

Bit-Wise Unequal Error Protection for Variable-Length Block Codes With Feedback

Bariş Nakiboğlu, Siva K. Gorantla, Lizhong Zheng, and Todd P. Coleman, *Senior Member, IEEE*

Abstract—The bit-wise unequal error protection problem, for the case when the number of groups of bits ℓ is fixed, is considered for variable-length block codes with feedback. An encoding scheme based on fixed-length block codes with erasures is used to establish inner bounds to the achievable performance for finite expected decoding time. A new technique for bounding the performance of variable-length block codes is used to establish outer bounds to the performance for a given expected decoding time. The inner and the outer bounds match one another asymptotically and characterize the achievable region of rate-exponent vectors, completely. The single-message message-wise unequal error protection problem for variable-length block codes with feedback is also solved as a necessary step on the way.

Index Terms—Block codes, Burnashev’s exponent, discrete memoryless channels (DMCs), error exponents, errors-and-erasures decoding variable-length block coding, feedback, Kudryashov’s signaling, unequal error protection (UEP), variable-length communication, Yamamoto–Itoh scheme.

I. INTRODUCTION

IN the conventional formulation of digital communication problem, the primary concern is the correct transmission of the message; hence, there is no distinction between different error events. In other words, there is a tacit assumption that all error events are equally undesirable; incorrectly decoding to a message \bar{m} when a message \tilde{m} is transmitted is as undesirable

as incorrectly decoding to a message \bar{m} when a message \tilde{m} is transmitted, for any \bar{m} other than \tilde{m} and \bar{m} other than \tilde{m} . Therefore, the performance criteria used in the conventional formulation (minimum distance between codewords, maximum conditional error probability among messages, average error probability, etc.) are oblivious to any precedence order that might exist among the error events.

In many applications, however, there is a clear order of precedence among the error events. For example in Internet communication, packet headers are more important than the actual payload data. Hence, a code used for Internet communication can enhance the protection against the erroneous transmission of the packet headers at the expense of the protection against the erroneous transmission of payload data. In order to appreciate such a coding scheme, one may analyze error probability of the packet headers and error probability of payload data separately, instead of analyzing the error probability of the overall message composed of packet header and payload data. Such a formulation for Internet communication is an unequal error protection (UEP) problem, because of the separate calculation of the error probabilities of the parts of the messages.

Problems capturing the disparity of undesirability among various classes of error events, by assigning and analyzing distinct performance criteria for different classes of error events, are called UEP problems. UEP problems have already been studied widely by researchers in communication theory, coding theory, and computer networks from the perspectives of their respective fields. In this paper, we enhance the information theoretic perspective on UEP problems [2], [5] for variable-length block codes by generalizing the results of [2] to the rates below capacity.

In information theoretic UEP, error events are grouped into different classes and the probabilities associated with these different classes of error events are analyzed separately. In order to prioritize protection against one or the other class of error events, corresponding error exponent is increased at the expense of the other error exponents. There are various ways to choose the error event classes but two specific choices of error event classes stand out because of their intuitive familiarity and practical relevance; they correspond to the message-wise UEP and the bit-wise UEP. In the following, we first describe these two types of UEP and then specify the UEP problems we are interested in this paper.

In the message-wise UEP, the message set \mathcal{M} is assumed to be the union of ℓ disjoint sets for some fixed ℓ , i.e., $\mathcal{M} = \cup_{j=1}^{\ell} \mathcal{M}_j$ where $\mathcal{M}_i \cap \mathcal{M}_j = \emptyset$ for all $i \neq j$. For

Manuscript received January 10, 2011; revised April 23, 2012; accepted October 08, 2012. Date of publication November 16, 2012; date of current version February 12, 2013. This work was supported in part by the National Science Foundation (NSF) through the NSF Cyberphysical Systems Program NSF0932410, in part by the NSF Science and Technology Center under Grant CCF-0939370, and in part by the Defense Advanced Research Projects Agency’s Information Theory of Mobile Ad-Hoc Networks program under Grant W911NF-07-1-0029.

B. Nakiboğlu is with the Department of Electrical Engineering and Computer Sciences, University of California Berkeley, Berkeley, CA 94720-1770 USA (e-mail: nakib@eecs.berkeley.edu).

S. K. Gorantla was with the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA. He is now with Adchemy Inc., Foster City, CA 94402 USA (e-mail: sgorant2@illinois.edu).

L. Zheng is with the Department of Electrical Engineering and Computer Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (e-mail: lizhong@mit.edu).

T. P. Coleman was with the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA. He is now with the Department of Bioengineering, University of California at San Diego, La Jolla, CA 92093 USA (e-mail: tpcoleman@ucsd.edu).

Communicated by D. Guo, Associate Editor for Shannon Theory.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIT.2012.2227671

each set \mathcal{M}_j , the maximum error probability¹ $P_e\{j\}$, the rate $R_{\{j\}}$, and the error exponent $E_{\{j\}}$ are defined as the corresponding quantities defined in the conventional problem, i.e., $P_e\{j\} = \max_{m \in \mathcal{M}_j} \mathbf{P}[\widehat{M} \neq m \mid M = m]$, $R_{\{j\}} = \frac{\ln |\mathcal{M}_j|}{n}$, $E_{\{j\}} = \frac{-\ln P_e\{j\}}{n}$, for all j in $\{1, 2, \dots, \ell\}$ where n is the length of the code. The ultimate aim is calculating the achievable region of rate vector error exponent vector pairs, $(R_{\{\cdot\}}, E_{\{\cdot\}})$'s where² $R_{\{\cdot\}} = (R_{\{1\}}, R_{\{2\}}, \dots, R_{\{\ell\}})$ and $E_{\{\cdot\}} = (E_{\{1\}}, E_{\{2\}}, \dots, E_{\{\ell\}})$. The *message-wise* UEP problem was the first information theoretic UEP problem to be considered; it was considered by Csiszár in his work on joint source channel coding [5]. Csiszár showed that for any integer ℓ , block length n , and ℓ -dimensional rate vector $R_{\{\cdot\}}$ such that $0 \leq R_{\{j\}} \leq C$ for $j = 1, 2, \dots, \ell$, there exists a length n block code with message set $\mathcal{M} = \cup_{j=1}^{\ell} \mathcal{M}_j$ where $|\mathcal{M}_j| = e^{n(R_{\{j\}} - \varepsilon_n)}$ such that the conditional error probability of each message in each \mathcal{M}_j is less than $e^{-n(E_r(R_{\{j\}}) - \varepsilon_n)}$ where $E_r(\cdot)$ is the random coding exponent and ε_n converges to zero as n diverges.³

The *bit-wise* UEP problem is the other canonical form of the information theoretic UEP problems. In the *bit-wise* UEP problem, the message set \mathcal{M} is assumed to be the Cartesian product of $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_\ell$ for some fixed ℓ , i.e., $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_\ell$. Thus, the transmitted message M and the decoded message \widehat{M} are given by $M = (M_1, M_2, \dots, M_\ell)$ and $\widehat{M} = (\widehat{M}_1, \widehat{M}_2, \dots, \widehat{M}_\ell)$, respectively. Furthermore, M_j 's and \widehat{M}_j 's are called the transmitted and decoded submessages, respectively. The error events of interest in the *bit-wise* UEP problem are the ones corresponding to the erroneous transmission of the submessages. The error probability $P_e(j)$, rate R_j , and the error exponent E_j of submessages are given by $P_e(j) = \mathbf{P}[\widehat{M}_j \neq M_j]$, $R_j = \frac{\ln |\mathcal{M}_j|}{n}$, $E_j = \frac{-\ln P_e(j)}{n}$ for all j in $\{1, 2, \dots, \ell\}$ where n is the block length. As was the case in the *message-wise* UEP problem, the ultimate aim in the *bit-wise* UEP problem is determining the achievable region of the rate vector error exponent vector pairs⁴ (\bar{R}, \bar{E}) . The formulation of Internet communication problem we have considered above, with packet header and payload data, is a *bit-wise* UEP problem with two submessages, i.e., with $\ell = 2$.

¹This formulation is called the missed detection formulation of the *message-wise* UEP problem in [2]. If $\mathbf{P}[\widehat{M} \neq m \mid M = m]$ is replaced with $\mathbf{P}[\widehat{M} = m \mid M \neq m]$, we get the false alarm formulation of the *message-wise* UEP problem. In this paper, we restrict our discussion to the missed detection problem and use *message-wise* UEP without any qualifications to refer to the missed detection formulation of the *message-wise* UEP problem.

²Here, ℓ is assumed to be a fixed integer. All rate-exponent vectors, achievable or not, are in the region of $\mathbf{R}^{2\ell}$ in which $R_{\{j\}} \geq 0$ and $E_{\{j\}} \geq 0$ for all $1 \leq j \leq \ell$. $\mathbf{R}^{2\ell}$ is the 2ℓ -dimensional real vector space with the norm $\|\vec{x}\| = \sup_j |x_j|$.

³Csiszár proved the aforementioned result not only for the case when ℓ is constant for all n but also for the case when ℓ_n is a sequence such that $\lim_{n \rightarrow \infty} \frac{\ln \ell_n}{n} = 0$. See [5, Th. 5].

⁴Similar to the *message-wise* UEP problem discussed previously, in the current formulation of *bit-wise* UEP problem, we assume ℓ to be fixed. Thus, all rate-exponent vectors, achievable or not, are in region of $\mathbf{R}^{2\ell}$ in which $R_j \geq 0$ and $E_j \geq 0$ for all $1 \leq j \leq \ell$, by definition.

There is some resemblance in the definitions of *message-wise* and *bit-wise* UEP problems, but they have very different behavior in many problems. For example, consider the *message-wise* UEP problem and the *bit-wise* UEP problem with $\ell = 2$, $\mathcal{M}_1 = \{1, 2\}$, and $\mathcal{M}_2 = \{3, 4, \dots, e^{n(C - \varepsilon_n)}\}$ for some ε_n that goes to zero as n diverges. It is shown in [2, Th. 1] that if $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2$ and $\mathbf{P}[M_2 \neq \widehat{M}_2] \leq \tilde{\varepsilon}_n$ for some $\tilde{\varepsilon}_n$ that goes to zero as n diverges, then⁵ $E_1 = 0$. Thus, in the *bit-wise* UEP problem, even a bit cannot have a positive error exponent. As result of [5, Th. 5], on the other hand, if $\mathcal{M} = \mathcal{M}_1 \cup \mathcal{M}_2$, we know that \mathcal{M}_1 can have an error exponent $E_{\{1\}}$ as high as $E_r(0) > 0$, while having a small error probability for \mathcal{M}_2 , i.e., $\max_{m \in \mathcal{M}_2} \mathbf{P}[\widehat{M} \neq m \mid M = m] \leq \tilde{\varepsilon}_n$ for some $\tilde{\varepsilon}_n$ that goes to zero as n diverges. Thus, in the *message-wise* UEP problem, it is possible to give an error exponent as high as $E_r(0)$ to \mathcal{M}_1 .

The *message-wise* and the *bit-wise* UEP problems cover a wide range of problems of practical interest. Yet, as noted in [2], there are many UEP problems of practical importance that are neither *message-wise* nor *bit-wise* UEP problems. One of our aims in studying the *message-wise* and the *bit-wise* UEP problems is gaining insights and devising tools for the analysis of those more complicated problems.

In the aforementioned discussion, the UEP problems are described for fixed-length block codes for the sake of simplicity. One can, however, easily define the corresponding problems for various families of codes: with or without feedback, fixed or variable length, by modifying the definitions of the error probability, the rate, and the error exponent appropriately. Furthermore, parameter ℓ representing the number of groups of bits or messages is assumed to be fixed in the aforementioned discussion for simplicity. However, both the *message-wise* and the *bit-wise* UEP problems can be defined for ℓ 's that are increasing with block length n in fixed-length block codes and for ℓ 's that are increasing with expected block length $\mathbf{E}[T]$ in variable-length block codes. In fact, Csiszár's result discussed previously [5, Th. 5] is proved not only for constant ℓ but also for any ℓ_n sequence satisfying $\lim_{n \rightarrow \infty} \frac{\ln \ell_n}{n} = 0$.

In this paper, we consider two closely related UEP problems for variable-length block codes over a discrete memoryless channels (DMCs) with noiseless feedback: the *bit-wise* UEP problem and the single-message *message-wise* UEP problem.

- 1) In the *bit-wise* UEP problem, there are ℓ submessages each with different priority and rate. For all fixed values of ℓ , we characterize the tradeoff between the rates and the error exponents of these submessages by revealing the region of achievable rate vector, exponent vector pairs. For fixed ℓ , this problem is simply the variable-length code version of the previously described *bit-wise* UEP problem.
- 2) In the single-message *message-wise* UEP problem, we characterize the tradeoff between the exponents of the minimum and the average conditional error probability. Thus, this problem is similar to the previously described *message-wise* UEP problem for the case $\ell = 2$ and $\mathcal{M}_1 = \{1\}$. But unlike that problem, we work with variable-length codes and average conditional error prob-

⁵The channel is assumed to have no zero probability transition.

ability rather than fixed-length codes and the maximum error probability.

The *bit-wise* UEP problem for fixed number of groups of bits, i.e., fixed ℓ , and the single-message *message-wise* UEP problem were first considered in [2], for the case when the rate is (very close to) the channel capacity; we solve both of these problems for all achievable rates.

In fact, in [2], single-message *message-wise* UEP problem is solved not only at capacity, but also for all the rates below capacity both for fixed-length block codes without feedback and for variable-length block codes with feedback, but only for case when overall error exponent is zero (see [2, Appendix D]). Recently, Wang *et al.* [11] put forward a new proof based on method of types for the same problem.⁶ Nazer *et al.* [9], on the other hand, investigated the problem for fixed-length block codes on additive white Gaussian noise channels at zero error exponent and derived the exact analytical expression in terms of rate and power constraints.

Before starting our presentation, let us give a brief outline of this paper. In Section II, we specify the channel model and make a brief overview of stopping times and variable-length block codes. In Section III, we first present the single-message *message-wise* UEP problem and fixed ℓ version of the *bit-wise* UEP problem for variable-length block codes; then, we state the solutions of these two UEP problems. In Section IV, we present inner bounds for both the single-message *message-wise* UEP problem and the *bit-wise* UEP problem. In Section V, we introduce a new technique, Lemma 5, for deriving outer bounds for variable-length block codes and apply it to the two UEP problems we are interested in. Finally, in Section VI, we discuss the qualitative ramifications of our results in terms of the design of communication systems with UEP and the limitations of our analysis. The proofs of the propositions in Sections III–V are deferred to the Appendices.

II. PRELIMINARIES

As it is customary, we use upper case letters, e.g., M, X, Y, T for random variables and lower case letters, e.g., m, x, y, t for their sample values.

We denote discrete sets by capital letters with calligraphic fonts, e.g., $\mathcal{M}, \mathcal{X}, \mathcal{Y}$ and power sets of discrete sets by $\wp(\cdot)$, e.g., $\wp(\mathcal{M}), \wp(\mathcal{X}), \wp(\mathcal{Y})$. In order to denote the set of all probability distributions on a discrete set, we use $P(\cdot)$, e.g., $P(\mathcal{M}), P(\mathcal{X}), P(\mathcal{Y})$.

Definition 1 (Total Variation): For any discrete set \mathcal{Z} and for any $\mu_1, \mu_2 \in P(\mathcal{Z})$, the total variation $\Delta(\mu_1, \mu_2)$ is defined as

$$\Delta(\mu_1, \mu_2) = \frac{1}{2} \sum_{z \in \mathcal{Z}} |\mu_1(z) - \mu_2(z)|. \quad (1)$$

We denote the indicator function by $\mathbb{1}_{\{\cdot\}}$, i.e., $\mathbb{1}_{\{\Gamma\}} = 1$ when event Γ happens $\mathbb{1}_{\{\Gamma\}} = 0$ otherwise. We denote the binary entropy function by $h(\cdot)$, i.e.,

$$h(s) \triangleq -s \ln s - (1-s) \ln(1-s) \quad \forall s \in [0, 1]. \quad (2)$$

⁶In addition to their new proof in missed-detection problem [11, Th. 1], Wang *et al.* present a completely new result on the false-alarm formulation of the problem [11, Th. 5].

A. Channel Model

We consider a DMC with input alphabet \mathcal{X} , output alphabet \mathcal{Y} , and $|\mathcal{X}| - \text{by} - |\mathcal{Y}|$ transition probability matrix W . Each row of W corresponds to a probability distribution on \mathcal{Y} , i.e., $W_x \in P(\mathcal{Y})$ for all $x \in \mathcal{X}$. For the reasons that will become clear shortly, in Section II-D, we assume that $W_x(y) > 0$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ and denote the smallest transitions probability by λ

$$\lambda \triangleq \min_{x,y} W_x(y) > 0. \quad (3)$$

The input and output letters at time τ , up to time τ , and between time τ_1 and τ_2 are denoted by $X_\tau, Y_\tau, X^\tau, Y^\tau, X_{\tau_1}^{\tau_2}$, and $Y_{\tau_1}^{\tau_2}$ respectively. DMCs are both memoryless and stationary; hence, the conditional probability of $Y_\tau = y$ given $(X^\tau, Y^{\tau-1})$ is given by

$$\mathbf{P}[Y_\tau = y | X^\tau, Y^{\tau-1}] = W_{X_\tau}(y).$$

Definition 2 (Empirical Distribution): For any $\tau_2 \geq \tau_1$ and any sequence $z_{\tau_1}^{\tau_2}$ such that $z_j \in \mathcal{Z}$ for all $j \in [\tau_1, \tau_2]$, the empirical distribution $Q_{\{z_{\tau_1}^{\tau_2}\}}$ is given by

$$Q_{\{z_{\tau_1}^{\tau_2}\}}(z) = \frac{1}{\tau_2 - \tau_1 + 1} \sum_{\tau=\tau_1}^{\tau_2} \mathbb{1}_{\{z_\tau=z\}} \quad \forall z \in \mathcal{Z}. \quad (4)$$

Note that we replace $z_{\tau_1}^{\tau_2}$ by $Z_{\tau_1}^{\tau_2}$ when the empirical distribution $Q_{\{z_{\tau_1}^{\tau_2}\}}(z)$ becomes a random variable for each $z \in \mathcal{Z}$.

B. Stopping Times

Stopping times are central in the formal treatment of variable-length codes; it is not possible to define or comprehend variable-length codes without a solid understanding of stopping times. For those readers who are not already familiar with the concept of the stopping times, we present a brief overview in this section.

In order to make our presentation more accessible, we use the concept of power sets, rather than sigma fields in the definitions. We can do that only because the random variables we use to define stopping times are discrete random variables. In the general case, when the underlying variables are not necessarily discrete, one needs to use the concept of sigma fields instead of power set.

Let us start with introducing the concept of Markov times. For an infinite sequence of random variables Z_1, Z_2, \dots , a positive, *integer**-valued⁷ function T defined on \mathcal{Z}^∞ is a Markov time, if for all positive integers τ , it is possible to determine whether $T = \tau$ or not by considering Z^τ only, i.e., if $\mathbb{1}_{\{T=\tau\}}$ is not only a function of Z^∞ but also a function of Z^τ for all positive integers τ . The formal definition is given in the following.

Definition 3 (Markov Time): Let Z_1^∞ be an infinite sequence of \mathcal{Z} -valued random variables Z_τ for $\tau \in \{1, 2, \dots\}$ and T be a function of Z^∞ which takes values from the set $\{1, 2, \dots, \infty\}$. Then, the random variable T is a Markov time with respect to Z^τ if

$$\{z^\infty : T = \tau \text{ if } Z^\infty = z^\infty\} \in \wp(\mathcal{Z}^\tau) \times \{z_{\tau+1}^\infty\} \quad \forall \tau \in \{1, 2, \dots\} \quad (5)$$

⁷*Integer** is the set of all integers together with two infinities, i.e., $\{-\infty, \dots, -1, 0, 1, \dots, \infty\}$.

where $\wp(\mathcal{Z}^\tau) \times \{\mathcal{Z}_{\tau+1}^\infty\}$ is the Cartesian product of the power set of \mathcal{Z}^τ and the one element set $\{\mathcal{Z}_{\tau+1}^\infty\}$.

