

# Maximum Likelihood Upper Bounds on the Capacities of Discrete Information Stable Channels

Tongxin Li, *Student Member, IEEE*

**Abstract**—Motivated by a greedy approach for generating *information stable* processes, we prove a universal maximum likelihood (ML) upper bound on the capacities of discrete information stable channels, including the binary erasure channel (BEC), the binary symmetric channel (BSC) and the binary deletion channel (BDC). The bound is derived leveraging a system of equations obtained via the Karush-Kuhn-Tucker conditions. Intriguingly, for some memoryless channels, *e.g.*, the BEC and BSC, the resulting upper bounds are *tight* and equal to their capacities. For the BDC, the universal upper bound is related to a function counting the number of possible ways that a length- $m$  binary subsequence can be obtained by deleting  $n - m$  bits (with  $n - m$  close to  $nd$  and  $d$  denotes the *deletion probability*) of a length- $n$  binary sequence. To get explicit upper bounds from the universal upper bound, it requires to compute a maximization of the matching functions over a Hamming cube containing all length- $n$  binary sequences. Calculating the maximization exactly is hard. Instead, we provide a combinatorial formula approximating it. Under certain assumptions, several approximations and an *explicit* upper bound for deletion probability  $d \geq 1/2$  are derived.

**Index Terms**—Information Stable Channels; Channel Capacity

## I. INTRODUCTION

THE *information stable* channels were introduced by Dubrushin in [1]. Under the *information stability* condition, sufficiently, their capacities can be expressed as

$$C = \liminf_{n \rightarrow \infty} \frac{1}{n} \sup_{\mathbf{X}} I(\mathbf{X}; \mathbf{Y}(\mathbf{X})). \quad (1)$$

Essentially, a channel satisfies information stability is equivalent to having the capacity expression above [2]. Preceding works have considered a variety of more general frameworks, *e.g.*, a formula for channel capacity [3] based on the information-spectrum method; a general capacity expression for channels with feedback [4]; general capacity formulas for classical-quantum channels [5], to list just a few.

Despite the simplicity of the formula in (1), for some channels with memory, explicitly computing the capacities directly using the general formula is often not trivial. A famous example is the binary deletion channel (BDC), which was introduced by Levenshtein in [6] more than fifty years ago to model synchronization errors. In his model, a transmitter sends an infinite stream of bits representing messages over a communication channel. Before reaching at a receiver, the bits are deleted independently and identically with some *deletion probability*  $d \in (0, 1)$ . The receiver wishes to recover

the original message based on the deleted bits, with an asymptotically zero (in the length of the stream) probability of error. The BDC satisfies the information stability [7]. Thus, the channel capacity denoted by  $C(d)$  can be expressed via the formula in (1). However, a precise characterization of  $C(d)$  is still unknown; it does not seem even possible to accurately compute the capacity numerically relying on the existing methods, for instance, the Blahut-Arimoto Algorithm (BAA) [8–10].

In this work, we consider discrete channels with finite alphabets and derive a general upper bound (called the maximum likelihood (ML) upper bound in Section III) on the capacities of information stable channels by analyzing a system of equations derived from the general formula in (1). We demonstrate that for channels without memory, *e.g.*, the binary erasure channel (BEC) and the binary symmetric channel (BSC). The corresponding upper bounds are tight for the BEC and BSC and equal to their channel capacities. For channels with memory, as a case study, we apply the ML upper bound to derive (*implicit* and *explicit*) approximations for  $C(d)$ , under certain assumptions.

### A. Background

A discrete channel with a finite alphabet can be regarded as a stochastic matrix from an input space of all infinite-length sequences to an output space containing all sequences that can be obtained via the channel law. Formally, we follow the approach in [11, 12] and define the transmitted and received bit-streams via infinite processes. For each fixed *block-length*  $n$ , there is a sequence of elements  $\mathbf{X}_1^{(n)} \dots \mathbf{X}_n^{(n)}$  selected from a finite set  $\mathcal{X}$ , and there is a probability distribution  $P_{\mathbf{X}^n}$  over this sequence. Let  $\mathbb{X}$  denote an input process in terms of finite-dimensional sequences such that  $\mathbb{X} := \{\mathbf{X}^n = (\mathbf{X}_1^{(n)}, \dots, \mathbf{X}_n^{(n)})\}_{n \geq 1}$ . Similarly, denote by  $\mathbb{Y} := \{\mathbf{Y}^m = (\mathbf{Y}_1^{(m)}, \dots, \mathbf{Y}_m^{(m)})\}_{m \geq 1}$  with each  $\mathbf{Y}_i^{(m)}$  in a finite set  $\mathcal{Y}$  the corresponding output process of finite-dimensional sequences induced by  $\mathbb{X}$  via the channel law  $\mathbb{W} := \{\mathbf{W}^n(\cdot|\cdot) : \mathcal{X}^n \rightarrow \mathcal{Y}^m\}_{n, m \geq 1}$ . So that

$$\mathbb{P}(\mathbf{Y}^m = \mathbf{y}^m | \mathbf{X}^n = \mathbf{x}^n) := \mathbf{W}^n(\mathbf{y}^m | \mathbf{x}^n).$$

Note that the block-length of received codewords  $m$  is not necessarily equal to  $n$ , the block-length of the transmitted codeword. Moreover, the output block-length  $m$  is allowed to be *flexible*, meaning that it can be regarded as a random variable with distribution specified by the channel law<sup>1</sup>. In the

Li is with the Computing + Mathematical Sciences Department, California Institute of Technology, Pasadena, CA 91125 USA (e-mail: tongxin@caltech.edu).

<sup>1</sup>Flexible output length allows us to apply this general framework to the BDC later in Section IV.

remaining part of this paper, we often omit the superscript  $m$  in  $\mathbf{y}^m$  and  $\mathbf{Y}^m$ , to avoid confusion. This indicates the length of the output sequence  $\mathbf{y}$  is not fixed.

### B. Outline of the Paper

The remaining content of the paper is organized as follows. In Section II we give a simplified version of  $C(d)$ , which is derived from the capacity formula in (2) for information stable channels. Based on it, we prove a general upper bound (the ML upper bound in Theorem 1) on information stable channels in Section III. Section III-C follows by verifying the tightness of the ML upper bound for the BEC and the BSC. Next, in Section IV, several approximations for the capacity of the BDC are reported.

## II. PRELIMINARIES

### A. Notational Convention

We use  $\log(\cdot)$  to denote logarithms over base 2, unless stated otherwise. Let  $\mathcal{X}^n$  and  $\overline{\mathcal{Y}}$  denote the set of all possible length- $n$  sequences and the set of all induced output sequences (having flexible lengths). Let  $N := |\mathcal{X}^n|$  and  $M := |\overline{\mathcal{Y}}|$ . We use the lowercase letter  $j$  to index the  $j$ -th length- $n$  input sequence  $\mathbf{x}_j^n$ , and the letter  $i$  to index the length- $m$  output sequence  $\mathbf{y}_i^m$  with  $j = 1, \dots, N$  and  $i = 1, \dots, M$  respectively. To distinguish between random variables and their realizations, we denote the former by capital letters and the latter by lower case letters, at most of the places throughout this work<sup>2</sup>.

### B. Capacity Proxies

For a fixed dimension  $n$ , we maximize the mutual information between  $\mathbf{X}^n \in \mathcal{X}^n$  and  $\mathbf{Y} \in \overline{\mathcal{Y}}$  in a way similar to defining the ‘‘information capacity’’ for discrete memoryless channels (DMCs) over the binary alphabet, to obtain the quantity:

$$C_n(\mathbf{W}^n) := \frac{1}{n} \sup_{\mathbf{X}^n} I(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n)) \quad (2)$$

where the supremum is taken over all  $\mathbf{X}^n \in \mathcal{X}^n$  with distributions in the set

$$\mathbb{P}^N := \left\{ P_{\mathbf{X}^n} \in \mathbb{R}^N : p_j \geq 0 \forall j = 1, \dots, N; \sum_{j=1}^N p_j = 1 \right\}. \quad (3)$$

### C. Information-stability

It turns out the quantity  $C_n(\mathbf{W}^n)$  is asymptotically (in  $n$ ) the same as the operational capacity under the following condition on channels, which is called *information stability*<sup>3</sup>.

<sup>2</sup>Except for the output length  $m$ , which is a random variable dependent on the channel law.

<sup>3</sup>The way of classifying the channels that have an operational meaning with the capacity expressions in (1.1) using a condition called *information stability* was first introduced by Dobrushin and Guoding Hu [1, 2]. It was restated and studied in many equivalent forms. For instance, in [11], information stability was proved to be insufficient to classify whether a source-channel separation holds or not. In [13], the expressions for optimistic channel capacity and optimal source coding rate are given for the class of information stable channels and similarly ‘‘information stable’’ sources respectively.

**Definition 1** (Information Stability for Channels [1–3]). *A channel  $\mathbb{W}$  is said to be information stable, if there exists an input process  $\mathbb{X}$  such that  $C_n(\mathbf{W}^n) < \infty$  for all sufficiently large  $n$  and*

$$\limsup_{n \rightarrow \infty} \mathbb{P} \left\{ \left| \frac{i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n))}{nC_n(\mathbf{W}^n)} - 1 \right| > \gamma \right\} = 0 \quad \forall \gamma > 0$$

where  $i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{x}^n; \mathbf{y}) := \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}^n)}{P_{\mathbf{Y}}(\mathbf{y})}$  denotes the information density for all  $\mathbf{x}^n \in \mathcal{X}^n$  and  $\mathbf{y} \in \overline{\mathcal{Y}}$ . In other words, the normalized information density  $\frac{1}{n} i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{x}^n; \mathbf{y})$  converges in probability to  $C_n(\mathbf{W}^n)$ .

Intuitively, information stability characterizes the types of channels in a manner similar to that the asymptotic equipartition property (AEP) characterizes stochastic sources. In fact, information stability for a channel  $\mathbb{W}$  implies the existence of a class of corresponding input processes  $\mathbb{X}$  such that  $\mathbb{X}$ , on being input to  $\mathbb{W}$ , results in a near-optimal code. For the case of discrete memoryless channels (DMCs), the operational meaning of the single-letter quantity in Eq. (4) below appears as a natural consequence of the law of large numbers. For general channels, by considering the asymptotic behavior of the information density taking on an optimal input process  $\mathbb{X}$  (which maximizes the mutual information), information stability provides (with sufficient generality, for a broader class of channels) an analogue of the law of large numbers. Such an optimal input process  $\mathbb{X}$  may be understood to be equivalent to a sequence of codes that are capacity-achieving asymptotically in the block-length  $n$ . The key idea relies on classical achievability bounds (for instance, Feinstein’s lemma [14], Shannon’s achievability bound [15]).

Dobrushin in [7] proved that BDCs are information stable as defined in Definition 1. For an arbitrary fixed  $n$ , maximizing  $I(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n))$  in Eq. (2) gives an optimal input distribution  $P_{\mathbf{X}^n}^*$ . Through appropriate achievability results ([14, 15]), it is possible to construct an  $(n, M, \lambda)$ -code whose error probability vanishes as  $n$  goes to infinity. In addition, the rate  $\log M/n$  approaches  $C_n(\mathbf{W}^n) < \infty$  for sufficiently large  $n$ . Hence for information stable channels, the capacities exist and can be written as<sup>4</sup>

$$C = \liminf_{n \rightarrow \infty} C_n(\mathbf{W}^n) < \infty. \quad (4)$$

### D. System of Equations for Optimality

Recall  $N := |\mathcal{X}^n|$  and  $M := |\overline{\mathcal{Y}}|^m$ .

