

Expurgated Random-Coding Ensembles: Exponents, Refinements and Connections

Jonathan Scarlett, Li Peng, Neri Merhav, Alfonso Martinez and Albert Guillén i Fàbregas

Abstract

This paper studies expurgated random-coding bounds and exponents for channel coding with a given (possibly suboptimal) decoding rule. Variations of Gallager's analysis are presented, yielding several asymptotic and non-asymptotic bounds on the error probability for an arbitrary codeword distribution. A simple non-asymptotic bound is shown to attain an exponent of Csiszár and Körner under constant-composition coding. Using Lagrange duality, this exponent is expressed in several forms, one of which is shown to permit a direct derivation via cost-constrained coding which extends to infinite and continuous alphabets. The method of type class enumeration is studied, and it is shown that this approach can yield improved exponents and better tightness guarantees for some codeword distributions. A generalization of this approach is shown to provide a multi-letter exponent which extends immediately to channels with memory. Finally, a refined analysis expurgated i.i.d. random coding is shown to yield a $O(\frac{1}{\sqrt{n}})$ prefactor, thus improving on the standard $O(1)$ prefactor. Moreover, the implied constant is explicitly characterized.

Index Terms

Expurgated error exponents, reliability function, random coding, mismatched decoding, maximum-likelihood decoding, type class enumeration

I. INTRODUCTION

Achievable performance bounds for channel coding are typically obtained by analyzing the average error probability of an ensemble of codebooks with independently generated codewords. For memoryless channels, random codes with independent and identically distributed (i.i.d.) symbols achieve the channel capacity [1], characterize the error exponent of the best code at sufficiently high rates [2, Ch. 5], and provide tight bounds on the finite-length performance [3].

J. Scarlett and L. Peng are with the Department of Engineering, University of Cambridge, Cambridge, CB2 1PZ, U.K. (e-mails: jmscarlett@gmail.com, lp327@cam.ac.uk). N. Merhav is with the Department of Electrical Engineering, Technion - Israel Institute of Technology, Technion City, Haifa 32000, Israel. (e-mail: merhav@ee.technion.ac.il). A. Martinez is with the Department of Information and Communication Technologies, Universitat Pompeu Fabra, 08018 Barcelona, Spain (e-mail: alfonso.martinez@ieec.org). A. Guillén i Fàbregas is with the Institució Catalana de Recerca i Estudis Avançats (ICREA), the Department of Information and Communication Technologies, Universitat Pompeu Fabra, 08018 Barcelona, Spain, and also with the Department of Engineering, University of Cambridge, Cambridge, CB2 1PZ, U.K. (e-mail: guillen@ieec.org).

This work has been funded in part by the European Research Council under ERC grant agreement 259663, by the European Union's 7th Framework Programme (PEOPLE-2011-CIG) under grant agreement 303633 and by the Spanish Ministry of Economy and Competitiveness under grants RYC-2011-08150 and TEC2012-38800-C03-03. The work of N. Merhav was partially supported by the Israel Science Foundation (ISF), grant no. 412/12.

This is the extended version of a paper which was accepted to *IEEE Transactions on Information Theory* (April 2014). A shortened version was presented at the 2014 International Zurich Seminar on Communications.

At low rates, the error probability of the best code in the random-coding ensemble can be significantly smaller than the average. In such cases, better performance bounds are obtained by considering an ensemble in which a subset of the randomly generated codewords are expurgated from the codebook. In particular, the error exponents resulting from such techniques are generally higher than the random-coding error exponent at low rates. Existing works exploring such techniques include those of Gallager [2, Sec. 5.7], Csiszár-Körner-Martón [4], [5, Ex. 10.18] and Csiszár-Körner [6]. The advantages of Gallager's approach include its simplicity and the fact that the analysis is not restricted to finite alphabets. On the other hand, as we will see in Section III, the exponents of [4]–[6] can improve on that of Gallager for a given input distribution or decoding rule.

In this paper, we provide techniques that attain the best of each of the above approaches. Using variations of Gallager's analysis, we obtain several asymptotic and non-asymptotic bounds for an arbitrary codeword distribution. Using these bounds, we provide derivations of both new and existing expurgated exponents, each yielding various advantages such as simplicity, generality, and guarantees of exponential tightness. We explore the method of type class enumeration (e.g. see [7]–[9]) for both discrete and continuous channels, and show that it can yield improved exponents and tightness guarantees, as well as providing a multi-letter exponent which extends immediately to channels with memory.

A. System Setup

The input and output alphabets are denoted by \mathcal{X} and \mathcal{Y} respectively. The channel is assumed to be memoryless, yielding an n -letter transition law given by $W^n(\mathbf{y}|\mathbf{x}) \triangleq \prod_{i=1}^n W(y_i|x_i)$ for some conditional distribution $W(y|x)$. In the case that both \mathcal{X} and \mathcal{Y} are finite, the channel is a discrete memoryless channel (DMC), but we do not assume this to be the case in general. The encoder takes as input a message m equiprobable on the set $\{1, \dots, M\}$, and transmits the corresponding codeword $\mathbf{x}^{(m)}$ from a codebook $\mathcal{C} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$. The decoder receives the vector \mathbf{y} at the output of the channel, and forms the estimate

$$\hat{m} = \arg \max_{j \in \{1, \dots, M\}} q^n(\mathbf{x}^{(j)}, \mathbf{y}), \quad (1)$$

where $q^n(\mathbf{x}, \mathbf{y}) \triangleq \prod_{i=1}^n q(x_i, y_i)$, and $q(x, y)$ is a non-negative function called the *decoding metric*. An error is said to have occurred if $\hat{m} \neq m$, and we assume that ties are broken as errors. We let $p_{e,m}(\mathcal{C})$ be the error probability induced by \mathcal{C} given a particular message m , and we denote the maximal error probability by $p_e(\mathcal{C}) \triangleq \max_m p_{e,m}(\mathcal{C})$.

When $q(x, y) = W(y|x)$, (1) is the optimal maximum-likelihood (ML) decoding rule. For other decoding metrics, this setting is that of *mismatched decoding* [10]–[13], which is relevant when ML decoding is not feasible, e.g. due to channel uncertainty or implementation constraints.

Throughout the paper, we consider channels with both constrained and unconstrained inputs. In the former setting, each codeword \mathbf{x} must satisfy a constraint of the form

$$\frac{1}{n} \sum_{i=1}^n c(x_i) \leq \Gamma, \quad (2)$$

where $c(\cdot)$ is referred to as a cost function, and Γ is a constant. Except where stated otherwise, it will be assumed that the input is unconstrained, which corresponds to $\Gamma = \infty$.

For a given rate R , an error exponent $E(R)$ is said to be achievable if there exists a sequence of codebooks \mathcal{C}_n of length n and rate R whose error probability $p_e(\mathcal{C}_n)$ satisfies

$$\liminf_{n \rightarrow \infty} -\frac{1}{n} \log p_e(\mathcal{C}_n) \geq E(R). \quad (3)$$

We focus on the maximal error probability rather than the average error probability, but the two are equivalent for the purposes of studying error exponents.

B. Previous Work

Considering ML decoding, Gallager [2, Ch. 5] studied an ensemble in which $2M - 1$ codewords are generated at random, and a subset of M codewords forms the codebook. Roughly speaking, the codewords which are kept are those which have the lowest error probability among the original codewords. A different approach was taken by Csiszár, Körner and Marton [4] (see also [5, Ex. 10.18]), who began by proving the existence of a collection of constant-composition codewords such that any two codewords have a joint empirical distribution satisfying certain properties. By analyzing this collection of codewords using the method of types, an error exponent was obtained which coincides with that of Gallager after the optimization of the input distribution. An exponent for mismatched decoding was derived by Csiszár and Körner [6], and was shown to coincide with that of [4] when particularized to the case of ML decoding.

As stated in the introduction, the exponents of [4], [6] can in fact improve on that of Gallager for a given input distribution. However, the proofs rely heavily on techniques which are valid only when the input and output alphabets are finite. In particular, [4] uses the type packing lemma [5, Ch. 10], and [6] uses a combinatorial graph decomposition lemma. For other related works, see [14]–[17].

Overviews of the mismatched decoding problem can be found in [10]–[13]. Most of the literature has focused on achievable rates, whereas this paper is concerned with the performance at low rates. The mismatched decoding paper most relevant to this one is [13], which studies random-coding error exponents for various non-expurgated ensembles.

C. Contributions

The main contributions of this paper are as follows:

- In Section II, we present variations of Gallager’s analysis which yield several asymptotic and non-asymptotic bounds on the error probability. In particular, we consider the use of a logarithmic function in the expurgation argument in place of the power function used by Gallager [2, Sec. 7.3].
- In Section III, we present an overview of various expurgated exponents and the connections between them. Using the method of Lagrange duality [18], we relate the exponents given in [2], [4], [6]. Generalizations of the exponents in [2], [4] to the setting of mismatched decoding are given, and an alternative form of the exponent in [6] is given which extends readily to channels with infinite or continuous alphabets.
- In Section IV, we present several methods for deriving both new and existing exponents:
 - In Section IV-A, we present simple techniques for deriving exponents using a non-asymptotic bound from Section II. Applying constant-composition coding and the method of types recovers the exponent in [6], thus

providing a simple and concise proof. Furthermore, applying cost-constrained coding with multiple auxiliary costs [13] recovers the generalization of this exponent to more general alphabets.

- In Section IV-B, we study the method of type class enumeration (e.g. see [7]–[9]), which is shown to yield better exponents than the simpler approach for some codeword distributions, as well as better guarantees of exponential tightness.
- In Section IV-C, we extend the type class enumeration analysis to allow for infinite and continuous alphabets. This is not only of interest in itself, but also yields a multi-letter exponent which can be directly applied to channels with memory and more general decoding metrics.
- In Section V, we present a refined derivation of Gallager’s exponent for i.i.d. random coding (and its generalization to mismatched decoding) with a $O(\frac{1}{\sqrt{n}})$ prefactor, thus improving on the original $O(1)$ prefactor. Similar improvements for the non-expurgated random-coding error exponent have recently been obtained by Altuğ and Wagner [19] (see also [20]).

D. Notation

We use bold symbols for vectors (e.g. \mathbf{x}), and denote the corresponding i -th entry using a subscript (e.g. x_i).

The set of all probability distributions on an alphabet, say \mathcal{X} , is denoted by $\mathcal{P}(\mathcal{X})$, and the set of all empirical distributions on a vector in \mathcal{X}^n (i.e. types [5, Ch. 2]) is denoted by $\mathcal{P}_n(\mathcal{X})$. For a given type $Q \in \mathcal{P}_n(\mathcal{X})$, the type class $T^n(Q)$ is defined to be the set of all sequences in \mathcal{X}^n with type Q .

The probability of an event is denoted by $\mathbb{P}[\cdot]$, and the symbol \sim means “distributed as”. The marginals of a joint distribution $P_{XY}(x, y)$ are denoted by $P_X(x)$ and $P_Y(y)$. We write $P_X = \tilde{P}_X$ to denote element-wise equality between two probability distributions on the same alphabet. Expectation with respect to a joint distribution $P_{XY}(x, y)$ is denoted by $\mathbb{E}_P[\cdot]$, or simply $\mathbb{E}[\cdot]$ when the associated probability distribution is understood from the context. Similarly, the mutual information with respect to P_{XY} is written as $I_P(X; Y)$, or simply $I(X; Y)$ when the distribution is understood from the context. Given a distribution $Q(x)$ and conditional distribution $W(y|x)$, we write $Q \times W$ to denote the joint distribution defined by $Q(x)W(y|x)$.

For two positive sequences f_n and g_n , we write $f_n \doteq g_n$ if $\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{f_n}{g_n} = 0$, and we write $f_n \stackrel{\cdot}{\leq} g_n$ if $\limsup_{n \rightarrow \infty} \frac{1}{n} \log \frac{f_n}{g_n} \leq 0$, and analogously for $\stackrel{\cdot}{\geq}$. We write $f_n = O(g_n)$ if $|f_n| \leq c|g_n|$ for some c and sufficiently large n . All logarithms have base e , and all rates are in units of nats except in the examples, where bits are used. We define $[c]^+ = \max\{0, c\}$, and denote the indicator function by $\mathbb{1}\{\cdot\}$.

II. EXPURGATED BOUNDS

In this section, we present a number of variations of Gallager’s bounds and techniques which will provide the starting points of the derivations of the exponents in Section IV. We let $P_{\mathbf{X}}$ denote a codeword distribution, and we define the random variables $(\mathbf{X}, \mathbf{Y}, \overline{\mathbf{X}})$ distributed according to

$$(\mathbf{X}, \mathbf{Y}, \overline{\mathbf{X}}) \sim P_{\mathbf{X}}(\mathbf{x})W^n(\mathbf{y}|\mathbf{x})P_{\mathbf{X}}(\overline{\mathbf{x}}). \quad (4)$$

In the case that a cost constraint of the form (2) is present, we assume that $P_{\mathbf{X}}$ is chosen such that \mathbf{X} satisfies the constraint with probability one.

We let $\mathcal{C} = \{\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M')}\}$ be a random codebook of size M' with each codeword independently generated according to $P_{\mathbf{X}}$. The symbol \mathcal{C} is used to denote a fixed expurgated codebook containing $M \leq M'$ codewords.

We begin with the following straightforward generalization of [2, Lemma, p. 151].

Lemma 1. *Fix a function $f : [0, 1] \rightarrow \mathbb{R}$ and a codeword distribution $P_{\mathbf{X}}$ such that $f(p_{e,m}(\mathcal{C}))$ is non-negative for all m with probability one. For any $\eta > 0$, there exists a codebook \mathcal{C} of size M such that $M' \frac{\eta}{1+\eta} < M \leq M'$ and*

$$f(p_{e,m}(\mathcal{C})) \leq (1 + \eta)\mathbb{E}[f(p_{e,m}(\mathcal{C}))] \quad (5)$$

for $m = 1, \dots, M$.

Proof: The proof is identical to [2, Lemma, p. 151], with the assumption of $f(p_{e,m}(\mathcal{C}))$ being non-negative ensuring the validity of Markov's inequality. \blacksquare

While Lemma 1 is valid for any function $f(\cdot)$, it is primarily of interest when $f(\cdot)$ is monotonically increasing, so that (5) can be inverted in order to obtain an upper bound on $p_{e,m}(\mathcal{C})$. Gallager [2] presented the lemma with the choices $\eta = 1$ and $f(\cdot) = (\cdot)^{1/\rho}$, where $\rho > 0$, thus proving the existence of a codebook \mathcal{C} of size M such that

$$p_e(\mathcal{C}) \leq \left(2\mathbb{E}[p_{e,m}(\mathcal{C})^{1/\rho}]\right)^\rho, \quad (6)$$

where \mathcal{C} contains $M' = 2M - 1$ codewords. In the following theorem, we provide non-asymptotic bounds on the error probability which follow in a straightforward fashion from (6). The proof alters Gallager's arguments for the purpose of better characterizing the non-asymptotic performance, and also for dealing with suboptimal decoding rules.