We denote Z_τ 's from $\tau = 1$ to $\tau = T$ by Z^T and their sample values by z^t . The set of all sample values of Z^T such that $T = \tau$, on the other hand, is denoted by $\mathcal{Z}_{\{T=\tau\}}^\tau$. We denote union of all $\mathcal{Z}_{\{T=\tau\}}^\tau$'s for finite τ 's by \mathcal{Z}^{T*} and the union of all $\mathcal{Z}_{\{T=\tau\}}^\tau$'s by \mathcal{Z}^T , i.e.,

$$\mathcal{Z}_{\{T=\tau\}}^\tau \triangleq \{z^\tau : T = \tau \text{ if } Z^\tau = z^\tau\} \quad \tau \in \{1, \dots, \infty\} \quad (6a)$$

$$\mathcal{Z}^{T*} \triangleq \bigcup_{1 \leq \tau < \infty} \mathcal{Z}_{\{T=\tau\}}^\tau \quad (6b)$$

$$\mathcal{Z}^T \triangleq \mathcal{Z}^{T*} \bigcup \mathcal{Z}_{\{T=\infty\}}^\infty. \quad (6c)$$

For an arbitrary, positive, *integer**-valued function T of Z^∞ , however, one cannot talk about Z^T , because the value of T can in principle depend on Z_{T+1}^∞ . For a Markov time T , however, the value of T does not depend on Z_{T+1}^∞ . That is why we can define Z^T , $\mathcal{Z}_{\{T=\tau\}}^\tau$, \mathcal{Z}^{T*} , and \mathcal{Z}^T for any Markov time T .

Given an infinite sequence of z_τ 's, i.e., z^∞ , either $z^\infty \in \mathcal{Z}_{\{T=\infty\}}^\infty$ or z^∞ has a unique subsequence z^τ that is in \mathcal{Z}^{T*} .

In most practical situations, one is interested in Markov times that are guaranteed to have a finite value; those Markov times are called stopping times.

Definition 4 (Stopping Time): A Markov time T with respect to Z^τ is a stopping time iff $\mathbf{P}[T < \infty]$.

Note that if T is a stopping time, then $\mathbf{P}[Z^T \in \mathcal{Z}^{T*}] = 1$. Furthermore, unlike \mathcal{Z}^T , \mathcal{Z}^{T*} is a countable set for all stopping times T because $|\mathcal{Z}|$ is finite.⁸

C. Variable-Length Block Codes

A variable-length block code on a DMC is given by a random decoding time T , an encoding scheme Φ , and a decoding rule Ψ satisfying $\mathbf{P}[T < \infty] = 1$.

1) *Decoding time* T is a Markov time with respect to the receiver's observation Y^τ , i.e., given Y^τ , receiver knows whether $T = \tau$ or not. Hence, T is a random quantity rather than a constant; thus, neither the decoder nor the receiver knows the value of T *a priori*. But as time passes, both the decoder and the encoder (because of feedback link) will be able to decide whether T has been reached or not, just by considering the current and past channel outputs.

2) *Encoding scheme* Φ is a collection of mappings that determines the input letter at time $(\tau + 1)$ for each message in the finite message set \mathcal{M} , for each $y^\tau \in \mathcal{Y}^\tau$ such that $T > \tau$

$$\Phi(\cdot, y^\tau) : \mathcal{M} \rightarrow \mathcal{X} \quad \forall y^\tau : T > \tau.$$

3) *Decoding rule* is a mapping from the set of output sequences y^τ such that $T = \tau$ to the finite message set \mathcal{M} which determines the decoded message, \hat{M} . With a slight

⁸ \mathcal{Z}^{T*} is a countable set even when $|\mathcal{Z}|$ is countably infinite.

abuse of notation, we denote the set of all, possibly infinite, output sequences y^τ such that $\{T = \tau \text{ if } Y^\tau = y^\tau\}$ by⁹ \mathcal{Y}^T and write the decoding rule Ψ as

$$\Psi(\cdot) : \mathcal{Y}^T \rightarrow \mathcal{M}.$$

4) Note that because of the condition $\mathbf{P}[T < \infty] = 1$, decoding time is not only a Markov time, but also a stopping time.¹⁰

At time zero, the message M chosen uniformly at random from \mathcal{M} is given to the transmitter; the transmitter uses the codeword associated M , i.e., $\Phi(M, \cdot)$, to convey the message M until the decoding time T . Then, the receiver chooses the decoded message \hat{M} using its observation Y^T and the decoding rule Ψ , i.e., $\hat{M} = \Psi(Y^T)$. The error probability, the rate, and the error exponent of a variable-length block code are given by

$$P_e = \mathbf{P}[\hat{M} \neq M] \quad (7a)$$

$$R = \frac{\ln |\mathcal{M}|}{\mathbf{E}[T]} \quad (7b)$$

$$E = \frac{-\ln P_e}{\mathbf{E}[T]}. \quad (7c)$$

Indeed, one can interpret the variable-length block codes on DMCs as trees; for a more detailed discussion of this interpretation, readers may go over [1, Sec. II].

D. Reliable Sequences for Variable-Length Block Codes

In order to suppress the secondary terms while discussing the main results, we use the concept of reliable sequences. In a sequence of codes, we denote the error probability and the message set of the κ th code of the sequence by $P_e^{(\kappa)}$ and $\mathcal{M}^{(\kappa)}$, respectively.

Definition 5 (Reliable Sequence): A sequence of variable-length block codes \mathbb{Q} is reliable if the error probabilities of the codes vanish and the size of the message sets of the codes diverge¹¹

$$\lim_{\kappa \rightarrow \infty} \left(P_e^{(\kappa)} + \frac{1}{|\mathcal{M}^{(\kappa)}|} \right) = 0$$

where $P_e^{(\kappa)}$ and $\mathcal{M}^{(\kappa)}$ are the error probability and the message set for the κ th code of the reliable sequence, respectively.

Note that in a sequence of codes, each code has an associated probability space. We denote the random variables in these probability spaces together with a superscript corresponding to the code. For example, the decoding time of the κ th code in the sequence is denoted by $T^{(\kappa)}$. The expected value of random variables in the probability space associated with the κ th code in the sequence is denoted¹² by $\mathbf{E}^{(\kappa)}[\cdot]$.

⁹See (6).

¹⁰Having a finite decoding time with probability one, i.e., $\mathbf{P}[T < \infty] = 1$, does not imply having a finite expected value for the decoding time, i.e., $\mathbf{E}[T] < \infty$. Thus, a variable-length code can, in principle, have an infinite expected decoding time.

¹¹Recall that the decoding time of a variable-length block code is finite with probability one. Thus, $\mathbf{P}^{(\kappa)}[T^{(\kappa)} < \infty] = 1$ for all κ for a reliable sequence.

¹²Evidently, it is possible to come up with a probability space that includes all of the codes in a reliable sequence and invoke independence between random quantities associated with different codes. We choose the current convention to emphasize independence explicitly in the notation we use.

Definition 6 (Rate of a Reliable Sequence): The rate of a reliable sequence \mathbb{Q} is the limit infimum of the rates of the individual codes

$$R_{\mathbb{Q}} \triangleq \liminf_{\kappa \rightarrow \infty} \frac{\ln |\mathcal{M}^{(\kappa)}|}{\mathbf{E}^{(\kappa)} [\mathsf{T}^{(\kappa)}]}.$$

Definition 7 (Capacity): The capacity of a channel for variable-length block codes is the supremum of the rates of the all reliable sequences

$$C \triangleq \sup_{\mathbb{Q}} R_{\mathbb{Q}}.$$

The capacity of a DMC for variable-length block codes is identical to the usual channel capacity [3]. Hence

$$C = \max_{\mu \in \mathcal{P}(\mathcal{X})} \sum_{x,y} \mu(x) W_x(y) \ln \frac{W_x(y)}{\bar{\mu}(y)} \quad (8)$$

where $\bar{\mu}(y) = \sum_x \mu(x) W_x(y)$.

Definition 8 (Error Exponent of a Reliable Sequence): The error exponent of a reliable sequence \mathbb{Q} is the limit infimum of the error exponents of the individual codes

$$E_{\mathbb{Q}} \triangleq \liminf_{\kappa \rightarrow \infty} \frac{-\ln P_e^{(\kappa)}}{\mathbf{E}^{(\kappa)} [\mathsf{T}^{(\kappa)}]}.$$

Definition 9 (Reliability Function): The reliability function of a channel for variable-length block codes at rate $R \in [0, C]$ is the supremum of the exponents of all reliable sequences whose rate is R or higher

$$E(R) \triangleq \sup_{\mathbb{Q}: R_{\mathbb{Q}} \geq R} E_{\mathbb{Q}}.$$

Burnashev [3] analyzed the performance of variable-length block codes with feedback and established inner and outer bounds to their performance. Results of [3] determine the reliability function of variable-length block codes on DMCs for all rates. According to [3], we have the following.

1) If all entries of W are positive, then¹³

$$E(R) = \left(1 - \frac{R}{C}\right) D \quad \forall R \in [0, C]$$

where D is maximum Kullback–Leibler divergence between the output distributions of any two input letters

$$D \triangleq \max_{x, \tilde{x} \in \mathcal{X}} D(W_x \| W_{\tilde{x}}). \quad (9)$$

¹³Problem is formulated somewhat differently in [3]; as a result, the work in [3] did not deal with the case $\mathbf{E}[\mathsf{T}] = \infty$. The bounds in [3] does not guarantee that the error probability of a variable-length code with infinite expected decoding time is greater than zero; however, this is the case if all the transition probabilities are positive. To see that, consider a channel with positive minimum transition probability λ , i.e., $\lambda = \min_{x,y} W_x(y) > 0$. In such a channel, any variable-length code satisfies $P_e \geq \frac{|\mathcal{M}|-1}{|\mathcal{M}|} \mathbf{E} \left[\left(\frac{\lambda}{1-\lambda} \right)^{\mathsf{T}} \right]$; then, $P_e > 0$ as $\lambda > 0$ and $\mathbf{P}[\mathsf{T} < \infty] = 1$. Consequently, both the rate and the error exponent are zero for variable-length block codes with infinite expected decoding time. A more detailed discussion of this fact can be found in Section H1 in the Appendix.

2) If there are one or more zero entries¹⁴ in W , i.e., if there are two input letters x, \tilde{x} and an output letter y such that, $W_x(y) = 0$ and $W_{\tilde{x}}(y) > 0$, then for all $R < C$, for large enough $\mathbf{E}[\mathsf{T}]$, there are rate R variable-length block codes that are error free, i.e., $P_e = 0$.

When $P_e = 0$, all error events can have zero probability at the same time. Consequently, all the UEP problems are answered trivially when there is a zero probability transition. This is why we have assumed that $W_x(y) > 0$ for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

We denote the input letters that get this maximum value of Kullback–Leibler divergence by¹⁵ \mathbf{a} and \mathbf{r}

$$D = D(W_{\mathbf{a}} \| W_{\mathbf{r}}). \quad (10)$$

III. PROBLEM STATEMENT AND MAIN RESULTS

A. Problem Statement

For each $m \in \mathcal{M}$, the conditional error probability is defined as¹⁶

$$P_{e|m} \triangleq \mathbf{P} \left[\widehat{\mathsf{M}} \neq \mathsf{M} \mid \mathsf{M} = m \right]. \quad (11)$$

In the conventional setting, we are interested in either the average or the maximum of the conditional error probability of the messages. The behavior of the minimum conditional error probability is scarcely investigated. Single-message UEP problem attempts to answer that question by determining the tradeoff between exponential decay rates of P_e and $\min_{m \in \mathcal{M}} P_{e|m}$. The operational definition of the problem in terms of reliable sequences is as follows.

Definition 10 (Single-Message Message-Wise UEP Problem): For any reliable sequence \mathbb{Q} , the missed detection exponent of the reliable sequence \mathbb{Q} is defined as

$$E_{\text{md}, \mathbb{Q}} = \liminf_{\kappa \rightarrow \infty} \frac{-\ln \min_{m \in \mathcal{M}^{(\kappa)}} P_{e|m}^{(\kappa)}}{\mathbf{E}^{(\kappa)} [\mathsf{T}^{(\kappa)}]} \quad (12)$$

where $P_{e|m}^{(\kappa)}$ is the conditional error probability of the message m for the κ th code of the reliable sequence \mathbb{Q} .

For any rate $R \in [0, C]$ and error exponent¹⁷ $E \in [0, (1 - \frac{R}{C})D]$, the missed detection exponent $E_{\text{md}}(R, E)$ is defined as

$$E_{\text{md}}(R, E) \triangleq \sup_{\substack{\mathbb{Q}: R_{\mathbb{Q}} \geq R \\ E_{\mathbb{Q}} \geq E}} E_{\text{md}, \mathbb{Q}}. \quad (13)$$

In variable-length block codes with feedback, the single-message message-wise UEP problem not only answers a curious

¹⁴Note that, in this situation, $D = \infty$.

¹⁵This particular naming of letters is reminiscent of the use of these letters in Yamamoto–Itoh scheme [12]. Although they are named differently in [12], \mathbf{a} is used for accepting and \mathbf{r} is used for rejecting the tentative decision in Yamamoto–Itoh scheme.

¹⁶Later, in the paper, we consider block codes with erasures. The conditional error probabilities, $P_{e|m}$ for $m \in \mathcal{M}$, are defined slightly differently for them, see (24).

¹⁷Burnashev’s expression for error exponent of variable-length block codes is used explicitly in the definition because we know, as a result of [3], that the error exponents of all reliable sequences are upper bounded by Burnashev’s exponent. An alternative definition oblivious to Burnashev’s result can simply define $E_{\text{md}}(R, E)$ for all rate-exponent vectors that are achievable. That definition is equivalent to Definition 10, because of [3].

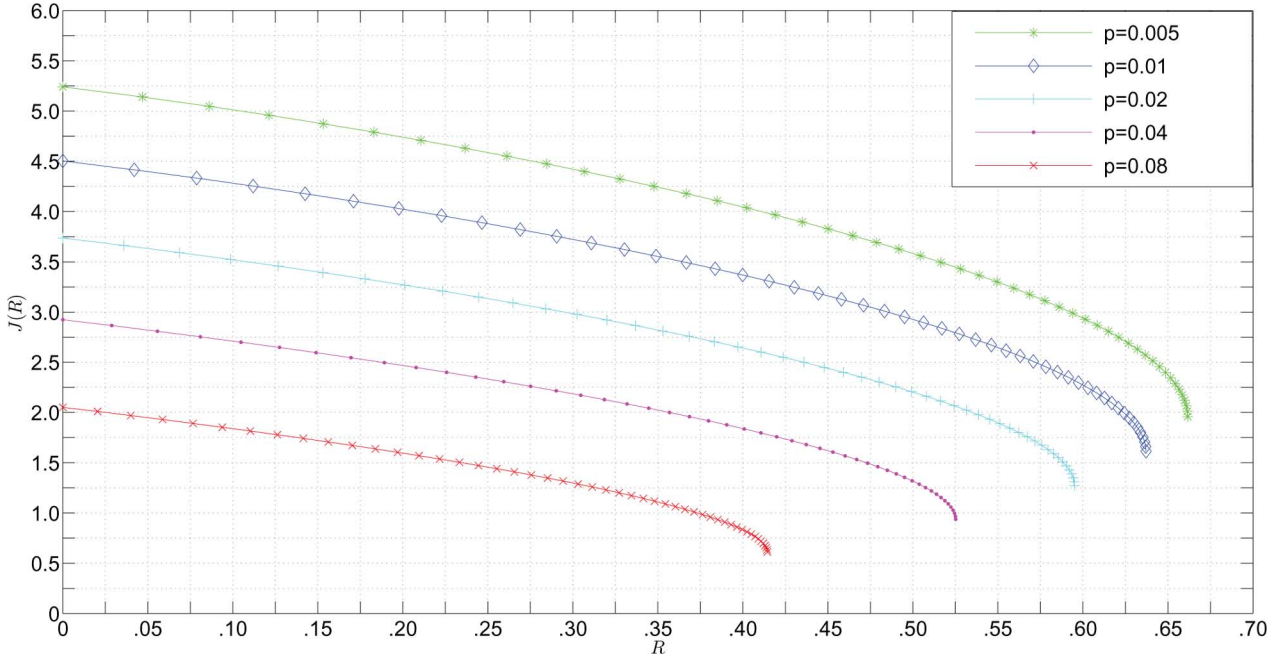


Fig. 1. $J(R)$ function is drawn for BSCs with crossover probabilities $p \in \{0.005, 0.01, 0.02, 0.04, 0.08\}$.

question about the decay rate of the minimum conditional error probability of a code, but also plays a key role in the bit-wise UEP problem, which is our main focus in this paper.

Though they are central in the message-wise UEP problems, the conditional error probabilities of the messages are not relevant in the bit-wise UEP problems. In the bit-wise UEP problems, we analyze the error probabilities of groups of submessages. In order to do that, consider a code with a message set \mathcal{M} of the form

$$\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2 \times \cdots \times \mathcal{M}_\ell.$$

Then, the transmitted message \mathbf{M} and decoded message $\widehat{\mathbf{M}}$ of the code are of the form

$$\begin{aligned} \mathbf{M} &= (M_1, M_2, \dots, M_\ell) \\ \widehat{\mathbf{M}} &= (\widehat{M}_1, \widehat{M}_2, \dots, \widehat{M}_\ell) \end{aligned}$$

where $M_j, \widehat{M}_j \in \mathcal{M}_j$ for all $j = 1, 2, \dots, \ell$. Furthermore, M_j and \widehat{M}_j are called j th transmitted submessage and j th decoded submessage, respectively.

The error probabilities we are interested in correspond to erroneous transmission of certain parts of the message. In order to define them succinctly, let us define \mathcal{M}^j, M^j , and \widehat{M}^j for all j between one and ℓ as follows:

$$\begin{aligned} \mathcal{M}^j &\triangleq \mathcal{M}_1 \times \mathcal{M}_2 \times \cdots \times \mathcal{M}_j \\ M^j &\triangleq (M_1, M_2, \dots, M_j) \\ \widehat{M}^j &\triangleq (\widehat{M}_1, \widehat{M}_2, \dots, \widehat{M}_j). \end{aligned}$$

Then, $P_e(j)$ is defined¹⁸ as the probability of the event that $\widehat{M}^j \neq M^j$

$$P_e(j) \triangleq \mathbf{P}[\widehat{M}^j \neq M^j] \quad \text{for } j = 1, 2, \dots, \ell. \quad (14)$$

¹⁸Similar to the conditional error probabilities, $P_{e|\mathbf{m}}$'s for $\mathbf{m} \in \mathcal{M}$, error probabilities of submessages, $P_e(j)$'s for $j = 1, 2, \dots, \ell$, are defined slightly differently for codes with erasures, see (30).

Note that if $\widehat{M}^j \neq M^j$, then $\widehat{M}^i \neq M^i$ for all i greater than j . Thus

$$P_e(1) \leq P_e(2) \leq P_e(3) \leq \cdots \leq P_e(\ell). \quad (15)$$

Definition 11 (Bit-Wise UEP Problem for Fixed ℓ): For any positive integer ℓ , let \mathbb{Q} be a reliable sequence whose message sets $\mathcal{M}^{(\kappa)}$ are of the form $\mathcal{M}^{(\kappa)} = \mathcal{M}_1^{(\kappa)} \times \mathcal{M}_2^{(\kappa)} \times \cdots \times \mathcal{M}_\ell^{(\kappa)}$. Then, the entries of the rate vector $\vec{R}_{\mathbb{Q}}$ and the error exponent vector $\vec{E}_{\mathbb{Q}}$ are defined as

$$\begin{aligned} R_{\mathbb{Q},j} &\triangleq \liminf_{\kappa \rightarrow \infty} \frac{\ln |\mathcal{M}_j^{(\kappa)}|}{\mathbf{E}^{(\kappa)}[\mathbf{T}^{(\kappa)}]} \quad \forall j \in \{1, 2, \dots, \ell\} \\ E_{\mathbb{Q},j} &\triangleq \liminf_{\kappa \rightarrow \infty} \frac{-\ln P_e(j)^{(\kappa)}}{\mathbf{E}^{(\kappa)}[\mathbf{T}^{(\kappa)}]} \quad \forall j \in \{1, 2, \dots, \ell\}. \end{aligned}$$

A rate-exponent vector (\vec{R}, \vec{E}) is achievable if and only if there exists a reliable sequence \mathbb{Q} such that $(\vec{R}, \vec{E}) = (\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}})$.

