in [n]$ , such that  $\psi_i : \{0, 1\}^r \times \prod_{k \in A_i} \{0, 1\}^{t_k} \rightarrow \{0, 1\}^{t_i}$  maps the received sequence  $y$  and the side information  $\mathbf{x}_{A_i}$  to  $\hat{x}_i$ .

We say that a rate tuple  $\mathbf{R} = (R_i, i \in [n])$  is achievable for every  $\epsilon > 0$ , there exist a  $(\mathbf{t}, r)$  index code such that

$$R_i \leq \frac{t_i}{r}, \quad i \in [n], \quad (1)$$

and the probability of error satisfies

$$\mathbb{P}\{(\hat{X}_1, \dots, \hat{X}_n) \neq (X_1, \dots, X_n)\} \leq \epsilon. \quad (2)$$

The capacity region  $\mathcal{C}$  of this index coding problem is the closure of the set of all achievable rate tuples  $\mathbf{R}$ . The symmetric capacity is defined as

$$C_{\text{sym}} = \max\{R_{\text{sym}} : (R_{\text{sym}}, \dots, R_{\text{sym}}) \in \mathcal{C}\}. \quad (3)$$

We will compactly represent an index coding instance by a sequence  $(i|j \in A_i)$ ,  $i \in [n]$ . For example, for  $A_1 = \emptyset$ ,  $A_2 = \{3\}$ , and  $A_3 = \{2\}$ , we write  $(1|-, (2|3), (3|2)$ .

### III. REVIEW OF FLAT CODING AND COMPOSITE CODING

Let  $r \in \mathbb{N}$  and  $t_i = \lceil rR_i \rceil$ ,  $i \in [n]$ , where  $R_i$  is the rate of message  $i$ . Flat coding is a single-layer random coding scheme for index coding. For each realization of messages  $\mathbf{x}_{[n]}$  a codeword  $y(\mathbf{x}_{[n]})$  is drawn uniformly at random from  $[2^r]$ . The codebook is revealed to all parties. To communicate messages  $\mathbf{x}_{[n]}$ , the server transmits  $y(\mathbf{x}_{[n]})$ . Since the encoding is flat, receiver  $i$  must decode all messages it does not know and then discard the unwanted or interfering messages  $\mathbf{x}_{B_i}$  to obtain  $x_i$ , where  $B_i = (A_i \cup \{i\})^c$ . The achievable rate region of flat coding is

$$R_i + \sum_{j \in B_i} R_j < 1, \quad i \in [n]. \quad (4)$$

For the problem  $(1|4), (2|3, 4), (3|1, 2), (4|2, 3)$ , flat coding for the original problem or time-sharing of flat coding over subproblems gives a suboptimal rate region as  $R_1 + R_2 + R_3 < 1$ ,  $R_1 + R_4 < 1$ , and  $R_3 + R_4 < 1$ . However, using a two-layer CC scheme, described below,  $R_1 + R_2 < 1$ ,  $R_1 + R_3 < 1$ ,  $R_1 + R_4 < 1$ , and  $R_3 + R_4 < 1$  is achievable.

Let  $s_K = \lceil rS_K \rceil$ ,  $K \subseteq [n]$ , where  $S_K$  is the rate of composite index  $K$  (to be defined shortly). We review the original version of the CC scheme [4], in which each receiver uses a fixed decoding configuration.

**Codebook generation.** *Step 1.* For each  $K \subseteq [n]$ , generate a composite index  $w_K(\mathbf{x}_K)$  uniformly at random from  $[2^{s_K}]$ . *Step 2.* Generate the codeword  $y(w_K, K \subseteq [n])$  uniformly at random from  $[2^r]$ . Reveal the codebook to all parties.

**Encoding.** To communicate messages  $\mathbf{x}_{[n]}$ , the server transmits  $y(w_K(\mathbf{x}_K), K \subseteq [n])$ .

Receiver  $i$  decodes for a subset of messages indexed by  $D_i \subseteq [n] \setminus A_i$ , such that  $i \in D_i$ . The tuple of decoding message sets is denoted by  $\mathbf{D} = (D_i, i \in [n])$  and referred to as the decoding configuration for composite coding.

**Decoding.** *Step 1.* Receiver  $i$  finds the unique composite index tuple  $(\hat{w}_K, K \subseteq [n])$  such that  $y = y(\hat{w}_K, K \subseteq [n])$ . If there is more than one such tuple it declares an error. *Step 2.* Assuming  $(\hat{w}_K, K \subseteq [n])$  is correct, receiver  $i$  finds the unique message tuple  $\hat{\mathbf{x}}_{D_i}$  such that  $\hat{w}_K = w_K(\hat{\mathbf{x}}_K)$  for all  $K \subseteq D_i \cup A_i$ . If there is more than one such tuple it declares an error.

The achievable rate region of this coding scheme [4], which will be referred to as the original CC or just the CC later, is summarized below. The proof can be found in [3], [7].