Theorem 1. *For any pair (n, M) , codeword distribution $P_{\mathbf{X}}$, and parameters $\rho \geq 1$ and $s \geq 0$, there exists a codebook \mathcal{C}_n with M codewords of length n whose maximal error probability satisfies*

$$p_e(\mathcal{C}_n) \leq \text{rcux}_\rho(n, M) \leq \text{rcux}_{\rho,s}(n, M) \quad (7)$$

where

$$\text{rcux}_\rho(n, M) \triangleq \left(4(M-1)\mathbb{E}\left[\mathbb{P}\left[q^n(\bar{\mathbf{X}}, \mathbf{Y}) \geq q^n(\mathbf{X}, \mathbf{Y}) \mid \mathbf{X}, \bar{\mathbf{X}}\right]^{1/\rho}\right]\right)^\rho \quad (8)$$

$$\text{rcux}_{\rho,s}(n, M) \triangleq \left(4(M-1)\mathbb{E}\left[\mathbb{E}\left[\left(\frac{q^n(\bar{\mathbf{X}}, \mathbf{Y})}{q^n(\mathbf{X}, \mathbf{Y})}\right)^s \mid \mathbf{X}, \bar{\mathbf{X}}\right]^{1/\rho}\right]\right)^\rho. \quad (9)$$

Proof: We obtain (8) from (6) by weakening the expectation as follows:

$$\mathbb{E}[p_{e,m}(\mathcal{C})^{1/\rho}] \leq \mathbb{E}\left[\left(\sum_{\bar{m} \neq m} \mathbb{P}\left[q^n(\mathbf{X}^{(\bar{m})}, \mathbf{Y}) \geq q^n(\mathbf{X}^{(m)}, \mathbf{Y}) \mid \mathbf{X}^{(m)}, \mathbf{X}^{(\bar{m})}\right]\right)^{1/\rho}\right] \quad (10)$$

$$\leq \mathbb{E}\left[2(M-1)\mathbb{P}\left[q^n(\bar{\mathbf{X}}, \mathbf{Y}) \geq q^n(\mathbf{X}, \mathbf{Y}) \mid \mathbf{X}, \bar{\mathbf{X}}\right]^{1/\rho}\right], \quad (11)$$

where (10) follows from the union bound, and (11) follows using $M' = 2M - 1$ along with the inequality

$$\left(\sum_i a_i\right)^{1/\rho} \leq \sum_i a_i^{1/\rho}, \quad (12)$$

which holds for any $\rho \geq 1$. We obtain (9) by applying Markov's inequality to the inner probability in (8). \blacksquare

Following the terminology of Polyanskiy *et al.* [3], we refer to the bounds in (8)–(9) as *expurgated random-coding union* (RCUX) bounds. These bounds are computable for sufficiently symmetric setups, and are thus of independent interest for characterizing the finite-length performance [3]. It should be noted that both rcux_ρ and $\text{rcux}_{\rho,s}$ extend immediately to channels with memory and general decoding metrics (not necessarily single-letter).

The bound $\text{rcux}_{\rho,s}$ was presented by Gallager [2] under ML decoding with $s = \frac{1}{2}$. For the random-coding ensembles we consider, it will be seen that this choice of s is optimal for ML decoding, at least in terms of the error exponent. However, for mismatched decoding it is important to allow for an arbitrary choice of $s \geq 0$.

The following theorem gives an asymptotic bound which follows by using Lemma 1 with a choice of $f(\cdot)$ which differs from that of Gallager.

Theorem 2. *Consider a sequence of codebooks \mathcal{C}_n containing $M'_n = e^{nR}$ codewords which are generated independently according to $P_{\mathbf{X}}$. Suppose that there exists a non-negative sequence $E(n)$ growing subexponentially in n (i.e. $E(n) \doteq 1$) such that*

$$\mathbb{P}[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \mid \mathbf{X} = \mathbf{x}] \geq e^{-E(n)} \quad (13)$$

for all \mathbf{x} and $\bar{\mathbf{x}}$ on the support of $P_{\mathbf{X}}$. Then there exists a sequence of codebooks \mathcal{C}_n with M_n codewords such that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log M_n = R \quad (14)$$

and

$$p_e(\mathcal{C}_n) \leq \exp\left(\mathbb{E}[\log p_{e,m}(\mathcal{C}_n)]\right) \quad (15)$$

$$\leq \exp\left(\rho \mathbb{E}\left[\log \mathbb{E}[p_{e,m}(\mathcal{C}_n)^{1/\rho} \mid \mathbf{X}^{(m)}]\right]\right), \quad (16)$$

where (16) holds for any $\rho > 0$.

Proof: The error probability associated with the transmitted codeword \mathbf{x} is lower bounded by the left-hand side of (13), where $\bar{\mathbf{x}}$ is any incorrect codeword. The assumption in (13) thus implies that the function $f(p_{e,m}(\mathcal{C})) = E(n) + \log p_{e,m}(\mathcal{C})$ is non-negative for $m = 1, \dots, M$. Applying Lemma 1, we obtain that for each n and any $\eta_n > 0$ there exists a codebook \mathcal{C}_n of size $M_n = e^{nR} \frac{\eta_n}{1+\eta_n}$ such that

$$E(n) + \log p_e(\mathcal{C}_n) \leq (1 + \eta_n)(E(n) + \mathbb{E}[\log p_{e,m}(\mathcal{C}_n)]). \quad (17)$$

Since $\log \alpha \leq 0$ for $\alpha \in (0, 1]$, it follows that

$$\log p_e(\mathcal{C}_n) \leq \eta_n E(n) + \mathbb{E}[\log p_{e,m}(\mathcal{C}_n)]. \quad (18)$$

Choosing $\eta_n = \frac{1}{E(n)}$, we obtain (15), and the assumption that $E(n) \doteq 1$ implies (14). We obtain (16) by writing $\log \alpha = \rho \log(\alpha^{1/\rho})$, writing $\mathbb{E}[\cdot] = \mathbb{E}[\mathbb{E}[\cdot \mid \mathbf{X}^{(m)}]]$, and applying Jensen's inequality. ■

The assumption of Theorem 2 is mild, allowing ensembles for which the error probability associated with any two permissible codewords decays nearly *double*-exponentially fast. However, it is a multi-letter condition, and may therefore be difficult to verify directly. A single-letter sufficient condition depending only on the channel, metric and cost constraint (2) is that

$$\lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \log \log \frac{1}{\pi(\gamma)} = 0, \quad (19)$$

where

$$\pi(\gamma) \triangleq \min_{(x, \bar{x}) : c(x) \leq \gamma, c(\bar{x}) \leq \gamma} \mathbb{P}[Y_x \in \mathcal{E}(x, \bar{x})] \quad (20)$$

$$\mathcal{E}(x, \bar{x}) \triangleq \{y : q(\bar{x}, y) \geq q(x, y)\}, \quad (21)$$

where in (20) we define $Y_x \sim W(\cdot|x)$. Under this assumption, the probability in (13) is lower bounded by the probability that $Y_i \in \mathcal{E}(X_i, \bar{X}_i)$ for $i = 1, \dots, n$, which in turn is lower bounded by $\pi(n\Gamma)^n$. Since n times a subexponential sequence is also subexponential, the condition of Theorem 2 follows from (19). Further discussion is given in Appendix A, along with some examples.

From (15), we can see the advantage of the expurgated ensemble over the non-expurgated one. The former yields the exponent corresponding to $-\frac{1}{n}\mathbb{E}[\log p_{e,m}(C_n)]$, which is higher in general than that of $-\frac{1}{n}\log \mathbb{E}[p_{e,m}(C_n)]$ due to Jensen's inequality.

Using L'Hôpital's rule, it is easily shown that $\lim_{\rho \rightarrow \infty} \rho \log \mathbb{E}[Z^{1/\rho}] = \mathbb{E}[\log Z]$ for any random variable Z . It follows that the inequality in (16) is actually an equality in the limit as $\rho \rightarrow \infty$. At first glance, it may appear that a similar argument can be used to show that (6) yields the same exponent as (15). However, there is an issue with the order of the limits of n and ρ . If we take $\rho \rightarrow \infty$ in (6), the factor 2^ρ makes the right-hand side equal ∞ . Letting ρ grow slowly with n is also potentially problematic, since the random variable $p_{e,m}(C)$ varies with n .

The bounds in Theorem 2 will prove useful for deriving improved exponents compared to Theorem 1 for some codeword distributions, and for extending the type class enumeration method beyond the finite-alphabet setting.

III. EXPURGATED ENSEMBLES AND EXPONENTS

In this section, we present an overview of various expurgated exponents and the connections between them. Our focus here is primarily on existing exponents or simple variations thereof, though we also provide a dual form of the exponent in [6] which is new to the best of our knowledge. Further exponents which appear for the first time in this paper are given in Theorems 5 and 7 in Section IV.

Throughout the paper, we consider three expurgated ensembles, each of which depends on an input distribution Q :

- 1) The i.i.d. ensemble is characterized by

$$P_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n Q(x_i). \quad (22)$$

This codeword distribution is valid for both discrete and continuous alphabets, but it is not suitable for channels with cost constraints, since in all non-trivial cases there is a non-zero probability of violating the constraint.

- 2) The constant-composition ensemble is characterized by

$$P_{\mathbf{X}}(\mathbf{x}) = \frac{1}{|T^n(Q_n)|} \mathbb{1}\{\mathbf{x} \in T^n(Q_n)\}, \quad (23)$$

where Q_n is a type with the same support as Q such that $|Q_n(x) - Q(x)| = O(\frac{1}{n})$ for all x . This codeword distribution relies on the input being finite. It is directly applicable to channels with cost constraints, since each codeword satisfies (2) provided that $\mathbb{E}_{Q_n}[c(X)] \leq \Gamma$, which in turn can be achieved provided that $\mathbb{E}_Q[c(X)] \leq \Gamma$.

3) The cost-constrained ensemble is characterized by

$$P_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\mu_n} \prod_{i=1}^n Q(x_i) \mathbb{1}\{\mathbf{x} \in \mathcal{D}_n\}, \quad (24)$$

where

$$\mathcal{D}_n \triangleq \left\{ \mathbf{x} : \frac{1}{n} \sum_{i=1}^n c(x_i) \leq \Gamma, \left| \frac{1}{n} \sum_{i=1}^n a_l(x_i) - \phi_l \right| \leq \frac{\delta}{n}, l = 1, \dots, L \right\}, \quad (25)$$

and where δ is a positive constant (independent of n), $\{a_l(\cdot)\}_{l=1}^L$ are functions with means $\phi_l \triangleq \mathbb{E}_Q[a_l(X)]$, and μ_n is a normalizing constant. This codeword distribution is valid for both discrete and continuous alphabets, and ensures that each codeword satisfies (2). Both $c(\cdot)$ and $\{a_l(\cdot)\}$ can be thought of as cost functions, and we will distinguish between the two by referring to them as the *system cost* and *auxiliary costs* respectively. In contrast to the system cost, the auxiliary costs are functions which can be optimized. That is, while the system cost is given as part of the problem statement, the auxiliary costs are introduced to improve the performance of the random-coding ensemble itself [12], [13], [21].

We proceed by stating and comparing the exponents obtained by the above ensembles; derivations will be given in Section IV. Except where stated otherwise, we assume that the channel is a DMC with unconstrained inputs.

A straightforward generalization of Gallager's i.i.d. exponent to the setting of mismatched decoding is as follows:

$$E_{\text{ex}}^{\text{iid}}(Q, R) \triangleq \sup_{\rho \geq 1} E_{\text{x}}^{\text{iid}}(Q, \rho) - \rho R, \quad (26)$$

where

$$E_{\text{x}}^{\text{iid}}(Q, \rho) \triangleq \sup_{s \geq 0} -\rho \log \sum_{x, \bar{x}} Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \right)^{1/\rho}. \quad (27)$$

The objective in (27) is concave in s , and under ML decoding (i.e. $q(x, y) = W(y|x)$), it is also unchanged when s is replaced by $1 - s$. From these properties, it follows that $s = \frac{1}{2}$ is optimal for ML decoding, and thus the exponent is the same as that of Gallager [2].

Csiszár and Körner [6] make use of the constant-composition codeword distribution in (23). The analysis is significantly different to that of Gallager, and yields an exponent in a different form, namely

$$E_{\text{ex}}^{\text{cc}}(Q, R) \triangleq \min_{\substack{P_{X\bar{X}Y} \in \mathcal{T}(Q) \\ I_P(X; \bar{X}) \leq R}} D(P_{X\bar{X}Y} \| Q \times Q \times W) - R, \quad (28)$$

where the notation $Q \times Q \times W$ denotes the distribution $Q(x)Q(\bar{x})W(y|x)$, and

$$\mathcal{T}(Q) \triangleq \left\{ P_{X\bar{X}Y} \in \mathcal{P}(\mathcal{X} \times \mathcal{X} \times \mathcal{Y}) : P_X = Q, P_{\bar{X}} = Q, \mathbb{E}_P[\log q(\bar{X}, Y)] \geq \mathbb{E}_P[\log q(X, Y)] \right\}. \quad (29)$$

The objective in (28) follows from [6, Eq. (32)] and the identity

$$D(P_{X\bar{X}Y} \| Q \times Q \times W) = D(P_{X\bar{X}Y} \| P_{X\bar{X}} \times W) + I_P(X; \bar{X}), \quad (30)$$

which holds for any $P_{X\bar{X}Y}$ such that $P_X = P_{\bar{X}} = Q$. Defining $P_Y(y) \triangleq \sum_x Q(x)W(y|x)$, we observe that $E_{\text{ex}}^{\text{cc}}$ is positive for sufficiently small R provided that $\mathbb{E}_{Q \times W}[\log q(X, Y)] > \mathbb{E}_{Q \times P_Y}[\log q(X, Y)]$. It was shown in [11] that the mismatched capacity is in fact zero unless this condition holds for some Q .

The following theorem provides the means for comparing the above two exponents, as well as that of [4].