This definition of the *bit-wise* UEP problem is slightly different than the one described in the introduction, because $P_e(j)$ is defined as $\mathbf{P}[\widehat{M}^j \neq M^j]$ rather than $\mathbf{P}[\widehat{M}_j \neq M_j]$. Note that if $\widehat{M}_j \neq M_j$, then $\widehat{M}^j \neq M^j$; consequently, $\mathbf{P}[\widehat{M}^j \neq M^j] \geq \mathbf{P}[\widehat{M}_j \neq M_j]$ for all j 's. In addition, if we assume without loss of generality that $\mathbf{P}[\widehat{M}_j \neq M_j] \geq \mathbf{P}[\widehat{M}_i \neq M_i]$ for all $j \geq i$, the union bound implies that $\mathbf{P}[\widehat{M}^j \neq M^j] \leq j \mathbf{P}[\widehat{M}_j \neq M_j]$. Thus, for the case when ℓ is fixed, both formulations of the problem result in exactly the same achievable region of rate-exponent vectors.

The achievable region of rate-exponent vectors could have been defined as the closure of the points of the form $(\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}})$ for some reliable sequence \mathbb{Q} . Using the definition of $(\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}})$'s, one can easily show that, in this case too both

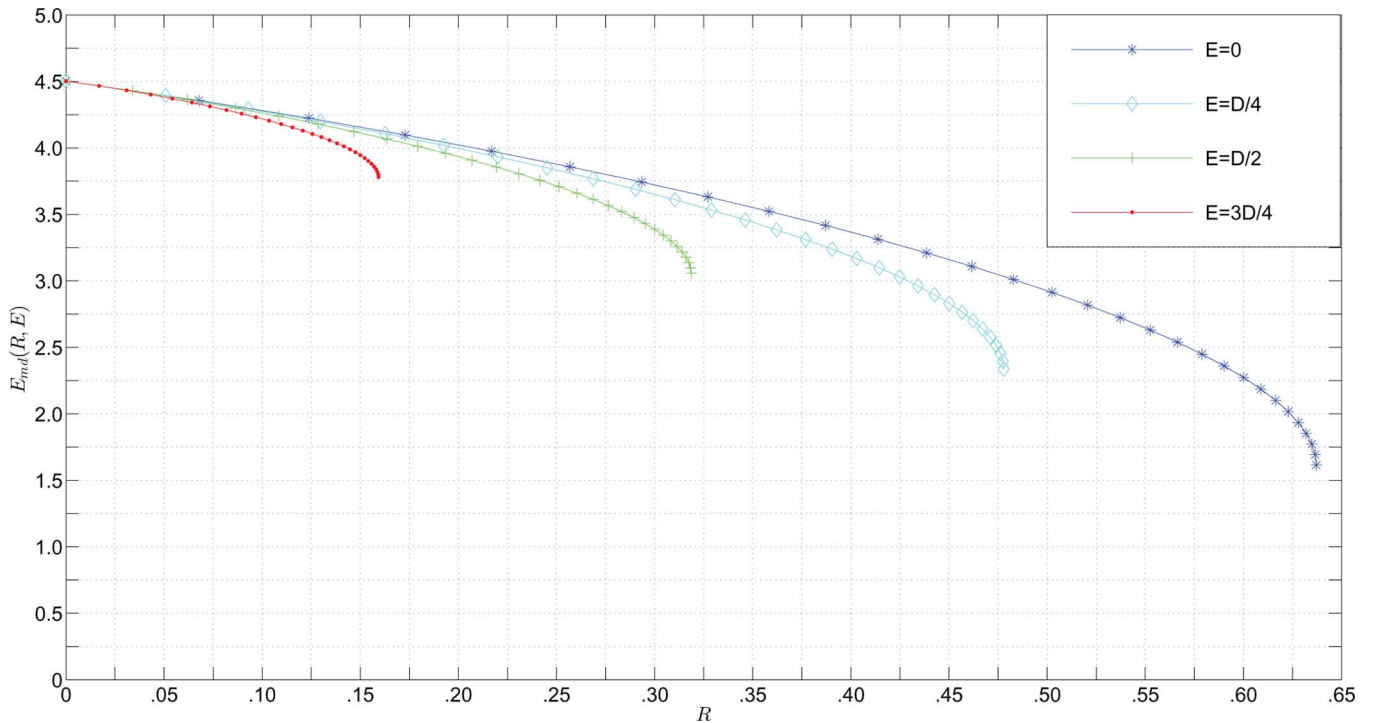


Fig. 2. $E_{\text{md}}(R, E)$ is drawn at various values of the error exponent E as a function of rate R for a BSC with crossover probability $p = 0.01$. Note that when $p = 0.01$, $C = 0.6371$ Nats per channel use and $D = 4.503$. As we increase the exponent of the average error probability, i.e., E , the value of $E_{\text{md}}(R, E)$ decreases, as one would expect.

definitions result in exactly the same achievable region of rate-exponent vectors.

B. Main Results

For variable-length block codes with feedback, the results of both the single-message *message-wise* UEP problem and the *bit-wise* UEP problem are given in terms of the $J(R)$ function defined in the following. The $J(R)$ function is first introduced by¹⁹ Kudryashov [7, eq. (2.6)], while describing the performance of nonblock variable-length codes with feedback and delay constraints. Later, the $J(R)$ function is used in [2] for describing the performance of block codes in single-message *message-wise* UEP problem. It is shown in [2, Appendix D] that for both fixed-length block codes without feedback and variable-length block codes with feedback on DMCs satisfy

$$E_{\text{md}}(R, 0) = J(R). \quad (16)$$

Recently Nazer *et al.* obtained the closed-form expression for $E_{\text{md}}(R, 0)$ for fixed-length block codes on the additive white Gaussian noise channel, under certain average and peak power constraints [9, Th. 1]. Curiously, equality given in (16) holds in that case too.²⁰

Definition 12: For any $R \in [-\infty, C]$, $J(R)$ is defined as

$$J(R) \triangleq \max_{\substack{0 \leq \alpha \leq 1 \\ x_1, x_2 \in \mathcal{X} \\ \alpha, x_1, x_2, \mu_1, \mu_2: \\ \mu_1, \mu_2 \in \mathcal{P}(\mathcal{X}) \\ \alpha I(\mu_1, W) + (1-\alpha)I(\mu_2, W) \geq R}} \alpha D(\bar{\mu}_1 \| W_{x_1}) + (1-\alpha)D(\bar{\mu}_2 \| W_{x_2}) \quad (17)$$

¹⁹In [7, eq. (2.6)], there is no optimization over the parameter α . Thus, strictly speaking, what is introduced in [7, eq. (2.6)] is $j(R)$ given in (64) rather than $J(R)$ given in (17).

²⁰Unlike DMC for these channels, it is possible to obtain a closed-form expression in terms of the rate and the power constraints.

where $\bar{\mu}_i(y) = \sum_x W_x(y) \mu_i(x)$ for $i = 1, 2$.

We have plotted the $J(R)$ function for binary symmetric channels²¹ (BSCs) with various crossover probabilities in Fig. 1. Note that as the channel becomes noisier, i.e., as the crossover probability becomes closer to $1/2$, the value of $J(R)$ function decreases at all values of rate where it is positive. Furthermore, the highest value of rate where it is positive, i.e., the channel capacity, decreases.

Lemma 1: The function $J(R)$ defined in (17) is a concave, decreasing function such that $J(R) = D$ for $R \leq 0$.

Proof of Lemma 1 is given in Section A in the Appendix.

Now, let us consider the single-message *message-wise* UEP problem given in Definition 10.

Theorem 1: For any rate $0 \leq R \leq C$ and error exponent $E \leq (1 - \frac{R}{C})D$, the missed detection exponent $E_{\text{md}}(R, E)$ defined in (13) is equal to²²

$$E_{\text{md}}(R, E) = E + \left(1 - \frac{E}{D}\right) J\left(\frac{R - \frac{E}{D}}{1 - \frac{E}{D}}\right) \quad (18)$$

where C , D , and $J(\cdot)$ are given in (8), (9), and (17), respectively. Furthermore $E_{\text{md}}(R, E)$ is jointly concave in (R, E) pairs.

We have plotted $E_{\text{md}}(R, E)$ as a function of rate, for various values of E in Fig. 2. When rate is zero, the exponent of the average error probability can be made as high as D . Thus, all

²¹Recall that in a BSC with crossover probability p , $\mathcal{X} = \{0, 1\}$, $\mathcal{Y} = \{0, 1\}$, and $W_x(y) = (1-p)\mathbb{1}_{\{x=y\}} + p\mathbb{1}_{\{x \neq y\}}$.

²²For the case when $R = 0$ and $E = D$, the $\left(1 - \frac{E}{D}\right) J\left(\frac{R - \frac{E}{D}}{1 - \frac{E}{D}}\right)$ term should be interpreted as 0, i.e., $\left(1 - \frac{E}{D}\right) J\left(\frac{R - \frac{E}{D}}{1 - \frac{E}{D}}\right) \Big|_{\substack{R=0 \\ E=D}} = 0$.

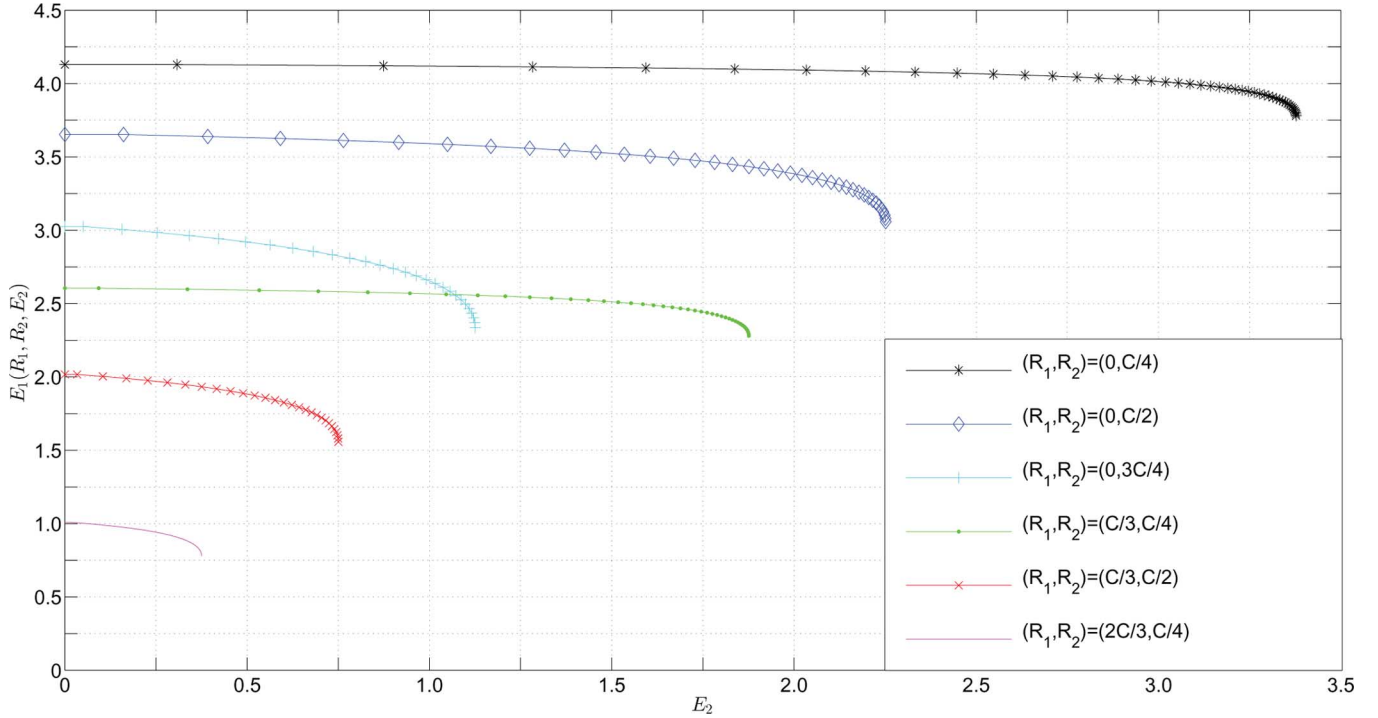


Fig. 3. $E_1(R_1, R_2, E_2)$ is drawn for various values of rate pairs (R_1, R_2) as a function of error exponent E_2 for a BSC with crossover probability $p = 0.01$. Recall that when $p = 0.01$, $C = 0.6371$ Nats per channel use and $D = 4.503$.

the curves meet at $(0, D)$ point. But for all positive rates, the exponent of the average error probability makes a difference; as E increases, $E_{\text{md}}(R, E)$ decreases. Furthermore, for any given rate R , the exponent of the average error probability can only be as high as $(1 - \frac{R}{C})D$. This is why the curves corresponding to higher values of E have smaller support on rate axis.

Proof of Theorem 1 is presented in Section I in the Appendix.

Similar to the single-message *message-wise* UEP problem, the solution of the *bit-wise* UEP problem is given in terms of the $J(R)$ function.

Theorem 2: A rate-exponent vector (\vec{R}, \vec{E}) is achievable if and only if there exists a $\vec{\eta}$ such that²³

$$E_i \leq (1 - \sum_{j=1}^{\ell} \eta_j)D + \sum_{j=i+1}^{\ell} \eta_j J\left(\frac{R_j}{\eta_j}\right) \quad \forall i \in \{1, 2, \dots, \ell\} \quad (19a)$$

$$R_i \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (19b)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (19c)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1 \quad (19d)$$

where C , D , and $J(\cdot)$ are given in (8), (9), and (17), respectively. Furthermore, the set of all achievable rate-exponent vectors is convex.

Proof of Theorem 2 is presented in Section J in the Appendix.

For the special case when there are only two submessages, the condition given in Theorem 2 for the achievability of a rate

²³For the case when $R_j = 0$ and $\eta_j = 0$, the $\eta_j J\left(\frac{R_j}{\eta_j}\right)$ term should be interpreted as 0, i.e., $\eta_j J\left(\frac{R_j}{\eta_j}\right)\Big|_{\substack{R_j=0 \\ \eta_j=0}} = 0$.

vector error exponent vector pair can be turned into an analytical expression for the optimal E_1 in terms of R_1 , R_2 , and E_2 . In order to see why, note that revealing the region of achievable (R_1, R_2, E_1, E_2) vectors is equivalent to revealing the region of achievable (R_1, R_2, E_2) 's and the value of the maximum achievable E_1 for all the (R_1, R_2, E_2) 's in the achievable region.

Corollary 1: For any rate pair (R_1, R_2) such that $R_1 + R_2 \leq C$ and error exponent E_2 such that $E_2 \leq (1 - \frac{R_1 + R_2}{C})D$, the optimal value of E_1 is given by²⁴

$$E_1(R_1, R_2, E_2) = E_2 + \left(1 - \frac{R_1}{C} - \frac{E_2}{D}\right) J\left(\frac{R_2}{1 - \frac{R_1}{C} - \frac{E_2}{D}}\right) \quad (20)$$

where C , D , and $J(\cdot)$ are given in (8), (9), and (17), respectively. Furthermore $E_1(R_1, R_2, E_2)$ is concave in (R_1, R_2, E_2) .

Note that for the $E_1(R_1, R_2, E_2)$ given in (20), $E_1(R_1, R_2, E_2) \geq E_2$ for all (R_1, R_2, E_2) triples such that $R_1 + R_2 \leq C$ and $E_2 \leq (1 - \frac{R_1 + R_2}{C})D$. Furthermore, inequality is strict as long as $R_2 > 0$. We have drawn $E_1(R_1, R_2, E_2)$ for various (R_1, R_2) pairs as a function of E_2 in Fig. 3.

IV. ACHIEVABILITY

In both the single-message *message-wise* UEP problem and the *bit-wise* UEP problem, the codes that achieve the optimal performance employ a number of different ideas at the same time. In order to avoid introducing all of those ideas at once,

²⁴For the case when $R_2 = 0$ and $E_2 = (1 - \frac{R_1}{C})D$, the second term on the right-hand side of (20) should be interpreted as zero, i.e., $\left(1 - \frac{R_1}{C} - \frac{E_2}{D}\right) J\left(\frac{R_2}{1 - \frac{R_1}{C} - \frac{E_2}{D}}\right)\Big|_{R_2=0, E_2=(1-\frac{R_1}{C})D} = 0$

we first describe two families of codes and analyze the probabilities of various error events in those two families of codes. Later, we use those two families of codes as the building blocks for the codes that achieve the optimal performance in the UEP problems we are interested in. Before going into a more detailed description and analysis of those codes, let us first give a birds eye plan for this section.

- 1) *A Single-Message Message-wise UEP Scheme without Feedback*: First, in Section IV-A, we consider a family of fixed-length codes without feedback. We prove that these codes can achieve any rate R less than channel capacity, with vanishing²⁵ error probability P_e while having a minimum conditional error probability, $\min_m P_{e|m}$, as low as $e^{-nJ(R)}$. The main drawback of this family of codes is that the decay rate of the average error probability P_e has to be subexponential in this family of codes.
- 2) *Control Phase and Error-Erasure Decoding*: In Section IV-B, in order to obtain nonzero exponential decay for the average error probability, we use a method introduced by Yamamoto and Itoh in [12]. We append the fixed-length codes described in Section IV-A with a control phase and use an error-erasure decoder. This new family of codes with control phase and error-erasure decoding is shown, in Section IV-B, to achieve any rate R less than the channel capacity C with exponentially decaying average error probability P_e , exponentially decaying minimum conditional error probability $\min_m P_{e|m}$, and vanishing erasure probability P_x .
- 3) *Single-Message Message-wise UEP for Variable-Length Codes*: In Section IV-C, we obtain variable-length codes for single-message *message-wise* UEP problem using the codes described in Section IV-B. In order to do that, we use the fixed-length codes with feedback and erasures described in Section IV-B, repetitively until a nonerasure decoding happens. This idea too was employed by Yamamoto and Itoh in [12].
- 4) *Bit-wise UEP for Variable-Length Codes*: In Section IV-D, we first use the codes described in Section IV-A and the control phase discussed in Section IV-B to obtain a family of fixed-length codes with feedback and erasures which has *bit-wise* UEP, i.e., which has different bounds on error probabilities for different submessages. While using the codes described in Section IV-A, we employ an implicit acceptance explicit rejection scheme first introduced in [7] by Kudryashov. Once we obtain a fixed-length code with erasures and *bit-wise* UEP, we use a repeat at erasures scheme like the one described in Section IV-C to obtain a variable-length code with *bit-wise* UEP.

²⁵Vanishing with increasing block length.

The achievability results we derive in this section are revealed to be the optimal ones, in terms of the decay rates of error probabilities with expected decoding time $\mathbf{E}[T]$, as a result of the outer bounds we derive in Section V.

A. Single-Message Message-Wise UEP Scheme Without Feedback

In this section, we describe a family of fixed-length block codes without feedback that achieves any rate R less than capacity with small error probability while having an exponentially small $\min_m P_{e|m}$, for sufficiently large block length n . We describe these codes in terms of a time-sharing constant $\alpha \in [0, 1]$, two input letters $x_1, x_2 \in \mathcal{X}$, and two probability distributions on the input alphabet, $\mu_1, \mu_2 \in P(\mathcal{X})$.

In order to point out that certain sequence of input letters is a codeword or part of a codeword for message m , we put (m) after it. Hence, we denote the codeword for m by $x^n(m)$ in a given code and by $X^n(m)$ in a code ensemble, as a random quantity.

Let us start with describing the encoding scheme. The codeword of the first message, i.e., $x^n(1)$, is x_1 in first $n_\alpha = \lfloor \alpha n \rfloor$ time instances and x_2 in the rest, i.e., $x_\tau(1) = x_1$ for $\tau = 1, \dots, n_\alpha$ and $x_\tau(1) = x_2$ for $\tau = (n_\alpha + 1), \dots, n$. The codewords of the other messages are described via a random coding argument. In the ensemble of codes, we are considering all entries of all codewords other than the first codeword, i.e., $X_\tau(m) \forall \tau \in [1, n], \forall m \neq 1$, are generated independently of other codewords and other entries of the same codeword. In the first n_α time instances, $X_\tau(m)$ is generated using μ_1 , in the rest using μ_2 , i.e., $\mathbf{P}[X_\tau(m) = x] = \mu_1(x)$ for $\tau = 1, \dots, n_\alpha$ and $\mathbf{P}[X_\tau(m) = x] = \mu_2(x)$ for $\tau = (n_\alpha + 1), \dots, n$.