*Proposition 1:* A rate tuple  $\mathbf{R}$  is achievable for the index coding problem  $(i|A_i)$ ,  $i \in [n]$ , under a given  $\mathbf{D}$  if

$$\sum_{K \subseteq [n], K \not\subseteq A_i} S_K < 1, \quad \forall i \in [n], \quad (5)$$

$$\sum_{j \in L} R_j < \sum_{\substack{K \subseteq D_i \cup A_i, \\ K \cap L \neq \emptyset}} S_K, \quad \forall L \subseteq D_i, i \in [n], \quad (6)$$

for some  $S_K \geq 0$ ,  $K \subseteq [n]$ .

To compute an explicit achievable rate or rate region from (5) and (6), one has to eliminate the intermediate variables  $(S_K, K \subseteq [n])$  using an optimization tool.

Due to space limitations, the more general enhanced composite coding [5], or enhanced CC, will not be presented here, but will be used for our numerical analysis. Enhancement is due to better convexification of message rates and composite index rates over different decoding configurations.

As shown by the following example, even the best composite coding, namely, time sharing of enhanced CC over all possible subproblems (including the problem itself), is not optimal in general.

*Example 1:* Consider the achievable symmetric rate  $R_{\text{sym}} = R_i, \forall i \in [7]$  for the 7-message problem  $(1|5), (2|3, 5, 6), (3|4, 6, 7), (4|1, 2, 7), (5|2, 3, 4, 7), (6|3, 4, 7), (7|1, 2, 4)$ . Time sharing of enhanced CC over subproblems gives  $R_{\text{sym}} < \frac{1}{3.25}$ . However, by simple linear encoding scheme such as sending three codewords,  $x_1 \oplus x_5, x_2 \oplus x_4 \oplus x_5 \oplus x_7$ , and  $x_3 \oplus x_4 \oplus x_6 \oplus x_7$ , a symmetric rate of  $R_{\text{sym}} = \frac{1}{3}$  can be achieved for the problem, which indeed reaches the symmetric capacity.

*Simplifications:* In Algorithm 1 we repeat a recent result from [8], which can greatly reduce the number of decoding configurations to consider in the CC scheme without any performance loss in achievable rates. Only supersets of  $\underline{\mathbf{D}}$  (i.e.,  $D_i \supseteq \underline{D}_i$  for all  $i \in [n]$ ) need to be considered.

In [8] we also proposed another way of complexity reduction by removing composite index rate variables,  $S_K$  that are guaranteed to be dominated by other variables  $S_{K'}$ . We will use a similar idea later in Section V for reducing the number of doubly composite index rate variables in the TLCC.

---

#### Algorithm 1: Natural decoding configuration

---

**Input :** Index coding problem  $(i|A_i)$ ,  $i \in [n]$ .

**Output:** Natural decoding configuration

$$\underline{\mathbf{D}} = (\underline{D}_i, i \in [n]).$$

1 Initialize  $\underline{D}_i = \{i\}$ ,  $i \in [n]$ .

2 As long as there exists  $i, j \in [n]$  such that  $A_j \subseteq A_i \cup \underline{D}_i$  and  $\underline{D}_j \not\subseteq A_i \cup \underline{D}_i$ , update  $\underline{D}_i \leftarrow \underline{D}_i \cup (\underline{D}_j \setminus A_i)$ . If no such  $i, j$  exist, terminate the algorithm.

---

### IV. THREE-LAYER COMPOSITE CODING

In this section, we present the main result of this paper, the three-layer composite coding scheme and its corresponding achievable rate region. Due to the space limitations, we only present results for a fixed decoding configuration. Extension to enhanced three-layer composite coding is straightforward.

**Codebook generation.** *Step 1.* For each  $K \subseteq [n]$ , generate a composite index  $w_K(\mathbf{x}_K)$  drawn uniformly at random from  $[2^{s_K}]$ . *Step 2.* For each  $M \subseteq N$ , generate a doubly composite index  $v_M(w_K, K \in M)$  drawn uniformly at random from  $[2^{z_M}]$ , where  $z_M = \lceil rZ_M \rceil$  and  $Z_M$  is the rate for the doubly composite index  $v_M$ . *Step 3.* Generate the codeword  $y(v_M, M \subseteq N)$  drawn uniformly at random from  $[2^r]$ . The codebook  $\{(w_K(\mathbf{x}_K), K \subseteq [n]), (v_M(w_K, K \in M), M \subseteq N), y(v_M, M \subseteq N)\}$  is revealed to all parties.

For brevity, when we say (composite) index  $K$  and (doubly composite) index  $M$ , we mean composite index  $w_K(\mathbf{x}_K)$  and doubly composite index  $v_M(w_K, K \in M)$ .

**Encoding.** To communicate messages  $\mathbf{x}_{[n]}$ , the server transmits  $y(v_M(w_K(\mathbf{x}_K), K \in M), M \subseteq N)$ .