Our approach focuses on bounding  $C_n(d)$ . Expressing the mutual information in terms of the channel law, the capacity-proxy  $C_n(\mathbf{W}^n)$  defined in (2) equals to

$$C_n(\mathbf{W}^n) = \frac{1}{n} \sup_{P_{\mathbf{X}^n}} \sum_{\mathbf{x}, \mathbf{y}} p(\mathbf{x}) \mathbf{W}^n(\mathbf{y}|\mathbf{x}) \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{\sum_{\mathbf{x}} p(\mathbf{x}) \mathbf{W}^n(\mathbf{y}|\mathbf{x})}. \quad (5)$$

<sup>4</sup>Note that this limiting expression does not always hold for general channels. For instance, consider one example in [3]: a binary channel with output codewords equal to the input codeword with probability 1/2 and changed independently of the input codewords with probability 1/2. The capacity of this channel is 0 since the error probability is always strictly positive and hence not vanishes. However, the formula in (4) gives 1/2.

Here, the supremum is taken over all distributions  $\{p(\mathbf{x})\}_{\mathbf{x} \in \mathcal{X}^n}$  in the set  $\mathbb{P}^N$  (defined in (3)) and the summation is taken over all length- $n$  input sequences  $\mathbf{x}_j \in \mathcal{X}^n$  and all output sequences  $\mathbf{y}_i \in \cup_m \mathcal{Y}^m$ .

From an optimization perspective, the asymptotic behavior of  $C_n(\mathbf{W}^n)$  can be captured by establishing a sequence of capacity-achieving distributions  $\{P_{\mathbf{X}^n}^*\}_n$  maximizing the following quantity for each  $n \geq 1$ :

$$\sum_{j=1}^N \sum_{i=1}^M p_j \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}{\sum_{j=1}^N p_j \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}. \quad (6)$$

Derived from the Karush-Kuhn-Tucker conditions, the following lemma generalizes Theorem 4.5.1. in [16] (cf. [17]), which was established to find channel capacities of DMCs with non-binary input/output alphabets. The lemma states a necessary and sufficient condition of the existence of  $\{P_{\mathbf{X}^n}^*\}_n$  maximizing (6) and it can be proved along the same line as in [16]. The only difference is that for general channels, the summation is taken over all sequences in  $\mathcal{Y}^m$  (this is in general exponential in  $m$ ). While for DMCs, the summation can be decomposed and taken over the alphabet set of each individual coordinate of the sequence, thus the number of summations is linear in  $m$ . For brevity the proof is omitted.

**Lemma 1** ([8, 9, 16, 17]). *Fix a block-length  $n \geq 1$ . There exists an optimal probability vector  $P_{\mathbf{X}^n}^* = (p_1^*, p_2^*, \dots, p_N^*)$  such that the quantity in (6) is maximized if and only if there exists  $\lambda_n \geq 0$  and for all  $j = 1, \dots, N$ ,*

$$\frac{1}{n} \sum_{i=1}^M \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)} \begin{cases} = \lambda_n & \text{if } p_j^* \neq 0 \\ \leq \lambda_n & \text{if } p_j^* = 0 \end{cases}. \quad (7)$$

Moreover the capacity  $C = \lim_{n \rightarrow \infty} \lambda_n$  if the limit exists.

Indeed, (see [16, 17]) a probability distribution for an information stable source  $\mathbb{X}$  satisfying (7) always exists as  $n$  grows. Thus, the capacity-achieving distribution with fixed block-length  $n$  can be attained by solving the system (7). Finding such an optimal  $P_{\mathbf{X}^n}^*$  for the system of equations (7) is equivalent to solving a non-linear system of equations that consists of exponentially (in  $n$ ) many variables. As introduced in Section I, the BAA is one of the algorithms that can be applied to search for numerical solutions of (7).

However, this approach has several limitations. On the one hand, in direct implementation of the BAA, as  $n$  grows, it becomes computationally intractable even to store the variables to be computed. On the other hand, as the BAA is itself an iterative algorithm attempting to solve the non-convex optimization problem (7), and to the best of our knowledge for general channels, there are no guarantees on how quickly the numerical solution converges as a function of the number of iterations. Therefore, instead of looking for numerical answers, we concentrate on finding a universal upper bound on the capacities of general channels. This motivates the next section.

### III. MAXIMUM LIKELIHOOD UPPER BOUND

In the sequel, we present some definitions. First, motivated by the notion of information stability defined in Definition 1, we characterize a subset of the joint set  $\mathcal{X}^n \times \overline{\mathcal{Y}}$  consisting of all possible combinations of input and output sequences. This subset satisfies two vital properties. First, it behaves as a ‘‘typical set’’ and contains nearly all pairs of  $(\mathbf{x}^n, \mathbf{y})$  randomly generated according to an arbitrary distributions  $P_{\mathbf{X}^n}$  for every large  $n$ . Second, conditioned on the pair  $(\mathbf{X}^n, \mathbf{Y})$  belongs to the subset, the conditional mutual information does not differ too much from  $C_n(\mathbf{W}^n)$ . Note that the concentration of information densities is stronger than that for information stable sources in two perspectives – the concentration is in expectation; and it is required to hold for every source  $\mathbb{X}$ .

**Definition 2.** *For information stable channels with any source  $\mathbb{X}$ , a subset  $\mathcal{A}$  of  $\mathcal{X}^n \times \overline{\mathcal{Y}}$  is called a concentration set if it satisfies*

$$\begin{aligned} \liminf_{n \rightarrow \infty} \mathbb{P}((\mathbf{X}^n, \mathbf{Y}) \in \mathcal{A}) &= 1, & (8) \\ \limsup_{n \rightarrow \infty} \mathbb{E} \left[ \left| \frac{i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n))}{nC_n(\mathbf{W}^n)} - 1 \right| \middle| (\mathbf{X}^n, \mathbf{Y}) \in \mathcal{A} \right] &= 0 & (9) \end{aligned}$$

where the randomness is over the source  $\mathbb{X}$  and channel law  $\mathbb{W}$ .

Later in Section III-C and Section IV-B1, we provide concrete and nontrivial examples of the concentration sets for the BEC, BSC and BDC respectively.

Since  $\mathcal{A} \subseteq \mathcal{X}^n \times \overline{\mathcal{Y}}$ , the next lemma is straightforward.

**Lemma 2.** *For each block-length  $n$ , there exists a subset  $\mathcal{B}$  of  $\overline{\mathcal{Y}}$  such that*

$$\mathcal{A} \subseteq \mathcal{X}^n \times \mathcal{B}.$$

It is useful to introduce the following ‘‘constant’’ version of the stochastic matrix  $\mathbf{W}$ , called the *stochastic factors* for convenience. Again, we will carefully construct them in Section III-C for both the BEC and the BSC, and in Section IV-B for the BDC.

**Definition 3.** *We call a set of functions  $f_k(\cdot|\cdot) : \mathcal{Y}^m \times \mathcal{X}^n \mapsto [0, 1]$  stochastic factors if there exists a decomposition of  $\mathcal{B} = \bigcup_{k \in \mathcal{K}} \mathcal{B}_k$  ( $\mathcal{K}$  is a discrete set) such that*

$$\sum_{\mathbf{y} \in \mathcal{B}_k} f_k(\mathbf{y}|\mathbf{x}) = 1, \quad \forall \mathbf{x} \in \mathcal{X}^n, k \in \mathcal{K}, \quad (10)$$

$$\sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k(\mathbf{y}|\mathbf{x})} \leq 1 \quad (11)$$

where  $\mathcal{A}_k := \mathcal{X}^n \times \mathcal{B}_k$ .

Based on the concentration set and the stochastic factor defined above, we obtain the following upper bound on the capacity of an information stable channel:

**Theorem 1** (Maximum Likelihood Upper Bound<sup>5</sup>). *For a discrete information stable channel defined in Section I-A,*

<sup>5</sup>For intuition on why we call a maximum likelihood (ML) upper bound, see Section III-A1.

assume there exist a concentration set  $\mathcal{A}$  and stochastic factors  $f$  defined above. The following upper bound on the channel capacity holds for any  $\mathcal{A} = \bigcup_{k \in \mathcal{K}} \mathcal{X}^n \times \mathcal{B}_k$  and  $\{f_k\}_{k \in \mathcal{K}}$ :

$$C \leq \liminf_{n \rightarrow \infty} \bar{C}_n(\mathbf{W}^n) \quad (12)$$

where  $\bar{C}_n(\mathbf{W}^n)$  denotes the following quantity:

$$\bar{C}_n(\mathbf{W}^n) := \frac{1}{n} \max_{k \in \mathcal{K}} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_k} \max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x}) \right). \quad (13)$$

An intuitive derivation of the bound (12) is described below by formulating a simplified system in a greedy approach. The formal proof using Jensen's inequality is provided in Section III-B.

### A. Intuition

Recall that the system of equations in (7) gives, for every fixed dimension  $n$ , an optimizing probability distribution  $P_{\mathbf{X}^n}^*$  for the capacity proxy  $C_n(\mathbf{W}^n)$  in (5). Since actually solving the system of equations in (7) is computationally intractable for large  $n$ , it is desirable to relax this system to a computationally tractable system that nonetheless provides a good outer bound to (7). Our starting point is the observation that an *information stable* input process  $\hat{\mathbb{X}} = \{\hat{\mathbf{X}}^n = (\hat{\mathbf{X}}_1^{(n)}, \dots, \hat{\mathbf{X}}_n^{(n)})\}_{n \geq 1}$  (with corresponding sequence of probability distributions  $\{P_{\hat{\mathbf{X}}^n} = \hat{p}_1, \dots, \hat{p}_N\}_{n \geq 1}$ ) satisfies all but an asymptotically (in  $n$ ) vanishing fraction of the constraints in (7). To see this, one may notice that by the definition of information stability (in Definition 1), for any fixed  $\gamma > 0$  it holds that there is a sufficiently large  $N_\gamma$  such that for all  $n > N_\gamma$ , in probability (over the input process  $\hat{\mathbb{X}}$  and the channel law  $\mathbb{W}$  it holds that the ratio of the information density  $\log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}^n)}{P_{\mathbf{Y}}(\mathbf{y})}$  over  $nC_n(\mathbf{W}^n)$  converges to 1. Moreover, for an arbitrary but fixed integer  $n$  and for all  $\mathbf{x} \in \mathcal{X}^n$ ,

$$\sum_{\mathbf{y} \in \mathcal{Y}} \mathbf{W}^n(\mathbf{y}|\mathbf{x}) = 1. \quad (14)$$

Therefore, for all sufficiently large  $n$ , w.h.p. using the distribution of the information stable process  $\hat{\mathbb{X}}$ , the quantity on the LHS of (7) is approximately equal to  $C_n(\mathbf{W}^n)$ . Thus any information stable input process  $\hat{\mathbb{X}}$ , for sufficiently large  $n$ , becomes a reasonable approximation of the input process  $\mathbb{X}^*$  optimizing (7). This encourages us to construct a new input process  $\bar{\mathbb{X}}$  by maximizing, for every integer  $n$ , the probability in (15) below (using a greedy approach):

$$\mathbb{P} \left\{ \left| \frac{i_{\bar{\mathbf{X}}^n, \mathbf{W}^n}(\bar{\mathbf{X}}^n; \mathbf{Y}(\mathbf{X}^n))}{nC_n(d)} \right| = 1 \right\}. \quad (15)$$

Through this process we are able to introduce such a process  $\bar{\mathbb{X}}$  (with corresponding distribution  $P_{\bar{\mathbf{X}}^n}$ ) that mimics the one for information stability in Definition 1.