Let us begin the description of the decoding scheme, by specifying the decoding region of the first message $\mathcal{G}[1]$: it is the set of all output sequences y^n whose empirical distribution is not typical with $(\alpha, \bar{\mu}_1, \bar{\mu}_2)$. More precisely, the decoding region of the first message $\mathcal{G}[1]$ is given by (21) shown at the bottom of the page where Δ is the total variation distance defined in (1), $Q_{\{y_1^{n_\alpha}\}}$ and $Q_{\{y_{n_\alpha+1}^n\}}$ are the empirical distributions of $y_1^{n_\alpha}$ and $y_{n_\alpha+1}^n$ defined in (4), and $\bar{\mu}_1$ and $\bar{\mu}_2$ are probability distributions on \mathcal{Y} , i.e., $\bar{\mu}_1, \bar{\mu}_2 \in P(\mathcal{Y})$, such that $\bar{\mu}_i(y) = \sum_x \mu_i(x) W_x(y)$.

For other messages, $m \neq 1$, decoding regions $\mathcal{G}[m]$ are the set of all output sequences for which $Q_{\{x^n(m), y^n\}}$ is typical with $(\alpha, \mu_1 W, \mu_2 W)$ and $Q_{\{x^n(\tilde{m}), y^n\}}$ is not typical with $(\alpha, \mu_1 W, \mu_2 W)$ for any $\tilde{m} \neq m$. To be precise, the decoding region of the messages other than the first message are

$$\mathcal{G}[m] = \mathcal{B}[x^n(m)] \cap \left(\bigcap_{\tilde{m} \neq m} \overline{\mathcal{B}[x^n(\tilde{m})]} \right) \quad \forall m \in \{2, 3, \dots, |\mathcal{M}|\} \quad (22)$$

where for all $x^n \in \mathcal{X}^n$, $\mathcal{B}[x^n]$ is the set of all y^n 's for which (x^n, y^n) is typical with $(\alpha, \mu_1 W, \mu_2 W)$ shown at the bottom

$$\mathcal{G}[1] = \left\{ y^n : n_\alpha \Delta \left(Q_{\{y_1^{n_\alpha}\}}, \bar{\mu}_1 \right) + (n - n_\alpha) \Delta \left(Q_{\{y_{n_\alpha+1}^n\}}, \bar{\mu}_2 \right) \geq |\mathcal{X}| |\mathcal{Y}| \sqrt{n \ln(1+n)} \right\} \quad (21)$$

of the page where Δ is the total variation distance defined in (1), $Q_{\{x_1^{n_\alpha}, y_1^{n_\alpha}\}}$ and $Q_{\{x_{n_\alpha+1}^{n_\alpha}, y_{n_\alpha+1}^{n_\alpha}\}}$ are the empirical distributions of $(x_1^{n_\alpha}, y_1^{n_\alpha})$ and $(x_{n_\alpha+1}^{n_\alpha}, y_{n_\alpha+1}^{n_\alpha})$ defined in (4), and $\mu_1 W$ and $\mu_2 W$ are probability distributions on $\mathcal{X} \times \mathcal{Y}$, i.e., $\mu_1 W \in P(\mathcal{X} \times \mathcal{Y})$ and $\mu_2 W \in P(\mathcal{X} \times \mathcal{Y})$.

In Section B in the Appendix, we have analyzed the conditional error probabilities $P_{e|m}$ for the previously described code and proved Lemma 2 given in the following.

Lemma 2: For any block length n , time-sharing constant $\alpha \in [0, 1]$, input letters $x_1, x_2 \in \mathcal{X}$, and input distributions $\mu_1, \mu_2 \in P(\mathcal{X})$, there exists a length n block code such that

$$\begin{aligned} |\mathcal{M}| &\geq e^{n(\alpha I(\mu_1, W) + (1-\alpha)I(\mu_2, W) - \varepsilon_n)} \\ P_{e|1} &\leq e^{-n(\alpha D(\bar{\mu}_1 \| W_{x_1}) + (1-\alpha)D(\bar{\mu}_2 \| W_{x_2}) - \varepsilon_n)} \\ P_{e|m} &\leq \varepsilon_n \quad m = 2, 3, \dots, |\mathcal{M}| \end{aligned}$$

where $\bar{\mu}_1(y) = \sum_x W_x(y) \mu_1(x)$, $\bar{\mu}_2(y) = \sum_x W_x(y) \mu_2(x)$ and $\varepsilon_n = \frac{10|\mathcal{X}||\mathcal{Y}|(\ln \frac{e}{\alpha})\sqrt{\ln(1+n)}}{\sqrt{n}}$.

Given the channel W , if we discard the error terms ε_n , for a given value of rate, $0 \leq R \leq C$, we can optimize exponent of $P_{e|1}$ over the time-sharing constant α , the input letters x_1, x_2 , and input distributions μ_1, μ_2 . Evidently, the optimization problem we get is the one given for the definition of $J(R)$, in (17). Thus, Lemma 2 implies that for any $R \in [0, C]$ and block length n , there exists a length n code such that $|\mathcal{M}| \geq e^{n(R-\varepsilon_n)}$, $P_{e|m} \leq \varepsilon_n$ for $m = 2, 3, \dots, |\mathcal{M}|$ and $P_{e|1} \leq e^{-n(J(R)-\varepsilon_n)}$.

One curious question is whether or not the exponent of $P_{e|1}$ can be increased by including more than two phases. Carathéodory's theorem answers that question negatively, i.e., to obtain the largest value of $J(R)$ one does not need to do time sharing between more than two input-letter-input-distribution pairs.

B. Control Phase and Error-Erasure Decoding

The family of codes described in Lemma 2 has a large exponent for the conditional error probability of the first message, i.e., $P_{e|1}$. But the conditional error probabilities of other messages, $P_{e|m}$ for $m \neq 1$, decay subexponentially. In order to facilitate an exponential decay of $P_{e|m}$ for $m \neq 1$ with block length, we append the codes described in Lemma 2 with a control phase and allow erasures. The idea of using a control phase and an error-erasure decoding, in establishing achievability results for variable-length code, was first employed by Yamamoto and Itoh in [12].

In order to explain what we mean by the control phase, let us describe our encoding scheme and decoding rule briefly. First, a code from the family of codes described in Section IV-A is used to transmit M and the receiver makes a tentative decision \hat{M} using the decoder of the very same code. The transmitter

knows \hat{M} because of the feedback link. In the remaining time instances, i.e., in the control phase, the transmitter sends the input letter a if $\hat{M} = M$, the input letter r if $\hat{M} \neq M$. The input letters a and r are described in (10). At the end of the control phase, the receiver checks whether or not the output sequence in the control phase is typical with W_a , if it is then $\hat{M} = \hat{M}$ otherwise an erasure is declared.

Lemma 3 given in the following states the results of the performance analysis of the previously described code. In order to understand what is stated in Lemma 3 accurately, let us make a brief digression and elaborate on the codes with erasure. We have assumed in our models until now that $\hat{M} \in \mathcal{M}$. However, there are many interesting problems in which this might not hold. In codes with erasures, for example, we replace $\hat{M} \in \mathcal{M}$ with $\hat{M} \in \hat{\mathcal{M}}$ where $\hat{\mathcal{M}} = \mathcal{M} \cup \{x\}$ and x is the erasure symbol. Furthermore, in codes with erasures for each $m \in \mathcal{M}$, the conditional error probability $P_{e|m}$ and conditional erasure probability $P_{x|m}$ are defined as follows:

$$P_{e|m} = \mathbf{P} \left[\hat{M} \notin \{m, x\} \mid M = m \right] \quad \forall m \in \mathcal{M} \quad (24a)$$

$$P_{x|m} = \mathbf{P} \left[\hat{M} = x \mid M = m \right] \quad \forall m \in \mathcal{M}. \quad (24b)$$

Note that definitions of $P_{e|m}$ and $P_{x|m}$ given previously can be seen as the generalizations of the corresponding definitions in block codes without erasures. In erasure free codes, aforementioned definitions are equivalent to corresponding definitions there.

Lemma 3: For any block length n , rate $0 \leq R \leq C$, and error exponent $0 \leq E \leq (1 - \frac{R}{C})D$, there exists a length n block code with erasures such that

$$\begin{aligned} |\mathcal{M}| &\geq e^{n(R-\varepsilon_n)} \\ P_{e|1} &\leq e^{-n(E + (1 - \frac{R}{D})J(\frac{R}{1-E/D}) - \varepsilon_n)} \\ P_{e|m} &\leq \varepsilon_n \min\{1, e^{-n(E-\varepsilon_n)}\} \quad m = 2, 3, \dots, |\mathcal{M}| \\ P_{x|m} &\leq \varepsilon_n + e^{-n((1 - \frac{R}{D})J(\frac{R}{1-E/D}) - \varepsilon_n)} \quad m = 1, 2, \dots, |\mathcal{M}| \end{aligned}$$

where $\varepsilon_n = \frac{10|\mathcal{X}||\mathcal{Y}|(\ln \frac{e}{\alpha})\sqrt{\ln(1+n)}}{\sqrt{n}}$.

Proof of Lemma 3 is given in Section C in the Appendix.

Note that in Lemma 3, unlike $P_{e|m}$'s which decrease exponentially with n , $P_{x|m}$'s decays as $\frac{\ln n}{\sqrt{n}}$. It is possible to tweak the proof so as to have a nonzero exponent for $P_{x|m}$'s, see [8]. But this can only be done at the expense of $P_{e|m}$'s. Our aim, however, is to achieve the optimal performance in variable-length block codes. As we will see in the following section, for that what matters is exponents of error probabilities and having vanishing erasure probabilities. The rate at which erasure probability decays does not effect the performance of variable-length block codes in terms of error exponents.

$$\mathcal{B}[x^n] = \{y^n : n_\alpha \Delta(Q_{\{x_1^{n_\alpha}, y_1^{n_\alpha}\}}, \mu_1 W) + (n - n_\alpha) \Delta(Q_{\{x_{n_\alpha+1}^{n_\alpha}, y_{n_\alpha+1}^{n_\alpha}\}}, \mu_2 W) \leq |\mathcal{X}||\mathcal{Y}| \sqrt{n \ln(1+n)}\} \quad (23)$$

C. Single-Message Message-Wise UEP Achievability

In this section, we construct variable-length block codes for the single-message *message-wise* UEP problem using Lemma 3. In first n time units, the variable-length encoding scheme uses a fixed-length block code with erasures which has the performance described in Lemma 3. If the decoded message of the fixed-length code is in the message set, i.e., if $\widehat{M} \in \mathcal{M}$, then decoded message of the fixed-length code becomes the decoded message of the variable-length code. If the decoded message of the fixed-length code is the erasure symbol, i.e., if $\widehat{M} = \mathbf{x}$, then the encoder uses the fixed-length code again in the second n time units. By repeating this scheme until the decoded message of the fixed-length code is in \mathcal{M} , i.e., $\widehat{M} \in \mathcal{M}$, we obtain a variable-length code.

Let L be the number of times the fixed-length code is used until a $\widehat{M} \in \mathcal{M}$ is observed. Then, given the message M , L is a geometrically distributed random variable with success probability $(1 - P_{\mathbf{x}|M})$ where $P_{\mathbf{x}|M}$ is the conditional erasure probability of the fixed-length code given the message M . Then, the conditional probability distribution and the conditional expected value of L given M are

$$\mathbf{P}[L = l | M] = (1 - P_{\mathbf{x}|M})(P_{\mathbf{x}|M})^{l-1} \quad l = 1, 2, \dots \quad (25a)$$

$$\mathbf{E}[L | M] = (1 - P_{\mathbf{x}|M})^{-1}. \quad (25b)$$

Furthermore, the conditional expected value of decoding time and the conditional error probability given the message M are

$$\mathbf{E}[T | M] = n\mathbf{E}[L | M] \quad (26a)$$

$$\mathbf{P}[\widehat{M} \neq M | M] = P_{e|M}\mathbf{E}[L | M] \quad (26b)$$

where n is the block length of the fixed-length code and $P_{e|M}$ is the conditional error probability given the message M for the fixed-length code.

Thus, as result of (25b) and (26) and Lemma 3, we know that for any rate $R \in [0, C]$, error exponent $E \in [0, (1 - \frac{R}{C})D]$, there exists a reliable sequence \mathbb{Q} such that $R_{\mathbb{Q}} = R$, $E_{\mathbb{Q}} = E$, and

$$E_{\text{md}\mathbb{Q}} = E + (1 - \frac{E}{D})J\left(\frac{R}{1-E/D}\right). \quad (27)$$

We show in Section V-C that for any reliable sequence \mathbb{Q} with rate $R_{\mathbb{Q}} = R$ and error exponent $E_{\mathbb{Q}} = E$, $E_{\text{md}\mathbb{Q}}$ is upper bounded by the expression on the right-hand side of (27).

D. Bit-Wise UEP Achievability

In this section, we first use the family of codes described in Section IV-A and the control phase idea described in Section IV-B to construct fixed-length block codes with erasures which have *bit-wise* UEP. Then, we use them with a repeat until nonerasure decoding scheme, similar to the one described in Section IV-C, to obtain variable-length block codes with *bit-wise* UEP.

Let us start with describing the encoding scheme for the fixed-length block code with *bit-wise* UEP. If there are ℓ submessages, i.e., if $\mathcal{M} = (\mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_\ell)$, then the encoding scheme has $\ell + 1$ phases with lengths $n_1, n_2, \dots, n_{\ell+1}$ such that $n_1 + n_2 + \dots + n_{\ell+1} = n$.

1) In the first phase, a length n_1 code from the family of codes described in Section IV-A is used. The message set of the code ${}_{\tau}\mathcal{M}_1$ is $\mathcal{M}_1 \cup \{|\mathcal{M}_1| + 1\}$ and the message ${}_{\tau}M_1$ of the code is determined by the first submessage ${}_{\tau}M_1 = M_1 + 1$. At the end of first phase, receiver uses the decoder of the length n_1 code to get a tentative decision ${}_{\tau}\widehat{M}_1$ which is known by the transmitter at the beginning of the second phase because of the feedback link.

2) In the second phase, a length n_2 code from the family of codes described in Section IV-A, with the message set ${}_{\tau}\mathcal{M}_2 = \mathcal{M}_2 \cup \{|\mathcal{M}_2| + 1\}$, is used. If ${}_{\tau}\widehat{M}$ is decoded correctly at the end of the first phase, then the message ${}_{\tau}M_2$ of the code used in the second phase is determined by the second submessage as ${}_{\tau}M_2 = M_2 + 1$, else ${}_{\tau}M_2 = 1$. At the end of the second phase, the receiver uses the decoder of the second phase code to get the tentative decision ${}_{\tau}\widehat{M}_2$ which is known by the transmitter at the beginning of the third phase because of the feedback link.

3) In phases 3 to ℓ , the previously described scheme is used. In phase i , a length n_i code, with the message set ${}_{\tau}\mathcal{M}_i = \mathcal{M}_i \cup \{|\mathcal{M}_i| + 1\}$, from the family of codes described in Section IV-A is used. The message of the length n_i code ${}_{\tau}M_i$ is $M_i + 1$ if ${}_{\tau}\widehat{M}_{i-1} = {}_{\tau}M_{i-1}$, 1 otherwise for $i = 3, 4, \dots, \ell$.

4) The last phase is a $n_{\ell+1}$ long control phase, i.e., a $n_{\ell+1}$ long code with the message set ${}_{\tau}\mathcal{M}_{\ell+1} = \{1, 2\}$ is used in the last phase. The codewords for the first and second messages are $n_{\ell+1}$ long sequences of input letters \mathbf{r} and \mathbf{a} , respectively, where \mathbf{r} and \mathbf{a} are described in (10). The tentative decision in the last phase ${}_{\tau}\widehat{M}_{\ell+1}$ is equal to the first message if the output sequence in the last phase is not typical with $W_{\mathbf{a}}$, the second message otherwise. The message of the $n_{\ell+1}$ long code ${}_{\tau}M_{\ell+1}$ is equal to 2 if ${}_{\tau}\widehat{M}_{\ell} = {}_{\tau}M_{\ell}$, 1 otherwise.

Note that if we define ${}_{\tau}\widehat{M}_0, {}_{\tau}M_0$ and $M_{\ell+1}$ all to be 1, i.e., ${}_{\tau}\widehat{M}_0 = {}_{\tau}M_0 = M_{\ell+1} = 1$, we can write the following rule for determining the ${}_{\tau}M_i$'s for $i = 1$ to $\ell + 1$

$${}_{\tau}M_i = 1 + \mathbb{1}_{\{\widehat{M}_{i-1} = {}_{\tau}M_{i-1}\}} M_i \quad i = 1, 2, \dots, (\ell + 1). \quad (28)$$

It is important however to keep in mind that the last phase is a control phase and the codes in the first ℓ phases are from the family of codes described in Section IV-A.

Note that during the phases $i = 2$ to ℓ , erroneous transmission of ${}_{\tau}M_{i-1}$ is conveyed using ${}_{\tau}M_i = 1$; hence, the transmission of M_i through ${}_{\tau}M_i$, i.e., ${}_{\tau}M_i = 1 + M_i$, is a tacit approval of the tentative decision ${}_{\tau}\widehat{M}_{i-1}$. Because of this, the aforementioned encoding scheme is said to have an implicit acceptance explicit rejection property. The idea of implicit acceptance explicit rejection was first introduced by Kudryashov in [7] in the context of nonblock variable-length codes with feedback and delay constraints.

After finishing the description of the encoding scheme, we are ready to describe the decoding scheme. The receiver determines the decoded message using the tentative decisions, ${}_{\tau}\widehat{M}_i$ for $i = 1$ to $\ell + 1$. If one or more of the tentative decisions are equal to 1, then an erasure is declared. If all $\ell + 1$ tentative decision are

different from 1, then $\widehat{M}_i = \widehat{\tau}M_i - 1$ for all $i = 1, 2, \dots, \ell$. Hence, the decoding rule is

$$(\widehat{M}_1, \widehat{M}_2, \dots, \widehat{M}_\ell) = \begin{cases} (\widehat{\tau}M_1 - 1, \widehat{\tau}M_2 - 1, \dots, \widehat{\tau}M_\ell - 1) & \text{if } \prod_{i=1}^{\ell+1} (\widehat{\tau}M_i - 1) > 0 \\ \mathbf{x} & \text{if } \prod_{i=1}^{\ell+1} (\widehat{\tau}M_i - 1) = 0. \end{cases} \quad (29)$$

For *bit-wise* UEP codes with erasure, the definition of $P_e(i)$ is slightly different from the original one given in (14)

$$P_e(i) = \mathbf{P} \left[\{ \widehat{M}^i \neq m^i, \widehat{M} \neq \mathbf{x} \} \right]. \quad (30)$$

With this alternative definition in mind, let us define $P_{e|m}(i)$ as the conditional probability of the erroneous transmission of any one of the first i submessage when $M = m$

$$P_{e|m}(i) = \mathbf{P} \left[\{ \widehat{M}^i \neq m^i, \widehat{M} \neq \mathbf{x} \} \mid M = m \right]. \quad (31)$$

The error analysis of the previously described fixed-length codes, presented in Section D in the Appendix, leads to Lemma 4 given in the following.