Receiver  $i$  decodes for a subset of messages indexed by  $D_i \subseteq [n] \setminus A_i$ , such that  $i \in D_i$  as well as a subset of composite indices indexed by  $P_i \subseteq N \setminus 2^{A_i}$ . The tuple of decoding message sets and decoding (composite) index sets is denoted by  $(\mathbf{D}, \mathbf{P}) = ((D_i, i \in [n]), (P_i, i \in [n]))$  and referred to as the decoding configuration for the TLCC.

**Decoding.** *Step 1.* Receiver  $i$  finds the unique doubly composite index tuple  $(\hat{v}_M, M \subseteq N)$  such that  $y = y(\hat{v}_M, M \subseteq N)$ . If there is more than one such tuple, it declares an error. *Step 2.* Assuming that  $(\hat{v}_M, M \subseteq N)$  is correct, receiver  $i$  finds the unique composite index tuple  $(\hat{w}_K, K \in P_i)$  such that  $\hat{v}_M = v_M(\hat{w}_K, K \in M)$  for every  $M \subseteq 2^{A_i} \cup P_i$ . If there is more than one such tuple, it declares an error. *Step 3.* Assuming that  $(\hat{w}_K, K \in P_i)$  is correct, receiver  $i$  finds the unique message tuple  $\hat{\mathbf{x}}_{D_i}$  such that  $\hat{w}_K = w_K(\hat{\mathbf{x}}_K)$  for every  $K \subseteq D_i \cup A_i, K \in P_i$ . If there is more than one such tuple, it declares an error.

We have the following theorem.

*Theorem 1:* A rate tuple  $\mathbf{R}$  is achievable for the indexing problem  $(i|A_i, i \in [n])$ , under a given  $(\mathbf{D}, \mathbf{P})$  if

$$\sum_{M \subseteq N, M \not\subseteq 2^{A_i}} Z_M < 1, \quad \forall i \in [n], \quad (7)$$

$$\sum_{K \in Q} S_K < \sum_{\substack{M \subseteq 2^{A_i} \cup P_i, \\ M \cap Q \neq \emptyset}} Z_M, \quad \forall Q \subseteq P_i, i \in [n], \quad (8)$$

$$\sum_{j \in L} R_j < \sum_{\substack{K \subseteq D_i \cup A_i, \\ K \cap L \neq \emptyset, \\ K \in P_i}} S_K, \quad \forall L \subseteq D_i, i \in [n], \quad (9)$$

for some  $Z_M \geq 0, M \subseteq N$  and  $S_K \geq 0, K \subseteq [n]$ .

The proof for Theorem 1 is presented in Appendix A. The inequalities in (7), (8), and (9) signify the first-step, second-step, and third-step decoding constraints for the TLCC. It can be shown that for the TLCC,  $\underline{\mathbf{D}}$  found in Algorithm 1 still serves as the baseline or natural decoding message set, such that it suffices to only consider the supersets of  $\underline{\mathbf{D}}$ .

It will be shown later in the paper that the TLCC scheme subsumes the CC scheme in general. For some problems, it can give strictly better results.

*Example 2:* Consider the problem in Example 1. While time sharing of enhanced CC over subproblems gives  $R_{\text{sym}} < \frac{1}{3.25}$ , according to Theorem 1, any  $R_{\text{sym}} < \frac{1}{3}$  is achievable as follows. Set  $(\mathbf{D}, \mathbf{P})$  such that  $D_i = \{i\}, \forall i \in [n]$ , and  $P_1 = \{\{1\}\}, P_2 = \{\{2\}, \{4, 7\}\}, P_5 = \{\{5\}\}, P_3 = P_6 = \{\{3, 6\}\}, P_4 = P_7 = \{\{5\}, \{4, 7\}\}$ . Then set  $Z_M = 0, \forall M \subseteq N$  except for  $Z_{\{\{1\}, \{5\}\}}, Z_{\{\{2\}, \{4, 7\}, \{5\}\}}$  and  $Z_{\{\{3, 6\}, \{4, 7\}\}}$ . Also set  $S_K = 0, \forall K \subseteq [n]$  except for  $S_{\{1\}}, S_{\{2\}}, S_{\{5\}}, S_{\{3, 6\}}$ , and  $S_{\{4, 7\}}$ . Writing all the active

decoding inequalities of Theorem 1 yields

$$\begin{aligned} Z_{\{\{1\}, \{5\}\}} + Z_{\{\{2\}, \{4, 7\}, \{5\}\}} + Z_{\{\{3, 6\}, \{4, 7\}\}} &< 1, \\ S_{\{1\}} &< Z_{\{\{1\}, \{5\}\}}, \\ S_{\{2\}} + S_{\{4, 7\}} &< Z_{\{\{2\}, \{4, 7\}, \{5\}\}} + Z_{\{\{3, 6\}, \{4, 7\}\}}, \\ S_{\{2\}} &< Z_{\{\{2\}, \{4, 7\}, \{5\}\}}, \\ S_{\{3, 6\}} &< Z_{\{\{3, 6\}, \{4, 7\}\}}, \\ S_{\{5\}} + S_{\{4, 7\}} &< Z_{\{\{1\}, \{5\}\}} + Z_{\{\{2\}, \{4, 7\}, \{5\}\}}, \\ S_{\{4, 7\}} &< Z_{\{\{2\}, \{4, 7\}, \{5\}\}}, \\ S_{\{5\}} &< Z_{\{\{2\}, \{4, 7\}, \{5\}\}}, \\ R_{\text{sym}} &< S_K, \forall K \in \{\{1\}, \{2\}, \{5\}, \{3, 6\}, \{4, 7\}\}. \end{aligned}$$