1) *Approximate Information Stable Processes*: To find a system that obtains such a sub-optimal input distribution  $P_{\bar{\mathbf{X}}^n}$  efficiently, one simple heuristic method is to maximize the probability in (15) greedily.

For fixed input block-length  $n$ , we consider the set of all output sequences  $\mathbf{y}$  in the concentration set  $\mathcal{B}$ . For each  $\mathbf{y}$  in  $\mathcal{B}$ , we *greedily* choose the corresponding  $\mathbf{x} \in \mathcal{X}^n$  that maximizes the *a posteriori* probability of an instance  $\mathbf{x}$  being transmitted under the channel law  $\mathbb{W}$  (this is intuitively where the term  $\max_{\mathbf{x}} f_k(\mathbf{y}|\mathbf{x})$  comes from) that shows up in Eq. (13).<sup>6</sup>

Now, guided by the intuition in the previous paragraph about the LHS of (7) being approximately equal to  $C_n(\mathbf{W}^n)$  for many  $j$ , we *fix* the information density

$$i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{x}_j; \mathbf{y}_j) = \log \frac{\max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}_j|\mathbf{x})}{\sum_{j=1}^n \bar{p}_j^n \mathbf{W}^n(\mathbf{y}_j|\mathbf{x}_j)}$$

for each such  $(\mathbf{x}_j, \mathbf{y}_j)$  drawn in above to equal a certain constant  $\bar{\lambda}_n$  (which shows up later, in Eq. (16)).<sup>7</sup> The fixing of the information density  $i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{x}_j; \mathbf{y}_j)$  is done in a manner such that, using Bayes' rule, a probability distribution  $P_{\bar{\mathbf{X}}^n} = (\bar{p}_j^n)_{j=1, \dots, N}$  is induced on  $\mathbf{x}_j$ . In particular, the value of  $\bar{\lambda}_n$  is chosen so that the summation of  $\bar{p}_j^n$  over all  $\mathbf{x}_j$  equals 1.

2) *Simplified System*: Formally, we describe the new (as a simplified version of (7)) system as follows.

For all  $\mathbf{y} \in \mathcal{B}_k$ ,  $k \in \mathcal{K}$  and some  $\bar{\lambda}_n \geq 1$ , we let

$$\frac{1}{n} \log \frac{\max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x})}{\sum_{j=1}^n \bar{p}_j^n \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} = \bar{\lambda}_n. \quad (16)$$

As explained above, by exhausting the set  $\mathcal{B}_k$ , the constraints in (16) suggest a *greedy* approach for finding the sub-optimal distribution  $P_{\bar{\mathbf{X}}^n}$ .

Recall that (Definition 3 in Section III)  $\sum_{\mathbf{y} \in \mathcal{B}_k} f_k(\mathbf{y}|\mathbf{x}) = 1$  for all  $\mathbf{x} \in \mathcal{X}^n$  and  $k \in \mathcal{K}$ . Given an input process  $\bar{\mathbb{X}}$  satisfying Eqs. (16) for each integer  $n$ , we can rewrite the constraints in (16) as

$$\frac{\max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x})}{\sum_{j=1}^n \bar{p}_j^n \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} = 2^{n\bar{\lambda}_n}, \quad \forall \mathbf{y} \in \mathcal{A}_k.$$

Summing both sides over all  $\mathbf{y} \in \mathcal{B}_k$ ,

$$\begin{aligned} \sum_{\mathbf{y} \in \mathcal{B}_k} \frac{\max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x})}{2^{n\bar{\lambda}_n}} &= \sum_{\mathbf{y} \in \mathcal{B}_k} \sum_{j=1}^N \bar{p}_j^n \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) \\ &= \sum_{j=1}^N \bar{p}_j^n \sum_{\mathbf{y} \in \mathcal{B}_k} \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) \\ &= \sum_{j=1}^N \bar{p}_j^n = 1. \end{aligned}$$

<sup>6</sup>This procedure coincides with the greedy decoding suggested in [18] for the BDC, which can be used to derive lower bounds on  $C(d)$ . However, making use of the system (7), the greedy selection is also capable to give upper bounds.

<sup>7</sup>We do not claim to have an efficient computational process for determining this constant  $\bar{\lambda}_n$ . However, this  $\bar{\lambda}_n(d)$  has a strong operational meaning – it provides an outer bound on the capacity  $C(d)$  of the deletion channel, as discussed in Eq. (17).

Multiplying both sides with  $2^{n\bar{\lambda}_n}$  and taking logarithms,

$$\bar{C}_n(\mathbf{W}^n) := \frac{1}{n} \max_{k \in \mathcal{K}} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_k} \max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x}) \right) = \bar{\lambda}_n. \quad (17)$$

The number of constraints in (16) is much smaller than in (7), suggesting  $\bar{\lambda}_n \geq \lambda_n$  (without a proof) for each  $n$ . This indicates that the ML upper bound defined in Theorem 3 makes sense. Next we prove Theorem 1.

### B. Proof of Theorem 1

Denote by  $P_{\mathbf{X}}^*$  the optimizing probability distribution maximizing the quantity in (6). Based on Definition 2, Lemma 2 and Definition 3, we prove Theorem 1.

Considering the constraints in (7), it follows that

$$C_n(\mathbf{W}^n) = \frac{1}{n} \sum_{i=1}^M \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)} \quad (20)$$

for all  $j \in \{1, \dots, N\}$  with  $p_j^* \neq 0$ .

Now we introduce an auxiliary probability distribution  $Q_{\mathbf{X}}^n := (q_1, \dots, q_N)$  with  $q_j = 0$  once  $p_j^* = 0$  in the set  $\mathbb{P}^N$ . Multiplying both sides of (20) by  $q_j$  and summing over all  $j$ ,

$$C_n(\mathbf{W}^n) \leq \frac{1}{n} \sum_{j=1}^N \sum_{i=1}^M q_j \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}_i|\mathbf{x}_j)}.$$

Making use of the concentration set  $\mathcal{A}$  in Definition 2, we get (18) where  $\gamma_n \rightarrow 0$  as  $n \rightarrow \infty$ . Moreover, the decomposition  $\mathcal{A} = \bigcup_{k \in \mathcal{K}} \mathcal{A}_k$  (with  $\mathcal{A}_k := \mathcal{X}^n \times \mathcal{B}_k$ ) yields (19). Since logarithmic functions are concave and Eq. (10) implies

$$\sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}_k} q_j f_k(\mathbf{y}|\mathbf{x}_j) = \sum_{\mathbf{y} \in \mathcal{B}_k} \sum_{j=1}^N q_j f_k(\mathbf{y}|\mathbf{x}_j) = \sum_{j=1}^N q_j = 1,$$

applying Jensen's inequality to (19), it follows that

$$\begin{aligned} & \sum_{k \in \mathcal{K}} \sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}_k} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{f_k(\mathbf{y}|\mathbf{x}_j)} q_j f_k(\mathbf{y}|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} \\ & \leq \max_{k \in \mathcal{K}} \log \left( \sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}_k} \frac{q_j \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) f_k(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} \right) \end{aligned} \quad (21)$$

where the last inequality holds since Eq. (10) guarantees that

$$\sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k(\mathbf{y}|\mathbf{x})} \leq 1.$$

Next, we set  $q_j = p_j^*$  for all  $j \in \{1, \dots, N\}$ . The quantity inside the logarithm of (21) becomes

$$\sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}_k} \frac{p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) f_k(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} \quad (22)$$

$$= \sum_{\mathbf{y} \in \mathcal{B}_k} \frac{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) f_k(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} \quad (23)$$

$$\leq \sum_{\mathbf{y} \in \mathcal{B}_k} \max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x}). \quad (24)$$

Putting (21) and (22) into (18), for any concentration set  $\mathcal{A} = \bigcup_{k \in \mathcal{K}} \mathcal{X}^n \times \mathcal{B}_k$  and stochastic factors  $\{f_k\}_k$ ,

$$\begin{aligned} C_n(\mathbf{W}^n) & \leq \bar{C}_n(\mathbf{W}^n) \\ & := \frac{1}{n} \max_{k \in \mathcal{K}} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_k} \max_{\mathbf{x} \in \mathcal{X}^n} f_k(\mathbf{y}|\mathbf{x}) \right) + \gamma_n. \end{aligned} \quad (25)$$

Note that the term  $\gamma_n$  is vanishing (in  $n$ ). Hence, for information stable channels, the general formula in (4) implies that

$$C \leq \liminf_{n \rightarrow \infty} \bar{C}_n(\mathbf{W}^n).$$

This completes the proof of Theorem 1.  $\square$

### C. Verification of the Tightness for the BEC and BSC

This section is devoted to verifying the tightness of the upper bound in Theorem 1 on the BEC and the BSC. Denote by  $p \in (0, 1)$  the erasure/bit-flip probability. Note that under the settings of the BEC( $p$ ) and BSC( $p$ ), we have the following realizations of the input and output spaces:

$$\mathcal{X}_{\text{BEC}}^n = \mathcal{X}_{\text{BSC}}^n := \{0, 1\}^n$$

and

$$\begin{aligned} \bar{\mathcal{Y}}_{\text{BEC}} & := \{0, 1, E\}^n, \\ \bar{\mathcal{Y}}_{\text{BSC}} & := \{0, 1\}^n. \end{aligned}$$

Denote by  $d_E(\cdot)$  the number of erasures (marked as  $E$ ) in a given the Hamming distance and  $d_H(\cdot, \cdot)$  the Hamming distance. We select the following concentration sets (and the corresponding decompositions) for the two types of channels respectively.

**Definition 4.** We consider the following concentration sets denoted by  $\mathcal{A}^{\text{BEC}}$  and  $\mathcal{A}^{\text{BSC}}$  for the BEC and BSC respectively according to Definition 2:

$$\mathcal{A}^{\text{BEC}} := \bigcup_{k \in \mathcal{K}_\varepsilon} \mathcal{A}_k^{\text{BEC}}, \quad (26)$$

$$\mathcal{A}^{\text{BSC}} := \bigcup_{k \in \mathcal{K}_\varepsilon} \mathcal{A}_k^{\text{BSC}} \quad (27)$$

where  $\mathcal{A}_k^{\text{BEC}}$  and  $\mathcal{A}_k^{\text{BSC}}$  are defined as follows in agreement with Definition 3:

$$\mathcal{K}_\varepsilon := \lfloor n(p - \varepsilon), n(p + \varepsilon) \rfloor,$$

$$\mathcal{A}_k^{\text{BEC}} := \left\{ (\mathbf{x}, \mathbf{y}) \in \{0, 1\}^n \times \{0, 1, E\}^n : d_E(\mathbf{y}) = k \right\}, \quad (28)$$

$$\mathcal{A}_k^{\text{BSC}} := \left\{ (\mathbf{x}, \mathbf{y}) \in \{0, 1\}^n \times \{0, 1\}^n : d_H(\mathbf{x}, \mathbf{y}) = k \right\} \quad (29)$$

and the sets  $\mathcal{B}_k^{\text{BEC}}$  and  $\mathcal{B}_k^{\text{BSC}}$  are defined as

$$\mathcal{B}_k^{\text{BEC}} := \left\{ \mathbf{y} \in \{0, 1, E\}^n : d_E(\mathbf{y}) = k \right\}, \quad (30)$$

$$\mathcal{B}_k^{\text{BSC}} := \{0, 1\}^n \quad (31)$$

with some  $\varepsilon$  satisfying  $\min\{p, 1 - p\} > \varepsilon > 0$ .