Theorem 3. For any input distribution Q and rate R , we have

$$E_{\text{ex}}^{\text{cc}}(Q, R) = \sup_{s \geq 0} \min_{\substack{P_{X\bar{X}}: P_X=Q, P_{\bar{X}}=Q, \\ I_P(X; \bar{X}) \leq R}} \mathbb{E}_P[d_s(X, \bar{X})] + I_P(X; \bar{X}) - R \quad (31)$$

$$= \sup_{\rho \geq 1} E_x^{\text{cc}}(Q, \rho) - \rho R, \quad (32)$$

where

$$d_s(x, \bar{x}) \triangleq -\log \sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \quad (33)$$

$$E_x^{\text{cc}}(Q, \rho) \triangleq \sup_{s \geq 0, a(\cdot)} -\rho \sum_x Q(x) \log \sum_{\bar{x}} Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \frac{e^{a(\bar{x})}}{e^{a(x)}} \right)^{1/\rho}. \quad (34)$$

Proof: See Appendix B. ■

Equations (32) and (34) strongly resemble (26)–(27). Equation (31) is a generalization of the exponent in [4], which is recovered by setting $q(x, y) = W(y|x)$ and $s = \frac{1}{2}$. Using the same argument as the one following (27), it can be shown that the latter choice is optimal. From the proof of Theorem 3, this implies the optimality of $s = \frac{1}{2}$ in (34) under ML decoding, though the optimal choice of $a(\cdot)$ is unclear in general. To our knowledge, the expression in (34) has not appeared previously even for ML decoding.

As noted in [6], [16], we can write (31) in the language of rate-distortion theory [22, Ch. 10]. Fix $s \geq 0$ and define

$$D_s(Q, R) \triangleq \min_{\substack{P_{X\bar{X}}: P_X=Q, P_{\bar{X}}=Q, \\ I_P(X; \bar{X}) \leq R}} \mathbb{E}_P[d_s(X, \bar{X})]. \quad (35)$$

This can be interpreted as the distortion-rate function of a source X with a reproduction variable \bar{X} , subject to the additional constraint that each reproduction codeword \bar{x} has empirical distribution Q . For any $s \geq 0$, the constraint on the mutual information in (31) is active for sufficiently small R . The supremum of all such rates is given by

$$R_s(Q) \triangleq I_{P^*}(X; \bar{X}), \quad (36)$$

where

$$P_{X\bar{X}}^* \triangleq \arg \min_{P_{X\bar{X}}: P_X=Q, P_{\bar{X}}=Q} \mathbb{E}_P[d_s(X, \bar{X})] + I_P(X; \bar{X}). \quad (37)$$

For $R \leq R_s$ we have $I_P(X; \bar{X}) = R$ under the minimizing $P_{X\bar{X}Y}$, whereas for $R \geq R_s$ the minimum in (31) decreases linearly with R for any fixed s . It follows that

$$E_{\text{ex}}^{\text{cc}}(Q, R) = \sup_{s \geq 0} E_{\text{ex}}^{\text{cc}}(Q, R, s), \quad (38)$$

where

$$E_{\text{ex}}^{\text{cc}}(Q, R, s) \triangleq \begin{cases} D_s(Q, R) & R \leq R_s(Q) \\ D_s(Q, R_s) + R_s(Q) - R & R > R_s(Q). \end{cases} \quad (39)$$

By applying Jensen's inequality to (34) and setting $a(x) = 0$, we immediately obtain

$$E_{\text{ex}}^{\text{cc}}(Q, R) \geq E_{\text{ex}}^{\text{iid}}(Q, R). \quad (40)$$

It was shown in [5, Ex. 10.18] that (40) holds with equality under ML decoding with an optimized input distribution Q . However, when either the decoding rule or input distribution is fixed, the inequality in (40) can be strict; an example is given at the end of this section. In Section IV-A, we show that the stronger exponent $E_{\text{ex}}^{\text{cc}}$, in the form given in (32), remains achievable in the case of continuous alphabets, with the summations in (34) replaced by integrals. This is proved using the cost-constrained ensemble in (24).

The following proposition generalizes Gallager's expression for the expurgated exponent as $R \rightarrow 0^+$ for channels whose zero-error capacity [23] is zero, and shows that the inequality in (40) becomes an equality in the limit.

Proposition 1. *Fix any input distribution Q such that all pairs (x, \bar{x}) with $Q(x)Q(\bar{x}) > 0$ share a common output, i.e. $W(y|x)W(y|\bar{x}) > 0$ for some y . Then*

$$\lim_{R \rightarrow 0^+} E_{\text{ex}}^{\text{cc}}(Q, R) = \lim_{R \rightarrow 0^+} E_{\text{ex}}^{\text{iid}}(Q, R) = \sup_{s \geq 0} \mathbb{E}[d_s(X, \bar{X})], \quad (41)$$

where d_s is defined in (33), and the expectation is taken with respect to $Q(x)Q(\bar{x})$.

Proof: See Appendix C. ■

We conclude this section with a numerical example. The channel is defined by the entries of the $|\mathcal{X}| \times |\mathcal{Y}|$ matrix

$$\begin{bmatrix} 1 - 2\delta_0 & \delta_0 & \delta_0 \\ \delta_1 & 1 - 2\delta_1 & \delta_1 \\ \delta_2 & \delta_2 & 1 - 2\delta_2 \end{bmatrix}, \quad (42)$$

and the decoding metric is defined similarly with a fixed $\delta \in (0, \frac{1}{3})$ in place of each δ_i ($i = 1, 2, 3$), yielding a minimum Hamming distance rule. Figure 1 plots the exponents in the case that $\delta_0 = 0.01$, $\delta_1 = 0.05$, $\delta_2 = 0.25$ and $Q = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. We observe that $E_{\text{ex}}^{\text{cc}} > E_{\text{ex}}^{\text{iid}}$ at all positive rates, and the gap is particularly significant in the mismatched case. However, consistent with Proposition 1, the two coincide in the limit as $R \rightarrow 0$.

As noted in [6], if Q is optimized, then the two exponents coincide for ML decoding. However, the strict inequality $E_{\text{ex}}^{\text{cc}} > E_{\text{ex}}^{\text{iid}}$ remains possible for other decoding rules.

IV. DERIVATIONS OF THE EXPURGATED EXPONENTS

In this section, we provide several techniques for deriving the expurgated exponents, including those introduced in Section III and a further two in Theorems 5 and 7 below. The approaches given here have various advantages which were outlined in Section I-C, and which are discussed further in Section IV-D. Throughout the section, expectations are written using summations for notational simplicity (e.g. $\mathbb{E}_Q[f(X)] = \sum_x Q(x)f(x)$). However, we will highlight that certain results apply in the case of continuous alphabets upon replacing the summations by integrals.

A. Derivations Using Theorem 1

1) *i.i.d. ensemble:* We immediately obtain the exponent in (26), as well as its generalization to continuous alphabets, by substituting the i.i.d. distribution in (22) into $\text{rcux}_{\rho, s}$ in (9).

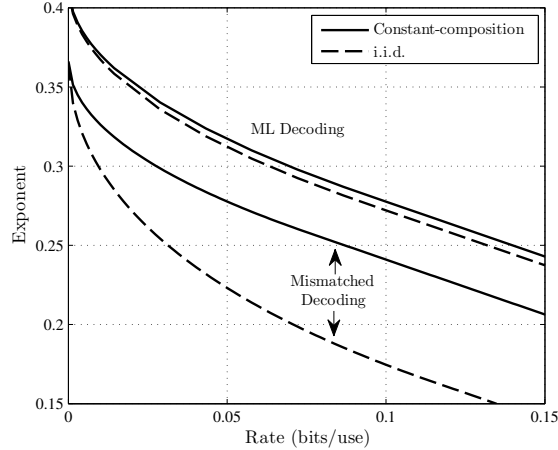


Figure 1. Expurgated exponents for the channel described in (42) with minimum Hamming distance decoding and ML decoding. The parameters are $\delta_0 = 0.01$, $\delta_1 = 0.05$, $\delta_2 = 0.25$ and $Q = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

2) *Constant-composition Ensemble*: In the case of finite alphabets, the method of types [5, Ch. 2] can be used to obtain the exact exponents corresponding to rcux_ρ and $\text{rcux}_{\rho,s}$ for each of the ensembles defined in (22)–(24). The analysis is similar for each of these, so we focus on the constant-composition ensemble described by (23). We define

$$\mathcal{S}(Q) \triangleq \left\{ \tilde{P}_{X\bar{X}} \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) : \tilde{P}_X = Q, \tilde{P}_{\bar{X}} = Q \right\} \quad (43)$$

$$\mathcal{T}(\tilde{P}_{X\bar{X}}) \triangleq \left\{ P_{X\bar{X}Y} \in \mathcal{P}(\mathcal{X} \times \mathcal{X} \times \mathcal{Y}) : P_{X\bar{X}} = \tilde{P}_{X\bar{X}}, \mathbb{E}_P[\log q(\bar{X}, Y)] \geq \mathbb{E}_P[\log q(X, Y)] \right\} \quad (44)$$

$$\mathcal{S}_n(Q) \triangleq \mathcal{S}(Q) \cap \mathcal{P}_n(\mathcal{X} \times \mathcal{X}) \quad (45)$$

$$\mathcal{T}_n(\tilde{P}_{X\bar{X}}) \triangleq \mathcal{T}(\tilde{P}_{X\bar{X}}) \cap \mathcal{P}_n(\mathcal{X} \times \mathcal{X} \times \mathcal{Y}). \quad (46)$$

where we overload the symbol \mathcal{T} (see (29)). It follows that $P_{X\bar{X}Y} \in \mathcal{T}(Q)$ (defined in (29)) if and only if $P_{X\bar{X}Y} \in \mathcal{T}(\tilde{P}_{X\bar{X}})$ (defined in (44)) for some $\tilde{P}_{X\bar{X}} \in \mathcal{S}(Q)$. We note the following properties of types [5, Ch. 2]:

1) For any $\tilde{P}_{X\bar{X}} \in \mathcal{S}_n(Q_n)$,

$$\mathbb{P}[(\mathbf{X}, \bar{\mathbf{X}}) \in T^n(\tilde{P}_{X\bar{X}})] \doteq e^{-nI_{\tilde{P}}(X; \bar{X})}. \quad (47)$$

2) If $(\mathbf{x}, \bar{\mathbf{x}}) \in T^n(\tilde{P}_{X\bar{X}})$, then for any $P_{X\bar{X}Y} \in \mathcal{T}_n(\tilde{P}_{X\bar{X}})$,

$$\mathbb{P}[(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}) \in T^n(P_{X\bar{X}Y}) \mid \mathbf{X} = \mathbf{x}] \doteq e^{-nD(P_{X\bar{X}Y} \parallel \tilde{P}_{X\bar{X}} \times W)}. \quad (48)$$

Theorem 4. Consider a discrete memoryless channel, and let the codeword distribution $P_{\mathbf{X}}$ be the constant-composition distribution in (23) for some input distribution Q . The bound rcux_ρ in (8) satisfies the following for any rate $R > 0$:

$$\inf_{\rho \geq 1} \text{rcux}_\rho(n, e^{nR}) \doteq e^{-nE_{\text{ex}}^{\text{cc}}(Q, R)}. \quad (49)$$

Proof: Using the codeword distribution in (23) and expanding (8) in terms of types, we obtain

$$\text{rcux}_\rho(n, M)^{1/\rho} = 4(M-1) \sum_{\tilde{P}_{X\bar{X}} \in \mathcal{S}_n(Q_n)} \mathbb{P}[(\mathbf{X}, \bar{\mathbf{X}}) \in T^n(P_{X\bar{X}})] \sum_{P_{X\bar{X}Y} \in \mathcal{T}_n(\tilde{P}_{X\bar{X}})} \mathbb{P}[(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}) \in T^n(P_{X\bar{X}Y}) \mid \mathbf{X} = \mathbf{x}]^{1/\rho} \quad (50)$$

$$\doteq M \max_{\tilde{P}_{X\bar{X}} \in \mathcal{S}_n(Q_n)} \max_{P_{X\bar{X}Y} \in \mathcal{T}_n(\tilde{P}_{X\bar{X}})} \exp\left(-nI_{\tilde{P}}(X; \bar{X})\right) \exp\left(-n \cdot \frac{1}{\rho} D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W)\right) \quad (51)$$

$$\doteq M \max_{P_{X\bar{X}Y} \in \mathcal{T}(Q)} \exp\left(-n\left(\frac{1}{\rho} D(P_{X\bar{X}Y} \| P_{X\bar{X}} \times W) + I_P(X; \bar{X})\right)\right), \quad (52)$$

where in (50) we define $(\mathbf{x}, \bar{\mathbf{x}})$ to be an arbitrary pair with joint type $\tilde{P}_{X\bar{X}}$, (51) follows from the properties of types in (47)–(48) and the fact that the number of joint types is polynomial in n , and (52) follows from the definitions of \mathcal{S}_n , \mathcal{T}_n and \mathcal{T} , and by using a standard continuity argument to expand the maximization from types to general distributions (e.g. see [24]). We thus obtain the exponent

$$\sup_{\rho \geq 1} \min_{P_{X\bar{X}Y} \in \mathcal{T}(Q)} D(P_{X\bar{X}Y} \| P_{X\bar{X}} \times W) + \rho(I_P(X; \bar{X}) - R) \quad (53)$$

$$= \min_{P_{X\bar{X}Y} \in \mathcal{T}(Q)} \sup_{\rho \geq 1} D(P_{X\bar{X}Y} \| P_{X\bar{X}} \times W) + \rho(I_P(X; \bar{X}) - R), \quad (54)$$

where (54) follows from Fan's minimax theorem [25], the conditions of which are satisfied here since the objective is linear in ρ and convex in $P_{X\bar{X}Y}$. Using

$$\sup_{\rho \geq 1} \rho \alpha = \begin{cases} \infty & \alpha > 0 \\ \alpha & \alpha \leq 0 \end{cases} \quad (55)$$

and the identity in (30), it follows that (54) coincides with (28). \blacksquare

The preceding derivation of $E_{\text{ex}}^{\text{cc}}$ provides a simple alternative to that of Csiszár and Körner [6], while yielding the exponent in the same form.

3) *Cost-constrained Ensemble:* Here we provide a derivation of $E_{\text{ex}}^{\text{cc}}$ in the form given in (34), as well as its generalization to continuous alphabets, using the cost-constrained ensemble in (24). We allow for a system cost constraint of the form given in (2). A key property of the ensemble which will prove useful in the derivations is

$$\mathbf{x} \in \mathcal{D}_n \implies e^{r\left(\sum_{i=1}^n a(x_i) - n\phi_a\right)} e^{|\mathbf{r}|\delta} \geq 1, \quad (56)$$

which holds for any real number r , and follows immediately from (25). Furthermore, we have the following.