For an arbitrary  $\epsilon \in (0, \frac{1}{3})$ , assigning  $R_{\text{sym}} = \frac{1}{3} - \epsilon, S_{\{1\}} = S_{\{2\}} = S_{\{5\}} = S_{\{3, 6\}} = S_{\{4, 7\}} = \frac{1}{3} - \frac{\epsilon}{2}, Z_{\{\{1\}, \{5\}\}} = Z_{\{\{2\}, \{4, 7\}, \{5\}\}} = Z_{\{\{3, 6\}, \{4, 7\}\}} = \frac{1}{3} - \frac{\epsilon}{4}$  satisfies all the inequalities. Hence any  $R_{\text{sym}} < \frac{1}{3}$  is achievable by the TLCC.

## V. SIMPLIFICATIONS FOR THE TLCC

The main challenges for the TLCC's computation are the overwhelming number of  $Z_M$  variables, which is  $2^{2^n - 1} - 1$ , as well as the choice of  $\mathbf{P}$ . To circumvent these, we now present a series of simplifications.

### A. Limiting the Choice of $\mathbf{P}$

First, we propose a heuristic baseline decoding index set  $\underline{\mathbf{P}}^{\mathbf{D}}$  for a given  $\mathbf{D}$ . The idea of the heuristic is as follows. For a given  $D_i$ , the starting set  $P_i$  for receiver  $i$  contains all subsets of  $[n]$  that intersect with  $D_i$  and do not intersect with interfering messages  $(A_i \cup D_i)^c$ . Then, following similar lines of thought as in Algorithm 1, we iteratively add missing elements from  $P_j$  to  $P_i$  if  $2^{A_j} \subseteq 2^{A_i} \cup P_i^{\mathbf{D}}$  and  $P_j^{\mathbf{D}} \not\subseteq 2^{A_i} \cup P_i^{\mathbf{D}}$ . This is summarized in Algorithm 2, which makes use of the following notation. For any message subsets  $K, L \subseteq [n]$ :

$$T_{K, \bar{L}} = \{J \in N : J \cap K \neq \emptyset, J \cap L = \emptyset\}.$$

After running Algorithm 2, one can only consider the supersets of  $\underline{\mathbf{P}}^{\mathbf{D}}$  in the second-step decoding, which can lead to great reduction in the number of possible decoding configurations, albeit with possible performance loss.

For any collection of composite indices  $M \subseteq N$ , we set

$$\Gamma_*(M) = \bigcup_{K \in M} \{L \in N : L \subseteq K\} = \bigcup_{K \in M} 2^K$$

to be the subset completion of  $M$ . Note  $2^K = \Gamma_*(\{K\})$ . The set  $M \subseteq N$  is subset-complete if  $M = \Gamma_*(M)$ .

For any  $\mathbf{P}$ , if  $2^{A_i} \cup P_i$  is subset-complete for any  $i \in [n]$ , we simply say that  $\mathbf{P}$  is subset-complete. Then we have the following lemma.

*Lemma 1:* For any  $\mathbf{D}$ , its corresponding  $\mathbf{P}^{\mathbf{D}}$  given by Algorithm 2 is subset-complete.

*Proof:* For any  $i \in [n]$ , consider the initial set  $P_i^{\mathbf{D}} = T_{D_i, \overline{(A_i \cup D_i)^c}}$ . As  $2^{A_i} \cup P_i^{\mathbf{D}} = 2^{A_i \cup D_i}$ , we have  $\Gamma_*(2^{A_i} \cup P_i^{\mathbf{D}}) = \Gamma_*(2^{A_i \cup D_i}) = 2^{A_i \cup D_i} = 2^{A_i} \cup P_i^{\mathbf{D}}$ . Since the union of any two subset-complete sets is also subset-complete, for the final  $P_i^{\mathbf{D}}, 2^{A_i} \cup P_i^{\mathbf{D}}$  must be subset-complete as well. ■

For a given  $\mathbf{D}$ , despite possible performance loss, we may further narrow down the range of  $\mathbf{P}$  to consider by enforcing that  $\mathbf{P}$  is a superset of  $\mathbf{P}^{\mathbf{D}}$  such that  $\mathbf{P}$  is subset-complete.

### B. Reducing Doubly Composite Indices

In the following proposition, we adopt and modify the composite index rate transfer and removal technique proposed in [8, Theorem 1] to doubly composite index rates  $Z_M, M \subseteq N$ .

*Proposition 2:* For a given decoding configuration and arbitrary  $M, M' \subseteq N, M \neq M'$ , compare the relative presence for  $Z_M, Z_{M'}$  in the inequalities identified by Theorem 1.