**Lemma 3.** The concentration sets  $\mathcal{A}^{\text{BEC}}$  and  $\mathcal{A}^{\text{BSC}}$  defined in above satisfy the conditions (8)-(9) in Definition 2.

$$C_n(\mathbf{W}^n) \leq \frac{1}{n} \sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}} q_j \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} + \gamma_n. \quad (18)$$

$$\sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}} q_j \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)} \leq \sum_{k \in \mathcal{K}} \sum_{(\mathbf{x}_j, \mathbf{y}) \in \mathcal{A}_k} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{f_k(\mathbf{y}|\mathbf{x}_j)} q_j f_k(\mathbf{y}|\mathbf{x}_j) \log \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}{\sum_{j=1}^N p_j^* \mathbf{W}^n(\mathbf{y}|\mathbf{x}_j)}. \quad (19)$$

*Proof.* For the BEC( $p$ ), since each bit of the length- $n$  input sequence is erased i.i.d., the number of erased bits is distributed according to Bernoulli( $n, p$ ). Therefore, fix input length  $n$ , by the Chernoff bound (see [19, 20]), source  $\mathbb{X}$ , the probability of the output sequence being inside the concentration set  $\mathcal{A}^{\text{BEC}}$  defined in (26) is always bounded from below by  $1 - 2 \exp(-np\varepsilon^2/3)$ . Moreover, the information densities corresponding to the outliers must be bounded. Thus, the concentration set  $\mathcal{A}^{\text{BEC}}$  satisfies the condition in (9). For the BSC, we have a trivial decomposition  $\mathcal{A}^{\text{BSC}} = \{0, 1\}^n \times \bar{\mathcal{Y}}^{\text{BSC}}$ , automatically guarantees the condition in (9).  $\square$

Furthermore, we associate the following stochastic factors to both channels.

**Definition 5.** The stochastic factors for the BEC and BSC are defined as follows.

$$f_k^{\text{BEC}}(\mathbf{x}, \mathbf{y}) := \begin{cases} 1/\binom{n}{k} & \text{if } d_E(\mathbf{y}) = k \\ 0 & \text{otherwise} \end{cases} \quad (32)$$

and

$$f_k^{\text{BSC}}(\mathbf{x}, \mathbf{y}) := \begin{cases} 1/\binom{n}{k} & \text{if } d_H(\mathbf{x}, \mathbf{y}) = k \\ 0 & \text{otherwise} \end{cases} \quad (33)$$

with  $1 \leq k \leq n$  satisfying

$$\lfloor n(p - \varepsilon) \rfloor \leq k \leq \lceil n(p + \varepsilon) \rceil \quad (34)$$

for some  $\varepsilon \in (0, \min\{p, 1 - p\})$ .

**Lemma 4.** The stochastic factors  $f_k^{\text{BEC}}(\mathbf{x}, \mathbf{y})$  and  $f_k^{\text{BSC}}(\mathbf{x}, \mathbf{y})$  defined in above satisfy the conditions (10)-(11) in Definition 3.

*Proof.* We check the stochastic factors defined above satisfy the conditions in (10)-(11). For the BEC( $p$ ), plugging in (28) and (32),

$$\sum_{\mathbf{y} \in \mathcal{B}_k^{\text{BEC}}} f_k^{\text{BEC}}(\mathbf{y}|\mathbf{x}) = \frac{|\mathcal{B}_k^{\text{BEC}}|}{2^{n-k}} \cdot \frac{1}{\binom{n}{k}} = 1, \quad \forall \mathbf{x} \in \mathcal{X}^n, k \in \mathcal{K}.$$

Since  $\mathbf{W}^n(\mathbf{y}|\mathbf{x}) = p^k (1-p)^{n-k}$  for all  $(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BEC}}$ ,

$$\begin{aligned} & \sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BEC}}} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k^{\text{BEC}}(\mathbf{y}|\mathbf{x})} \\ &= \sum_{k=\lfloor n(p-\varepsilon) \rfloor}^{\lceil n(p+\varepsilon) \rceil} \binom{n}{k} p^k (1-p)^{n-k} \in [0, 1] \end{aligned}$$

showing that  $f_k^{\text{BEC}}(\mathbf{x}, \mathbf{y})$  are stochastic factors.

For the the BSC( $p$ ), similarly, based on the definitions in (33), since for each fixed  $\mathbf{x} \in \{0, 1\}^n$ , there are in total  $\binom{n}{k}$

many  $\mathbf{y} \in \mathcal{B}_k^{\text{BSC}} = \{0, 1\}^n$  satisfying  $d_H(\mathbf{x}, \mathbf{y}) = k$ , it follows that

$$\sum_{\mathbf{y} \in \mathcal{B}_k^{\text{BSC}}} f_k^{\text{BSC}}(\mathbf{y}|\mathbf{x}) = \binom{n}{k} \cdot \frac{1}{\binom{n}{k}} = 1, \quad \forall \mathbf{x} \in \mathcal{X}^n, k \in \mathcal{K}.$$

Moreover,  $\mathbf{W}^n(\mathbf{y}|\mathbf{x}) = p^k (1-p)^{n-k}$  for all  $(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BSC}}$ , clearly,

$$0 \leq \sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BSC}}} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k^{\text{BSC}}(\mathbf{y}|\mathbf{x})}. \quad (35)$$

From the definition of the set  $\mathcal{A}_k^{\text{BSC}}$ ,

$$\begin{aligned} & \sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BSC}}} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k^{\text{BSC}}(\mathbf{y}|\mathbf{x})} \\ &= \sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}): d_H(\mathbf{x}, \mathbf{y})=k} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k^{\text{BSC}}(\mathbf{y}|\mathbf{x})} \\ &= \sum_{k=\lfloor n(p-\varepsilon) \rfloor}^{\lceil n(p+\varepsilon) \rceil} \binom{n}{k} p^k (1-p)^{n-k} \in [0, 1] \end{aligned}$$

verifying that the stochastic factors in (32)-(33) satisfy the conditions in Definition 2 and Definition 3.  $\square$

Based on the concentration sets  $\mathcal{A}^{\text{BEC}}$  and  $\mathcal{A}^{\text{BSC}}$  defined in (26)-(31) and the stochastic factors in (32)-(33), the ML upper bound in Theorem 1 is tight, as stated in the following theorem.

**Theorem 2** (Tightness for the BEC and BSC). Let  $p \in (0, 1)$  be the erasure/bit-flip probability. The ML upper bound in Theorem 1 is tight for the BEC and BSC, i.e.,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \max_{k \in \mathcal{K}_\varepsilon} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_k^{\text{BEC}}} \max_{\mathbf{x} \in \mathcal{X}_{\text{BEC}}^n} f_k^{\text{BEC}}(\mathbf{y}|\mathbf{x}) \right) = 1 - p,$$

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \max_{k \in \mathcal{K}_\varepsilon} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_k^{\text{BSC}}} \max_{\mathbf{x} \in \mathcal{X}_{\text{BSC}}^n} f_k^{\text{BSC}}(\mathbf{y}|\mathbf{x}) \right) = 1 - h(p)$$

where  $h(p) := -p \log p - (1-p) \log(1-p)$  denotes the binary entropy.

*Proof.* Note that the parameter  $\varepsilon > 0$  can be arbitrarily small, taking  $k = \lceil np \rceil$  and applying Theorem 1, the capacity of the BEC denoted by  $C_{\text{BEC}}(p)$  satisfies

$$\begin{aligned} C_{\text{BEC}}(p) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_{\lceil np \rceil}^{\text{BEC}}} \max_{\mathbf{x} \in \mathcal{X}_{\text{BEC}}^n} f_{\lceil np \rceil}^{\text{BEC}}(\mathbf{y}|\mathbf{x}) \right) \\ &= \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \left| \mathcal{B}_{\lceil np \rceil}^{\text{BEC}} \right| / \binom{n}{\lceil np \rceil} \right). \end{aligned}$$

Putting  $|\mathcal{B}_{\lceil np \rceil}^{\text{BEC}}| = 2^{n(1-p)} \binom{n}{\lceil np \rceil}$  into above,

$$\begin{aligned} C_{\text{BEC}}(p) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( 2^{n(1-p)} \binom{n}{\lceil np \rceil} / \binom{n}{\lceil np \rceil} \right) \\ &= 1 - p. \end{aligned} \quad (36)$$

Furthermore, for the BSC, the capacity  $C_{\text{BSC}}(p)$  is bounded from above as

$$\begin{aligned} C_{\text{BSC}}(p) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_{\lceil np \rceil}^{\text{BSC}}} \max_{\mathbf{x} \in \mathcal{X}_{\text{BSC}}^n} f_{\lceil np \rceil}^{\text{BSC}}(\mathbf{y}|\mathbf{x}) \right) \\ &= \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( |\mathcal{B}_{\lceil np \rceil}^{\text{BSC}}| / \binom{n}{\lceil np \rceil} \right). \end{aligned}$$

Since  $|\mathcal{B}_{\lceil np \rceil}^{\text{BSC}}| = 2^n$ , it follows that

$$C_{\text{BSC}}(p) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( 2^n / \binom{n}{\lceil np \rceil} \right) = 1 - h(p). \quad (37)$$

□

The theorem above indicates that the ML upper bound in Theorem 1 is actually tight for some memoryless channels, e.g., the BEC( $p$ ) and BSC( $p$ ). In next section, we analyze the binary deletion channel (BDC) as an example for channels with memory, and show that the ML upper bound is capable of providing a nontrivial and explicit approximation for the capacity  $C(d)$  ( $d$  denotes the deletion probability) assuming Hypothesis 2 in Section IV-C.

#### IV. BINARY DELETION CHANNEL

For brevity, we consider specifically the binary deletion channel (BDC), though the approach generalizes to arbitrary alphabet sizes. The following section is devoted to summarizing related work on finding the capacity  $C(d)$  of the BDC. In particular, we focus on the existing upper bounds on  $C(d)$  and the known asymptotic results. The survey by Mitzenmacher [18] elucidates critical problems, useful techniques and further applications in a more comprehensive way. The recent paper by Cheraghchi also provides a decent summary of the state-of-the-art literature [21]. Before moving to the contexts, we first give a brief summary of known bounds on  $C(d)$ , the capacity of the BDC, at the risk of missing much of the literature.

##### A. Previous Work

1) *Existing Upper Bounds*: Recently, Cheraghchi in [21] gave an explicit and concise upper bound on  $C(d)$  such that  $C(d) \leq 1 - d \log(4/\phi)$  for  $d < 1/2$  and  $C(d) \leq (1-d) \log \phi$  for  $d \geq 1/2$  where  $\phi := (1 + \sqrt{5})/2$  is the golden ratio. The bound was obtained by first deriving an upper bound on  $C(1/2)$ ; then applying the fact that  $C(d)$  is convex as showed in [22].

Running the BAA (cf. [8, 9]) up to  $n = 17$ , tighter numerical upper bounds were provided in [10] improving the upper bounds in [23] for a wide range of the deletion probability  $d$ . They proved that increasing the dimension  $n$  in the BAA always provides a better upper bound on  $C(d)$ . The convexity of  $C(d)$  in [22] can be used to tighten the bounds in [10]. However, the space complexity of BAA is exponential in  $n$ , prohibiting

obtaining better bounds by trying larger dimensions. Following previous literature ([24]), we sometimes replace the maximized finite-length mutual information  $C_n(\mathbf{W}^n)$  with  $C_n(d)$  since the quantity is determined entirely by  $d$  and  $n$  in particular for deletion channels. Denote by  $C_{n,T}^{\text{BAA}}(d)$  the approximation of  $C_n(d)$  using the BAA with  $T$  iterations. *A priori*, if one could use the BAA to obtain the value of  $C_n(d)$  precisely, one would be able to get a  $\log(n+1)/n$  additive approximation to the capacity  $C(d)$  (for  $n \leq 17$ ). In [10], essentially this approach is followed to obtain numerical solutions for (7), for  $n \leq 17$ .