Lemma 4: For block length n , any integer $\ell \leq \frac{n}{\ln(1+n)}$, rate vector \vec{R} , and time-sharing vector $\vec{\eta}$ such that

$$R_i \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (32a)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (32b)$$

$$\sum_{i=1}^{\ell} \eta_i \leq 1 \quad (32c)$$

there exists a length n block code such that

$$|\mathcal{M}_i| \geq e^{n(R_i - \varepsilon_{n,\ell})} \quad \forall i \in \{1, 2, \dots, \ell\}$$

$$|\mathcal{M}^i| \geq e^{n(-\varepsilon_{n,\ell} + \sum_{j=1}^i R_j)} \quad \forall i \in \{1, 2, \dots, \ell\}$$

$$P_{e|m}(i) \leq \varepsilon_{n,\ell} e^{-n \left(-\varepsilon_{n,\ell} + \sum_{j=i+1}^{\ell+1} \eta_j J \left(\frac{R_j}{\eta_j} \right) \right)} \quad \forall m \in \mathcal{M}, i \in \{1, 2, \dots, \ell\}$$

$$P_{\mathbf{x}|m} \leq \varepsilon_{n,\ell} \quad \forall m \in \mathcal{M}$$

where

$$\eta_{\ell+1} = 1 - \sum_{i=1}^{\ell} \eta_i, R_{\ell+1} = 0, \varepsilon_{n,\ell} = \frac{10|\mathcal{X}||\mathcal{Y}| \ln(\frac{e}{\lambda}) \sqrt{\ln(1+n)}}{\sqrt{n}} \sqrt{1+\ell}.$$

Recall the repeat at erasures scheme described in Section IV-C. If we use that scheme to obtain a variable-length code from the fixed-length *bit-wise* UEP code described in Lemma 4, we obtain a variable-length code with UEP such that

$$\mathbf{E}[T \mid M] = \frac{n}{1 - P_{\mathbf{x}|M}} \quad (33a)$$

$$\mathbf{P} \left[\widehat{M}^i \neq M^i \mid M \right] \leq \frac{P_{e|M}(i)}{1 - P_{\mathbf{x}|M}} \quad i = 1, 2, \dots, \ell. \quad (33b)$$

As a result of (33) and Lemma 4, we know that for any rate vector \vec{R} , error exponent vector \vec{E} , and time-sharing vector $\vec{\eta}$ such that

$$E_i \leq (1 - \sum_{j=1}^{\ell} \eta_j) D + \sum_{j=i+1}^{\ell} \eta_j J \left(\frac{R_j}{\eta_j} \right) \quad \forall i \in \{1, 2, \dots, \ell\} \quad (34a)$$

$$R_i \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (34b)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (34c)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1 \quad (34d)$$

there exists a reliable sequence \mathbb{Q} such that $(\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}}) = (\vec{R}, \vec{E})$. Thus, the existence of the time-sharing vector $\vec{\eta}$ satisfying the constraints given in (34) is a sufficient condition for the achievability of a rate-exponent vector (\vec{R}, \vec{E}) . We show in Section V-D that the existence of a time-sharing vector $\vec{\eta}$ satisfying the constraints given in (34) is also a necessary condition for the achievability of a rate-exponent vector (\vec{R}, \vec{E}) .

V. CONVERSE

Berlin *et al.* [1] used the error probability of a random binary query posed at a stopping time for bounding the error probability of a variable-length block code. Later, similar techniques have been applied in [2] for establishing outer bounds in UEP problems. Our approach is similar to that of [1] and [2]; we, too, use error probabilities of random queries posed at stopping times for establishing outer bounds. Our approach, nevertheless, is novel because of the error events we choose to analyze and the bounding techniques we use. Furthermore, the relation we establish in Lemma 5 between the error probabilities and the decay rate of the conditional entropy of the messages with time is a brand new tool for UEP problems.

For rigorously and unambiguously generalizing the technique used in [1] and [2], we introduce the concept of anticipative list decoders (ALDs) in Section V-A. Then, in Section V-B, we bound the probabilities of certain error events associated with ALDs from the following. This bound, i.e., Lemma 5, is used in Sections V-C and V-D to derive tight outer bounds for the performance of variable-length block codes in the single-message *message-wise* UEP problem and in the *bit-wise* UEP problem, respectively.

A. Anticipative List Decoders (ALDs)

In this section, we first introduce the concepts of ALDs and nontrivial ALDs. After that, we show that for a given variable-length code, any nontrivial ALD (\tilde{T}, \mathcal{A}) can be used to define a probability distribution, $P_{\{\mathcal{A}\}}$, on $\mathcal{M} \times \mathcal{Y}^{T^*}$. Finally, we use $P_{\{\mathcal{A}\}}$ to define the probability measure $\mathbf{P}_{\{\mathcal{A}\}}[\cdot]$ for the events in $\wp(\mathcal{M} \times \mathcal{Y}^T)$. Both the nontrivial ALDs (\tilde{T}, \mathcal{A}) and the probability measures $\mathbf{P}_{\{\mathcal{A}\}}[\cdot]$ associated with them play key roles in Lemma 5 of Section V-B.

An ALD for a variable-length code is a list decoder \mathcal{A} that decodes at a stopping time \tilde{T} that is always less than or equal to the decoding time of the code T . The ALDs are used to formulate

questions about the transmitted message or the decoded message, in terms of a subset of the message set \mathcal{M} that is chosen at a stopping time \tilde{T} . For example, let \mathcal{A} be the set of all $m \in \mathcal{M}$ whose posterior probability at time one is larger than $1/|\mathcal{M}|$. Evidently, for all values of Y_1 , \mathcal{A} is a subset of \mathcal{M} , but it is not necessarily the same subset for all values of Y_1 . Indeed, \mathcal{A} is a function from \mathcal{Y}_1 to the power set of \mathcal{M} and (\tilde{T}, \mathcal{A}) is an ALD, for which $\tilde{T} = 1$. Formal definition, for ALDs, is given in the following. In order to avoid separate treatment in certain special cases, we include the case when $\tilde{T} = 0$ and \mathcal{A} is fixed subset of \mathcal{M} , in the definition.

Definition 13 (ALD): For a variable-length code with decoding time T , a pair (\tilde{T}, \mathcal{A}) is called an ALD if

- 1) either \tilde{T} is the constant random variable 0 and \mathcal{A} is a fixed subset of \mathcal{M} , i.e.,

$$\begin{aligned} \tilde{T} &= 0 \\ \mathcal{A} &\in \wp(\mathcal{M}) \end{aligned}$$

- 2) or \tilde{T} is a stopping time, which is smaller than T with probability one, and \mathcal{A} is a $\wp(\mathcal{M})$ -valued function defined on $\mathcal{Y}^{\tilde{T}}$, i.e.,

$$\begin{aligned} \mathbf{P}[\tilde{T} \leq T] &= 1 \\ \mathcal{A} : \mathcal{Y}^{\tilde{T}} &\rightarrow \wp(\mathcal{M}). \end{aligned}$$

Definition of ALD does not require \mathcal{A} to be of some fixed size, nor it requires \mathcal{A} to include more likely or less likely messages. Thus, for certain values of $Y^{\tilde{T}}$, \mathcal{A} might not include any $m \in \mathcal{M}$ with positive posterior probability. In other words, for some values of $Y^{\tilde{T}}$, we might have

$$\mathbf{P}[M \in \mathcal{A}(Y^{\tilde{T}}) \mid Y^{\tilde{T}} = y^{\tilde{i}}] = 0.$$

The ALD's in which such $y^{\tilde{i}}$'s have zero probability are called nontrivial ALD's.

Definition 14 (Nontrivial ALD): An ALD (\tilde{T}, \mathcal{A}) is called a nontrivial anticipative list decoder (NALD) if $\mathbf{P}[M \in \mathcal{A}(Y^{\tilde{T}}) \mid Y^{\tilde{T}}] > 0$ with probability one, i.e.,

$$\mathbf{P}[\mathbf{P}[M \in \mathcal{A}(Y^{\tilde{T}}) \mid Y^{\tilde{T}}] > 0] = 1. \quad (35)$$

In the following, for any variable-length code and an associated NALD (\tilde{T}, \mathcal{A}) , we define a probability distribution $P_{\{\mathcal{A}\}}$ on $\mathcal{M} \times \mathcal{Y}^{T^*}$ and a probability measure $\mathbf{P}_{\{\mathcal{A}\}}[\cdot]$ for the events in $\wp(\mathcal{M} \times \mathcal{Y}^T)$. For doing that, first note that the probability measure generated by the code, i.e., $\mathbf{P}[\cdot]$, can be used to define a probability distribution P on $\mathcal{M} \times \mathcal{Y}^{T^*}$ as follows:

$$P(m, y^t) \triangleq \mathbf{P}[M = m, Y^T = y^t] \quad \forall m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*} \quad (36)$$

where \mathcal{Y}^{T^*} is a countable set for any stopping time, given in (6b).

As T is a stopping time, the probability of any event Γ in $\wp(\mathcal{M} \times \mathcal{Y}^T)$ under $\mathbf{P}[\cdot]$, i.e., $\mathbf{P}[\Gamma]$, is equal to

$$\mathbf{P}[\Gamma] = \sum_{(m, y^t) \in \Gamma \cap (\mathcal{M} \times \mathcal{Y}^{T^*})} P(m, y^t). \quad (37)$$

Evidently, we can extend the definition of P and assume that P is zero whenever y^t is in $\mathcal{Y}_{\{T=\infty\}}^\infty$, i.e.,

$$P(m, y^t) \triangleq 0 \quad \forall m \in \mathcal{M}, y^t \in \mathcal{Y}_{\{T=\infty\}}^\infty. \quad (38)$$

This extension is neither necessary nor relevant for calculating the probabilities of the events in $\wp(\mathcal{M} \times \mathcal{Y}^T)$, because T is a stopping time, i.e., $\mathbf{P}[T < \infty] = 1$.

Definition 15: Given a variable-length code with decoding time T , for any NALD (\tilde{T}, \mathcal{A}) let $P_{\{\mathcal{A}\}}$ be²⁶

$$P_{\{\mathcal{A}\}}(m, y^t) \triangleq \mathbf{P}(y^{\tilde{i}}) \frac{P(m|y^{\tilde{i}}) \mathbb{1}_{\{m \in \mathcal{A}(y^{\tilde{i}})\}}}{\sum_{\tilde{m} \in \mathcal{M}} P(\tilde{m}|y^{\tilde{i}}) \mathbb{1}_{\{\tilde{m} \in \mathcal{A}(y^{\tilde{i}})\}}} P(y_{\tilde{i}+1}^t | y^{\tilde{i}}, m) \quad \forall m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*}. \quad (39)$$

Note that Definition 15 is a parametric definition in the sense that it assigns a $P_{\{\mathcal{A}\}}$ for all NALDs (\tilde{T}, \mathcal{A}) . While proving outer bounds, we will employ not one but multiple NALD's and use them in conjunction with our new result, i.e., Lemma 5. But before introducing Lemma 5, let us elaborate on the relations between marginal and conditional distributions of $P_{\{\mathcal{A}\}}$ and P .

For $P_{\{\mathcal{A}\}}$ defined in (39), we have

$$\sum_{m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*}} P_{\{\mathcal{A}\}}(m, y^t) = 1.$$

Hence, $P_{\{\mathcal{A}\}}$ is a probability distribution on $\mathcal{M} \times \mathcal{Y}^{T^*}$, i.e., $P_{\{\mathcal{A}\}} \in P(\mathcal{M} \times \mathcal{Y}^{T^*})$.

Note that the marginal distributions of $P_{\{\mathcal{A}\}}$ and P are the same on $\mathcal{Y}^{\tilde{T}^*}$. Furthermore, for all $y^{\tilde{i}} \in \mathcal{Y}^{\tilde{T}^*}$ and $m \in \mathcal{M}$, the conditional distributions of $P_{\{\mathcal{A}\}}$ and P are the same on $\mathcal{Y}_{\{\tilde{T}^*\}}^{\tilde{T}^*}$. The probability distributions $P_{\{\mathcal{A}\}}$ and P differ only in their conditional distributions on \mathcal{M} given $y^{\tilde{i}}$. More specifically

$$P_{\{\mathcal{A}\}}(y^{\tilde{i}}) = P(y^{\tilde{i}}) \quad \forall y^{\tilde{i}} \in \mathcal{Y}^{\tilde{T}^*} \quad (40a)$$

$$P_{\{\mathcal{A}\}}(m|y^{\tilde{i}}) = \frac{P(m|y^{\tilde{i}}) \mathbb{1}_{\{m \in \mathcal{A}(y^{\tilde{i}})\}}}{\sum_{\tilde{m} \in \mathcal{M}} P(\tilde{m}|y^{\tilde{i}}) \mathbb{1}_{\{\tilde{m} \in \mathcal{A}(y^{\tilde{i}})\}}} \quad \forall y^{\tilde{i}} \in \mathcal{Y}^{\tilde{T}^*}, \forall m \in \mathcal{M} \quad (40b)$$

$$P_{\{\mathcal{A}\}}(y_{\tilde{i}+1}^t | y^{\tilde{i}}, m) = P(y_{\tilde{i}+1}^t | y^{\tilde{i}}, m) \quad \forall y^t \in \mathcal{Y}^{T^*}, \forall m \in \mathcal{M}. \quad (40c)$$

Using the parametric definition of probability distribution $P_{\{\mathcal{A}\}}$ on $\mathcal{M} \times \mathcal{Y}^{T^*}$, we define a probability measure $\mathbf{P}_{\{\mathcal{A}\}}[\cdot]$ for the events in $\wp(\mathcal{M} \times \mathcal{Y}^T)$ as follows:

$$\mathbf{P}_{\{\mathcal{A}\}}[\Gamma] \triangleq \sum_{(m, y^t) \in \Gamma \cap (\mathcal{M} \times \mathcal{Y}^{T^*})} P_{\{\mathcal{A}\}}(m, y^t) \quad \forall \Gamma \in \wp(\mathcal{M} \times \mathcal{Y}^T). \quad (41)$$

²⁶There is a slight abuse of notation in (39); if \tilde{T} is not a stopping time but rather a constant random variable $\tilde{T} = 0$, $P_{\{\mathcal{A}\}}(m, y^t)$ should be interpreted as

$$P_{\{\mathcal{A}\}}(m, y^t) \triangleq \frac{\mathbb{1}_{\{m \in \mathcal{A}\}}}{|\mathcal{A}|} P(y^t | m) \quad \forall m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*}.$$

Evidently, we can extend the definition of $\mathbf{P}_{\{\mathcal{A}\}}$ to $\mathcal{M} \times \mathcal{Y}^T$ by defining it to be zero on $\mathcal{M} \times \mathcal{Y}_{\{T=\infty\}}^\infty$, i.e.,

$$\mathbf{P}_{\{\mathcal{A}\}}(m, y^t) \stackrel{\Delta}{=} 0 \quad \forall m \in \mathcal{M}, y^t \in \mathcal{Y}_{\{T=\infty\}}^\infty. \quad (42)$$

As in the case of \mathbf{P} , this extension is neither necessary nor relevant for calculating the probabilities $\mathbf{P}_{\{\mathcal{A}\}}[\Gamma]$ given in (41).

B. Error Probability and Decay Rate of Entropy

In this section, we lower bound the probability of the event that the decoded message \hat{M} is in \mathcal{A} under the probability measure $\mathbf{P}_{\{\mathcal{A}\}}[\cdot]$, i.e., $\mathbf{P}_{\{\mathcal{A}\}}[\hat{M} \notin \mathcal{A}(Y^T)]$. The bounds we derive depend on the decay rate of the conditional entropy of the messages in the interval between \tilde{T} and T .

Before even stating our bound, we need to specify what we mean by the conditional entropy of the messages. While defining the conditional entropy, many authors do take an average over the sample values of the conditioned random variable and obtain a constant. We, however, do not take an average over the conditioned random variable and define conditional entropy as a random variable itself, which is a function of the random variable that is conditioned on²⁷

$$H(M|Y^T) \stackrel{\Delta}{=} \sum_{m \in \mathcal{M}} \mathbf{P}[M = m | Y^T] \ln \frac{1}{\mathbf{P}[M = m | Y^T]}. \quad (43)$$

Using the probability distribution \mathbf{P} defined in (36), we see that the conditional entropy defined in (43) is equal to

$$H(M|Y^T) = \mathbf{E}\left[\ln \frac{1}{\mathbf{P}(M|Y^T)} \middle| Y^T\right]. \quad (44)$$

Lemma 5: For any variable-length block code with finite expected decoding time, $\mathbf{E}[T] < \infty$, let $(T_1, \mathcal{A}_1), (T_2, \mathcal{A}_2), \dots, (T_k, \mathcal{A}_k)$ be k NALDs²⁸ such that

$$\mathbf{P}\{0 \leq T_1 \leq T_2 \leq \dots \leq T_k \leq T\} = 1. \quad (45)$$

²⁷Recall the standard notation in probability theory about the conditional expectations and conditional probabilities. Let H be a real-valued random variable and G be a random quantity that takes values from a finite set \mathcal{G} , such that $\mathbf{P}[G = g] > 0$ for all $g \in \mathcal{G}$. Then, unlike $\mathbf{E}[H]$, which is constant, $\mathbf{E}[H|G]$ is a random variable. Thus, an equation of the form $Z = \mathbf{E}[H|G]$ implies not the equality of two constants but the equality of two random variables, i.e., it means that $Z = \mathbf{E}[H|G = g]$ for all $g \in \mathcal{G}$. Similarly, let \mathcal{H}_1 be a set of sample values of the random variable H ; then, unlike $\mathbf{P}[H \in \mathcal{H}_1]$, which is a constant, $\mathbf{P}[H \in \mathcal{H}_1|G]$ is a random variable. Equations (43) and (44) are such equations. Explaining conditional expectations and conditional probabilities are beyond the scope of this paper; readers who are not sufficiently fluent with these concepts are encouraged to read [10, Ch. I, Sec. VIII], which deal the case where random variables can take finitely many values. Appropriately generalized formal treatment of the subject in terms of sigma fields is presented in [10, Ch. II, Sec. VII].

²⁸Recall ALD's and NALD's are defined in Definitions 13 and 14, respectively.

Then, for all i in $\{1, 2, \dots, k\}$ such that $(\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})] + P_e) \leq 1/2$ we have (46) shown at the bottom of the page where $T_0 = 0, T_{k+1} = T$, and for all j in $\{1, 2, \dots, (k+1)\}$, r_j 's are given by

$$r_j = \begin{cases} 0 & \text{if } \mathbf{P}[T_j = T_{j-1}] = 1 \\ \frac{\mathbf{E}[H(M|Y^{T_{j-1}}) - H(M|Y^{T_j})]}{\mathbf{E}[T_j - T_{j-1}]} & \text{if } \mathbf{P}[T_j = T_{j-1}] < 1 \end{cases}. \quad (47)$$

Proof of Lemma 5 is presented in Section E in the Appendix.

Before presenting the application of Lemma 5 in UEP problems, let us elaborate on its hypothesis and ramifications. We assumed that (T_i, \mathcal{A}_i) are all NALDs. Thus, for each (T_i, \mathcal{A}_i) , the set of all $y^{T_i} \in \mathcal{Y}^{T_i}$ such that the transmitted message is guaranteed to be outside $\mathcal{A}_i(y^{T_i})$ has zero probability and there is an associated probability measure $\mathbf{P}_{\{\mathcal{A}_i\}}[\cdot]$ given in (41). Furthermore, $\mathbf{P}_{\{\mathcal{A}_i\}}[\hat{M} \notin \mathcal{A}_i(Y^{T_i})]$ is the probability of the event that decoded message \hat{M} is not in \mathcal{A}_i under the probability measure $\mathbf{P}_{\{\mathcal{A}_i\}}[\cdot]$.

Condition given in (45) ensures that the decoding times of the k NALD's we are considering, T_1, T_2, \dots, T_k , are reached in their indexing order and before the decoding time of the variable-length code T . Any T_1, T_2, \dots, T_k satisfying (45) divides the time interval between 0 and T into $k+1$ disjoint intervals. The duration of these intervals as well as the decrease of the conditional entropy during them is random. For the j th interval, the expected values of the duration and the decrease in the conditional entropy are given by $\mathbf{E}[T_j - T_{j-1}]$ and $\mathbf{E}[H(M|Y^{T_{j-1}}) - H(M|Y^{T_j})]$, respectively. Hence, r_j 's defined in (47) are rate of decrease of the conditional entropy of the messages per unit time in different intervals.

Lemma 5 bounds the probability of \hat{M} being outside \mathcal{A}_i under $\mathbf{P}_{\{\mathcal{A}_i\}}[\cdot]$ from below in terms of r_j 's and $\mathbf{E}[T_j - T_{j-1}]$'s for $j > i$. The bound on $\mathbf{P}_{\{\mathcal{A}_i\}}[\hat{M} \notin \mathcal{A}_i(Y^{T_i})]$ also depends on $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$ and P_e . But the particular choice of \mathcal{A}_j 's for $j \neq i$ has no effect on the bound. This feature of the bound is its main merit over bounds resulting from the previously suggested techniques.