- 1) If  $Z_{M'}$  appears in any first-step decoding inequality then so does  $Z_M$ , AND
- 2) If  $Z_M$  appears in any second-step decoding inequality then so does  $Z_{M'}$ ,

then  $Z_M$  can be removed from the rate expressions without affecting the resulting rate region.

The proof is similar to that in [8] for removing  $S_K$  rates.

The number of  $Z_M$  variables remaining after applying the above proposition to every  $M, M'$  pairs can be much smaller than its original value  $2^{2^n-1} - 1$ . This leads to considerably lower complexity in computation. However, as the original number of  $Z_M$  variables is extremely large for large  $n$ , even the computational complexity for applying Proposition 2 for every possible  $M, M'$  can be a problem. Hence, we propose to systematically exclude some  $Z_M$  variables even before applying Proposition 2. This can be done if  $\mathbf{P}$  is subset-complete, such as  $\mathbf{P}^D$  at the output of Algorithm 2.

*Corollary 1:* For any decoding configuration  $(\mathbf{D}, \mathbf{P})$  such that  $\mathbf{P}$  is subset-complete, it suffices to only consider  $Z_M$  such that  $M$  is subset-complete.

*Proof:* Consider an arbitrary  $M \subseteq N$  such that  $M \neq \Gamma_*(M)$ , set  $M' = \Gamma_*(M)$ .  $M'$  is subset-complete and  $M \subset M'$ . Since  $\bigcup_{K \in M} K = \bigcup_{K \in M'} K$ , whenever  $Z_M$  appears in a first-step decoding inequality in (7), so does  $Z_{M'}$  and vice versa. Consider the relative presence for  $Z_M, Z_{M'}$  in the second-stage inequalities. For any  $i \in [n]$ , since  $2^{A_i} \cup P_i$  is subset-complete, if  $M \subseteq 2^{A_i} \cup P_i$ , we must have  $M' \subseteq 2^{A_i} \cup P_i$ . Also, as  $M \subset M'$ , if  $M \cap L \neq \emptyset$  for any  $L \subseteq D_i$ , we must have  $M' \cap L \neq \emptyset$ . Therefore, whenever  $Z_M$  appears in a second-step decoding inequality in (8) so does  $Z_{M'}$ . According to Proposition 2 any such  $Z_M$  can be removed. ■

### C. Reducing Doubly Composite Indices Based on Already Reduced Composite Indices

For a given problem and decoding message set  $\mathbf{D}$ , use  $N_{\mathcal{K}}(\mathbf{D})$  and  $N'_{\mathcal{K}}(\mathbf{D})$  to denote the collection of  $K$  such that  $S_K$  remains after applying [8, Theorem 1] or [8, Algorithm 1], which remove some  $S_K$  variables from the expressions of the rate region (5)-(6) of the CC scheme. The difference is that [8, Theorem 1] is guaranteed not to reduce the CC achievable rate region, whereas [8, Algorithm 1] is a possibly suboptimal heuristic. When the context is clear, we simply use the shorthand notation  $N_{\mathcal{K}}$  and  $N'_{\mathcal{K}}$ . Note that  $N'_{\mathcal{K}} \subseteq N_{\mathcal{K}} \subseteq N$ . For the 7-message problem in Examples 1 and 2, setting  $\mathbf{D}$  to be  $\underline{\mathbf{D}}$  and applying [8, Theorem 1] or [8, Algorithm 1], we find  $|N_{\mathcal{K}}| = 16$  and  $|N'_{\mathcal{K}}| = 7$ , respectively. Note that the original number of  $S_K$  variables is  $2^7 - 1 = 127$ .

For the TLCC, one can remove any composite index  $w_K, K \notin N_{\mathcal{K}}$  or  $K \notin N'_{\mathcal{K}}$  from the coding scheme. Since doubly composite indices  $v_M$  are generated from composite

---

### Algorithm 2: Heuristic baseline decoding index set.

---

**Input :** Index coding problem  $(i|A_i), i \in [n]$ , and decoding message set  $\mathbf{D} = (D_i, i \in [n])$

**Output:** Heuristic baseline decoding index set  $\mathbf{P}^D = (P_i^D, i \in [n])$ .

---

- 1 Initialize  $P_i^D = T_{D_i, (\overline{A_i \cup D_i})^c}, i \in [n]$ .
  - 2 If there exists  $i, j \in [n]$  such that  $2^{A_j} \subseteq 2^{A_i} \cup P_i^D$  and  $P_j^D \not\subseteq 2^{A_i} \cup P_i^D$ , update  $P_i^D \leftarrow P_i^D \cup (P_j^D \setminus 2^{A_i})$ . If no such  $i, j$  exist, terminate the algorithm.
- 

indices  $w_K$ , this will naturally narrow down the range of  $v_M$  from  $M \subseteq N$  to  $M \subseteq N_{\mathcal{K}}$  or  $M \subseteq N'_{\mathcal{K}}$ , which can lead to a huge reduction in the number of  $Z_M$  variables from  $2^{|N|} - 1 = 2^{2^n-1} - 1$  to  $2^{|N_{\mathcal{K}}|} - 1$  or  $2^{|N'_{\mathcal{K}}|} - 1$ . Note that such reduction may lead to performance loss.