Proposition 2. [13, Prop. 1] *Fix any input distribution Q and set of cost functions $\{a_l\}_{l=1}^L$ such that $E_Q[c(X)] \leq \Gamma$, $E_Q[c(X)^2] < \infty$ and $E_Q[a_l(X)^2] < \infty$ for $l = 1, \dots, L$. Then the normalizing constant μ_n in (24) satisfies*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \mu_n = 0. \quad (57)$$

The following theorem gives an achievable error exponent for a fixed set of auxiliary costs.

Theorem 5. *Consider a memoryless (possibly continuous) channel, and fix any input distribution Q and functions $\{a_l\}$ satisfying the assumptions of Proposition 2. Under the cost-constrained distribution in (24), we have*

$$\inf_{\rho \geq 1, s \geq 0} \text{rcux}_{\rho, s}(n, e^{nR}) \stackrel{\leq}{\leq} e^{-nE_{\text{ex}}^{\text{cost}}(Q, R, \{a_l\})} \quad (58)$$

for any rate $R > 0$, where

$$E_{\text{ex}}^{\text{cost}}(Q, R, \{a_l\}) \triangleq \sup_{\rho \geq 1} E_{\text{x}}^{\text{cost}}(Q, \rho, \{a_l\}) - \rho R, \quad (59)$$

and¹

$$E_{\text{x}}^{\text{cost}}(Q, R, \{a_l\}) \triangleq \sup_{s \geq 0, \{r_l\}, \{\bar{r}_l\}} -\rho \log \sum_{x, \bar{x}} Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \frac{e^{\sum_{l=1}^L \bar{r}_l (a_l(\bar{x}) - \phi_l)}}{e^{\sum_{l=1}^L r_l (a_l(x) - \phi_l)}} \right)^{1/\rho}. \quad (60)$$

Proof: Let $a_l^n(\mathbf{x}) \triangleq \sum_{i=1}^n a_l(x_i)$ and $Q^n(\mathbf{x}) \triangleq \prod_{i=1}^n Q(x_i)$. We start with (9), and write

$$\text{rcux}_{\rho, s}(n, M)^{1/\rho} = 4(M-1) \sum_{\mathbf{x}, \bar{\mathbf{x}}} P_{\mathbf{X}}(\mathbf{x})P_{\mathbf{X}}(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{y})}{q^n(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho} \quad (61)$$

$$\leq M \sum_{\mathbf{x}, \bar{\mathbf{x}}} P_{\mathbf{X}}(\mathbf{x})P_{\mathbf{X}}(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{y})}{q^n(\mathbf{x}, \mathbf{y})} \right)^s \frac{e^{\sum_{l=1}^L \bar{r}_l (a_l^n(\bar{\mathbf{x}}) - n\phi_l)}}{e^{\sum_{l=1}^L r_l (a_l^n(\mathbf{x}) - n\phi_l)}} \right)^{1/\rho} \quad (62)$$

$$\leq M \sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x})Q^n(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{y})}{q^n(\mathbf{x}, \mathbf{y})} \right)^s \frac{e^{\sum_{l=1}^L \bar{r}_l (a_l^n(\bar{\mathbf{x}}) - n\phi_l)}}{e^{\sum_{l=1}^L r_l (a_l^n(\mathbf{x}) - n\phi_l)}} \right)^{1/\rho}, \quad (63)$$

where (62) holds for any $\{r_l\}$ and $\{\bar{r}_l\}$ from (56), and (63) follows from (24) and Proposition 2. The proof is concluded by expanding each term in (63) as a product from 1 to n and optimizing ρ , s , $\{r_l\}$ and $\{\bar{r}_l\}$. ■

We now show that we can recover $E_{\text{ex}}^{\text{cc}}$ from $E_{\text{ex}}^{\text{cost}}$ upon setting $L = 2$ and optimizing the auxiliary costs; an analogous statement was shown to be true for the random-coding exponent in [13]. Setting $\bar{r}_1 = r_2 = 1$ and $\bar{r}_2 = r_1 = 0$, and optimizing $a_1(\cdot)$ and $a_2(\cdot)$, we obtain

$$E_{\text{x}}^{\text{cost}}(Q, \rho) = \sup_{s \geq 0, a_1(\cdot), a_2(\cdot)} -\rho \log \sum_{x, \bar{x}} Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \frac{e^{a_1(\bar{x}) - \phi_1}}{e^{a_2(x) - \phi_2}} \right)^{1/\rho} \quad (64)$$

$$\leq \sup_{s \geq 0, a_1(\cdot), a_2(\cdot)} -\rho \sum_x Q(x) \log \sum_{\bar{x}} Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \frac{e^{a_1(\bar{x}) - \phi_1}}{e^{a_2(x) - \phi_2}} \right)^{1/\rho}, \quad (65)$$

where (65) follows from Jensen's inequality. For any s and $a_1(\cdot)$, there exists a choice of $a_2(\cdot)$ that makes Jensen's inequality hold with equality in (65), and hence the same is true after taking the supremum. Hence, and by writing

$$-\sum_x Q(x) \log \left(\frac{e^{-\phi_1}}{e^{a_2(x) - \phi_2}} \right)^{1/\rho} = -\sum_x Q(x) \log (e^{-a_1(x)})^{1/\rho} = \frac{\phi_1}{\rho}, \quad (66)$$

we see that the $a_2(\cdot)$ achieving the supremum in (64) is the one yielding equality in (65). Renaming $a_1(\cdot)$ as $a(\cdot)$ and using the first equality in (66), we obtain (34).

It should be noted that, in accordance with Proposition 2, the supremum over s and $a(\cdot)$ in (34) is restricted to choices such that $E_Q[a(X)^2] < \infty$, and such that $E_Q[a_2(X)^2] < \infty$ for the choice of $a_2(\cdot)$ which makes Jensen's inequality hold with equality in (65) (expressed in terms of s and $a(\cdot)$). This may rule out some parameters in the case of infinite or continuous alphabets.

While the parameters $\{r_l\}$ and $\{\bar{r}_l\}$ are not necessary for obtaining (64), they can improve the exponent for a given set of auxiliary costs [13]. That is, the more general exponent of Theorem 5 serves as an indicator of the performance when the auxiliary costs are chosen suboptimally. Using a similar argument to that of (64)–(66), it is easily shown that $E_{\text{ex}}^{\text{cost}} \leq E_{\text{ex}}^{\text{cc}}$, and hence one cannot improve on the exponent obtained using $L = 2$ optimally chosen auxiliary costs.

¹In the case of continuous alphabets, the summations over sequences should be replaced by integrals.

B. Derivation Using Type Class Enumerators

In the proof of Theorem 4, we gave an exponentially tight analysis of rcux_ρ . In this subsection, we show that an exponentially tight analysis can be provided starting from an earlier step using the method of type class enumeration (e.g. see [7]–[9]). Once again, the analysis is similar for each of the ensembles in (22)–(24), so we focus on the constant-composition ensemble described by (23).

Substituting (10) into (6) and defining

$$d_q(\mathbf{x}, \bar{\mathbf{x}}) \triangleq -\log \mathbb{P}[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \mid \mathbf{X} = \mathbf{x}], \quad (67)$$

we obtain the bound

$$p_e(\mathcal{C}) \leq (2A_n(R, \rho))^\rho, \quad (68)$$

where

$$A_n(R, \rho) \triangleq \mathbb{E} \left[\left(\sum_{\bar{m} \neq m} e^{-d_q(\mathbf{X}^{(m)}, \mathbf{X}^{(\bar{m})})} \right)^{1/\rho} \right]. \quad (69)$$

This bound provides the starting point for our analysis. Note that since we have not used the inequality in (12), we may allow for $\rho \geq 0$ rather than just $\rho \geq 1$.

Theorem 6. *Consider a discrete memoryless channel, and let the codeword distribution $P_{\mathbf{X}}$ be the constant-composition distribution in (23) for some input distribution Q . Then the following holds for any rate $R > 0$:*

$$\inf_{\rho \geq 0} (2A_n(R, \rho))^\rho \doteq e^{-nE_{\text{ex}}^{\text{cc}}(Q, R)}. \quad (70)$$

Proof: For $m = 1, \dots, M$ and each joint type $\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}$, we define the random variable

$$N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}) \triangleq \sum_{\bar{m} \neq m} \mathbb{1}\{(\mathbf{X}^{(m)}, \mathbf{X}^{(\bar{m})}) \in T^n(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})\}. \quad (71)$$

Under the random-coding distribution in (23), we have $N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}) = 0$ with probability one if $\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}} \notin \mathcal{S}_n(Q_n)$. That is, the marginal distribution of each codeword must agree with Q . Since d_q depends only on the joint type of its arguments, we define $d_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}) \triangleq \frac{1}{n} d_q(\mathbf{x}, \bar{\mathbf{x}})$, where $(\mathbf{x}, \bar{\mathbf{x}}) \in T^n(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})$.

Making repeated use of the fact that the number of joint types is polynomial in n , we have the following:

$$A_n(R, \rho) = \mathbb{E} \left[\left(\sum_{\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}} N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}) e^{-nd_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})} \right)^{1/\rho} \right] \quad (72)$$

$$\doteq \mathbb{E} \left[\max_{\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}} N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})^{1/\rho} e^{-nd_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})/\rho} \right] \quad (73)$$

$$\doteq \mathbb{E} \left[\sum_{\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}} N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})^{1/\rho} e^{-nd_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})/\rho} \right] \quad (74)$$

$$\doteq \max_{\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}} \mathbb{E} \left[N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})^{1/\rho} \right] e^{-nd_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})/\rho}, \quad (75)$$

where (75) follows by first taking the summation outside the expectation. It follows from (75) that

$$(2A_n(R, \rho))^\rho \doteq \max_{\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}}} \left(\mathbb{E} \left[N_m(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})^{1/\rho} \right] \right)^\rho e^{-nd_q(\tilde{P}_{\mathbf{X}\bar{\mathbf{X}}})}. \quad (76)$$

Similarly to [7, Eq. (34)], we have for all $\tilde{P}_{X\bar{X}} \in \mathcal{S}_n(Q_n)$ that

$$\mathbb{E}\left[N_m(\tilde{P}_{X\bar{X}})^{1/\rho}\right] \doteq \begin{cases} \exp(n(R - I_{\tilde{P}}(X; \bar{X}))) & R < I_{\tilde{P}}(X; \bar{X}) \\ \exp(n(R - I_{\tilde{P}}(X; \bar{X}))/\rho) & R \geq I_{\tilde{P}}(X; \bar{X}). \end{cases} \quad (77)$$

This follows from the fact that given $\mathbf{X}^{(m)} = \mathbf{x}$, $N_m(\tilde{P}_{X\bar{X}})$ is the sum of $e^{nR} - 1$ binary independent random variables,

$$U_{\bar{m}} \triangleq \mathbf{1}\left\{(\mathbf{x}, \mathbf{X}^{(\bar{m})}) \in T^n(\tilde{P}_{X\bar{X}})\right\}, \quad (78)$$

whose expectations are of the exponential order of $e^{-nI_{\tilde{P}}(X; \bar{X})}$ (see (47)). Furthermore, expanding (67) in terms of types and using the property in (48), we obtain

$$e^{-nd_q(\tilde{P}_{X\bar{X}})} \doteq \exp\left(-n \min_{P_{X\bar{X}Y} \in \mathcal{T}(\tilde{P}_{X\bar{X}})} D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W)\right) \quad (79)$$

$$\triangleq e^{-nD_q(\tilde{P}_{X\bar{X}})}. \quad (80)$$

Upon taking into account all the possible empirical distributions $\{\tilde{P}_{X\bar{X}}\}$ in (78), we obtain

$$(2A_n(R, \rho))^\rho \doteq e^{-n \min\{E_1(R, \rho), E_2(R)\}}, \quad (81)$$

where

$$E_1(R, \rho) \triangleq \min_{\substack{\tilde{P}_{X\bar{X}} \in \mathcal{S}(Q) \\ I_{\tilde{P}}(X; \bar{X}) \geq R}} D_q(\tilde{P}_{X\bar{X}}) + \rho(I_{\tilde{P}}(X; \bar{X}) - R) \quad (82)$$

and

$$E_2(R) = \min_{\substack{\tilde{P}_{X\bar{X}} \in \mathcal{S}(Q) \\ I_{\tilde{P}}(X; \bar{X}) \leq R}} D_q(\tilde{P}_{X\bar{X}}) + I_{\tilde{P}}(X; \bar{X}) - R. \quad (83)$$

Combining (30), (80) and (83), we see that $E_2(R)$ coincides with $E_{\text{ex}}^{\text{cc}}$ in the form given in (28). It remains to show that $E_1(R, \rho)$, for the optimum choice of ρ , is never smaller than $E_2(R)$. This can be seen by noting that since (82) contains the constraint $I_{\tilde{P}}(X; \bar{X}) \geq R$, the term multiplying ρ in (82) is non-negative. Thus, the best choice of ρ is to take the limit as $\rho \rightarrow \infty$, and hence the minimum in (82) is achieved by some $\tilde{P}_{X\bar{X}}$ satisfying $I_{\tilde{P}}(X; \bar{X}) = R$. Since this joint distribution also satisfies the constraints in (83), we conclude that $E_1 \geq E_2$, thus completing the proof. ■

While the exponents of Theorems 4 and 6 coincide for the constant-composition ensemble, the type enumeration approach can yield strictly higher exponents for other codeword distributions; see Section IV-D for details.

C. Derivation Using Distance Enumerators

In this subsection, we extend the preceding type enumeration analysis to channels with infinite or continuous alphabets, and then discuss the further extension to channels with memory. We make use of Theorem 2, and we assume that the technical assumption therein is satisfied (see Appendix A for discussion). We fix $s \geq 0$ and make use of d_s in (33) (or its counterpart for continuous outputs with an integral in place of the summation), as well as its multi-letter extension

$$d_s^n(\mathbf{x}, \bar{\mathbf{x}}) \triangleq \sum_{i=1}^n d_s(x_i, \bar{x}_i). \quad (84)$$

Theorem 7. Consider a memoryless (possibly continuous) channel, and fix any codeword distribution $P_{\mathbf{X}}$ satisfying the assumption of Theorem 2. The exponent

$$E_{\text{ex}}(R) \triangleq \mathbb{E} \left[\inf_{D: R(D, \mathbf{X}) \leq R} D + R(D, \mathbf{X}) - R \right] \quad (85)$$

is achievable for any function $R(D, \mathbf{x})$ such that $\mathbb{P}[d_s^n(\mathbf{x}, \bar{\mathbf{X}}) < nD] \leq e^{-nR(D, \mathbf{x})}$ uniformly in \mathbf{x} , and such that $R(\cdot, \mathbf{x})$ is continuous for any given \mathbf{x} .