C. Single-Message Message-Wise UEP Converse

In this section, we bound the conditional error probabilities of the messages, i.e., $P_{e|m}$'s, from below uniformly over the message set \mathcal{M} in a variable-length block code with average error probability P_e , using Lemma 5. Resulting outer bound reveals that the inner bound we obtained in Section IV-C for the single-message *message-wise UEP* problem is tight.

Consider a variable-length block code with finite expected decoding time, i.e., $\mathbf{E}[T] < \infty$. In order to bound $P_{e|m}$, defined

$$\mathbf{P}_{\{\mathcal{A}_i\}}[\hat{M} \notin \mathcal{A}_i(Y^{T_i})] \geq \exp\left(\frac{-h(P_e + \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]) - \sum_{j=i+1}^{k+1} \mathbf{E}[T_j - T_{j-1}] J(r_j)}{1 - P_e - \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]}\right) \quad (46)$$

in (11), from the following, we apply Lemma 5 for $k = 2$ with $(T_1, \mathcal{A}_1), (T_2, \mathcal{A}_2)$ given in the following.

- 1) Let T_1 be zero and \mathcal{A}_1 be $\{m\}$, i.e.,

$$T_1 = 0 \quad (48)$$

$$\mathcal{A}_1 = \{m\}. \quad (49)$$

- 2) Let T_2 be the first time instance before T such that one message, not necessarily the one chosen for \mathcal{A}_1 , i.e., m , has a posterior probability $1 - \delta$ or higher

$$T_2 \triangleq \min\{\tau : \max_m \mathbf{P}[M = \tilde{m} | Y^\tau] \geq (1 - \delta) \text{ or } \tau = T\}. \quad (50)$$

Let \mathcal{A}_2 be the set of all messages whose posterior probability at time T_2 is less than $(1 - \delta)$

$$\mathcal{A}_2(Y^{T_2}) \triangleq \{\tilde{m} \in \mathcal{M} : \mathbf{P}[M = \tilde{m} | Y^{T_2}] < (1 - \delta)\}. \quad (51)$$

We apply Lemma 5 for (T_1, \mathcal{A}_1) and (T_2, \mathcal{A}_2) given in (48)–(51). Then, using the fact that $J(\cdot) \leq D$, we get

$$\ln P_{e|m} \geq \frac{-h(P_e + |\mathcal{M}|^{-1}) - \mathbf{E}[T_2] J\left(\frac{\mathbf{E}[H(M) - H(M|Y^{T_2})]}{\mathbf{E}[T_2]}\right) - \mathbf{E}[T - T_2] D}{1 - P_e - |\mathcal{M}|^{-1}} \quad (52a)$$

$$\ln \mathbf{P}_{\{\mathcal{A}_2\}}[\widehat{M} \notin \mathcal{A}_2(Y^{T_2})] \geq \frac{-h(P_e + \mathbf{P}[M \in \mathcal{A}_2(Y^{T_2})]) - \mathbf{E}[T - T_2] D}{1 - P_e - \mathbf{P}[M \in \mathcal{A}_2(Y^{T_2})]}. \quad (52b)$$

If $\delta < 1/2$, one can show $\mathbf{P}_{\{\mathcal{A}_2\}}[\widehat{M} \notin \mathcal{A}_2(Y^{T_2})]$ is roughly equal to P_e/δ . Thus, inequality in (52b) becomes a lower bound on $\mathbf{E}[T - T_2]$ in terms of P_e . It can be shown that the lower bound (52a) takes its smallest value for the smallest value of $\mathbf{E}[T - T_2]$. Then, using Fano's inequality for $\mathbf{E}[H(M|Y^{T_2})]$, we obtain Lemma 6 given in the following.

A complete proof of Lemma 6 for variable-length block codes with finite expected decoding time is presented in Section F in the Appendix. For variable-length block codes with infinite expected decoding time, Lemma 6 follows from the lower bounds on P_e and $P_{e|m}$ derived in Sections H1 and H2 in the Appendix.

Lemma 6: For any variable-length block code and positive δ such that $P_e + \delta + \frac{P_e}{\delta} + |\mathcal{M}|^{-1} \leq 1/2$

$$-\frac{\ln P_{e|m}}{\mathbf{E}[T]} \leq \mathbf{E} + \left(1 - \frac{\mathbf{E} - \tilde{\epsilon}}{D}\right) J\left(\frac{R - \tilde{\epsilon}}{1 - \frac{\mathbf{E} - \tilde{\epsilon}}{D}}\right) \quad \forall m \in \mathcal{M} \quad (53)$$

where $R = \frac{|\mathcal{M}|}{\mathbf{E}[T]}$, $\mathbf{E} = \frac{-\ln P_e}{\mathbf{E}[T]}$, $\tilde{\epsilon} = \frac{\tilde{\epsilon}_1 D + \tilde{\epsilon}_2}{1 - \tilde{\epsilon}_1}$, $\tilde{\epsilon}_1 = P_e + \delta + \frac{P_e}{\delta} + |\mathcal{M}|^{-1}$ and $\tilde{\epsilon}_2 = \frac{h(\tilde{\epsilon}_1) - \ln \lambda \delta}{\mathbf{E}[T]}$.

Lemma 6 is a generalization of [2, Th. 8] and [2, Lemma 1]. While deriving bounds given in [2, Th. 8] and [2, Lemma 1], no attention is paid to the fact that the rate of decrease of the conditional entropy of the messages can be different in different time intervals. As a result, both [2, Th. 8] and [2, Lemma 1] are tight only when the error exponent is very close to zero. While deriving the bound given in Lemma 6, on the other hand, the variation in the rate the conditional entropy decreases in different intervals is taken into account. Hence, the outer bound given in Lemma 6 matches the inner bound given in Section IV-C for all achievable values of error exponent, $0 \leq \mathbf{E} \leq (1 - \frac{R}{C})D$.

Consider a reliable sequence of codes \mathbb{Q} with rate $R_{\mathbb{Q}}$ and error exponent $E_{\mathbb{Q}}$. Then, if we apply Lemma 6 with $\delta = \frac{1}{\ln(1/P_e)}$, we get

$$E_{\text{md}, \mathbb{Q}} \leq E_{\mathbb{Q}} + \left(1 - \frac{E_{\mathbb{Q}}}{D}\right) J\left(\frac{R_{\mathbb{Q}}}{1 - E_{\mathbb{Q}}/D}\right). \quad (54)$$

Note that the upper bound on $E_{\text{md}, \mathbb{Q}}$'s given in (54) is achievable by at least one \mathbb{Q} described in Section IV-C.

D. Bit-Wise UEP Converse

In this section, we apply Lemma 5 to a variable-length block code with a message set \mathcal{M} of the form $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_\ell$, in order to obtain lower bounds on $P_e(i)$'s for $i = 1, 2, \dots, \ell$ in terms of the sizes of the submessage sets $|\mathcal{M}_1|, |\mathcal{M}_2|, \dots, |\mathcal{M}_\ell|$ and the expected decoding time $\mathbf{E}[T]$. When applied to reliable code sequences, these bounds on $P_e(i)$'s in terms of $|\mathcal{M}_i|$'s and $\mathbf{E}[T]$ give a necessary condition for the achievability of a rate vector and error exponent vector pair (\vec{R}, \vec{E}) that matches the sufficient condition for the achievability derived in Section IV-D.

In order to bound $P_e(i)$'s, we use Lemma 5 with ℓ NALD's, $(T_1, \mathcal{A}_1), \dots, (T_\ell, \mathcal{A}_\ell)$. Let us start with defining T_i 's and $\mathcal{A}_i(Y^{T_i})$'s.

- 1) For any i in $\{1, 2, \dots, \ell\}$, let T_i be the first time instance that a member of \mathcal{M}^i gains a posterior probability larger than or equal to $(1 - \delta)$ if it happens before T , T otherwise

$$T_i \triangleq \min\{\tau : \max_{m^i} \mathbf{P}[M^i = m^i | Y^\tau] \geq 1 - \delta \text{ or } \tau = T\}. \quad (55)$$

- 2) For any i in $\{1, 2, \dots, \ell\}$, let $\mathcal{A}_i(Y^{T_i})$ be the set of all messages of the form $m = (m^i, m_{i+1}, \dots, m_\ell)$ for which posterior probability of m^i is less than $(1 - \delta)$ at T_i

$$\mathcal{A}_i(Y^{T_i}) \triangleq \{(m^i, m_{i+1}, \dots, m_\ell) \in \mathcal{M} : \mathbf{P}[M^i = m^i | Y^{T_i}] < 1 - \delta\}. \quad (56)$$

If we apply Lemma 5 for $(T_1, \mathcal{A}_1), \dots, (T_\ell, \mathcal{A}_\ell)$ defined in (55) and (56), we obtain lower bounds on $\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})]$'s in terms of $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$'s and r_j 's and $\mathbf{E}[T_j - T_{j+1}]$'s for $j > i$. In order to turn these bounds into bounds on $P_e(i)$'s, we bound $\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})]$'s and $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$'s from above.

- 1) The posterior probability of a message at time $\tau + 1$ cannot be smaller than λ times its value at time τ because $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) = \lambda$. Thus, if $\delta < 1/2$, one can bound $\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})]$'s from above

$$\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})] < \frac{1}{\lambda \delta} P_e(i) \quad \forall i \in \{1, 2, \dots, \ell\}. \quad (57)$$

- 2) Note that if at T_i there is a m^i with posterior probability $(1 - \delta)$, then $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | Y^{T_i}] \leq \delta$. If at T_i there is no m^i with posterior probability $(1 - \delta)$, then $\mathbf{P}[\widehat{M}^i \neq M^i | Y^{T_i}] \geq \delta$. Using these facts, one can bound $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$ from above

$$\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})] \leq \frac{P_e}{\delta} + \delta \quad \forall i \in \{1, 2, \dots, \ell\}. \quad (58)$$

More detailed derivations of the inequalities given in (57) and (58) can be found in Section G in the Appendix.

Using (57) and (58) together with Lemma 5, we can conclude that

$$\ln P_e(i) \geq \ln(\lambda\delta) + \frac{-h(P_e + \delta + P_e/\delta) - \sum_{j=i+1}^{\ell+1} \mathbf{E}[\mathbb{T}_j - \mathbb{T}_{j-1}]J(r_j)}{1 - P_e - \delta - P_e/\delta} \quad (59)$$

$$\forall i \in \{1, 2, \dots, \ell\}$$

provided that $P_e + \delta + P_e/\delta \leq 1/2$, where r_j 's are defined in (47).

Note that the lower bound on $P_e(i)$'s given in (59) takes different values depending on the rate of decrease of the conditional entropy of the messages in different intervals, i.e., r_j 's, and the expected duration of different intervals, i.e., $\mathbf{E}[\mathbb{T}_j - \mathbb{T}_{j-1}]$'s. Making a worst case assumption on the rate of decrease of entropy and the durations of the intervals, one can obtain Lemma 7 given in the following.

A complete proof of Lemma 7 for variable-length block codes with finite expected decoding time is presented in Section G in the Appendix. For variable-length block codes with infinite expected decoding time, Lemma 7 follows from the lower bounds on P_e and $P_e(i)$'s derived in Sections H1 and H3 in the Appendix.

Lemma 7: For any variable-length block code with feedback with a message set \mathcal{M} of the form²⁹ $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_\ell$ and for any positive δ such that $P_e + \delta + \frac{P_e}{\delta} \leq \frac{1}{5}$, we have

$$(1 - \tilde{\epsilon}_3)E_i - \tilde{\epsilon}_5 \leq (1 - \sum_{j=1}^{\ell} \eta_j)D + \sum_{j=i+1}^{\ell} \eta_j J\left(\frac{(1 - \tilde{\epsilon}_3)R_i}{\eta_j}\right) \quad (60a)$$

$$i = 1, 2, \dots, \ell$$

$$(1 - \tilde{\epsilon}_3)R_i - \tilde{\epsilon}_4 \mathbb{1}_{\{i=1\}} \leq C\eta_i \quad i = 1, 2, \dots, \ell \quad (60b)$$

for some time-sharing vector $\vec{\eta}$ such that

$$\eta_i \geq 0 \quad i = 1, 2, \dots, \ell \quad (61a)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1 \quad (61b)$$

where $R_i = \frac{|\ln \mathcal{M}_i|}{\mathbf{E}[\mathbb{T}]}$, $E_i = \frac{-\ln P_e(i)}{\mathbf{E}[\mathbb{T}]}$, $\tilde{\epsilon}_3 = P_e + \delta + \frac{P_e}{\delta}$, $\tilde{\epsilon}_4 = \frac{h(\tilde{\epsilon}_3)}{\mathbf{E}[\mathbb{T}]}$, and $\tilde{\epsilon}_5 = \frac{h(\tilde{\epsilon}_3) - \ln \lambda\delta}{\mathbf{E}[\mathbb{T}]}$.

For any reliable sequence \mathbb{Q} whose message sets $\mathcal{M}^{(\kappa)}$ are of the form $\mathcal{M}^{(\kappa)} = \mathcal{M}_1^{(\kappa)} \times \mathcal{M}_2^{(\kappa)} \times \dots \times \mathcal{M}_\ell^{(\kappa)}$, if we set δ to $\delta = \frac{1}{-\ln P_e}$, Lemma 7 implies that there exists a $\vec{\eta}$ such that³⁰

$$E_{\mathbb{Q},i} \leq (1 - \sum_{j=1}^k \eta_j)D + \sum_{j=i+1}^{\ell} \eta_j J\left(\frac{R_{\mathbb{Q},i}}{\eta_j}\right) \quad (62a)$$

$$\forall i \in \{1, 2, \dots, \ell\}$$

$$R_{\mathbb{Q},i} \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (62b)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (62c)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1. \quad (62d)$$

²⁹We tacitly assume, without loss of generality, that $|\mathcal{M}_i| \geq 2$.

³⁰This fact is far from trivial; yet, it is intuitive to all who has worked with sequences of vectors in a bounded subset of \mathbf{R}^ℓ where \mathbf{R}^ℓ is the ℓ -dimensional real vector space with the norm $\|\vec{X}\| = \sup_j |x_j|$. For details, see Section J in the Appendix.

Recall that a rate-exponent vector (\vec{R}, \vec{E}) is achievable only if there exists a reliable code sequence \mathbb{Q} such that $(\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}}) = (\vec{R}, \vec{E})$. Thus, a rate-exponent vector (\vec{R}, \vec{E}) is achievable only if there exists a time-sharing vector $\vec{\eta}$ satisfying (34). In other words, the sufficient condition for the achievability of (\vec{R}, \vec{E}) we have derived in Section IV-D is also a necessary condition.

VI. CONCLUSION

We have considered the single-message *message-wise* and the fixed ℓ *bit-wise* UEP problems and characterized the achievable rate error exponent regions completely for both of the problems.

In the *bit-wise* UEP problem, we have observed that encoding schemes decoupling the communication and bulk of the error correction both at the transmitter and at the receiver can achieve optimal performance. This result is extending the similar observations made for conventional variable-length block coding schemes without UEP. However, for doing that one needs to go beyond the idea of communication phase and control phase introduced in [12], and harness the implicit confirmation explicit rejection schemes, introduced by Kudryashov in [7].

For the converses results, we have introduced a new technique for establishing outer bounds to the performance of the variable-length block codes, which can be used in both *message-wise* and *bit-wise* UEP problems.³¹

We were only interested in *bit-wise* UEP problem in this paper. We have analyzed single-message *message-wise* UEP problem, because it is closely related to *bit-wise* UEP problem and its analysis allowed us to introduce the ideas we use for *bit-wise* UEP, gradually. However, it seems using the technique employed in [2, Th. 9] on the achievability side and Lemma 5 on the converse side, one might be able to determine the achievable region of rate-exponent vectors for variable-length block codes in *message-wise* UEP problem. Such a work would allow us to determine the gains of feedback and variable-length decoding, because Csiszár [5] had already solved the problem for fixed-length block codes.

Arguably, the most important shortcoming of our *bit-wise* UEP result is that it only addresses the case when the number of groups of bits ℓ is a fixed integer. However, this has more to do with the formal definition of the problem we have chosen in Section III than our analysis and nonasymptotic results given in Sections IV and V, i.e., Lemmas 4 and 7.

Using the rate-exponent vectors for representing the performance of a reliable sequence with *bit-wise* UEP is apt only when the number of groups of bits are fixed or bounded. When the number of groups of bits ℓ in a reliable sequence diverge with increasing κ , i.e., when $\lim_{\kappa \rightarrow \infty} \ell_\kappa = \infty$, the rate-exponent vector formulation becomes fundamentally inapt. Consider, for example, a reliable sequence in which $|\mathcal{M}_i^{(\kappa)}| = \lceil e^{\mathbf{E}[\mathbb{T}^{(\kappa)}] \frac{R}{\ell_\kappa}} \rceil$. The rate of this reliable sequence is R ; yet, the rate of all of the submessages are zero. Thus, when ℓ_κ diverges, the rate vector does not have the same operational relevance or meaning it has when ℓ_κ is fixed or bounded. In order to characterize the change of error performance among submessages in the case when ℓ_κ

³¹We have not employed the bound in any hybrid problem but it seems that the result is abstract enough to be employed even in those problems with judicious choice of NALD's.

diverges, one needs to come up with an alternative formulation of the problem, in terms of cumulative rate of submessages.

Our nonasymptotic results are useful to some extent even when ℓ_κ diverges. Although infinite dimensional rate-exponent vectors fall short of representing all achievable performances, one can still use Lemma 4 of Section IV and Lemma 7 of Section V to characterize the set of achievable rate vector error exponent vector pairs.

1) As a result of Lemma 7, the necessary condition given in (19) is still a necessary condition for the achievability of rate-exponent vector.

2) Using Lemma 4, we see that the sufficient condition given in (19) is still a sufficient condition as long as the number of submessages in the reliable sequence satisfy $\limsup_{n \rightarrow \infty} \frac{\ell_n}{n/\ln n} = 0$.

Thus, for the case when $\ell \sim o\left(\frac{\mathbb{E}[\ell]}{\ln \mathbb{E}[\ell]}\right)$, i.e., $\limsup_{\kappa \rightarrow \infty} \frac{\ell_\kappa}{\mathbb{E}^{(\kappa)}[\ell]/\ln \mathbb{E}^{(\kappa)}[\ell]} = 0$, the condition given in (19) is still a necessary and sufficient condition for the achievability of a rate-exponent vector.

APPENDIX

A. Proof of Lemma 1

Proof: Note that $J(\mathbf{R})$ defined in (17) is also equal to

$$\begin{aligned} J(\mathbf{R}) &= \max_{\substack{0 \leq \alpha \leq 1 \\ x_1, x_2 \in \mathcal{X} \\ \mu_1, \mu_2 \in \mathcal{P}(\mathcal{X}) \\ \alpha, x_1, x_2, \mu_1, \mu_2, \mathbf{R}_1, \mathbf{R}_2: \\ \mathbf{R}_1, \mathbf{R}_2 \in [0, C] \\ I(\mu_1, W) \geq \mathbf{R}_1 \\ I(\mu_2, W) \geq \mathbf{R}_2 \\ \alpha \mathbf{R}_1 + (1 - \alpha) \mathbf{R}_2 = \mathbf{R}}} \alpha D(\bar{\mu}_1 \| W_{x_1}) + (1 - \alpha) D(\bar{\mu}_2 \| W_{x_2}) \\ &= \max_{\substack{0 \leq \alpha \leq 1 \\ \alpha, \mathbf{R}_1, \mathbf{R}_2: \mathbf{R}_1, \mathbf{R}_2 \in [0, C] \\ \alpha \mathbf{R}_1 + (1 - \alpha) \mathbf{R}_2 = \mathbf{R}}} \alpha j(\mathbf{R}_1) + (1 - \alpha) j(\mathbf{R}_2) \end{aligned} \quad (63)$$

where $j(\mathbf{R})$ is given by

$$j(\mathbf{R}) \triangleq \max_{\substack{x \in \mathcal{X} \\ \alpha, x, \mu: \mu \in \mathcal{P}(\mathcal{X}) \\ I(\mu, W) \geq \mathbf{R}}} D(\bar{\mu} \| W_x) \quad \forall \mathbf{R} \in C. \quad (64)$$

Note that $j(\mathbf{R})$ is a bounded real-valued function of a real variable. Therefore, Carathéodory's Theorem implies that consid-

ering two point convex combinations suffices in order make $j(\mathbf{R})$ a concave function. In other words, for any k , we have

$$\begin{aligned} & \max_{\substack{0 \leq \alpha \leq 1 \\ \alpha, \mathbf{R}_1, \mathbf{R}_2: \mathbf{R}_1, \mathbf{R}_2 \in [0, C] \\ \alpha \mathbf{R}_1 + (1 - \alpha) \mathbf{R}_2 = \mathbf{R}}} \alpha j(\mathbf{R}_1) + (1 - \alpha) j(\mathbf{R}_2) \\ &= \max_{\substack{0 \leq \alpha_i \leq 1 \forall i \\ 0 \leq \mathbf{R}_i \leq C \forall i \\ \alpha_1, \dots, \alpha_k, \mathbf{R}_1, \dots, \mathbf{R}_k: \\ \sum_i \alpha_i = 1 \\ \sum_i \alpha_i \mathbf{R}_i = \mathbf{R}}} \sum_{i=1}^k \alpha_i j(\mathbf{R}_i). \end{aligned} \quad (65)$$

Then, the concavity of $J(\mathbf{R})$ follows from (63)–(65).