The simplification techniques are summarized in Table I.

### D. Simplified TLCC Subsumes Composite Coding

In this subsection, we prove that even with some of the simplifications discussed so far, the TLCC scheme is guaranteed to perform at least as well as the original CC.

Consider an arbitrary index coding problem and an arbitrary  $\mathbf{D}$  that is a superset of  $\underline{\mathbf{D}}$ , let  $\mathcal{R}_{CC}(\mathbf{D})$  denote the achievable rate region of the CC with  $\mathbf{D}$ .

With slight abuse of notation, for any  $K \in N_{\mathcal{K}}$  and  $M \subseteq N_{\mathcal{K}}$ , let  $2^K$  denote the collection of subsets of  $K$  with respect to  $N_{\mathcal{K}}$ ,  $2^K = \{L \in N_{\mathcal{K}} : L \subseteq K\}$ , and let  $\Gamma_*(M)$  denote the subset completion of  $M$  with respect to  $N_{\mathcal{K}}$ ,  $\Gamma_*(M) = \bigcup_{K \in M} \{L \in N_{\mathcal{K}} : L \subseteq K\} = \bigcup_{K \in M} 2^K$ . Also, set  $T_{K,L} = \{J \in N_{\mathcal{K}} : J \cap K \neq \emptyset, J \cap L = \emptyset\}$ , and hence the baseline decoding index set  $\mathbf{P}^D = (P_i^D, i \in [n])$  given by Algorithm 2 is within the range of  $N_{\mathcal{K}}$ . Let  $\mathcal{R}_{TLCC}(\mathbf{D}, \mathbf{P})$  denote the achievable rate region of the TLCC with decoding configuration  $(\mathbf{D}, \mathbf{P})$  with the following simplifications:

- 1) any  $S_K, K \notin N_{\mathcal{K}}$  and  $Z_M, M \not\subseteq N_{\mathcal{K}}$  are removed from the decoding inequalities (7)-(9), and
- 2)  $\mathbf{P} = (P_i, i \in [n])$  is a decoding index set such that for any  $i \in [n]$ ,  $P_i^D \subseteq P_i \subseteq N_{\mathcal{K}} \setminus 2^{A_i}$  and that  $2^{A_i} \cup P_i$  is subset-complete with respect to  $N_{\mathcal{K}}$ .

*Proposition 3:*  $\mathcal{R}_{CC}(\mathbf{D}) \subseteq \mathcal{R}_{TLCC}(\mathbf{D}, \mathbf{P})$ .

*Proof:* Let  $\mathbf{R} \in \mathcal{R}_{CC}(\mathbf{D})$  be an achievable rate tuple of the CC. According to [8], removing any  $S_K$  that  $K \notin N_{\mathcal{K}}$  does not affect the achievable rate region of the CC, hence there exists some  $(S_K, K \in N_{\mathcal{K}})$  such that  $\mathbf{R}$  and  $(S_K, K \in N_{\mathcal{K}})$  satisfy (5) and (6) with  $\mathbf{D}$ . Set  $Z_M, M \subseteq N_{\mathcal{K}}$  as

$$Z_M = \begin{cases} S_K, & \text{if } M = 2^K \text{ for some } K \in N_{\mathcal{K}}, \\ 0, & \text{otherwise.} \end{cases}$$

Now we show that  $\mathbf{R}, (S_K, K \in N_{\mathcal{K}})$  and  $(Z_M, M \subseteq N_{\mathcal{K}})$  satisfy the decoding inequalities of the TLCC, (7)-(9), with  $(\mathbf{D}, \mathbf{P})$  defined above.

First, for any  $i \in [n]$ , we have  $\sum_{M \subseteq N_{\mathcal{K}}: M \not\subseteq 2^{A_i}} Z_M = \sum_{2^K \subseteq N_{\mathcal{K}}: 2^K \not\subseteq 2^{A_i}} Z_{2^K} = \sum_{K \in N_{\mathcal{K}}: K \not\subseteq A_i} S_K < 1$ .