However, this approach has several limitations. On the one hand, in direct implementation of BAA, as  $n$  grows, it becomes computationally intractable even to store the variables to be computed. On the other hand, as BAA is itself an iterative algorithm attempting to solve the non-convex optimization problem (7), and to the best of our knowledge there are no guarantees on how quickly  $C_{n,T}^{\text{BAA}}(d)$  converges to  $C_n(d)$  as a function of the number of iterations  $t$ .

In fact, in [7], Dobrushin showed the following quantitative bound on  $C_n(d)$  providing a  $\log(n+1)/n$ -additive approximation to  $C(d)$  (see also, [24, 25]. A tighter bound can be found in [26]):

$$C_n(d) - \frac{\log(n+1)}{n} \leq C(d) \leq C_n(d). \quad (38)$$

2) *Asymptotic Results*: Besides the upper bounds, in [25] Kanoria and Montanari gave a polynomial expression of channel capacity which is optimal with a residual term  $O(d^{3-\varepsilon})$  through the optimality of sources with i.i.d. coordinates distributed as Bernoulli ( $1/2$ ) when  $d = 0$  and a perturbed version when  $d$  is slightly larger than 0. Meanwhile for the regime  $d \rightarrow 1$ , Dalai in [24] provided an asymptotic result  $\lim_{d \rightarrow 1} C(d)/(1-d) \leq 0.4143$  with constant surpassing the upper bound 0.49 given in [10] by Fertonani and Duman. Furthermore, [27] studied the mutual information for deletion channels concerning both general sources and memoryless sources. In particular, their results for memoryless sources coincide with Kanoria and Montanari's expansion of mutual information [25] as  $d \rightarrow 0$ .

In the sequel, we apply the ML upper bound in Theorem 1 to derive approximations for the capacity  $C(d)$  of the BDC.

##### B. Capacity Upper Bound via Theorem 1

1) *Concentration Set and Stochastic Factors*: The input and output spaces of the BDC are

$$\mathcal{X}_{\text{BDC}}^n := \{0, 1\}^n$$

and

$$\bar{\mathcal{Y}}_{\text{BDC}} := \bigcup_{1 \leq m \leq n} \{0, 1\}^m.$$

We consider the following concentration set and the corresponding decompositions for the BDC:

**Definition 6** (Concentration Set for the BDC). *For the BDC, we define the following concentration set according to Definition 2:*

$$\mathcal{A}^{\text{BDC}} := \bigcup_{k \in \mathcal{K}} \mathcal{A}_k^{\text{BDC}} \quad (39)$$

with the following decomposition in agreement with Definition 3:

$$\begin{aligned}\mathcal{K}_\varepsilon &:= \lfloor n(d - \varepsilon), n(d + \varepsilon) \rfloor, \\ \mathcal{A}_k^{\text{BDC}} &:= \{0, 1\}^n \times \{0, 1\}^{n-k}, \\ \mathcal{B}_k^{\text{BDC}} &:= \{0, 1\}^{n-k}.\end{aligned}\quad (40)$$

As  $n$  grows without bound, standard concentration inequalities (for instance, Chernoff bound) imply that the length of the output sequence  $m$  is tightly concentrated around the “typical length”  $n(1 - d)$  in probability. Using the bounds in (38), [24] showed that  $C(d)$  is continuous (the continuity can be verified via various approaches, for instance, the information spectrum method [28]). Leveraging the continuity, Theorem 1 in [24] (cf. [10, 25]) proved that it is sufficient to consider output sequences with lengths in  $\mathcal{K}_\varepsilon = \lfloor n(d - \varepsilon), n(d + \varepsilon) \rfloor$ . In detail, Dalai showed that the following lemma holds.

**Lemma 5** (Theorem 1 [24]). *For any  $0 < \varepsilon \leq \min\{d, 1 - d\}$ ,*

$$\begin{aligned}&\liminf_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[ i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n)) \mid (\mathbf{X}^n, \mathbf{Y}) \in \mathcal{A}^{\text{BDC}} \right] \\ &= \liminf_{n \rightarrow \infty} C_n(\mathbf{W}^n) = C(d).\end{aligned}$$

The lemma above validates that  $\mathcal{A}^{\text{BDC}}$  is a concentration set according to Definition 2.

**Lemma 6.** *The concentration set  $\mathcal{A}^{\text{BDC}}$  defined in above satisfies the conditions (8)-(9) in Definition 2.*

*Proof.* Since by definition the concentration set  $\mathcal{A}^{\text{BDC}}$  contains all received codewords with lengths in the range  $\lfloor n(d - \varepsilon), n(d + \varepsilon) \rfloor$ , standard concentration inequalities (e.g. Chernoff bound) guarantees the condition in (8), and

$$\liminf_{n \rightarrow \infty} \mathbb{P}((\mathbf{X}^n, \mathbf{Y}) \in \mathcal{A}^{\text{BDC}}) = 1.$$

Moreover, using Lemma 5 above, directly,

$$\limsup_{n \rightarrow \infty} \mathbb{E} \left[ \left| \frac{i_{\mathbf{X}^n, \mathbf{W}^n}(\mathbf{X}^n; \mathbf{Y}(\mathbf{X}^n))}{nC_n(\mathbf{W}^n)} - 1 \right| \mid (\mathbf{X}^n, \mathbf{Y}) \in \mathcal{A}^{\text{BDC}} \right] = 0.$$

□

Before proceeding to the corresponding stochastic factors, it is helpful to introduce a quantity pertinent to relationships between length- $n$  input sequences  $x$  and length- $m$  output sequences  $y$ .

**Definition 7** (Deletion Pattern). *A deletion pattern from a length- $n$  input sequence  $\mathbf{x}$  to a length- $m$  output sequence  $\mathbf{y}$  is a binary sequence denoted by  $\mathbf{d} \in \{0, 1\}^n$ . If the  $i$ -th coordinate  $d_i = 1$ , the corresponding  $i$ -th coordinate  $x_i$  in  $\mathbf{x}$  is deleted; otherwise  $d_i = 0$  implies that  $x_i$  is kept.*

Below we define a function computing the number of possible ways that a fixed length- $n$  sequence is deleted to form a fixed shorter length- $m$  sequence. Previous studies have been focusing on similar quantities, Drmota *et al.* defined a similar quantity as the number of occurrences of a shorter sequence in a longer sequence, Liron and Langberg characterized the number of subsequences obtained from a fixed length- $n$  sequence via deletions [27, 29], to name just a few.

**Definition 8** (Number of Deletion Patterns<sup>8</sup>). *We define the number of deletion patterns as a quantity  $\#(\mathbf{x}, \mathbf{y}) \in \{1, \dots, \binom{n}{m}\}$  counting the number of distinct deletion patterns from an input  $\mathbf{x} \in \{0, 1\}^n$  to an output  $\mathbf{y} \in \{0, 1\}^m$ .*

Over the years it has been repeatedly noted that the number of deletion patterns plays an important role in finding the capacity  $C(d)$ . Part of the reason is that the number of deletion patterns can be regarded as a “normalized version” of the transition probability  $\mathbf{W}^n(\mathbf{y}|\mathbf{x})$ , as the following remark explains.

**Remark 1.** *Denote by  $\mathbf{W}^n(\cdot|\cdot)$  the corresponding stochastic matrix for the BDC (with block length  $n$ ). Since the probability of a particular deletion pattern of weight  $n - m$  occurring equals  $(1 - d)^m d^{n-m}$ , hence*

$$\#(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{(1 - d)^m d^{n-m}}. \quad (41)$$

The fact that  $\#(\mathbf{y}, \mathbf{x})$  is a scaled version of  $\mathbf{W}^n(\mathbf{y}|\mathbf{x})$  (corresponding to a conditional probability distribution for the input  $\mathbf{x}$  being mapped to the fixed output  $\mathbf{y}$ ) takes on operational significance later. For instance, it ensures us to define stochastic functions. In Section IV-D, on the other hand, we utilize the operational meaning of  $\#(\mathbf{y}, \mathbf{x})$  to derive explicit bounds on  $C(d)$ .

**Definition 9.** *The corresponding stochastic factors for the BDC are set to be*

$$f_k^{\text{BDC}}(\mathbf{x}, \mathbf{y}) := \begin{cases} \#(\mathbf{x}, \mathbf{y}) / \binom{n}{k} & \text{if } \mathbf{y} \in \{0, 1\}^{n-k} \\ 0 & \text{otherwise} \end{cases}. \quad (42)$$

It remains to check the validity of the stochastic factors  $f_k^{\text{BDC}}(\mathbf{x}, \mathbf{y})$ . We first claim that they satisfy the definition of stochastic factors.

**Lemma 7.** *The stochastic factors  $f_k^{\text{BEC}}(\mathbf{x}, \mathbf{y})$  and  $f_k^{\text{BSC}}(\mathbf{x}, \mathbf{y})$  defined in above satisfy the conditions (10)-(11) in Definition 3.*

*Proof.* Plugging in (40) and (42),

$$\begin{aligned}\sum_{\mathbf{y} \in \mathcal{B}_k^{\text{BDC}}} f_k^{\text{BDC}}(\mathbf{y}|\mathbf{x}) &= \sum_{\mathbf{y} \in \{0, 1\}^{n-k}} \frac{\#(\mathbf{x}, \mathbf{y})}{\binom{n}{k}} = 1, \\ \forall \mathbf{x} \in \mathcal{X}^n, k \in \mathcal{K}_\varepsilon.\end{aligned}\quad (43)$$

Considering (41),

$$\sum_{k \in \mathcal{K}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{A}_k^{\text{BDC}}} \frac{\mathbf{W}^n(\mathbf{y}|\mathbf{x})}{f_k^{\text{BDC}}(\mathbf{y}|\mathbf{x})} \quad (44)$$

$$= \sum_{k=\lfloor n(p-\varepsilon) \rfloor}^{\lfloor n(p+\varepsilon) \rfloor} \binom{n}{k} p^k (1-p)^{n-k} \in [0, 1] \quad (45)$$

showing that  $f_k^{\text{BDC}}(\mathbf{x}, \mathbf{y})$  are stochastic factors. □

Making use of the concentration set and stochastic factors constructed in (40) and (42), and substituting them into the ML upper bound in Theorem 1, the following upper bound on  $C(d)$  holds.

<sup>8</sup>This counting function is alternatively called *hidden pattern matching function* in [30] using the terminology in statistics.