Evidently, if the constraint set in a maximization is curtailed, then resulting maximum value cannot increase. Hence, $J(\mathbf{R})$ function defined in (17) is a decreasing function of \mathbf{R} .

As a result of the definition of D given in (9) and the convexity of Kullback–Leibler divergence, we have $D \geq J(0)$. On the other hand, $D(\bar{\mu} \| W_x) = D$ and $I(\mu, W) \geq 0$ for $x = \mathbf{r}$ and $\mu(\cdot) = \mathbb{1}_{\{\cdot = \mathbf{a}\}}$ where \mathbf{a} and \mathbf{r} are described in (10). Therefore, we have $j(0) \geq D$. Using the fact that $J(\mathbf{R}) \geq j(\mathbf{R})$, we conclude that $J(0) = j(0) = D$. ■

B. Proof of Lemma 2

Proof: We prove the lemma for a slightly more general setting and establish a result that will be easier to make use of in the proofs of other achievability results. Let $\mathcal{G}_\gamma[1]$, $\mathcal{G}_\gamma[m]$, and $\mathcal{B}_\gamma[x^n]$ be (66a)–(66c) shown at the bottom of the page. Note that $\mathcal{G}[1]$, $\mathcal{G}[m]$, and $\mathcal{B}[x^n]$ given, (21)–(23) are simply the $\mathcal{G}_\gamma[1]$, $\mathcal{G}_\gamma[m]$, and $\mathcal{B}_\gamma[x^n]$ for $\gamma = |\mathcal{X}| |\mathcal{Y}| \sqrt{n \ln(1+n)}$.

For all $y^n \notin \mathcal{G}_\gamma[1]$, we have

$$\begin{aligned} & n_\alpha D(Q_{\{y_1^{n_\alpha}\}} \| W_{x_1}) + (n - n_\alpha) D(Q_{\{y_{n_\alpha+1}^n\}} \| W_{x_2}) \\ &= n_\alpha D(Q_{\{y_1^{n_\alpha}\}} \| \bar{\mu}_1) + (n - n_\alpha) D(Q_{\{y_{n_\alpha+1}^n\}} \| \bar{\mu}_2) \\ & \quad + n_\alpha \sum_y Q_{\{y_1^{n_\alpha}\}}(y) \ln \frac{\bar{\mu}_1(y)}{W_{x_1}(y)} \\ & \quad + (n - n_\alpha) \sum_y Q_{\{y_{n_\alpha+1}^n\}}(y) \ln \frac{\bar{\mu}_2(y)}{W_{x_2}(y)} \\ & \stackrel{(a)}{\geq} n_\alpha \sum_y Q_{\{y_1^{n_\alpha}\}}(y) \ln \frac{\bar{\mu}_1(y)}{W_{x_1}(y)} \\ & \quad + (n - n_\alpha) \sum_y Q_{\{y_{n_\alpha+1}^n\}}(y) \ln \frac{\bar{\mu}_2(y)}{W_{x_2}(y)} \\ & \stackrel{(b)}{\geq} n_\alpha D(\bar{\mu}_1 \| W_{x_1}) + (n - n_\alpha) D(\bar{\mu}_2 \| W_{x_2}) + 2\gamma \ln \lambda. \end{aligned}$$

$$\mathcal{G}_\gamma[1] = \left\{ y^n : n_\alpha \Delta(Q_{\{y_1^{n_\alpha}\}}, \bar{\mu}_1) + (n - n_\alpha) \Delta(Q_{\{y_{n_\alpha+1}^n\}}, \bar{\mu}_2) \geq \gamma \right\} \quad (66a)$$

$$\mathcal{G}_\gamma[m] = \mathcal{B}_\gamma[x^n(m)] \cap \left(\bigcap_{m \neq \tilde{m}} \overline{\mathcal{B}_\gamma[x^n(\tilde{m})]} \right) \quad \forall m \in \{2, 3, \dots, |\mathcal{M}|\} \quad (66b)$$

$$\mathcal{B}_\gamma[x^n] = \left\{ y^n : n_\alpha \Delta(Q_{\{x_1^{n_\alpha}, y_1^{n_\alpha}\}}, \mu_1 W) + (n - n_\alpha) \Delta(Q_{\{x_{n_\alpha+1}^n, y_{n_\alpha+1}^n\}}, \mu_2 W) < \gamma \right\} \quad (66c)$$

Inequality (a) follows from the nonnegativity of the Kullback–Leibler divergence. In order to see why (b) holds, first recall that $\min_{x,y} W_x(y) = \lambda$. Hence, $|\ln \frac{\bar{\mu}_1(y)}{W_{x_1}(y)}| \leq \ln \frac{1}{\lambda}$ and $|\ln \frac{\bar{\mu}_2(y)}{W_{x_2}(y)}| \leq \ln \frac{1}{\lambda}$. Then, the inequality (b) follows from the definitions of total variation Δ and $\mathcal{G}_\gamma[1]$, given in (1) and (66a) and the fact that $y^n \notin \mathcal{G}_\gamma[1]$.

Note that the conditional error probability of the first message is given by

$$\begin{aligned} P_{e|1} &= \mathbf{P}[\hat{M} \neq 1 | M = 1] \\ &= \sum_{y^n \notin \mathcal{G}_\gamma[1]} \mathbf{P}[Y^n = y^n | M = 1]. \end{aligned}$$

Recall that the codeword of the message $M = 1$ is the concatenation of n_α x_1 's and $(n - n_\alpha)$ x_2 's where $n_\alpha = \lfloor n\alpha \rfloor$. Hence, the probability of all y^n 's whose empirical distribution in first n_α times instances is $Q_{\{y_1^{n_\alpha}\}}$ and whose empirical distribution in $[(n_\alpha + 1), n]$ is $Q_{\{y_{n_\alpha+1}^n\}}$ is upper bounded by $e^{-n_\alpha D(Q_{\{y_1^{n_\alpha}\}} \| W_{x_1}) - (n - n_\alpha) D(Q_{\{y_{n_\alpha+1}^n\}} \| W_{x_2})}$. Furthermore, there are less than $(n_\alpha + 1)^{|\mathcal{Y}|}$ distinct empirical distributions in the first phase and there are less than $(n - n_\alpha + 1)^{|\mathcal{Y}|}$ distinct empirical distributions in the second phase. Thus, we have the equation shown at the bottom of the page where $\varepsilon_2(\gamma, n) = \frac{-2\gamma \ln \lambda + D + 2|\mathcal{Y}| \ln(n+1)}{n}$.

The codewords and the decoding regions of the remaining messages are specified using a random coding argument together with an empirical typicality decoder. Consider an ensemble of codes in which first n_α entries of all the codewords are independent and identically distributed (i.i.d.) with input distribution μ_1 and the rest of the entries are i.i.d. with the input distribution μ_2 .

For any message m other than the first one, i.e., $m \neq 1$, the decoding region is $\mathcal{G}_\gamma[m]$ given in (66b). In other words, for any message m , other than the first one, the decoding region is the set of output sequences for which $(x^n(m), y^n)$ is typical with $(\alpha, \mu_1 W, \mu_2 W)$, i.e., $y^n \in \mathcal{B}_\gamma[x^n(m)]$, and $(x^n(\tilde{m}), y^n)$ is not typical with $(\alpha, \mu_1 W, \mu_2 W)$, i.e., $y^n \in \mathcal{B}_\gamma[x^n(\tilde{m})]$, for any $\tilde{m} \neq m$.

Since the decoding regions of different messages are disjoint, the previously described code does not decode to more than one message. Disjointness of decoding regions of messages $2, 3, \dots, |\mathcal{M}|$ follows from the definitions of $\mathcal{G}_\gamma[2], \mathcal{G}_\gamma[3], \dots, \mathcal{G}_\gamma[|\mathcal{M}|]$, given in (66b). In order to see why $\mathcal{G}_\gamma[1] \cap (\cup_{m \neq 1} \mathcal{G}_\gamma[m]) = \emptyset$ holds, note that for any pair probability of distributions, the total variation between them is lower bounded by the total variation between their marginals. In particular

$$\begin{aligned} \Delta(Q_{\{x_1^{n_\alpha}(m), y_1^{n_\alpha}\}}, \mu_1 W) &\geq \Delta(Q_{\{y_1^{n_\alpha}\}}, \bar{\mu}_1) \\ \Delta(Q_{\{x_{n_\alpha+1}^n(m), y_{n_\alpha+1}^n\}}, \mu_2 W) &\geq \Delta(Q_{\{y_{n_\alpha+1}^n\}}, \bar{\mu}_2). \end{aligned}$$

Then, as results of definitions of $\mathcal{G}_\gamma[1]$, $\mathcal{B}_\gamma[x^n]$, and $\mathcal{G}_\gamma[m]$ for $m \neq 1$ given in (66a), (66c), and (66b), we have

$$\mathcal{G}_\gamma[1] \cap \mathcal{G}_\gamma[m] = \emptyset \quad m = 2, 3, \dots, |\mathcal{M}|.$$

Then, for $m \in \{2, 3, \dots, |\mathcal{M}|\}$, the average of the conditional error probability of m th message over the ensemble is upper bounded as

$$\begin{aligned} \mathbf{E}[P_{e|m}] &\leq \mathbf{P}[Y^n \notin \mathcal{B}_\gamma[X^n(m)] | M = m] \\ &\quad + \sum_{\tilde{m} \neq m} \mathbf{P}[Y^n \in \mathcal{B}_\gamma[X^n(\tilde{m})] | M = m]. \end{aligned} \quad (67)$$

Let us start with bounding $\mathbf{P}[Y^n \notin \mathcal{B}_\gamma[X^n(m)] | M = m]$. Let $S_1(x, y)$ and $S_2(x, y)$ be

$$S_1(x, y) \triangleq n_\alpha |Q_{\{x_1^{n_\alpha}(m), y_1^{n_\alpha}\}}(x, y) - \mu_1(x) W_x(y)| \quad (68a)$$

$$S_2(x, y) \triangleq (n - n_\alpha) |Q_{\{x_{n_\alpha+1}^n(m), y_{n_\alpha+1}^n\}}(x, y) - \mu_2(x) W_x(y)|. \quad (68b)$$

As a result of the definition of total variation distance given in (1) and aforementioned definitions, we have

$$\begin{aligned} n_\alpha \Delta(Q_{\{x_1^{n_\alpha}(m), y_1^{n_\alpha}\}}, \mu_1 W) + (n - n_\alpha) \Delta(Q_{\{x_{n_\alpha+1}^n(m), y_{n_\alpha+1}^n\}}, \mu_2 W) \\ = \frac{1}{2} \sum_{x,y} [S_1(x, y) + S_2(x, y)]. \end{aligned}$$

Thus, the definition of $\mathcal{B}_\gamma[x^n(m)]$ given in (66c) implies that

$$\begin{aligned} \mathbf{P}[Y \notin \mathcal{B}_\gamma[X^n(m)] | M = m] \\ = \mathbf{P}\left[\sum_{x,y} [S_1(x, y) + S_2(x, y)] \geq 2\gamma \mid M = m\right]. \end{aligned} \quad (69)$$

If for all $x \in \mathcal{X}$, $y \in \mathcal{Y}$, and $j \in \{1, 2\}$, $S_j(x, y) \leq \gamma |\mathcal{X}|^{-1} |\mathcal{Y}|^{-1}$, then $\sum_{x,y} [S_1(x, y) + S_2(x, y)] \leq 2\gamma$. Thus, if $Y \notin \mathcal{B}_\gamma[X^n(m)]$, then for at least one (x, y, j) triple $S_j(x, y) \geq \gamma |\mathcal{X}|^{-1} |\mathcal{Y}|^{-1}$. Using the union bound, we get

$$\begin{aligned} \mathbf{P}\left[\sum_{x,y} [S_1(x, y) + S_2(x, y)] \geq 2\gamma \mid M = m\right] \\ \leq \sum_{x,y,j} \mathbf{P}\left[S_j(x, y) \geq \frac{\gamma}{|\mathcal{X}||\mathcal{Y}|} \mid M = m\right]. \end{aligned} \quad (70)$$

For bounding $\mathbf{P}\left[S_j(x, y) \geq \frac{\gamma}{|\mathcal{X}||\mathcal{Y}|} \mid M = m\right]$, we can simply use Chebyshev's inequality; however, in order to get better error terms, we use a standard concentration result about the sums of bounded random variables, [4, Th. 5.3].

Lemma 8: Let Z_1, Z_2, \dots, Z_k be independent random variables satisfying $|Z_i - \mathbf{E}[Z_i]| \leq c_i$ for all $1 \leq i \leq k$. Then

$$\mathbf{P}\left[\left|\sum_{i=1}^k (Z_i - \mathbf{E}[Z_i])\right| > \gamma\right] \leq 2e^{-\frac{\gamma^2}{2 \sum_{i=1}^k c_i^2}}.$$

$$\begin{aligned} P_{e|1} &\leq (n_\alpha + 1)^{|\mathcal{Y}|} (n - n_\alpha + 1)^{|\mathcal{Y}|} e^{-n_\alpha D(\bar{\mu}_1 \| W_{x_1}) + (n - n_\alpha) D(\bar{\mu}_2 \| W_{x_2}) - 2\gamma \ln \lambda} \\ &\leq e^{-n(\alpha D(\bar{\mu}_1 \| W_{x_1}) + (1 - \alpha) D(\bar{\mu}_2 \| W_{x_2}) - \varepsilon_2(\gamma, n))} \end{aligned}$$

For all $\mu_1 \in P(\mathcal{X})$, $x \in \mathcal{X}$, and $y \in \mathcal{Y}$, we have $c_i = 1$ for all $i = 1, 2, \dots, n_\alpha$; thus

$$\begin{aligned} \mathbf{P}\left[S_1(x, y) \geq \frac{\gamma}{|\mathcal{X}||\mathcal{Y}|} \mid M = m\right] &\leq 2e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n_\alpha}} \\ &\leq 2e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}. \end{aligned} \quad (71)$$

Similarly

$$\mathbf{P}\left[S_2(x, y) \geq \frac{\gamma}{|\mathcal{X}||\mathcal{Y}|} \mid M = m\right] \leq 2e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}. \quad (72)$$

Using (69)–(72), we get

$$\mathbf{P}[Y \notin \mathcal{B}_\gamma[X^n(m)] \mid M = m] \leq 4|\mathcal{X}||\mathcal{Y}|e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}. \quad (73)$$

Now, we focus on $\mathbf{P}[Y^n \in \mathcal{B}_\gamma[X^n(\tilde{m})] \mid M = m]$ terms. Note that all y^n in $\mathcal{B}_\gamma[X^n(\tilde{m})]$ satisfy

$$\begin{aligned} n_\alpha \Delta\left(Q_{\{x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha}\}}\right) \\ + (n - n_\alpha) \Delta\left(Q_{\{x_{n_\alpha+1}^{n_\alpha}(\tilde{m}), y_{n_\alpha+1}^{n_\alpha}\}}\right) \leq \gamma. \end{aligned} \quad (74)$$

On the other hand, when $M = m$, $X^n(\tilde{m})$ and Y^n are independent and their distribution is given by (75) shown at the bottom of the page. Furthermore, the number of $(x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha})$ sequences with an empirical distribution $Q_{\{x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha}\}}$ is upper

bounded as $e^{n_\alpha H(Q_{\{x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha}\}})}$. In addition, there are at most $(n_\alpha + 1)^{|\mathcal{X}||\mathcal{Y}|}$ different empirical distributions. Using these two bounds and their counterparts for $(x_{n_\alpha+1}^{n_\alpha}(\tilde{m}), y_{n_\alpha+1}^{n_\alpha})$ together with (74) and (75) we get (76) shown at the bottom of the page. Hence, if (see the third equation at the bottom of the page) then

$$\sum_{\tilde{m} \neq m} \mathbf{P}[Y^n \in \mathcal{B}_\gamma[X^n(\tilde{m})] \mid M = m] \leq 4|\mathcal{X}||\mathcal{Y}|e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}. \quad (77)$$

Thus, the average P_e over the ensemble can be bounded using (67), (70), and (77) as

$$\mathbf{E}[P_e] \leq 8|\mathcal{X}||\mathcal{Y}|e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}.$$

But if the ensemble average of the error probability is upper bounded like this, there is at least one code that has this low error probability. Furthermore, half of its messages have conditional error probabilities less than twice this average. Thus, for any block length n , time-sharing constant $\alpha \in [0, 1]$, input letters $x_1, x_2 \in \mathcal{X}$, input distributions $\mu_1, \mu_2 \in P(\mathcal{X})$, there exists a length n code such that

$$|\mathcal{M} \setminus \{1\}| \geq e^{n(\alpha I(\mu_1, W) + (1-\alpha)I(\mu_2, W) - \varepsilon_1(\gamma, n))} \quad (78a)$$

$$P_{e|1} \leq e^{-n(\alpha D(\bar{\mu}_1 \| W_{x_1}) + (1-\alpha)D(\bar{\mu}_2 \| W_{x_2}) - \varepsilon_2(\gamma, n))} \quad (78b)$$

$$P_{e|m} \leq \varepsilon_3(\gamma, n) \quad m = 2, 3, \dots, |\mathcal{M}| \quad (78c)$$

where

$$\begin{aligned} \varepsilon_1(\gamma, n) &= \frac{C - \ln(2|\mathcal{X}||\mathcal{Y}|) + 2|\mathcal{X}||\mathcal{Y}| \ln(n+1) - 2\gamma \ln \lambda}{n} \\ &\quad + \frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n^2} \end{aligned} \quad (79a)$$

$$\varepsilon_2(\gamma, n) = \frac{D + 2|\mathcal{Y}| \ln(n+1) - 2\gamma \ln \lambda}{n} \quad (79b)$$

$$\varepsilon_3(\gamma, n) = 16|\mathcal{X}||\mathcal{Y}|e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}}. \quad (79c)$$