Proof: We claim that (16) implies the following analog of (68) for a sequence of codebook \mathcal{C}_n of rate approaching R :

$$p_e(\mathcal{C}_n) \leq \exp \left(\rho \mathbb{E} [\log A_n(R, \rho, \mathbf{X}^{(m)})] \right), \quad (86)$$

where

$$A_n(R, \rho, \mathbf{X}^{(m)}) \triangleq \mathbb{E} \left[\left(\sum_{\bar{m} \neq m} e^{[-d_s^n(\mathbf{X}^{(m)}, \mathbf{X}^{(\bar{m})})]^+} \right)^{1/\rho} \middle| \mathbf{X}^{(m)} \right]. \quad (87)$$

In the absence of the $[\cdot]^+$ function in the exponent, this follows directly from the union bound and Markov's inequality, similarly to the proof of Theorem 1. The introduction of the $[\cdot]^+$ function corresponds to instead taking the better of Markov's inequality and the trivial bound $\mathbb{P}[\cdot] \leq 1$.²

For a fixed transmitted codeword $\mathbf{X}^{(m)} = \mathbf{x}$, we analyze $A_n(R, \rho, \mathbf{x})$ using *distance enumerators*:

$$\sum_{\bar{m} \neq m} e^{-[d_s^n(\mathbf{x}, \mathbf{X}^{(\bar{m})})]^+} \leq \sum_{k=0}^{\infty} e^{-nk\delta} N_m(k, \mathbf{x}), \quad (88)$$

where $\delta > 0$ is arbitrary, and

$$N_m(0, \mathbf{x}) \triangleq \sum_{\bar{m} \neq m} \mathbb{1} \{ d_s^n(\mathbf{x}, \mathbf{X}^{(\bar{m})}) < n\delta \} \quad (89)$$

$$N_m(k, \mathbf{x}) \triangleq \sum_{\bar{m} \neq m} \mathbb{1} \{ nk\delta \leq d_s^n(\mathbf{x}, \mathbf{X}^{(\bar{m})}) < n(k+1)\delta \} \quad (k \geq 1). \quad (90)$$

Using Markov's inequality, we can upper-bound the left-hand side of (13) by $e^{-d_s^n(\mathbf{x}, \bar{\mathbf{x}})}$. It thus follows from the assumption of Theorem 2 that the highest value of k ,

$$k_{\max}(n) \triangleq \max_{\mathbf{x}: P_{\mathbf{X}}(\mathbf{x}) > 0} \max \{ k : \mathbb{P}[N_m(k, \mathbf{x}) > 0] \neq 0 \}, \quad (91)$$

grows subexponentially in n for all $s \geq 0$. Thus, analogously to (76), the quantity $A_n(R, \rho, \mathbf{x})$ defined in (87) satisfies

$$A_n(R, \rho, \mathbf{x})^\rho \leq \max_{k \geq 0} \left(\mathbb{E}[N_m(k, \mathbf{x})^{1/\rho}] \right)^\rho e^{-nk\delta}. \quad (92)$$

We further upper bound this expression by removing the lower inequality in the indicator function in (90). The key issue is now to assess the exponential rate of decay of the binary random variable

$$U_{\bar{m}}(\mathbf{x}) \triangleq \mathbb{1} \{ d_s^n(\mathbf{x}, \mathbf{X}^{(\bar{m})}) < n(k+1)\delta \} \quad (93)$$

²This analysis corrects an error in the conference version of this work [26], where the $[\cdot]^+$ function was omitted. This omission does not affect the analysis for ML decoding, since the Bhattacharyya distance is non-negative. However, in general, the function $d_s(\cdot, \cdot)$ may be negative.

for a given transmitted codeword \mathbf{x} , i.e. to find the exponent of $\mathbb{P}[d_s^n(\mathbf{x}, \bar{\mathbf{X}}) < nD]$. This can be done using standard large deviations techniques such as the Chernoff bound. Letting $R(D, \mathbf{x})$ be as defined in the theorem statement, we have similarly to (81) that

$$A_n(R, \rho, \mathbf{x})^\rho \stackrel{\dot{\leq}}{e^{-n \min\{E_1(R, \rho, \delta, \mathbf{x}), E_2(R, \delta, \mathbf{x})\}}}, \quad (94)$$

where

$$E_1(R, \rho, \delta, \mathbf{x}) \triangleq \min_{k: R((k+1)\delta, \mathbf{x}) \geq R} k\delta + \rho(R((k+1)\delta, \mathbf{x}) - R) \quad (95)$$

$$E_2(R, \delta, \mathbf{x}) \triangleq \min_{k: R((k+1)\delta, \mathbf{x}) \leq R} k\delta + R((k+1)\delta, \mathbf{x}) - R. \quad (96)$$

Upon taking the limit $\delta \rightarrow 0$ and using the assumption that $R(\cdot, \mathbf{x})$ is lower semicontinuous, these become

$$E_1(R, \rho, \mathbf{x}) \triangleq \inf_{D: R(D, \mathbf{x}) \geq R} D + \rho(R(D, \mathbf{x}) - R) \quad (97)$$

$$E_2(R, \mathbf{x}) \triangleq \inf_{D: R(D, \mathbf{x}) \leq R} D + R(D, \mathbf{x}) - R. \quad (98)$$

Analogously to Section IV-B, the optimal choice of ρ is in the limit as $\rho \rightarrow \infty$, and we obtain $E_2 \leq E_1$, and hence

$$\inf_{\rho \geq 0} A_n(R, \rho, \mathbf{x})^\rho \stackrel{\dot{\leq}}{e^{-nE_2(R, \mathbf{x})}}. \quad (99)$$

Substituting (99) into (86), we obtain $p_e(\mathcal{C}) \stackrel{\dot{\leq}}{e^{-n\mathbb{E}[E_2(R, \mathbf{X})]}}$, thus yielding (85). \blacksquare

After a suitable modification of the definition of $d_s^n(\mathbf{x}, \bar{\mathbf{x}})$, (85) extends immediately to more general channels and metrics (e.g. channels with memory). The ability to simplify the exponent (e.g. to a single-letter expression) depends on the form of $R(D, \mathbf{x})$, which in turn depends strongly on the codeword distribution $P_{\mathbf{X}}$. In some cases, $P_{\mathbf{X}}$ can be chosen in such a way that $R(D, \mathbf{x})$ is the same for all \mathbf{x} with $P_{\mathbf{X}}(\mathbf{x}) > 0$, thus greatly simplifying (85).

In Appendix D, we particularize Theorem 7 to the cost-constrained ensemble with a single auxiliary cost $a_1(x) = a(x)$, and show that after optimizing $a(\cdot)$, (85) yields the exponent $E_{\text{ex}}^{\text{cc}}(Q, R)$ in (32). In accordance with Proposition 2, we require the auxiliary cost $a(\cdot)$ to satisfy $\mathbb{E}_Q[a(X)^2] < \infty$.

D. Comparison of Techniques

For the constant-composition codeword distribution, the approaches of Sections IV-A and IV-B led to the same exponent, namely $E_{\text{ex}}^{\text{cc}}$. It should be noted, however, that the type enumeration approach can yield a strictly higher exponent than that of rcux_ρ in Theorem 1 for some codeword distributions. Here we discuss the simple example of the i.i.d. distribution in (22). Applying properties of types in the same way in Section IV-A, it is easily verified that the exponent of rcux_ρ is

$$\min_{\substack{P_{X\bar{X}Y}: D(P_{X\bar{X}}\|Q \times Q) \leq R, \\ \mathbb{E}_P[\log q(\bar{X}, Y)] \geq \mathbb{E}_P[\log q(X, Y)]}} D(P_{X\bar{X}Y}\|Q \times Q \times W) - R. \quad (100)$$

On the other hand, the analysis of Section IV-B yields an exponent of the same form as (100) with an additional constraint $P_X = Q$ in the minimization. To see this, we note that the quantity $N_m(\tilde{P}_{X\bar{X}})$ defined in (71) satisfies

$$\mathbb{E}\left[N_m(\tilde{P}_{X\bar{X}})^{1/\rho}\right] = \mathbb{P}[\mathbf{X}^{(m)} \in T^n(\tilde{P}_X)] \mathbb{E}\left[N_m(\tilde{P}_{X\bar{X}})^{1/\rho} \mid \mathbf{X}^{(m)} \in T^n(\tilde{P}_X)\right] \quad (101)$$

$$\doteq \begin{cases} \exp(-nD(\tilde{P}_X\|Q)) \cdot \exp(n(R - D(\tilde{P}_{X\bar{X}}\|\tilde{P}_X \times Q))) & R < I_{\tilde{P}}(X; \bar{X}) \\ \exp(-nD(\tilde{P}_X\|Q)) \cdot \exp(n(R - D(\tilde{P}_{X\bar{X}}\|\tilde{P}_X \times Q))/\rho) & R \geq I_{\tilde{P}}(X; \bar{X}). \end{cases} \quad (102)$$

The additional factor $\exp(-nD(\tilde{P}_X\|Q))$ leads to an additive $\rho D(\tilde{P}_X\|Q)$ term in the exponent E_2 in (83). The optimal choice of ρ is again in the limit as $\rho \rightarrow \infty$, and under this choice the minimizing $\tilde{P}_{X\bar{X}}$ must satisfy $\tilde{P}_X = Q$ so that the divergence is forced to zero.

Depending on the channel, metric and input distribution, adding the constraint $P_X = Q$ to (100) may yield a strict improvement in the exponent. Since both derivations are exponentially tight from the step at which they start, we conclude that the weakness of the simpler derivation is in the inequality in (11), or more precisely, the use of (12). While this step simplifies the derivations, the above example shows that it is not exponentially tight in general.

Another approach to recovering the constraint $\tilde{P}_X = Q$ in the above example is to follow the steps of Theorem 1 and Section IV-A starting with Theorem 2. Since the expectation of the transmitted codeword is outside the logarithm in (16), we obtain the constraint $\tilde{P}_X = Q$ in the final minimization using the fact that the empirical distribution of \mathbf{X} is close to Q with high probability. We conclude that the inequality in (12) is exponentially tight for the i.i.d. ensemble when we start with (16), even though it is not tight when we start with (6).

We have provided two derivations of $E_{\text{ex}}^{\text{cc}}$ using the cost-constrained ensemble, namely, those in Sections IV-A and IV-C (along with Appendix D). A notable difference between the derivations is the method for ensuring that the average over x is outside the logarithm in (34), which is desirable due to Jensen's inequality. In Theorem 5, the expectation is inside the logarithm, but the desired result is obtained by choosing $a_2(x)$ to make Jensen's inequality hold with equality. On the other hand, in Appendix D the expectation arises outside the logarithm even in the case that $L = 1$.

Provided that the assumption of Theorem 2 is met, we can combine the two approaches and apply the techniques of Theorem 1 and Section IV-A to (16), in which case E_x^{cost} in (60) is improved to

$$E_x^{\text{cost}*}(Q, R, \{a_l\}) \triangleq \sup_{s \geq 0, \{\bar{r}_l\}} -\rho \sum_x Q(x) \log \sum_{\bar{x}} Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s e^{\sum_{l=1}^L \bar{r}_l (a_l(\bar{x}) - \phi_l)} \right)^{1/\rho}, \quad (103)$$

where the outer-most summation arises using Proposition 3 in Appendix D. This exponent can also be derived by extending the analysis of Appendix D to include multiple auxiliary costs.

In the case that $L = 0$ (i.e. i.i.d. coding), the Lagrange duality techniques of Theorem 3 reveal that (103) is in fact identical to (100) with the added constraint $P_X = Q$. That is, the additional constraint $P_X = Q$ in the primal expression corresponds to an average over x outside the logarithm in the dual expression.

E. Connections with Statistical Mechanics

It is instructive to look at the analysis of Sections IV-B and IV-C from the statistical-mechanical perspective. Let us take another look at the expression

$$Z(\mathbf{x}) = \sum_{\bar{\mathbf{m}} \neq \mathbf{m}} e^{-d(\mathbf{x}, \mathbf{X}^{(\bar{\mathbf{m}})})}, \quad (104)$$

where d can represent either d_q (see (67)) or $[d_s^n]^+$ (see (84)). From the viewpoint of statistical physics, Z can be interpreted as the partition function of a physical system, where for a fixed $\mathbf{x}^{(m)} = \mathbf{x}$, the various configurations (microstates) are $\{\mathbf{x}^{(\bar{m})}\}_{\bar{m} \neq m}$ and the energy function (Hamiltonian) is given by $d(\mathbf{x}, \bar{\mathbf{x}})$. The various ‘‘configurational energies’’ $\{d(\mathbf{x}, \mathbf{X}^{(\bar{m})})\}$ are independent random variables, since the codewords are generated independently. As explained in [27, Ch. 5-6] (see also [9, Ch. 6-7] and references therein), this setting is analogous to the random energy

model (REM) in the literature of statistical physics of magnetic materials. The REM was invented by Derrida [28]–[30] as a model of extremely disordered spin glasses. This model is exactly solvable and exhibits a phase transition: Below a certain critical temperature, the partition function becomes dominated by a subexponential number of configurations in the ground-state energy, which means that the system freezes and its entropy vanishes in the thermodynamic limit. This combination of freezing and disorder resembles the behavior of a glass, so this low temperature phase of zero entropy is called the *glassy phase*. Above the critical temperature, the partition function is dominated by an exponential number of configurations, so its entropy is positive. This high temperature phase is called the *paramagnetic phase*.

In the case that $P_{\mathcal{X}}$ is the constant-composition distribution in (23) and $d(\cdot, \cdot)$ represents $[d_s^n(\cdot, \cdot)]^+$, we can link these phases to the exponent $E_{\text{ex}}^{\text{cc}}$ in the form given in (39). The graph of $E_{\text{ex}}^{\text{cc}}(Q, R, s)$ is curved at rates below R_s (see (36)), and is a straight line at rates above R_s . The curved part corresponds to the glassy phase of the REM associated with (104), because the dominant contribution to $\mathbb{E}[Z(\mathbf{x})^{1/\rho}]$ (see (104)) is due to a subexponential number of codewords whose “distance” from \mathbf{x} (i.e. their “energy”) is roughly $nD_s(Q, R)$. The straight-line part, on the other hand, corresponds to the paramagnetic phase, where roughly $e^{n(R-R_s)}$ incorrect codewords at distance $nD_s(Q, R_s)$ dominate the behavior. Thus, the passage between the curved part and the straight-line part at $R = R_s$ can be interpreted as a glassy phase transition. A similar discussion applies for the multi-letter distance d_q used in Section IV-B, with $D_s(Q, R)$ replaced by

$$D_q(Q, R) \triangleq \min_{\tilde{P}_{\mathcal{X}\bar{\mathcal{X}}} \in \mathcal{S}(Q) : I_{\tilde{P}}(X; \bar{X}) \leq R} D_q(\tilde{P}_{\mathcal{X}\bar{\mathcal{X}}}), \quad (105)$$

where $D_q(\tilde{P}_{\mathcal{X}\bar{\mathcal{X}}})$ is defined in (80).