**Theorem 3.** For all  $d \in (0, 1)$  and  $m := \lceil n(1-d) \rceil$ , the capacity  $C(d)$  of the BDC is bounded from above by

$$C(d) \leq \bar{C}(d) := \liminf_{n \rightarrow \infty} \bar{C}_n(d) - h(d) \quad (46)$$

where  $h(d) = -d \log d - (1-d) \log(1-d)$  denotes the binary entropy and

$$\bar{C}_n(d) := \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \max_{\mathbf{x} \in \mathcal{X}_{\text{BDC}}^n} \#(\mathbf{x}, \mathbf{y}) \right). \quad (47)$$

*Proof.* Since  $0 < \varepsilon \leq \min\{d, 1-d\}$  is arbitrary, taking the maximizing  $k = m := \lceil n(1-d) \rceil$  and applying Theorem 1, the capacity of the BDC satisfies

$$C(d) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \mathcal{B}_m^{\text{BDC}}} \max_{\mathbf{x} \in \mathcal{X}_{\text{BDC}}^n} f_m^{\text{BDC}}(\mathbf{y}|\mathbf{x}) \right).$$

Applying the definitions of the sets  $\mathcal{B}_m^{\text{BDC}}$  and the stochastic factors  $f_m^{\text{BDC}}$  (in Definition 6 and 9),

$$\begin{aligned} C(d) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \max_{\mathbf{x} \in \{0,1\}^n} \frac{\#(\mathbf{x}, \mathbf{y})}{\binom{n}{m}} \right) \quad (48) \\ &= \liminf_{n \rightarrow \infty} \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \max_{\mathbf{x} \in \{0,1\}^n} \#(\mathbf{x}, \mathbf{y}) \right) \\ &\quad + \liminf_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{\binom{n}{m}} \\ &= \bar{C}_n(d) - h(d), \end{aligned}$$

which gives the desired bound in (46).  $\square$

**Remark 2.** Note that for any  $\mathbf{x} \in \{0,1\}^n$  and  $\mathbf{y} \in \{0,1\}^m$ , it always holds that

$$0 \leq \#(\mathbf{x}, \mathbf{y}) \leq \binom{n}{m}.$$

Putting this into (48), we recover the trivial upper bound  $C(d) \leq 1-d$ .

2) *Experimental Results:* We implement the BAA up to  $n = 18$  to compare  $\bar{C}_n(d)$ ,  $C_n(d)$  and the best known numerical bounds on  $C(d)$ . See Fig. 1.

In next section we make use of the operational meaning of the number of deletion patterns and analyze the quantity  $\bar{C}_n(d)$  in a more careful way. This allows us to derive upper bounds on  $C(d)$  based on an approximation ratio of a combinatorial problem defined later in Section IV-B.

### C. Maximal Number of Deletion Patterns

The remaining context of this paper is dedicated to approximate the terms in (48), which is summarized as the following combinatorial problem.

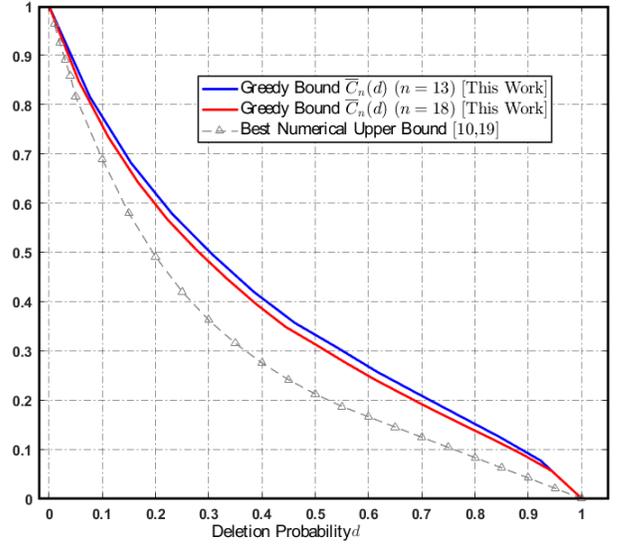


Fig. 1. The ML upper bounds (solid, blue and red)  $\bar{C}_n(d)$  for BDC with block-length  $n = 18$  and  $n = 13$ , together with the (convexified) numerical estimate of the capacity-proxy  $C_n(d)$  (dashed and marked black) for  $n = 17$ . The lower curve (dashed gray) is also known as the best numerical upper bounds provided in [22].

#### 1) The Maximum Deletion Matching problem:

**Definition 10** (MDM Problem). The maximum deletion matching (MDM) problem is to solve the following. Let  $n, m \in \mathbb{N}^+$  with  $1 \leq m \leq n$ . Given an arbitrary length- $m$  binary sequence  $\mathbf{y} \in \{0,1\}^m$ , the goal is to find the maximum corresponding length- $n$  binary sequence  $\mathbf{x}$  such that

$$\bar{\mathbf{x}}(\mathbf{y}) := \operatorname{argmax}_{\mathbf{x} \in \{0,1\}^n} \#(\mathbf{x}, \mathbf{y})$$

where  $\#(\mathbf{x}, \mathbf{y})$  denotes the number of deletion patterns of generating  $\mathbf{y}$  from  $\mathbf{x}$  defined in Definition 8. For notational convenience, write the maximal number of deletion patterns

$$\bar{\#}(\mathbf{y}) := \max_{\mathbf{x} \in \{0,1\}^n} \#(\mathbf{x}, \mathbf{y}) = \#(\bar{\mathbf{x}}(\mathbf{y}), \mathbf{y}).$$

2) *Run-length Representation:* One way to approximate  $\bar{\#}(\mathbf{y})$  and find approximation ratio of the MDM problem is to consider consecutive bits in  $\mathbf{y}$  as a “run” and jointly a distribution of run-lengths. Although encoding each sequence  $\mathbf{y}$  to the run-length representation suffers a loss of information (for instance, the ordering of runs is no longer kept in the run-length representation), it offers a concise approach to describe a binary sequence. Reprising the definitions from previous work [25, 29, 31], we consider the follows.

First, we associate each length- $m$  binary sequence  $\mathbf{y}$  with an integer sequence embedding the information of number of consecutive bits of  $\mathbf{y}$ . We call a subsequence  $y_i, \dots, y_{i+l-1}$  of  $\mathbf{y}$  an  $\ell$ -run if all the bits in the subsequence are the same and they differ from the bits next to the subsequence, i.e.,  $y_i \neq y_{i-1}, y_{i+l-1} = \dots = y_i$  and  $y_{i+l} \neq y_{i+l-1}$ . Thus, let  $R_{\mathbf{y}}^{(\ell)}$  be an integer counting the number of  $\ell$ -runs in  $\mathbf{y}$ . It

$\mathbf{y}$	Exact		Approximation		Duplication Ratio
	$\mathbf{x}$	$\#(\mathbf{y})$	$\mathbf{x}_{\text{dup}}$	$\#_{\text{dup}}(\mathbf{y})$	
0000	00000000	70	00000000	70	1
0001	00000011	40	00000011	40	1
0010	00001100	24	00001100	24	1
0011	00001111	36	00001111	36	1
0100	00110000	24	00110000	24	1
0101	<b>00101011</b>	16	<b>00110011</b>	16	1
0110	00111100	24	00111100	24	1
0111	00111111	40	00111111	40	1

TABLE I

A TABLE SHOWING THE CORRESPONDING RATIOS GIVEN  $m = 4$  AND  $n = 8$ . WHEN  $m \leq 4$ , THE RATIO IS ALWAYS 1.

$\mathbf{y}$	Exact		Approximation		Duplication Ratio
	$\mathbf{x}$	$\#(\mathbf{y})$	$\mathbf{x}_{\text{dup}}$	$\#_{\text{dup}}(\mathbf{y})$	
0000000	00000000000000	3432	00000000000000	3432	1
0000001	00000000000011	1848	00000000000011	1848	1
0000010	000000000001100	1008	000000000001100	1008	1
0000011	000000000001111	1512	000000000001111	1512	1
0000100	00000000110000	840	00000000110000	840	1
0000101	<b>00000000101011</b>	602	<b>00000000110011</b>	560	0.93023.
0000110	00000000111100	840	00000000111100	840	1
0000111	00000000111111	1400	00000000111111	1400	1
0000111	00000000111111	1400	00000000111111	1400	1
⋮	⋮	⋮	⋮	⋮	⋮
0001010	<b>00000000111111</b>	396	<b>00000011001100</b>	320	0.80808.
0001011	<b>00000010101111</b>	530	<b>00000011001111</b>	480	0.90566.
0010100	<b>00001010101000</b>	351	<b>00001100110000</b>	288	0.90566.
0010101	<b>00001010101011</b>	270	<b>00001100110011</b>	192	0.71111.
0010110	<b>00001010111100</b>	312	<b>00001100111100</b>	288	0.92308.
0010111	<b>00001010111111</b>	530	<b>00001100111111</b>	480	0.90566.
0011010	<b>00001111010100</b>	300	<b>00001111001100</b>	288	0.96
0100101	<b>00110000101011</b>	200	<b>00110000110011</b>	192	0.96
0101000	<b>01010101000000</b>	396	<b>00110011000000</b>	320	0.80808.
0101001	<b>01010101000011</b>	231	<b>00110011000011</b>	192	0.83117.
0101010	<b>00101010101010</b>	204	<b>00110011001100</b>	128	<b>0.62745.</b>
0101011	<b>00101010101111</b>	270	<b>00110011001111</b>	192	0.71111.
0101100	<b>00101011110000</b>	300	<b>00110011110000</b>	288	0.96
0101101	<b>00101011110011</b>	200	<b>00110011110011</b>	192	0.96
0101110	<b>00101011111100</b>	340	<b>00110011111100</b>	320	0.94118.
0101111	<b>00101011111111</b>	602	<b>00110011111111</b>	560	0.93023.
0110100	<b>00111101010000</b>	312	<b>00111100110000</b>	280	0.89743.
0110101	<b>00111101010101</b>	231	<b>00111100110011</b>	192	0.83117.
0111010	<b>00111110101010</b>	340	<b>00111111001100</b>	320	0.94118.

TABLE II

AN INCOMPLETE LIST COMPARING THE MAXIMAL NUMBER OF DELETION PATTERNS WITH THE DUPLICATION APPROXIMATIONS FOR  $m = 7$  AND  $n = 14$ . THE TABLE CONTAINS ALL RATIOS THAT ARE NOT EQUAL TO ONE. THE SMALLEST DUPLICATION RATIO  $\gamma(n, F)$  IS OBTAINED AT THE "FLIPPING SEQUENCE"  $\mathbf{y}_{\text{flip}} = 0101010$ .

follows that

$$\sum_{\ell=1}^m \ell R_{\mathbf{y}}^{(\ell)} = m \quad (49)$$

for all  $\mathbf{y} \in \{0, 1\}^m$ .

3) *Approximation and Duplication Ratio:* Figuring out the the maximal number of deletion patterns  $\#(\mathbf{y})$  as an explicit expression using the number of  $\ell$ -runs  $R_{\mathbf{y}}^{(1)}, \dots, R_{\mathbf{y}}^{(m)}$  is a nontrivial task. Instead, one might turn to consider approximations of the quantity  $\#(\mathbf{y})$ . Suppose  $F := n/m$

is an integer. Intuitively, duplicating  $F$  times each bit in  $\mathbf{y}$  may provide a decent estimate of  $\#(\mathbf{y})$ , which motivates the following definition.