Second, for any  $Q \subseteq P_i, i \in [n]$ , as  $2^{A_i} \cup P_i$  is subset-complete with respect to  $N_{\mathcal{K}}$ , we know that for any  $K \in Q$ ,

$K \in 2^{A_i} \cup P_i$  and thus  $2^K = \Gamma_*(K) \subseteq 2^{A_i} \cup P_i$ . Also, for any  $K \in Q$ ,  $2^K \cap Q \neq \emptyset$ . Hence, we have

$$\sum_{\substack{M \subseteq N_{\mathcal{K}}: \\ M \subseteq 2^{A_i} \cup P_i \\ M \cap Q \neq \emptyset}} Z_M = \sum_{\substack{2^K \subseteq N_{\mathcal{K}}: \\ 2^K \subseteq 2^{A_i} \cup P_i \\ 2^K \cap Q \neq \emptyset}} Z_{2^K} \geq \sum_{\substack{K \in N_{\mathcal{K}}: \\ K \in Q}} S_K. \quad (10)$$

Third, for any  $L \subseteq D_i, i \in [n]$ , note that  $T_{L, \overline{(A_i \cup D_i)^c}} \subseteq T_{D_i, \overline{(A_i \cup D_i)^c}} \subseteq P_i^{\mathbf{D}} \subseteq P_i$ . Therefore, if  $K \in N_{\mathcal{K}}, K \subseteq A_i \cup D_i$  and  $K \cap L \neq \emptyset$  then  $K$  must be in set  $P_i$ . Hence,

$$\sum_{j \in L} R_j < \sum_{K \in N_{\mathcal{K}}: K \subseteq A_i \cup D_i, K \cap L \neq \emptyset} S_K \quad (11)$$

$$= \sum_{K \in N_{\mathcal{K}}: K \subseteq A_i \cup D_i, K \cap L \neq \emptyset, K \in P_i} S_K. \quad (12)$$

We have proven that  $\mathbf{R}$  and  $(S_K, K \in N_{\mathcal{K}})$  and  $(Z_M, M \subseteq N_{\mathcal{K}})$  satisfy (7)-(9) with  $(\mathbf{D}, \mathbf{P})$ , which means that  $\mathbf{R} \in \mathcal{R}_{\text{TLCC}}(\mathbf{D}, \mathbf{P})$ . In summary,  $\mathcal{R}_{\text{CC}}(\mathbf{D}) \subseteq \mathcal{R}_{\text{TLCC}}(\mathbf{D}, \mathbf{P})$ . ■

$\mathcal{R}_{\text{TLCC}}(\mathbf{D}, \mathbf{P})$  can be strictly larger than  $\mathcal{R}_{\text{CC}}(\mathbf{D})$ .  $\mathcal{R}_{\text{TLCC}}(\mathbf{D}, \mathbf{P})$  can even strictly outperform the best composite coding for some cases, as shown below.

Consider Example 2 again. Time sharing of enhanced CC over subproblems gives  $R_{\text{sym}} < \frac{1}{3.25}$ . For the TLCC, we apply the following simplifications. Fix  $\mathbf{D}$  to be  $\underline{\mathbf{D}}$ . Apply [8, Theorem 1] to obtain  $N_{\mathcal{K}}$ , where  $|N_{\mathcal{K}}| = 16$ . Remove any  $S_K, K \notin N_{\mathcal{K}}$  and  $Z_M, M \not\subseteq N_{\mathcal{K}}$ . Thus the number of  $Z_M$  variables is hugely reduced from  $2^{2^7} - 1$  to  $2^{16} - 1$ . Find the output of Algorithm 2,  $\mathbf{P}^{\mathbf{D}}$  and set  $\mathbf{P}$  as  $P_i = P_i^{\mathbf{D}}, \forall i \in \{1, 3, 5, 6\}, P_2 = P_2^{\mathbf{D}} \cup (2^{\{4,7\}} \cap N_{\mathcal{K}})$  and  $P_i = P_i^{\mathbf{D}} \cup (2^{\{5\}} \cap N_{\mathcal{K}}), \forall i \in \{4, 7\}$ . Use Proposition 2 to further remove unnecessary  $Z_M$  variables. Applying Theorem 1 we obtain the optimal  $R_{\text{sym}} < \frac{1}{3}$ . Moreover,  $R_{\text{sym}} < \frac{1}{3}$  can be obtained with even much lower computational complexity via using  $N'_{\mathcal{K}}$  from [8, Algorithm 1], instead of  $N_{\mathcal{K}}$ , for the entire simplifying process, where  $|N'_{\mathcal{K}}| = 7$ . Thus the number of  $Z_M$  variables is only  $2^7 - 1$ . We can simply set  $\mathbf{P} = \mathbf{P}^{\mathbf{D}}$  (which is now within the range of  $N'_{\mathcal{K}}$ ). Applying Proposition 2, the final number of remaining  $Z_M$  variables is merely 6.

Table I

SIMPLIFICATION TECHNIQUES AND WHETHER THEY RETAIN OPTIMALITY

Simplification	References	Optimality
Limiting $\mathbf{D}$ to be supersets of $\underline{\mathbf{D}}$	Alg. 1, [8, Thm. 2]	Yes
Limiting $\mathbf{P}$ to be supersets of $\mathbf{P}^{\mathbf{D}}$ such that $\mathbf{P}$ is subset-complete	Alg. 2, Lem. 1	Unknown
Removing $Z$ variables by pairwise comparison	Prop. 2	Yes
Removing $Z_M, M \neq \Gamma_*(M)$ when $\mathbf{P}$ is subset-complete	Cor. 1	Yes
Removing $S_K$ and $Z_M$ not fully embedded in $N_{\mathcal{K}}$ or $N'_{\mathcal{K}}$	Sec. V-C, [8, Thm. 1, Alg. 1]	Unknown

It remains to compare the TLCC with other coding schemes in future works, especially the structured coding schemes such as those proposed in [9], [10]. Another fascinating

direction is to see whether the achievable rate region of such layered random coding scheme approaches the capacity region in general as the number of layers increases, and if not, to identify any fundamental gaps between them.