V. PREFACTOR TO THE I.I.D. EXPURGATED EXPONENT

Error exponents characterize the rate of decay of the error probability as the block length increases. At finite block lengths, the effect of the subexponential prefactor can be significant, and it is therefore of interest to characterize its behavior. There exist several works studying this prefactor for the random-coding exponent [13], [19], [31], [32] and the sphere-packing exponent [31]–[33]. In this section, we characterize the prefactor for the i.i.d. expurgated exponent. We will see that, under some technical conditions, the prefactor to rcux_ρ in (8) behaves as $O(\frac{1}{\sqrt{n}})$, thus improving on Gallager’s $O(1)$ prefactor. Our analysis builds on that of [13], [20].

A. Preliminary Definitions

We define the sets

$$\mathcal{Y}_1(x, \bar{x}) \triangleq \left\{ y : W(y|x)W(y|\bar{x}) > 0 \right\} \quad (106)$$

$$\mathcal{A}(Q) \triangleq \left\{ (x, \bar{x}) : Q(x)Q(\bar{x}) > 0, \frac{q(\bar{x}, y)}{q(x, y)} \neq \frac{q(\bar{x}, y')}{q(x, y')} \text{ for some } y, y' \in \mathcal{Y}_1(x, \bar{x}) \right\} \quad (107)$$

and make the following technical assumptions:

$$q(x, y) = 0 \iff W(y|x) = 0 \quad (108)$$

$$\mathcal{A}(Q) \neq \emptyset. \quad (109)$$

In the case that $q(x, y) = W(y|x)$ (i.e. ML decoding), (108) is trivial, and (109) reduces to the *non-singularity* assumption of [19]. A notable example where this condition fails is the binary erasure channel (BEC) with $Q = (\frac{1}{2}, \frac{1}{2})$.

We write

$$E_x^{\text{iid}}(Q, \rho, s) \triangleq -\rho \log \sum_{x, \bar{x}} Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \right)^{1/\rho} \quad (110)$$

to denote the objective in (27) with a fixed value of s . We define the tiled distribution

$$V_s(y|x, \bar{x}) \triangleq \frac{W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s}{\sum_{y'} W(y'|x) \left(\frac{q(\bar{x}, y')}{q(x, y')} \right)^s} \quad (111)$$

$$V_s^n(\mathbf{y}|\mathbf{x}, \bar{\mathbf{x}}) \triangleq \prod_{i=1}^n V_s(y_i|x_i, \bar{x}_i), \quad (112)$$

and the generalized information density

$$j_s(x, \bar{x}, y) \triangleq \log \frac{V_s(y|x, \bar{x})}{W(y|x)} \quad (113)$$

$$j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{y}) \triangleq \sum_{i=1}^n j_s(x_i, \bar{x}_i, y_i). \quad (114)$$

Furthermore, we define the joint tilted distribution

$$P_{\rho, s}^*(x, \bar{x}) = \frac{Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \right)^{1/\rho}}{\sum_{x', \bar{x}'} Q(x')Q(\bar{x}') \left(\sum_{y'} W(y'|x') \left(\frac{q(\bar{x}', y')}{q(x', y')} \right)^s \right)^{1/\rho}}, \quad (115)$$

and the conditional variance

$$c_0(Q, \rho, s) \triangleq \mathbb{E} \left[\text{Var} [j_s(X_s^*, \bar{X}_s^*, Y_s^*) | X_s^*, \bar{X}_s^*] \right], \quad (116)$$

where $(X_s^*, \bar{X}_s^*, Y_s^*) \sim P_{\rho, s}^*(x, \bar{x}) V_s(y|x, \bar{x})$. The arguments to c_0 will henceforth be omitted, since their values will be understood from the context.

Writing $Y_s \sim V_s(\cdot|x, \bar{x})$, the following arguments show that the assumptions in (108)–(109) imply that $c_0 > 0$ whenever $s > 0$:

$$\text{Var}[j_s(x, \bar{x}, Y_s)] = 0 \iff j_s(x, \bar{x}, y) \text{ is independent of } y \text{ wherever } V_s(y|x, \bar{x}) > 0 \quad (117)$$

$$\iff \frac{q(\bar{x}, y)}{q(x, y)} \text{ is independent of } y \text{ wherever } W(y|x)q(\bar{x}, y) > 0 \quad (118)$$

$$\iff (x, \bar{x}) \notin \mathcal{A}(Q), \quad (119)$$

where (118) follows from the definitions of j_s and V_s , and (119) follows from the assumption in (108) and the definition of $\mathcal{A}(Q)$. Using the assumption in (109), it follows that $c_0 > 0$.

Finally, we define the set

$$\mathcal{I}_s \triangleq \left\{ j_s(x, x, y) : W(y|x) > 0, (x, \bar{x}) \in \mathcal{A}(Q) \right\} \quad (120)$$

and the constant

$$\psi_s \triangleq \begin{cases} 1 & \mathcal{I}_s \text{ does not lie on a lattice} \\ \frac{\bar{h}}{1-e^{-\bar{h}}} & \mathcal{I}_s \text{ lies on a lattice with span } \bar{h}. \end{cases} \quad (121)$$

B. Statement of the Result

Theorem 8. Fix any DMC W , decoding metric q and input distribution Q satisfying (108)–(109). For any $R > 0$, $\rho \geq 1$ and $s > 0$, there exists a sequence of codebooks \mathcal{C}_n with $M \geq e^{nR}$ codewords whose maximal error probability satisfies

$$p_e(\mathcal{C}_n) \leq \frac{4^{\rho} \psi_s}{\sqrt{2\pi n c_0}} e^{-n(E_{\mathbf{x}}^{\text{iid}}(Q, \rho, s) - \rho R)} (1 + o(1)) \quad (122)$$

Proof: See Section V-C. ■

It is interesting to note that under ML coding and any rate where the expurgated exponent and random-coding exponent coincide (i.e. $\rho = 1$ in both cases), Theorem 8 gives the same prefactor growth rate as that of the random-coding exponent [13], [19]. There is an extra factor of four in (122), which can be attributed to the fact that Theorem 8 considers the maximal error rather than the average error. Of course, Theorem 8 is primarily of interest at low rates, where the expurgated exponent exceeds the random-coding exponent.

C. Proof of Theorem 8

The proof makes use of two technical lemmas. The first is a strong large deviations result which was proved in [13], building upon the analysis in the proof of [3, Lemma 47].

Lemma 2. [13, Lemma 1] Fix $K > 0$, and for each n , let (n_1, \dots, n_K) be integers such that $\sum_k n_k = n$. Fix the PMFs Q_1, \dots, Q_K on a finite subset of \mathbb{R} , and let $\sigma_1^2, \dots, \sigma_K^2$ be the corresponding variances. Let Z_1, \dots, Z_n be independent random variables, n_k of which are distributed according to Q_k for each k . Suppose that $\min_k \sigma_k > 0$ and $\min_k n_k = \Theta(n)$. Defining

$$\mathcal{I}_0 \triangleq \bigcup_{k: \sigma_k > 0} \{z : Q_k(z) > 0\} \quad (123)$$

$$\psi_0 \triangleq \begin{cases} 1 & \mathcal{I}_0 \text{ does not lie on a lattice} \\ \frac{h_0}{1 - e^{-h_0}} & \mathcal{I}_0 \text{ lies on a lattice with span } h_0, \end{cases} \quad (124)$$

the summation $S_n \triangleq \sum_i Z_i$ satisfies the following uniformly in t :

$$\mathbb{E} \left[e^{-S_n} \mathbf{1} \{S_n \geq t\} \right] \leq e^{-t} \left(\frac{\psi_0}{\sqrt{2\pi V_n}} + o\left(\frac{1}{\sqrt{n}}\right) \right), \quad (125)$$

where $V_n \triangleq \text{Var}[S_n]$.

The following lemma ensures the existence of a high probability set of $(\mathbf{x}, \bar{\mathbf{x}})$ pairs such that Lemma 2 can be applied to the inner probability in (8).

Lemma 3. For any $R > 0$, $\rho \geq 1$, $s > 0$ and (W, q, Q) satisfying (108)–(109), the sequence of sets

$$\mathcal{F}_{\rho, s}^n(\delta) \triangleq \left\{ (\mathbf{x}, \bar{\mathbf{x}}) : \max_{x, \bar{x}} \left| \hat{P}_{\mathbf{x}\bar{\mathbf{x}}}(x, \bar{x}) - P_{\rho, s}^*(x, \bar{x}) \right| \leq \delta \right\} \quad (126)$$

satisfies the following properties:

- 1) For any $\delta > 0$ and $(\mathbf{x}, \bar{\mathbf{x}}) \in \mathcal{F}_{\rho, s}^n(\delta)$, the random variable $\mathbf{Y}_s \sim V_s^n(\cdot | \mathbf{x}, \bar{\mathbf{x}})$ satisfies

$$\text{Var}[j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}_s)] \geq n(c_0 - r(\delta)), \quad (127)$$

where $r(\delta) \rightarrow 0$ as $\delta \rightarrow 0$.

2) For any $\delta > 0$, we have

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \frac{\sum_{(\mathbf{x}, \bar{\mathbf{x}}) \notin \mathcal{F}_{\rho, s}^n(\delta)} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho}}{\sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho}} > 0. \quad (128)$$

Proof of Theorem 8 Based on Lemma 3: Using the bound rcux_ρ in Theorem 1 with the i.i.d. codeword distribution $P_{\mathbf{X}} = Q^n$, we have for any $\delta > 0$ that

$$\frac{1}{4(M-1)} \text{rcux}_\rho(n, M)^{1/\rho} = \sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right]^{1/\rho} \quad (129)$$

$$\begin{aligned} &= \sum_{(\mathbf{x}, \bar{\mathbf{x}}) \in \mathcal{F}_{\rho, s}^n(\delta)} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right]^{1/\rho} \\ &\quad + \sum_{(\mathbf{x}, \bar{\mathbf{x}}) \notin \mathcal{F}_{\rho, s}^n(\delta)} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right]^{1/\rho}, \end{aligned} \quad (130)$$

where each probability is implicitly conditioned on $\mathbf{X} = \mathbf{x}$.

We first analyze the summation over $\mathcal{F}_{\rho, s}^n(\delta)$ in (130). In order to make the inner probability more amenable to an application of Lemma 2, we write it as

$$\mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right] = \mathbb{P} \left[\left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{Y})}{q^n(\mathbf{x}, \mathbf{Y})} \right)^s \geq 1 \right] \quad (131)$$

$$= \mathbb{P} \left[\frac{\left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{Y})}{q^n(\mathbf{x}, \mathbf{Y})} \right)^s}{\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{y})}{q^n(\mathbf{x}, \mathbf{y})} \right)^s} \geq \frac{1}{\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q^n(\bar{\mathbf{x}}, \mathbf{y})}{q^n(\mathbf{x}, \mathbf{y})} \right)^s} \right] \quad (132)$$

$$= \mathbb{P} \left[j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}) \geq -\log \sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right], \quad (133)$$

where j_s^n is defined in (114). Next, following [34, Sec. 3.4.5], we note that the following holds when $V_s^n(\mathbf{y}|\mathbf{x}, \bar{\mathbf{x}}) \neq 0$:

$$W^n(\mathbf{y}|\mathbf{x}) = W^n(\mathbf{y}|\mathbf{x}) \frac{V_s^n(\mathbf{y}|\mathbf{x}, \bar{\mathbf{x}})}{V_s^n(\mathbf{y}|\mathbf{x}, \bar{\mathbf{x}})} \quad (134)$$

$$= V_s^n(\mathbf{y}|\mathbf{x}, \bar{\mathbf{x}}) e^{-nj_s(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{y})}. \quad (135)$$

Summing (135) over all \mathbf{y} such that $j_s(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{y}) \geq t$, we obtain

$$\mathbb{P} \left[j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}) \geq t \right] = \mathbb{E} \left[e^{-j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}_s)} \mathbf{1} \{ j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}_s) \geq t \} \right], \quad (136)$$

where $\mathbf{Y}_s \sim V_s^n(\cdot|\mathbf{x}, \bar{\mathbf{x}})$. For any $(\mathbf{x}, \bar{\mathbf{x}}) \in \mathcal{F}_{\rho, s}^n(\delta)$, we obtain the following using Lemma 2, the first part of Lemma 3, and the fact that $c_0 > 0$ (see the arguments following (117)):

$$\mathbb{P} \left[j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}) \geq t \right] \leq \frac{\psi_s(1 + o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} e^{-t} \quad (137)$$

uniformly in t , provided that δ is sufficiently small so that $r(\delta) < c_0$. Substituting (137) into (133), we obtain

$$\mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right] \leq \frac{\psi_s(1 + o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} \sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s, \quad (138)$$

and hence

$$\begin{aligned} & \sum_{(\mathbf{x}, \bar{\mathbf{x}}) \in \mathcal{F}_{\rho, s}^n(\delta)} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \mathbb{P} \left[q^n(\bar{\mathbf{x}}, \mathbf{Y}) \geq q^n(\mathbf{x}, \mathbf{Y}) \right]^{1/\rho} \\ & \leq \sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \left(\frac{\psi_s(1+o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} \sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho}. \end{aligned} \quad (139)$$

We observe that the right-hand side of (139) has the same exponent as the denominator of (128). Using Markov's inequality, the summation over $\mathcal{F}_{\rho, s}^n(\delta)^c$ in (130) can be upper bounded by the numerator of (128), and thus the second part of Lemma 3 implies

$$\frac{1}{4(M-1)} \text{rcux}_{\rho, s}(n, M)^{1/\rho} \leq (1+o(1)) \sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \left(\frac{\psi_s(1+o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} \sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho}, \quad (140)$$

and hence

$$\text{rcux}_{\rho, s}(n, M) \leq \frac{4^p \psi_s(1+o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} M^p \left(\sum_{\mathbf{x}, \bar{\mathbf{x}}} Q^n(\mathbf{x}) Q^n(\bar{\mathbf{x}}) \left(\sum_{\mathbf{y}} W^n(\mathbf{y}|\mathbf{x}) \left(\frac{q(\bar{\mathbf{x}}, \mathbf{y})}{q(\mathbf{x}, \mathbf{y})} \right)^s \right)^{1/\rho} \right)^p \quad (141)$$

$$= \frac{4^p \psi_s(1+o(1))}{\sqrt{2\pi n}(c_0 - r(\delta))} e^{-n(E_x^{\text{iid}}(Q, \rho, s) - \rho R)}, \quad (142)$$

where (142) follows by expanding each term as a product from 1 to n and using the definition of E_x^{iid} . The proof is concluded by taking $\delta \rightarrow 0$ (and hence $r(\delta) \rightarrow 0$). \blacksquare