**Definition 11.** Suppose  $m$  divides  $n$ . Denoted by  $F := n/m \in \mathbb{N}^+$ . We define the duplication ratio  $\gamma(n, F, \mathbf{y}) \leq 1$  of the MDM problem to be the ratio of the approximated number of deletion patterns by duplicating each bit  $F$  times in  $\mathbf{y}$  and the

maximal number of deletion patterns  $\overline{\#}(\mathbf{y})$ :

$$\gamma(n, F, \mathbf{y}) := \frac{\overline{\#}_{\text{dup}}(\mathbf{y})}{\overline{\#}(\mathbf{y})}. \quad (50)$$

where  $\overline{\#}_{\text{dup}}(\mathbf{y})$  is given by

$$\overline{\#}_{\text{dup}}(\mathbf{y}) := \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\mathbf{y}}^{(\ell)}}. \quad (51)$$

Note that  $\overline{\#}_{\text{dup}}(\mathbf{y})$  in (51) equals to the number of deletion patterns of the length- $m$  binary sequence  $\mathbf{y}$  in the length- $n$  sequence  $\mathbf{x}_{\text{dup}}(\mathbf{y})$  by setting

$$\mathbf{x}_{\text{dup}}(\mathbf{y}) = \underbrace{y_1, y_1, \dots, y_1}_{F \text{ many}}, \dots, \underbrace{y_m, y_m, \dots, y_m}_{F \text{ many}}.$$

4) *Numerical Results:* We compute the duplication ratios with different block-lengths  $n$  and  $m$ . We exemplify part of the results in Table I and Table II. Furthermore, setting  $n = 18$ , we plot the following quantities:

$$\begin{aligned} \overline{C}_n(d) &:= \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^{\lceil n(1-d) \rceil}} \max_{\mathbf{x} \in \{0,1\}^n} \#(\mathbf{x}, \mathbf{y}) \right), \\ \tilde{C}_n(d) &:= \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^{\lceil n(1-d) \rceil}} \max_{\mathbf{x} \in \{0,1\}^n} \#_{\text{dup}}(\mathbf{x}, \mathbf{y}) \right). \end{aligned}$$

The output block-length  $m$  in above is set to be an integer between 1 and  $n = 18$ . In order to approximate the values of  $\overline{C}_n(d)$  and  $\tilde{C}_n(d)$  when  $m$  is not divided by  $n$ , we consider the following three different approaches:

*Approach 1 (Assign-to-the-last Approximation):*

First, duplicate  $F$  times each bit in  $\mathbf{y}$  where  $F$  is the largest integer satisfying  $mF \leq n$ ; then for the remaining  $n - mF$  bits, assign them proportionally to the last several runs in  $\mathbf{y}$ . For instance, suppose  $m = 6$ ,  $n = 15$  and  $\mathbf{y} = 010001$ . An approximation can be obtained by first constructing a length-12 sequence by duplicating the bits in  $\mathbf{y}$ ; then assigning 1 bit to the last run and 2 bits to the second last run.

*Approach 2 (Assign-by-the-length Approximation):*

The first step is the same as Approach 1. For the remaining bits, longer runs get more bits, i.e., assign them to the longest run (of length  $\ell$ ) until the length exceeds  $\ell n / \lfloor m \rfloor$ . For example, suppose  $m = 6$ ,  $n = 15$  and  $\mathbf{y} = 010001$ . Then the formed new length- $n$  sequence is 00110000000011.

*Approach 3 (Gamma Function Approximation):*

An alternative approximation is to substitute the binomial coefficients in Eq. (51) by Gamma functions, thus ensuring  $F = n/m$  taking non-integers.

The three approximations of  $\tilde{C}_n(d)$  and the ML upper bound  $\overline{C}_n(d)$  are depicted in Figure 2, together with the best known numerical upper bounds reported in [22].

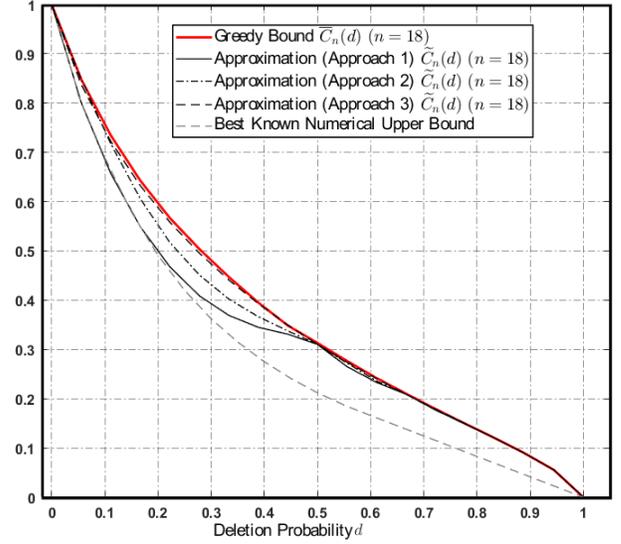


Fig. 2. The ML upper bound (solid, red)  $\overline{C}_n(d)$  for BDC with block-length  $n = 18$ , together approximations of  $\tilde{C}_n(d)$  obtained via the three approaches. The lower curve (dashed gray) corresponds to the best known numerical upper bounds provided in [22].

Some observations are summarized as hypotheses in sequel.

**Hypothesis 1.** For any block-lengths  $n, m \in \mathbb{N}^+$  and  $F = n/m \in \mathbb{N}^+$ ,

$$\min_{\mathbf{y} \in \{0,1\}^{n/F}} \gamma(n, F, \mathbf{y}) = \frac{F^{n/F}}{\#(\mathbf{y}_{\text{flip}})}$$

where  $\mathbf{y}_{\text{flip}}$  denotes a length- $m$  binary sequence with flipping bits, i.e.,  $R_{\mathbf{y}}^{(1)} = m$  and  $R_{\mathbf{y}}^{(\ell)} = 0$  for all  $\ell > 1$ .

Moreover, the approximations are tight, such that

**Hypothesis 2.** For any block-lengths  $n, m \in \mathbb{N}^+$  and  $F = n/m \in \mathbb{N}^+$ ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \gamma(n, F) = 0.$$

Furthermore, based on Hypothesis 1, Hypothesis 2 is true when  $d \rightarrow 1$ . Denote by  $\gamma(n, F) := \min_{\mathbf{y} \in \{0,1\}^{n/F}} \gamma(n, F, \mathbf{y})$  for notational convenience. We conclude the following asymptotic behavior of  $\gamma(n, F, \mathbf{y})$ :

**Lemma 8.**

$$\lim_{m/n \rightarrow 1} \lim_{n \rightarrow \infty} \frac{1}{n} \log \gamma(n, F) = 0.$$

*Proof.* Note that  $\gamma(n, F) = \frac{F^{n/F}}{\#(\mathbf{y}_{\text{flip}})} \geq \frac{F^{n/F}}{\binom{n}{m}}$ . Using Stirling's approximation (see (52)), we get

$$\binom{n}{m} \leq \frac{e}{2\pi} \cdot \frac{2^{nh(\frac{m}{n})}}{\sqrt{(1 - \frac{m}{n})m}}.$$

Since  $F = n/m$ ,

$$\gamma(n, F) \geq \frac{2\pi}{e} \cdot \left(\frac{n}{m}\right)^m \frac{\sqrt{(1 - \frac{m}{n})m}}{2^{nh(\frac{m}{n})}}.$$

$$\binom{\ell F}{\ell} = \frac{(\ell F)!}{\ell! (\ell F - \ell)!} \leq \frac{\exp(-\ell F + 1) (\ell F)^{\ell F + 1/2}}{(\sqrt{2\pi} \exp(-\ell) \ell^{\ell + 1/2}) \cdot (\sqrt{2\pi} \exp(-(\ell F - \ell)) (\ell F - \ell)^{(\ell F - \ell) + 1/2})} \quad (52)$$

Thus, taking logarithm and letting  $n \rightarrow \infty$ ,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \gamma(n, F) &\geq \lim_{n \rightarrow \infty} \frac{1}{n} \left( m \log \frac{n}{m} - nh \left( \frac{m}{n} \right) \right) \\ &= \lim_{n \rightarrow \infty} \left( 1 - \frac{m}{n} \right) \log \left( \frac{1}{1 - m/n} \right), \end{aligned}$$

which goes to 0 as  $m/n \rightarrow 1$ .  $\square$

The remaining part of this work is based on Hypothesis 2.

#### D. Explicit Approximation of $\bar{C}(d)$

Based on Hypothesis 2, we can further bound the capacity  $C(d)$  of the BDC when the deletion probability  $d \geq 1/2$ .

**Lemma 9.** *Suppose Hypothesis 2 is true. For any block-lengths  $n, m \in \mathbb{N}^+$  and  $F = n/m \in \mathbb{N}^+$ , the following bound on  $C(d)$  holds:*

$$C(d) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})} - h(d) \quad (53)$$

where  $m = \lceil n(1-d) \rceil$ .

*Proof.* We start with repeating Theorem 3:

$$C(d) \leq \bar{C}(d) = \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \bar{\#}(\mathbf{y}) \right) - h(d) \quad (54)$$

provided  $m = \lceil n(1-d) \rceil$ .

Suppose  $m$  divides  $n$ . Considering the definition of  $\gamma(n, F)$ , the logarithmic term in above can be bounded as

$$\begin{aligned} &\frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \bar{\#}(\mathbf{y}) \\ &\leq \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \frac{\prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})}}{\gamma(n, F)} \right) \\ &= \frac{1}{n} \log \left( \sum_{\mathbf{y} \in \{0,1\}^m} \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})} \right) - \frac{1}{n} \log \gamma(n, F) \quad (55) \end{aligned}$$

Taking the limits  $n \rightarrow \infty$  and  $d \rightarrow 1$ , Hypothesis 2 implies (53).  $\square$

We will then show the following lemma holds:

**Lemma 10.** *Suppose  $m = \lceil n(1-d) \rceil$ . For all deletion probability  $d \in [1/2, 1)$ ,*

$$\begin{aligned} &\frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})} \\ &= h(d) + 1 - d + \frac{1}{n} \log \left( \mathbb{E}[\exp(-\mu_d(\mathbf{Y}))] \right) \quad (56) \end{aligned}$$

where

$$\mu_d(\mathbf{y}) := \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right),$$

and the expectation is over a random length- $m$  sequence  $\mathbf{Y}$  wherein each bit is a random variable with distribution Bernoulli  $(1/2)$ .

*Proof.* Applying Stirling's approximation (inequalities) to the binomial coefficients, (53) follows. Therefore

$$\begin{aligned} \binom{\ell F}{\ell} &\leq \frac{e (\ell F)^{\ell F + 1/2}}{2\pi \ell^{\ell + 1/2} (\ell F - \ell)^{(\ell F - \ell) + 1/2}} \\ &= \frac{e}{2\pi} \cdot \frac{(\ell F)^{\ell F}}{\ell^{\ell} (\ell F - \ell)^{(\ell F - \ell)}} \cdot \sqrt{\frac{F}{\ell(F-1)}}. \quad (57) \end{aligned}$$

Since  $F = \frac{1}{1-d}$ , we have  $\frac{F}{F-1} = \frac{1}{d}$ . Thus,

$$\frac{(\ell F)^{\ell F}}{\ell^{\ell} (\ell F - \ell)^{(\ell F - \ell)}} = 2^{\ell F h((F-1)/F)} = 2^{\ell F h(d)}, \quad (58)$$

$$\sqrt{\frac{F}{\ell(F-1)}} = \sqrt{\frac{1}{d\ell}}. \quad (59)$$

Putting (58) and (59) into (57),

$$\binom{\ell F}{\ell} \leq \frac{e}{2\pi} \cdot \frac{2^{\ell F h(d)}}{\sqrt{d\ell}}.$$

Therefore,

$$\begin{aligned} &\ln \left( \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\mathbf{y}}^{(\ell)}} \right) \\ &= \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \binom{\ell F}{\ell} \\ &\leq \sum_{\ell=1}^m (\ln 2) \ell R_{\mathbf{y}}^{(\ell)} F h(d) \\ &\quad + \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \frac{e}{2\pi} - \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln(d\ell). \end{aligned}$$

According to (49),  $\sum_{\ell=1}^m \ell R_{\mathbf{y}}^{(\ell)} = m$ , implying that

$$\sum_{\ell=1}^m \ell R_{\mathbf{y}}^{(\ell)} F h(d) = m F h(d) = nh(d).$$

Continuing from above,

$$\begin{aligned} &\ln \left( \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\mathbf{y}}^{(\ell)}} \right) \\ &\leq (\ln 2) nh(d) + \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \frac{e}{2\pi} - \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln(d\ell) \quad (60) \end{aligned}$$

$$= (\ln 2) nh(d) - \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right). \quad (61)$$

Recall that

$$\mu_d(\mathbf{y}) := \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right).$$