## APPENDIX A

## PROOF OF THEOREM 1

We only present the error analysis for the second-step decoding of the TLCC, as the analysis for other steps is quite similar to that of the CC scheme. Assume that the doubly composite indices  $(\hat{v}_M, M \subseteq N)$  have been correctly decoded. For receiver  $i$ , we partition the error event according to the collection  $Q \subseteq P_i$  for erroneous composite indices. That is,  $\hat{w}_K \neq w_K$  iff  $K \in Q$ . Hence, for the second-step decoding error probability  $P_e$ , we have

$$P_e = P\{\hat{v}_M = v_M(\hat{w}_K, K \in M) \text{ for all } M \subseteq 2^{A_i} \cup P_i \text{ for some } \hat{w}_K \neq w_K, K \in P_i\} \quad (13)$$

$$\leq \sum_{Q \subseteq P_i} \sum_{\substack{(\hat{w}_K, K \in P_i): \\ \hat{w}_K \neq w_K, K \in Q \\ \hat{w}_K = w_K, K \notin Q}} P\left\{ \bigcap_{\substack{M \subseteq 2^{A_i} \cup P_i \\ M \cap Q \neq \emptyset}} \{v_M(\hat{w}_K, K \in M)\} \right\} \quad (14)$$

$$< \sum_{Q \subseteq P_i} 2^{\sum_{K \in Q} S_K} / 2^{\sum_{M \in \mathcal{M}} Z_M} \quad (15)$$

$$< \sum_{Q \subseteq P_i} 2^{r(\sum_{K \in Q} S_K - \sum_{M \in \mathcal{M}} Z_M) + \sum_{K \in Q}}, \quad (16)$$

where  $\mathcal{M} = \{M \subseteq N : M \subseteq 2^{A_i} \cup P_i, M \cap Q \neq \emptyset\}$ , and the first inequality is due to the union bound, and the second inequality holds since for each  $Q$ , the number of erroneous tuples is  $\prod_{K \in Q} (2^{S_K} - 1) \leq 2^{\sum_{K \in Q} S_K}$ , and the probability for any two distinct composite index tuples being mapped to the same doubly composite index  $v_M$  for all  $M \in \mathcal{M}$  is  $2^{-\sum_{M \in \mathcal{M}} Z_M}$ .  $P_e$  tends to zero as  $r \rightarrow \infty$  if (8) is satisfied.

## REFERENCES

- [1] Y. Birk and T. Kol, "Informed-source coding-on-demand (ISCOD) over broadcast channels," in *IEEE INFOCOM*, Mar. 1998, pp. 1257–1264.
- [2] L. Ong, "Optimal finite-length and asymptotic index codes for five or fewer receivers," *IEEE Transactions on Information Theory*, vol. 63, no. 11, pp. 7116–7130, Nov 2017.
- [3] F. Arbabjolfaei and Y.-H. Kim, *Fundamentals of Index Coding*. To appear in Foundations and Trends in Communications and Information Theory, 2018.
- [4] F. Arbabjolfaei, B. Bandemer, Y.-H. Kim, E. Sasoglu, and L. Wang, "On the capacity region for index coding," in *Proc. IEEE Int. Symp. on Information Theory (ISIT)*, 2013, pp. 962–966.
- [5] Y. Liu, P. Sadeghi, F. Arbabjolfaei, and Y.-H. Kim, "On the capacity for distributed index coding," in *Proc. IEEE Int. Symp. on Information Theory (ISIT)*, Aachen, Germany, Jun. 2017, pp. 3055–3059.
- [6] K. Wan, D. Tuninetti, and P. Piantanida, "A novel index coding scheme and its application to coded caching," in *2017 Information Theory and Applications Workshop (ITA)*, 2017, pp. 1–6.
- [7] F. Arbabjolfaei, B. Bandemer, and Y.-H. Kim, "Index coding via random coding," in *Iran Workshop on Communication and Information Theory (IWCIT)*, Tehran, Iran, May 2014.
- [8] Y. Liu, P. Sadeghi, F. Arbabjolfaei, and Y.-H. Kim, "Simplified composite coding for index coding," in *Proc. IEEE Int. Symp. on Information Theory (ISIT)*, Vail, CO, Jun. 2018, pp. 456–460.
- [9] C. Thapa, L. Ong, and S. J. Johnson, "Interlinked cycles for index coding: Generalizing cycles and cliques," *IEEE Transactions on Information Theory*, vol. 63, no. 6, pp. 3692–3711, 2017.
- [10] S. A. Jafar, "Topological interference management through index coding," *IEEE Transactions on Information Theory*, vol. 60, no. 1, pp. 529–568, 2014.