Proof of Lemma 3: We obtain (127) by expanding the variance as

$$\text{Var}[j_s^n(\mathbf{x}, \bar{\mathbf{x}}, \mathbf{Y}_s)] = \sum_{i=1}^n \text{Var}[j_s(x_i, \bar{x}_i, Y_{s,i})] \quad (143)$$

$$= \sum_{\mathbf{x}, \bar{\mathbf{x}}} n \hat{P}_{\mathbf{x}\bar{\mathbf{x}}}(x, \bar{x}) \text{Var}[j_s(x, \bar{x}, Y_s)] \quad (144)$$

and substituting the bound in the definition of $\mathcal{F}_{\rho, s}^n(\delta)$ in (126). To prove the second property, we note that a nearly identical argument to Section IV-A (based on types) reveals that the exponent of the denominator of (128) is equal to

$$\min_{P_{X\bar{X}}} D(P_{X\bar{X}} \| Q \times Q) + \frac{1}{\rho} \mathbb{E}_P[d_s(X, \bar{X})], \quad (145)$$

where d_s is defined in (33). Similarly, the exponent of the numerator of (128) is given by

$$\min_{P_{X\bar{X}} : \max_{\mathbf{x}, \bar{\mathbf{x}}} |P_{X\bar{X}}(\mathbf{x}, \bar{\mathbf{x}}) - P_{\rho, s}^*(\mathbf{x}, \bar{\mathbf{x}})| > \delta} D(P_{X\bar{X}} \| Q \times Q) + \frac{1}{\rho} \mathbb{E}_P[d_s(X, \bar{X})]. \quad (146)$$

A straightforward analysis of the Karush-Kuhn-Tucker (KKT) conditions [18, Sec. 5.5.3] reveals that (145) is uniquely minimized by $P_{\rho, s}^*$, defined in (115). On the other hand, $P_{\rho, s}^*$ does not satisfy the constraint in (146), and thus (146) is strictly higher than (145). \blacksquare

VI. DISCUSSION AND CONCLUSION

We have presented asymptotic and non-asymptotic expurgated bounds for channels with a given decoding rule. Several expurgated exponents have been derived, including that of Csiszár and Körner [6] and its generalization to continuous alphabets. The type class enumeration approach has been shown to provide better exponents for some codeword distributions, better guarantees of exponential tightness, and the opportunity for deriving expurgated exponents for

channels with memory. By refining the analysis of the i.i.d. ensemble, we have obtained a bound with a $O(\frac{1}{\sqrt{n}})$ prefactor, thus improving on Gallager's $O(1)$ prefactor.

APPENDIX

A. Technical Condition of Theorem 2

We begin by providing an example of a class of continuous channels and metrics satisfying the single-letter condition given in (19). Consider an additive noise channel $Y = X + Z$, and let $q(x, y)$ be any decreasing function of $|y - x|$. If the cost constraint is of the form $c(x) = |x|^\beta$ for some constant β , then $c(x) \leq \gamma$ if and only if $|x| \leq \gamma^{1/\beta}$. Thus, any two permissible points are separated by a distance of at most $2\gamma^{1/\beta}$, and the single-letter condition is satisfied if the additive noise satisfies $\mathbb{P}[Z > 2\gamma^{1/\beta}] \geq e^{-E'(\gamma)}$ and $\mathbb{P}[Z < -2\gamma^{1/\beta}] \geq e^{-E'(\gamma)}$ for some $E'(\gamma)$ growing subexponentially in γ . In particular, this holds for noise distributions with exponential tails (e.g. Gaussian). On the other hand, if the cost function is logarithmic, say $c(x) = \log(1 + |x|)$, then (19) fails for additive noise distributions with exponential tails, since in this case the limit on the left-hand side of (19) equals a positive constant.

For any DMC whose zero-error capacity [23] is zero, the condition of Theorem 2 is satisfied under ML decoding, since the error probability can only decay exponentially. On the other hand, the condition could fail for sufficiently “bad” metrics (e.g. one for which there exists a pair (x, \bar{x}) such that $q(x, y) > q(\bar{x}, y)$ for all y). Furthermore, the condition fails under ML decoding whenever the zero-error capacity is positive and Q has a support which includes two inputs not sharing a common output.

B. Proof of Theorem 3

Using the definitions of \mathcal{S} and \mathcal{T} in (43)–(44), we write (28) as

$$\hat{E}_{\text{ex}}^{\text{cc}}(Q, R) = \min_{\substack{\tilde{P}_{X\bar{X}} \in \mathcal{S}(Q) \\ I_{\tilde{P}}(X; \bar{X}) \leq R}} \min_{P_{X\bar{X}Y} \in \mathcal{T}(\tilde{P}_{X\bar{X}})} D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W) + I_{\tilde{P}}(X; \bar{X}) - R, \quad (147)$$

where the objective follows from (30). We will study (147) one minimization at a time.

Step 1: For a given $\tilde{P}_{X\bar{X}} \in \mathcal{S}(Q)$, $I_{\tilde{P}}(X; \bar{X}) - R$ is constant, and we thus consider the optimization problem

$$\min_{P_{X\bar{X}Y} \in \mathcal{T}(\tilde{P}_{X\bar{X}})} D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W). \quad (148)$$

The Lagrangian [18, Sec. 5.1.1] is given by

$$\begin{aligned} \mathsf{L}_1 = & \sum_{x, \bar{x}, y} P_{X\bar{X}Y}(x, \bar{x}, y) \log \frac{P_{X\bar{X}Y}(x, \bar{x}, y)}{\tilde{P}_{X\bar{X}}(x, \bar{x})W(y|x)} \\ & + s \left(\sum_{x, y} P_{XY}(x, y) \log q(x, y) - \sum_{\bar{x}, y} P_{\bar{X}Y}(\bar{x}, y) \log q(\bar{x}, y) \right) + \sum_{x, \bar{x}} \mu(x, \bar{x}) \left(\tilde{P}_{X\bar{X}}(x, \bar{x}) - P_{X\bar{X}}(x, \bar{x}) \right), \end{aligned} \quad (149)$$

where $s \geq 0$ and $\mu(\cdot, \cdot)$ are Lagrange multipliers. The optimization problem is convex with affine constraints, and thus the optimal value is equal to L_1 for some choice of $P_{X\bar{X}Y}$ and the Lagrange multipliers satisfying the Karush-Kuhn-Tucker (KKT) conditions [18, Sec. 5.5.3].

The simplification of (149) using the KKT conditions uses standard arguments, so we omit some details. Setting $\frac{\partial \mathcal{L}_1}{\partial P_{X\bar{X}Y}(x, \bar{x}, y)} = 0$, using the constraint $P_{X\bar{X}} = \tilde{P}_{X\bar{X}}$ to solve for $\mu(\cdot, \cdot)$, and substituting the resulting expressions back into (149), we obtain

$$\mathcal{L}_1 = - \sum_{x, \bar{x}} \tilde{P}_{X\bar{X}}(x, \bar{x}) \log \sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s. \quad (150)$$

Renaming $\tilde{P}_{X\bar{X}}$ as $P_{X\bar{X}}$, taking the supremum over $s \geq 0$, and adding $I_P(X; \bar{X}) - R$ (see (147)–(148)), we obtain the right-hand side of (31) with the minimum and supremum in the opposite order. Using Fan's minimax theorem [25], we can safely interchange the two.

Since we have taken the supremum over the parameter $s \geq 0$ without verifying that it satisfies the KKT conditions, we have only proved that (31) holds with the equality replaced by an inequality (\leq). To prove the reverse inequality, we use the log-sum inequality [22, Thm. 2.7.1] similarly to [10, Appendix A]. For any $P_{X\bar{X}Y} \in \mathcal{T}(\tilde{P}_{X\bar{X}})$, we have

$$D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W) \geq D(P_{X\bar{X}Y} \| \tilde{P}_{X\bar{X}} \times W) - s \sum_{x, \bar{x}, y} P_{X\bar{X}Y}(x, \bar{x}, y) \log \frac{q(\bar{x}, y)}{q(x, y)} \quad (151)$$

$$= \sum_{x, \bar{x}, y} P_{X\bar{X}Y}(x, \bar{x}, y) \log \frac{P_{X\bar{X}Y}(x, \bar{x}, y)}{\tilde{P}_{X\bar{X}}(x, \bar{x}) W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s} \quad (152)$$

$$\geq \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \frac{1}{\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s}, \quad (153)$$

where (151) holds for any $s \geq 0$ from the constraint $\mathbb{E}_P[\log q(\bar{X}, Y)] \geq \mathbb{E}_P[\log q(X, Y)]$ in (29), (152) follows from the definition of divergence, and (153) follows using the log-sum inequality [22, Thm. 2.7.1] and the constraint $P_{X\bar{X}} = \tilde{P}_{X\bar{X}}$. Equation (153) coincides with (150), thus completing the proof of (31).

Step 2: We now turn to the proof of (32). For any fixed $s \geq 0$, the Lagrangian corresponding to (31) is given by

$$\begin{aligned} \mathcal{L}_2 = & - \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s + (1 + \lambda) \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \frac{P_{X\bar{X}}(x, \bar{x})}{Q(x)Q(\bar{x})} - (1 + \lambda)R \\ & + \sum_x \nu_1(x) (Q(x) - P_X(x)) + \sum_{\bar{x}} \nu_2(\bar{x}) (Q(\bar{x}) - P_{\bar{X}}(\bar{x})), \end{aligned} \quad (154)$$

where $\lambda \geq 0$, $\nu_1(\cdot)$ and $\nu_2(\cdot)$ are Lagrange multipliers. Setting $\frac{\partial \mathcal{L}_2}{\partial P_{X\bar{X}}(x, \bar{x})} = 0$, using the constraint $P_X = Q$ to solve for $\nu_1(\cdot)$, and substituting the resulting expressions back into (154), we obtain

$$\mathcal{L}_2 = -(1 + \lambda) \sum_x Q(x) \log \sum_{\bar{x}} Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s \right)^{\frac{1}{1+\lambda}} e^{\frac{1}{1+\lambda}(\nu_2(\bar{x}) - \nu_2(x))} - (1 + \lambda)R. \quad (155)$$

Taking the supremum over $\nu_2(\cdot)$, $s \geq 0$ and $\lambda \geq 0$, we obtain the right-hand side of (32) after suitable renaming.

Once again, we have only proved that (32) holds with an inequality (\leq) in place of the equality, and we obtain a matching lower bound similarly to (151)–(153). For any $P_{X\bar{X}} \in \mathcal{S}(Q)$ with $I_{\bar{P}}(X; \bar{X}) \leq R$, we can lower bound the

objective in (31) as follows:

$$\begin{aligned} & - \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s + I_P(X; \bar{X}) - R \\ & \geq - \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s + \rho(I_P(X; \bar{X}) - R) \end{aligned} \quad (156)$$

$$= -\rho \sum_{x, \bar{x}} P_{X\bar{X}}(x, \bar{x}) \log \frac{Q(x)Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s e^{a(\bar{x}) - \phi_a} \right)^{1/\rho}}{P_{X\bar{X}}(x, \bar{x})} - \rho R \quad (157)$$

$$\geq -\rho \sum_x Q(x) \log \sum_{\bar{x}} Q(\bar{x}) \left(\sum_y W(y|x) \left(\frac{q(\bar{x}, y)}{q(x, y)} \right)^s e^{a(\bar{x}) - \phi_a} \right)^{1/\rho} - \rho R, \quad (158)$$

where (156) holds for any $\rho \geq 1$ from the constraint $I_{\bar{P}}(X; \bar{X}) \leq R$, (157) holds for any function $a(x)$ with mean $\phi_a = \mathbb{E}_Q[a(X)]$ by expanding the logarithm and applying simple manipulations, and (158) follows from the log-sum inequality [22, Thm. 2.7.1] and the constraint $P_X = Q$. Using the definition of ϕ_a and again expanding the logarithm, it is easily shown that (158) is unchanged when $e^{a(\bar{x}) - \phi_a}$ is replaced by $\frac{e^{a(\bar{x})}}{e^{a(\bar{x})}}$, thus completing the proof.

C. Proof of Proposition 1

The result for the i.i.d. exponent follows similarly to Gallager [2, Sec 5.7], so we only explain the differences. Let $E_x^{\text{iid}}(Q, \rho, s)$ be the function E_x^{iid} in (27), with a fixed value of s rather than a supremum. We claim that

$$\lim_{R \rightarrow 0^+} \sup_{\rho \geq 1, s \geq 0} E_x^{\text{iid}}(Q, \rho, s) - \rho R = \sup_{\rho \geq 1, s \geq 0} E_x^{\text{iid}}(Q, \rho, s). \quad (159)$$

It is easily seen that the left-hand side of (159) cannot exceed the right-hand side, since ρR is positive for any sequence of R values approaching zero from above. It remains to prove the converse. We have for all R that

$$\sup_{\rho \geq 1, s \geq 0} E_x^{\text{iid}}(Q, \rho, s) - \rho R \geq E_x^{\text{iid}}(Q, \rho, s) - \rho R. \quad (160)$$

Taking $R \rightarrow 0$ and then taking the supremum over $s \geq 0$ and $\rho \geq 1$ yields the desired result. The remainder of the proof follows using Gallager's argument: For any fixed s , the supremum over ρ is in the limit as $\rho \rightarrow \infty$, and this limit is easily evaluated using L'Hôpital's rule.

The result for the constant-composition exponent follows in the same way using the fact that $\sup_{s, a_1(\cdot), a_2(\cdot)} E_x^{\text{cost}}(Q, \rho, \{a_1, a_2\}) = E_x^{\text{cc}}(Q, \rho)$ (see Section IV-A; in particular, E_x^{cost} is defined in (60)). Once again, the supremum over ρ is in the limit as $\rho \rightarrow \infty$ when the remaining parameters are fixed.