Thus, summing over all  $\mathbf{y} \in \{0,1\}^m$ , (61) yields that

$$\begin{aligned} & \frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})} \\ & \leq \frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \exp \left( (\ln 2) n h(d) - \mu_d(\mathbf{y}) \right) \quad (62) \\ & = h(d) + \frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \exp \left( -\mu_d(\mathbf{y}) \right). \end{aligned}$$

The summation  $\sum_{\mathbf{y} \in \{0,1\}^m} \exp \left( -\mu_d(\mathbf{y}) \right)$  can be regarded as  $2^m$  times the expectation of  $\exp \left( -\mu_d(\mathbf{Y}) \right)$  given that each bit in  $\mathbf{Y}$  is selected Bernoulli  $(1/2)$ . Therefore,

$$\begin{aligned} & \frac{1}{n} \log \sum_{\mathbf{y} \in \{0,1\}^m} \prod_{\ell=1}^m \binom{\ell F}{\ell}^{R_{\ell}(\mathbf{y})} \\ & \leq h(d) + \frac{1}{n} \log \left( 2^m \mathbb{E} \left[ \exp \left( -\mu_d(\mathbf{Y}) \right) \right] \right) \\ & = h(d) + 1 - d + \frac{1}{n} \log \left( \mathbb{E} \left[ \exp \left( -\mu_d(\mathbf{Y}) \right) \right] \right). \quad (63) \end{aligned}$$

□

Taking the expectation outside, we derive the following approximation of the last term in (63):

$$\frac{1}{n} \mathbb{E} \left[ \log \left( \exp \left( -\mu_d(\mathbf{Y}) \right) \right) \right] = -\frac{1}{n} \mathbb{E} \left[ (\log e) \mu_d(\mathbf{Y}) \right]. \quad (64)$$

For a Bernoulli  $(1/2)$  process, the distribution of the number of  $\ell$ -runs  $R_{\mathbf{y}}^{(1)}, \dots, R_{\mathbf{y}}^{(m)}$  is proportional to the “run-length distribution” defined in [25], which is  $\{1/2^\ell\}_{\ell=1}^\infty$ . Hence, considering (49), the expectation of the number of runs is  $m/2^{\ell+1}$ , i.e.,

$$\mathbb{E} \left[ R_{\mathbf{Y}}^{(\ell)} \right] = \frac{m}{2^{\ell+1}}.$$

It follows that

$$\begin{aligned} \mathbb{E} \left[ \mu_d(\mathbf{Y}) \right] & = \mathbb{E} \left[ \frac{1}{2} \sum_{\ell=1}^m R_{\mathbf{y}}^{(\ell)} \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right) \right] \\ & = \frac{1}{2} \sum_{\ell=1}^m \mathbb{E} \left[ R_{\mathbf{Y}}^{(\ell)} \right] \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right) \\ & = \frac{m}{2} \sum_{\ell=1}^m \frac{1}{2^{\ell+1}} \ln \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right). \quad (65) \end{aligned}$$

Combining (53), (56), the approximation (65) and the identity (64) above, the approximation  $\tilde{C}(d)$  defined below holds:

$$\begin{aligned} \tilde{C}(d) & := 1 - d - \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[ (\log e) \mu_d(\mathbf{Y}) \right] \\ & = 1 - d - \frac{1}{2} \psi(1 - d) \quad (66) \end{aligned}$$

where

$$\begin{aligned} \psi & := \sum_{\ell=1}^{\infty} \frac{1}{2^{\ell+1}} \log \left( \left( \frac{2\pi}{e} \right)^2 d \ell \right) \\ & = \frac{1}{2} \log d + \sum_{\ell=1}^{\infty} \frac{\log \left( (2\pi/e) \sqrt{\ell} \right)}{2^\ell} \lesssim \frac{1}{2} \log d + 1.09179. \end{aligned}$$

A figure depicting the explicit approximation  $\tilde{C}(d)$  for  $d \geq 1/2$  is provided below.

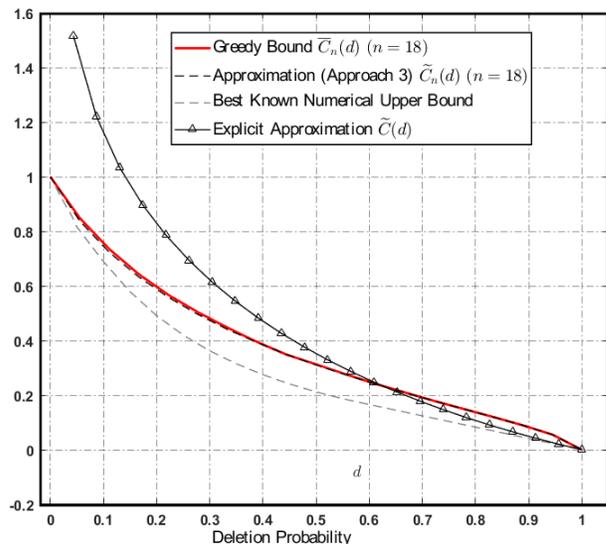


Fig. 3. The ML upper bound (solid, red)  $\bar{C}_n(d)$  for BDC with block-length  $n = 18$ , the approximation of  $\tilde{C}_n(d)$  obtained via the gamma function, and the *explicit approximation*  $\tilde{C}(d)$  derived in (66). The lower curve (dashed gray) corresponds to the best known numerical upper bounds provided in [22].

## V. CONCLUSION

We derive a general ML upper bound (See Theorem 1) for *information stable channels*. The corresponding bounds are shown to be tight for simple channels, e.g., the BEC and the BSC. Furthermore, we demonstrate the usage of the bound on the BDC, whose capacity remains unknown. The corresponding upper bound for the BDC derived from the general bound coincides with an intriguing combinatorial problem (defined as the MDM problem in Definition 10). Approximations for the derived upper bound are provided via three different approaches. Furthermore, analyzing and approximating the limiting behavior of the derived upper bound gives an explicit bound reported in (66) (and shown in Figure 3), validating that the general bound is capable of providing nontrivial results for sophisticated channels with memory. The next step is to validate the upper bounds on a variety types of channels and formalize a more general framework based on the main result (Theorem 1) stated in the paper.

## REFERENCES

- [1] R. Dobrushin, "General formulation of shannon's main theorem in information theory," *Amer. Math. Soc. Trans.*, vol. 33, pp. 323–438, 1963.
- [2] G. Hu, "On shannon theorem and its converse for sequence of communication schemes in the case of abstract random variables," in *Trans. 3rd Prague Conference on Information Theory, Statistical Decision Functions, Random Processes, Czechoslovak Academy of Sciences, Prague*, 1964, pp. 285–333.
- [3] S. Verdú and T. Han, "A general formula for channel capacity," *Information Theory, IEEE Transactions on*, vol. 40, no. 4, pp. 1147–1157, 1994.
- [4] S. Tatikonda and S. Mitter, "The capacity of channels with feedback," *IEEE Transactions on Information Theory*, vol. 55, no. 1, pp. 323–349, 2009.
- [5] M. Hayashi and H. Nagaoka, "General formulas for capacity of classical-quantum channels," *IEEE Transactions on Information Theory*, vol. 49, no. 7, pp. 1753–1768, 2003.
- [6] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions and reversals," in *Soviet physics doklady*, vol. 10, 1966, p. 707.
- [7] R. L. Dobrushin, "Shannon's theorems for channels with synchronization errors," *Problemy Peredachi Informatsii*, vol. 3, no. 4, pp. 18–36, 1967.
- [8] S. Arimoto, "An algorithm for computing the capacity of arbitrary discrete memoryless channels," *Information Theory, IEEE Transactions on*, vol. 18, no. 1, pp. 14–20, 1972.
- [9] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *Information Theory, IEEE Transactions on*, vol. 18, no. 4, pp. 460–473, 1972.
- [10] D. Fertoni and T. M. Duman, "Novel bounds on the capacity of the binary deletion channel," *Information Theory, IEEE Transactions on*, vol. 56, no. 6, pp. 2753–2765, 2010.
- [11] S. Vembu, S. Verdú, and Y. Steinberg, "The source-channel separation theorem revisited," *Information Theory, IEEE Transactions on*, vol. 41, no. 1, pp. 44–54, 1995.
- [12] T. Han and S. Verdú, "Approximation theory of output statistics," *Information Theory, IEEE Transactions on*, vol. 39, no. 3, pp. 752–772, 1993.
- [13] P.-N. Chen and F. Alajaji, "Optimistic shannon coding theorems for arbitrary single-user systems," *Information Theory, IEEE Transactions on*, vol. 45, no. 7, pp. 2623–2629, 1999.
- [14] A. Feinstein, *A new basic theorem of information theory*, 1954.
- [15] C. E. Shannon, "Certain results in coding theory for noisy channels," *Information and control*, vol. 1, no. 1, pp. 6–25, 1957.
- [16] R. G. Gallager, *Information theory and reliable communication*. Springer, 1968, vol. 2.
- [17] F. Jelinek, *Probabilistic information theory: discrete and memoryless models*. McGraw-Hill, 1968.
- [18] M. Mitzenmacher *et al.*, "A survey of results for deletion channels and related synchronization channels," *Probability Surveys*, vol. 6, pp. 1–33, 2009.
- [19] H. Chernoff, "A note on an inequality involving the normal distribution," *The Annals of Probability*, pp. 533–535, 1981.
- [20] M. Mitzenmacher and E. Upfal, *Probability and computing: Randomized algorithms and probabilistic analysis*. Cambridge University Press, 2005.
- [21] M. Cheraghchi, "Capacity upper bounds for deletion-type channels," *arXiv preprint arXiv:1711.01630*, 2017.
- [22] M. Rahmati and T. M. Duman, "Upper bounds on the capacity of deletion channels using channel fragmentation," *Information Theory, IEEE Transactions on*, vol. 61, no. 1, pp. 146–156, 2015.
- [23] S. Diggavi, M. Mitzenmacher, and H. Pfister, "Capacity upper bounds for deletion channels," in *Proceedings of the International Symposium on Information Theory*, 2007, pp. 1716–1720.
- [24] M. Dalai, "A new bound on the capacity of the binary deletion channel with high deletion probabilities," in *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*. IEEE, 2011, pp. 499–502.
- [25] Y. Kanoria and A. Montanari, "Optimal coding for the binary deletion channel with small deletion probability," *Information Theory, IEEE Transactions on*, vol. 59, no. 10, pp. 6192–6219, 2013.
- [26] A. Kalai, M. Mitzenmacher, and M. Sudan, "Tight asymptotic bounds for the deletion channel with small deletion probabilities," in *Information Theory Proceedings (ISIT), 2010 IEEE International Symposium on*. IEEE, 2010, pp. 997–1001.
- [27] M. Drmota, W. Szpankowski, and K. Viswanathan, "Mutual information for a deletion channel," in *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*. IEEE, 2012, pp. 2561–2565.
- [28] H. Koga *et al.*, *Information-spectrum methods in information theory*. Springer Science & Business Media, 2013, vol. 50.
- [29] Y. Liron and M. Langberg, "A characterization of the number of subsequences obtained via the deletion channel," *IEEE Transactions on Information Theory*, vol. 61, no. 5, pp. 2300–2312, 2015.
- [30] J. Bourdon and B. Vallée, "Generalized pattern matching statistics," in *Mathematics and Computer Science II*. Springer, 2002, pp. 249–265.
- [31] A. Kirsch and E. Drinea, "Directly lower bounding the information capacity for channels with iid deletions and duplications," *Information Theory, IEEE Transactions on*, vol. 56, no. 1, pp. 86–102, 2010.