D. Derivation of $E_{\text{ex}}^{\text{cc}}$ Using Theorem 7

Using similar arguments to those in Section IV-A, we can evaluate the lower tail probability of $d_s^n(\mathbf{x}, \bar{\mathbf{X}})$ as follows:

$$\sum_{\bar{\mathbf{x}}} P_{\mathbf{X}}(\bar{\mathbf{x}}) \mathbb{1}\{d_s^n(\mathbf{x}, \bar{\mathbf{x}}) \leq nD\} \leq \sum_{\bar{\mathbf{x}}} P_{\mathbf{X}}(\bar{\mathbf{x}}) e^{t(nD - d_s^n(\mathbf{x}, \bar{\mathbf{x}}))} \quad (161)$$

$$\leq \sum_{\bar{\mathbf{x}}} Q^n(\bar{\mathbf{x}}) e^{t(nD - d_s^n(\mathbf{x}, \bar{\mathbf{x}}))} e^{\bar{r}(a(\bar{\mathbf{x}}) - n\phi_a)} \quad (162)$$

$$= e^{n(tD - \bar{r}\phi_a)} \prod_{i=1}^n \sum_{\bar{x}} Q(x) e^{\bar{r}a(\bar{x}) - td_s(x_i, \bar{x})}, \quad (163)$$

where (161) holds or any $t \geq 0$ by upper bounding the indicator function, and (162) holds for any \bar{r} using (56) and (57). From (163), we may set

$$R(D, \mathbf{x}) = \sup_{t \geq 0, \bar{r}} \bar{r} \phi_a - tD - \frac{1}{n} \sum_{i=1}^n \theta(x_i, \bar{r}, t), \quad (164)$$

where

$$\theta(x, \bar{r}, t) \triangleq \log \mathbb{E}_Q [e^{\bar{r}a(\bar{X}) - td_s(x, \bar{X})}]. \quad (165)$$

Before proceeding, we present the following proposition.

Proposition 3. *Consider the cost-constrained distribution $P_{\mathbf{X}}$ in (24), and assume that the input distribution Q and auxiliary costs $\{a_l\}_{l=1}^L$ are such that assumptions of Proposition 2 are satisfied. For any function $f : \mathcal{X} \rightarrow \mathbb{R}$, we have*

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n f(X_i) \right] = \mathbb{E}_Q[f(X)] \quad (166)$$

provided that $\mathbb{E}_Q[f(X)]$ exists.

Proof: See Appendix E. ■

We can now simplify the exponent in (85) as follows:

$$\mathbb{E} \left[\inf_{D: R(D, \mathbf{X}) \leq R} D + R(D, \mathbf{X}) - R \right] \quad (167)$$

$$= \mathbb{E} \left[\inf_D \sup_{\rho \geq 1} D + \rho(R(D, \mathbf{X}) - R) \right] \quad (168)$$

$$\geq \sup_{\rho \geq 1} \mathbb{E} \left[\inf_D D + \rho(R(D, \mathbf{X}) - R) \right] \quad (169)$$

$$= \sup_{\rho \geq 1} \mathbb{E} \left[\inf_D \sup_{t \geq 0, \bar{r}} D(1 - \rho t) - \rho \left(-\bar{r} \phi_a + \frac{1}{n} \sum_{i=1}^n \theta(X_i, \bar{r}, t) + R \right) \right] \quad (170)$$

$$\geq \sup_{\rho \geq 1} \sup_{\bar{r}} -\rho \left(-\bar{r} \phi_a + \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \theta(X_i, \bar{r}, 1/\rho) \right] + R \right) \quad (171)$$

$$\rightarrow \sup_{\rho \geq 1} \sup_{\bar{r}} \rho \left(\bar{r} \phi_a - \mathbb{E}_Q[\theta(X, \bar{r}, 1/\rho)] - R \right), \quad (172)$$

where (168) follows from (55), (170) follows from (164), (171) follows by replacing the supremum over $t \geq 0$ by the choice $t = 1/\rho$, and (172) follows from Proposition 3.

Substituting (165) into (172) setting $\bar{r} = \frac{1}{\rho}$, and taking the supremum over $a(\cdot)$, we obtain (34), as desired.

E. Proof of Proposition 3

We first present the proof in the case that there is $L = 1$ auxiliary cost $a(\cdot)$ (with mean ϕ_a) and no system cost constraint, and then discuss the changes required to handle the general case. Throughout the proof, we define $a^n(\mathbf{x}) \triangleq \sum_{i=1}^n a(x_i)$ and $f^n(\mathbf{x}) \triangleq \sum_{i=1}^n f(x_i)$. We use summations to denote averaging with respect to Q , but the proof remains valid in the continuous case upon replacing these by integrals.

Let \mathbf{X} be the random cost-constrained codeword, and define $\mathbf{X}' \sim Q^n(x')$. From (24), we have

$$\frac{1}{n} \mathbb{E}[f^n(\mathbf{X})] = \frac{1}{n} \frac{1}{\mu_n} \mathbb{E} \left[f^n(\mathbf{X}') \mathbb{1} \{ |a^n(\mathbf{X}') - n\phi_a| \leq \delta \} \right]. \quad (173)$$

By a direct differentiation, this is equal to $\frac{d}{d\lambda} \left(\frac{1}{n} \log Z(\lambda) \right)$ evaluated at $\lambda = 0$, where

$$Z(\lambda) \triangleq \mathbb{E} \left[e^{\lambda f^n(\mathbf{X}')} \mathbb{1} \{ |a^n(\mathbf{X}') - n\phi_a| \leq \delta \} \right]. \quad (174)$$

Expanding the expectation and using the inverse Laplace transform relation

$$\mathbb{1} \{ z \geq 0 \} = \frac{1}{2\pi j} \int_{u-j\infty}^{u+j\infty} \frac{e^{tz}}{t} dt \quad (175)$$

for $u > 0$, we have the following:

$$Z(\lambda) = \sum_{\mathbf{x}'} Q^n(\mathbf{x}') e^{\lambda f^n(\mathbf{x}')} \left(\mathbb{1} \{ a^n(\mathbf{x}') \leq n\phi_a + \delta \} - \mathbb{1} \{ a^n(\mathbf{x}') \leq n\phi_a - \delta \} \right) \quad (176)$$

$$= \frac{1}{2\pi j} \sum_{\mathbf{x}'} Q^n(\mathbf{x}') e^{\lambda f^n(\mathbf{x}')} \int_{u-j\infty}^{u+j\infty} e^{t(n\phi_a - a^n(\mathbf{x}'))} \frac{e^{t\delta} - e^{-t\delta}}{t} dt \quad (177)$$

$$= \frac{1}{2\pi j} \int_{u-j\infty}^{u+j\infty} \frac{e^{t\delta} - e^{-t\delta}}{t} e^{n\phi_a t} \left(\sum_{\mathbf{x}'} Q(\mathbf{x}') e^{-ta(\mathbf{x}') + \lambda f(\mathbf{x}')} \right)^n dt. \quad (178)$$

Denoting the derivative of $Z(\cdot)$ by $Z'(\cdot)$, we have

$$Z'(0) = \frac{n}{2\pi j} \int_{u-j\infty}^{u+j\infty} \frac{e^{t\delta} - e^{-t\delta}}{t} e^{n\phi_a t} \left(\sum_{\mathbf{x}'} Q(\mathbf{x}') e^{-ta(\mathbf{x}')} \right)^{n-1} \sum_{\mathbf{x}'} Q(\mathbf{x}') f(\mathbf{x}') e^{-ta(\mathbf{x}')} dt \quad (179)$$

$$= \frac{n}{2\pi j} \int_{u-j\infty}^{u+j\infty} \frac{e^{t\delta} - e^{-t\delta}}{t} e^{n\phi_a t} \left(\sum_{\mathbf{x}'} Q(\mathbf{x}') e^{-ta(\mathbf{x}')} \right)^n \frac{\sum_{\mathbf{x}'} Q(\mathbf{x}') f(\mathbf{x}') e^{-ta(\mathbf{x}')}}{\sum_{\mathbf{x}'} Q(\mathbf{x}') e^{-ta(\mathbf{x}')}} dt. \quad (180)$$

Finally, using the assumption that $\mathbb{E}_Q[a(X)^2] < \infty$ and applying the saddlepoint method [35, Ch. 4-5] (see also [9, Sec. 4.2-4.3]), we obtain

$$\frac{d}{d\lambda} \left(\frac{1}{n} \log Z(\lambda) \right) \Big|_{\lambda=0} = \frac{Z'(0)}{Z(0)} \rightarrow \frac{\sum_{\mathbf{x}'} Q(\mathbf{x}') f(\mathbf{x}') e^{-t_0 a(\mathbf{x}')}}{\sum_{\mathbf{x}'} Q(\mathbf{x}') e^{-t_0 a(\mathbf{x}')}}, \quad (181)$$

where t_0 is the zero of the derivative (saddlepoint) of the function $h(t) = \phi_a t + \log \mathbb{E}_Q[e^{-ta(X)}]$. Since $\phi_a = \mathbb{E}_Q[a(X)]$ by definition, it is easily verified that $t_0 = 0$, and thus the right-hand side of (181) equals $\mathbb{E}_Q[f(X)]$, as desired.

In the case of multiple auxiliary costs, the argument is similar, but with $ta(\cdot)$ replaced by $\sum_l t_l a_l(\cdot)$. The system cost $c(x)$ in (25) can be handled similarly provided that $\mathbb{E}_Q[c(X)] \leq \Gamma$, which is an assumption of the proposition.

REFERENCES

- [1] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. Journal*, vol. 27, pp. 379–423, July and Oct. 1948.
- [2] R. Gallager, *Information Theory and Reliable Communication*. John Wiley & Sons, 1968.
- [3] Y. Polyanskiy, V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2307–2359, May 2010.
- [4] I. Csiszár, J. Körner, and K. Marton, "A new look at the error exponent of discrete memoryless channels," in *IEEE Int. Symp. Inf. Theory*, Ithaca, NY, 1977.
- [5] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, 2nd ed. Cambridge University Press, 2011.
- [6] —, "Graph decomposition: A new key to coding theorems," *IEEE Trans. Inf. Theory*, vol. 27, no. 1, pp. 5–12, Jan. 1981.
- [7] N. Merhav, "Error exponents of erasure/list decoding revisited via moments of distance enumerators," *IEEE Trans. Inf. Theory*, vol. 54, no. 10, pp. 4439–4447, Oct. 2008.
- [8] R. Etkin, N. Merhav, and E. Ordentlich, "Error exponents of optimum decoding for the interference channel," *IEEE Trans. Inf. Theory*, vol. 56, no. 1, pp. 40–56, 2010.
- [9] N. Merhav, "Statistical physics and information theory," *Foundations and Trends in Comms. and Inf. Theory*, vol. 6, no. 1-2, pp. 1–212, 2009.

- [10] N. Merhav, G. Kaplan, A. Lapidoth, and S. Shamai, "On information rates for mismatched decoders," *IEEE Trans. Inf. Theory*, vol. 40, no. 6, pp. 1953–1967, Nov. 1994.
- [11] I. Csiszár and P. Narayan, "Channel capacity for a given decoding metric," *IEEE Trans. Inf. Theory*, vol. 45, no. 1, pp. 35–43, Jan. 1995.
- [12] A. Ganti, A. Lapidoth, and E. Telatar, "Mismatched decoding revisited: General alphabets, channels with memory, and the wide-band limit," *IEEE Trans. Inf. Theory*, vol. 46, no. 7, pp. 2315–2328, Nov. 2000.
- [13] J. Scarlett, A. Martinez, and A. Guillén i Fàbregas, "Mismatched decoding: Error exponents, second-order rates and saddlepoint approximations," *IEEE Trans. Inf. Theory*, vol. 60, no. 5, pp. 2647–2666, May 2014.
- [14] F. Jelinek, "Evaluation of expurgated bound exponents," *IEEE Trans. Inf. Theory*, vol. 14, no. 3, pp. 501–505, 1968.
- [15] R. Blahut, "Composition bounds for channel block codes," *IEEE Trans. Inf. Theory*, vol. 23, no. 6, pp. 656–674, 1977.
- [16] J. K. Omura, "Expurgated bounds, Bhattacharyya distance, and rate distortion functions," *Inf. and Control*, vol. 24, no. 4, pp. 358 – 383, 1974.
- [17] I. Csiszár, "On the error exponent of source-channel transmission with a distortion threshold," *IEEE Trans. Inf. Theory*, vol. 28, no. 6, pp. 823–828, Nov. 1982.
- [18] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [19] Y. Altuğ and A. B. Wagner, "Refinement of the random coding bound," 2014, <http://arxiv.org/abs/1312.6875>.
- [20] J. Scarlett, A. Martinez, and A. Guillén i Fàbregas, "A derivation of the asymptotic random-coding prefactor," in *Allerton Conf. on Comm., Control and Comp.*, Monticello, IL, 2013.
- [21] J. Scarlett, A. Martinez, and A. Guillén i Fàbregas, "Cost-constrained random coding and applications," in *Inf. Theory and Apps. Workshop*, San Diego, CA, Feb. 2013.
- [22] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & Sons, Inc., 2001.
- [23] C. E. Shannon, "The zero error capacity of a noisy channel," *IRE Trans. Inf. Theory*, vol. 2, no. 3, pp. 8–19, Sept. 1956.
- [24] A. G. D'yachkov, "Bounds on the average error probability for a code ensemble with fixed composition," *Prob. Inf. Transm.*, vol. 16, no. 4, pp. 3–8, 1980.
- [25] K. Fan, "Minimax theorems," *Proc. Nat. Acad. Sci.*, vol. 39, pp. 42–47, 1953.
- [26] J. Scarlett, A. Martinez, and A. Guillén i Fàbregas, "Expurgated random-coding ensembles: Exponents, refinements and connections," in *Int. Zurich Sem. on Comms.*, Feb. 2014.
- [27] M. Mézard and A. Montanari, *Information, Physics and Computation*. Oxford University Press, 2009.
- [28] B. Derrida, "Random-energy model: Limit of a family of disordered models," *Phys. Rev. Lett.*, vol. 45, no. 2, pp. 79–82, 1980.
- [29] —, "The random energy model," *Physics Reports*, vol. 67, no. 1, pp. 29–35, 1980.
- [30] —, "Random-energy model: An exactly solvable model for disordered systems," *Phys. Rev. Lett.*, vol. 24, no. 5, pp. 2613–2626, 1981.
- [31] P. Elias, "Coding for two noisy channels," in *Third London Symp. Inf. Theory*, 1955.
- [32] R. L. Dobrushin, "Asymptotic estimates of the probability of error for transmission of messages over a discrete memoryless communication channel with a symmetric transition probability matrix," *Theory Prob. Apps.*, vol. 7, no. 3, pp. 270–300, 1962.
- [33] Y. Altuğ and A. B. Wagner, "Refinement of the sphere-packing bound: Asymmetric channels," *IEEE Trans. Inf. Theory*, vol. 60, no. 3, pp. 1592–1614, March 2013.
- [34] Y. Polyanskiy, "Channel coding: Non-asymptotic fundamental limits," Ph.D. dissertation, Princeton University, 2010.
- [35] N. G. de Bruijn, *Asymptotic Methods in Analysis*. Dover Publications, 1981.