Lemma 2 follows from (78) and the fact that

$$\varepsilon_i(\gamma, n) \Big|_{\gamma=|\mathcal{X}||\mathcal{Y}|\sqrt{n \ln(1+n)}} \leq \frac{9|\mathcal{X}||\mathcal{Y}|(1-\ln \lambda)\sqrt{\ln(1+n)}}{\sqrt{n}} \quad i = 1, 2, 3 \quad (80)$$

for $\varepsilon_1(\gamma, n)$, $\varepsilon_2(\gamma, n)$, and $\varepsilon_3(\gamma, n)$ given in (79). \blacksquare

$$\begin{aligned} \mathbf{P}[(X^n(\tilde{m}), Y^n) = (x^n(\tilde{m}), y^n) \mid M = m] &= \prod_{i=1}^{n_\alpha} \mu_1(x_i(\tilde{m}))\bar{\mu}_1(y_i) \prod_{j=n_\alpha+1}^n \mu_2(x_j(\tilde{m}))\bar{\mu}_2(y_j) \\ &= e^{-n_\alpha D(Q_{\{x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha}\}} \parallel \mu_1 \bar{\mu}_1)} e^{-n_\alpha H(Q_{\{x_1^{n_\alpha}(\tilde{m}), y_1^{n_\alpha}\}})} \\ &\quad e^{-(n-n_\alpha) D(Q_{\{x_{n_\alpha+1}^{n_\alpha}(\tilde{m}), y_{n_\alpha+1}^{n_\alpha}\}} \parallel \mu_2 \bar{\mu}_2)} e^{-(n-n_\alpha) H(Q_{\{x_{n_\alpha+1}^{n_\alpha}(\tilde{m}), y_{n_\alpha+1}^{n_\alpha}\}})} \end{aligned} \quad (75)$$

$$\begin{aligned} \mathbf{P}[Y^n \in \mathcal{B}_\gamma[X^n(\tilde{m})] \mid M = m] &\leq (n_\alpha + 1)^{|\mathcal{X}||\mathcal{Y}|} (n - n_\alpha + 1)^{|\mathcal{X}||\mathcal{Y}|} e^{-n_\alpha D(\mu_1 W \parallel \mu_1 \bar{\mu}_1) - (n-n_\alpha) D(\mu_2 W \parallel \mu_2 \bar{\mu}_2) - 2\gamma \ln \lambda} \\ &= (n_\alpha + 1)^{|\mathcal{X}||\mathcal{Y}|} (n - n_\alpha + 1)^{|\mathcal{X}||\mathcal{Y}|} e^{-n_\alpha I(\mu_1, W) - (n-n_\alpha) I(\mu_2, W) - 2\gamma \ln \lambda} \\ &\leq e^{-n(\alpha I(\mu_1, W) + (1-\alpha)I(\mu_2, W))} e^{C + 2|\mathcal{X}||\mathcal{Y}| \ln(n+1) - 2\gamma \ln \lambda} \end{aligned} \quad (76)$$

$$|\mathcal{M} \setminus \{1\}| = 4|\mathcal{X}||\mathcal{Y}|e^{-\frac{\gamma^2}{2|\mathcal{X}|^2|\mathcal{Y}|^2n}} e^{n(\alpha I(\mu_1, W) + (1-\alpha)I(\mu_2, W))} e^{-C - 2|\mathcal{X}||\mathcal{Y}| \ln(n+1) + 2\gamma \ln \lambda}$$

C. Proof of Lemma 3

Proof: Let n_1 be $n_1 = \lceil (1 - \frac{E}{D})n \rceil$. Recall that we have assumed $E \leq (1 - \frac{R}{C})D$; then, we have $\frac{R}{C} \leq 1 - \frac{E}{D}$. Consequently, $\frac{R}{C} \leq \frac{n_1}{n}$ and $\frac{n}{n_1}R \leq C$. On the other hand, as a result of (78) and the definition of $J(\cdot)$ given in (17), for any positive integer n_1 , positive real number γ_1 , rate $\tilde{R} \leq C$, there exists a length n_1 code such that

$$|\mathcal{M}| - 1 \geq e^{n_1[\tilde{R} - \varepsilon_1(\gamma_1, n_1)]} \quad (81a)$$

$$P_{e|1} \leq e^{-n_1[J(\tilde{R}) - \varepsilon_2(\gamma_1, n_1)]} \quad (81b)$$

$$P_{e|m} \leq \varepsilon_3(\gamma_1, n_1) \quad m = 2, 3, \dots, |\mathcal{M}| \quad (81c)$$

where $\varepsilon_1(\gamma_1, n_1)$, $\varepsilon_2(\gamma_1, n_1)$, and $\varepsilon_3(\gamma_1, n_1)$ are given in (79).

We use such a code in the first phase with $\tilde{R} = \frac{n}{n_1}R$ and call its decoded message $\hat{t}\hat{M}$, the tentative decision. Then, as a result of (81) and the fact that $^{32}n_1 J(\frac{n}{n_1}R) \geq n(1 - \frac{E}{D})J(\frac{R}{1 - E/D})$, we get

$$|\mathcal{M}| - 1 \geq e^{nR - n_1\varepsilon_1(\gamma_1, n_1)} \quad (82a)$$

$$P[\hat{t}\hat{M} \neq m | M = 1] \leq e^{-n(1 - \frac{E}{D})J(\frac{R}{1 - E/D}) + n_1\varepsilon_2(\gamma_1, n_1)} \quad (82b)$$

$$P[\hat{t}\hat{M} \neq m | M = m] \leq \varepsilon_3(\gamma_1, n_1) \quad m = 2, 3, \dots, |\mathcal{M}|. \quad (82c)$$

The transmitter knows what the tentative decision is and determines the channel inputs in the last $(n - n_1)$ time instances depending on its correctness. If $\hat{t}\hat{M} = M$, the channel inputs in the last $(n - n_1)$ time instances are all \mathbf{a} ; if $\hat{t}\hat{M} \neq M$, the channel inputs in the last $(n - n_1)$ time instances are all \mathbf{r} .

After observing Y^n , receiver checks whether the empirical distribution of the channel output in the last $(n - n_1)$ time units is typical with $W_{\mathbf{a}}$ if it is then $\hat{M} = \hat{t}\hat{M}$ otherwise $\hat{M} = \mathbf{x}$. Hence, the decoding region for erasures is given by

$$\mathcal{G}_\gamma[\mathbf{x}] = \{y^n : (n - n_1)\Delta(Q_{\{Y_{n_1+1}^n\}}, W_{\mathbf{a}}) \geq \gamma_2\}.$$

Let us start with bounding $P[\hat{M} = \mathbf{x} | \hat{t}\hat{M} = m, M = m]$, i.e., the probability of erasure for correct tentative decision. First note that

$$(n - n_1)\Delta(Q_{\{Y_{n_1+1}^n\}}, W_{\mathbf{a}}) = \frac{1}{2} \sum_y S(y)$$

where $S(y) = (n - n_1)|Q_{\{Y_{n_1+1}^n\}}(y) - W_{\mathbf{a}}(y)|$. Then, following an analysis similar to that one presented between (69) and (73), we get

$$P[\hat{M} = \mathbf{x} | \hat{t}\hat{M} = m, M = m] \leq 2|\mathcal{Y}|e^{-\frac{\gamma_2^2}{2|\mathcal{Y}|^2(n - n_1)}} = \varepsilon_3(\gamma_2, n - n_1) \quad \forall m \in \mathcal{M}. \quad (83)$$

³²Recall that $n_1 \geq (1 - E/D)n$ and $J(\cdot)$ is a nonincreasing and positive function.

In order to bound the probability of nonerasure decoding when tentative decision is incorrect, note that

$$\begin{aligned} & P[Y_{n_1+1}^n = y_{n_1+1}^n | \hat{t}\hat{M} \neq m, M = m] \\ &= \prod_{j=n_1+1}^n W_{\mathbf{r}}(y_j) \\ &= e^{-(n - n_1)D(Q_{\{Y_{n_1+1}^n\}} || \tilde{\mu}_2)} e^{-(n - n_1)H(Q_{\{Y_{n_1+1}^n\}})}. \end{aligned}$$

Then, following an analysis similar to the one between (75) and (76), we get

$$\begin{aligned} & P[\hat{M} \neq \mathbf{x} | \hat{t}\hat{M} \neq m, M = m] \\ & \leq \min\{(n - n_1 + 1)|\mathcal{Y}|e^{-(n - n_1)D - 2\gamma_2 \ln \lambda}, 1\} \\ & \leq \min\{e^{-nE + |\mathcal{Y}| \ln(n - n_1) + D - 2\gamma_2 \ln \lambda}, 1\} \\ & \leq \min\{e^{-nE + (n - n_1)\varepsilon_2(\gamma_2, n - n_1)}, 1\} \quad \forall m \in \mathcal{M}. \quad (84) \end{aligned}$$

Furthermore, the conditional error and erasure probabilities can be bounded in terms of $P[\hat{t}\hat{M} \neq m | M = m]$, $P[\hat{M} \neq \mathbf{x} | \hat{t}\hat{M} \neq m, M = m]$, and $P[\hat{M} = \mathbf{x} | \hat{t}\hat{M} = m, M = m]$ as follows:

$$P_{e|m} = P[\hat{t}\hat{M} \neq m | M = m] P[\hat{M} \neq \mathbf{x} | \hat{t}\hat{M} \neq m, M = m] \quad \forall m \in \mathcal{M} \quad (85a)$$

$$P_{\mathbf{x}|m} \leq P[\hat{t}\hat{M} \neq m | M = m] + P[\hat{M} = \mathbf{x} | \hat{t}\hat{M} = m, M = m] \quad \forall m \in \mathcal{M}. \quad (85b)$$

Using (82)–(85), we get

$$|\mathcal{M}| - 1 \geq e^{nR - n_1\varepsilon_1(\gamma_1, n_1)} \quad (86a)$$

$$P_{e|1} \leq e^{-n(1 - \frac{E}{D})J(\frac{R}{1 - E/D}) + n_1\varepsilon_2(\gamma_1, n_1)} \min\{e^{-nE + n_2\varepsilon_2(\gamma_2, n_2)}, 1\} \quad (86b)$$

$$P_{\mathbf{x}|1} \leq e^{-n(1 - \frac{E}{D})J(\frac{R}{1 - E/D}) + n_1\varepsilon_2(\gamma_1, n_1)} + \varepsilon_3(\gamma_2, n_2) \quad (86c)$$

$$P_{e|m} \leq \varepsilon_3(\gamma_1, n_1) \min\{e^{-nE + |\mathcal{Y}| \ln(n + 1) + n_2\varepsilon_2(\gamma_2, n_2)}, 1\} \quad m \neq 1 \quad (86d)$$

$$P_{\mathbf{x}|m} \leq \varepsilon_3(\gamma_1, n_1) + \varepsilon_3(\gamma_2, n_2) \quad m \neq 1 \quad (86e)$$

where $n_2 = n - n_1$.

We set $\gamma_j = |\mathcal{X}||\mathcal{Y}|\sqrt{5n_j \ln(1 + n)}$ for $j = 1, 2$ and obtain

$$n_j\varepsilon_1(\gamma_j, n_j) \leq 2|\mathcal{X}||\mathcal{Y}|(\ln(n + 1) - \sqrt{5n \ln(1 + n)}) \ln \lambda + (5/2) \ln(1 + n) \quad (87a)$$

$$n_j\varepsilon_2(\gamma_j, n_j) \leq 2|\mathcal{X}||\mathcal{Y}|(\ln(n + 1) - \sqrt{5n \ln(1 + n)}) \ln \lambda + D \quad (87b)$$

$$\varepsilon_3(\gamma_j, n_j) \leq 16|\mathcal{X}||\mathcal{Y}|/(1 + n)^{5/2}. \quad (87c)$$

Lemma 3 follows from the identities $|\mathcal{X}| \geq 2$, $|\mathcal{Y}| \geq 2$, $D \leq \ln(\frac{1}{\lambda})$, $n \geq 1$ and (86) and (87). \blacksquare

D. Proof of Lemma 4

Proof: Note that given the encoding scheme summarized in (28) and the decoding rule given in (29), if $\widehat{\mathbf{M}} = \mathbf{x}$, then there is a $i \leq \ell + 1$ such that $\widehat{M}_j = M_j$ for all $j < i$ and $\widehat{M}_i \neq M_i$. Thus, the conditional erasure probability $P_{\mathbf{x}|\mathbf{m}}$ is upper bounded as shown in (88) at the bottom of the page. Similarly, if $\widehat{\mathbf{M}} \neq \mathbf{x}$ and $\widehat{M}^i \neq M^i$, then for all $j > i$, $\widehat{M}_j = 1$ and $\widehat{M}_j \neq 1$; furthermore, there is a $k \leq i$ such that $\widehat{M}_j = M_j$ for all $j < k$ and $\widehat{M}_k \neq M_k$. Hence, one can bound $P_{\mathbf{e}|\mathbf{m}}(i)$ as shown in (89) at the bottom of the page.

In the first ℓ phases, we use $n_i = \lfloor \eta_i n_i \rfloor$ long codes with rate $\frac{R_i}{\eta_i}$ with the performance given in (81). Thus, for $1 \leq i \leq \ell$, we have

$$|\mathcal{M}_i| \geq 1 + e^{nR_i - C - n_i \varepsilon_1(\gamma_i, n_i)} \quad (90a)$$

$$\mathbf{P} \left[\widehat{M}_i \neq 1 \mid \mathcal{M}_i = 1 \right] \leq e^{-n\eta_i J\left(\frac{R_i}{\eta_i}\right) + D - n_i \varepsilon_2(\gamma_i, n_i)} \quad (90b)$$

$$\mathbf{P} \left[\widehat{M}_i \neq \mathcal{M}_i \mid \mathcal{M}_i = 1 + m_i \right] \leq \varepsilon_3(\gamma_i, n_i) \quad (90c)$$

$$m_i = 1, 2, 3, \dots, (|\mathcal{M}_i| - 1)$$

where $\varepsilon_1(\gamma_i, n_i)$, $\varepsilon_2(\gamma_i, n_i)$, $\varepsilon_3(\gamma_i, n_i)$ are given in (79).

In order to derive bounds corresponding to the ones given in (90) for the last phase, let us give the decoding regions for 1 and 2 for the length $n_{\ell+1}$ code employed between $(n + 1 - n_{\ell+1})$ and n

$$\mathcal{G}_\gamma[1] = \{y_{n+1-n_{\ell+1}}^n : n_{\ell+1} \Delta(Q_{\{y_{n+1-n_{\ell+1}}^n\}}, W_{\mathbf{a}}) \geq \gamma_{\ell+1}\}$$

$$\mathcal{G}_\gamma[2] = \{y_{n+1-n_{\ell+1}}^n : n_{\ell+1} \Delta(Q_{\{y_{n+1-n_{\ell+1}}^n\}}, W_{\mathbf{a}}) < \gamma_{\ell+1}\}.$$

Following an analysis similar to the one leading to (83) and (84), we get

$$\mathbf{P} \left[\widehat{M}_{\ell+1} \neq 1 \mid \mathcal{M}_{\ell+1} = 1 \right] \leq e^{-n_{\ell+1} D + n_{\ell+1} \varepsilon_2(\gamma_{\ell+1}, n_{\ell+1})} \quad (91a)$$

$$\mathbf{P} \left[\widehat{M}_{\ell+1} \neq 2 \mid \mathcal{M}_{\ell+1} = 2 \right] \leq \varepsilon_3(\gamma_{\ell+1}, n_{\ell+1}). \quad (91b)$$

Using (88)–(91), we obtain, see obtain (92a)–(92d) shown at the bottom of the page. If we set $\gamma_i = |\mathcal{X}||\mathcal{Y}| \sqrt{4n_i \ln(1+n)}$ for $i = 1, 2, \dots, (\ell + 1)$ for $\varepsilon_1(\gamma_i, n_i)$, $\varepsilon_2(\gamma_i, n_i)$, and $\varepsilon_3(\gamma_i, n_i)$ given in (79), we have

$$n_i \varepsilon_1(\gamma_i, n_i) \leq 2 \ln(1+n) + 2|\mathcal{X}||\mathcal{Y}| \ln(n_i + 1) - 2|\mathcal{X}||\mathcal{Y}| \sqrt{4n_i \ln(1+n)} \ln \lambda$$

$$n_i \varepsilon_2(\gamma_i, n_i) \leq D + 2|\mathcal{X}||\mathcal{Y}| \ln(n_i + 1) - 2|\mathcal{X}||\mathcal{Y}| \sqrt{4n_i \ln(1+n)} \ln \lambda$$

$$\varepsilon_3(\gamma_i, n_i) \leq 16|\mathcal{X}||\mathcal{Y}|/(1+n)^2.$$

Using the concavity of \sqrt{z} function, we can conclude that

$$\sum_{i=1}^{\ell+1} \frac{n_i \varepsilon_1(\gamma_i, n_i) + C}{n} \leq 2|\mathcal{X}||\mathcal{Y}| \left(\frac{(\ell+1) \ln(1+n)}{n} - \frac{(\ell+1)}{n} 2\sqrt{\frac{n}{\ell+1} \ln(1+n)} \ln \lambda \right) + \frac{2(\ell+1) \ln(1+n)}{n} + \frac{(\ell+1)}{n} C \quad (93a)$$

$$\sum_{i=1}^{\ell+1} \frac{n_i \varepsilon_2(\gamma_i, n_i) + D}{n} \leq 2|\mathcal{X}||\mathcal{Y}| \left(\frac{(\ell+1) \ln(1+n)}{n} - \frac{(\ell+1)}{n} 2\sqrt{\frac{n}{\ell+1} \ln(1+n)} \ln \lambda \right) + \frac{\ell+1}{n} 2D \quad (93b)$$

$$\sum_{i=1}^{\ell+1} \varepsilon_3(\gamma_i, n_i) \leq 8|\mathcal{X}||\mathcal{Y}| \frac{\ell+1}{1+n}. \quad (93c)$$

$$P_{\mathbf{x}|\mathbf{m}} \leq \sum_{i=1}^{\ell+1} \mathbf{P} \left[\widehat{M}_i \neq (1 + m_i) \mid \mathcal{M} = \mathbf{m}, \widehat{M}_1 = \mathcal{M}_1, \dots, \widehat{M}_{i-1} = \mathcal{M}_{i-1} \right]$$

$$= \sum_{i=1}^{\ell+1} \mathbf{P} \left[\widehat{M}_i \neq \mathcal{M}_i \mid \mathcal{M}_i = 1 + m_i \right] \quad (88)$$

$$P_{\mathbf{e}|\mathbf{m}}(i) \leq \left[\sum_{j=1}^i \mathbf{P} \left[\widehat{M}_j \neq \mathcal{M}_j \mid \mathcal{M}_j = 1 + m_j \right] \right] \prod_{j=i+1}^{\ell+1} \mathbf{P} \left[\widehat{M}_j \neq 1 \mid \mathcal{M}_j = 1 \right] \quad (89)$$

$$|\mathcal{M}_i| \geq e^{nR_i - n_i \varepsilon_1(\gamma_i, n_i) - C} \quad \forall i = 1, 2, \dots, \ell \quad (92a)$$

$$|\mathcal{M}^i| \geq e^{n \sum_{j=1}^i R_j} e^{-\sum_{j=1}^i (n_j \varepsilon_1(\gamma_j, n_j) + C)} \quad \forall i = 1, 2, \dots, \ell \quad (92b)$$

$$P_{\mathbf{e}|\mathbf{m}}(i) \leq \sum_{j=1}^i \varepsilon_3(\gamma_j, n_j) \min \left\{ 1, e^{-n \sum_{j=i+1}^{\ell+1} \eta_j J\left(\frac{R_j}{\eta_j}\right)} e^{n \sum_{j=i+1}^{\ell+1} n_j \varepsilon_2(\gamma_j, n_j) + D} \right\} \quad \forall i = 1, 2, \dots, \ell, \forall \mathbf{m} \in \mathcal{M} \quad (92c)$$

$$P_{\mathbf{x}|\mathbf{m}} \leq \sum_{j=1}^{\ell+1} \varepsilon_3(\gamma_j, n_j) \quad \forall \mathbf{m} \in \mathcal{M} \quad (92d)$$

Then, Lemma 4 follows from (92) and (93) for any $\ell \leq \frac{n}{\ln(n+1)}$. ■

E. Proof of Lemma 5

Proof: For \mathbf{P} defined in (36) as a result of (37), we have

$$\mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] = \sum_{y^t \in \{y^t: \widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\} \cap \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \quad (94a)$$

$$\mathbf{P}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] = \sum_{y^t \in \{y^t: \widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\} \cap \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t). \quad (94b)$$

For $\mathbf{P}_{\{\mathcal{A}_i\}}$ defined in (39) as a result of (41), we have

$$\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] = \sum_{y^t \in \{y^t: \widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\} \cap \mathcal{Y}^{\tau^*}} \mathbf{P}_{\{\mathcal{A}_i\}}(y^t) \quad (95a)$$

$$\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] = \sum_{y^t \in \{y^t: \widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\} \cap \mathcal{Y}^{\tau^*}} \mathbf{P}_{\{\mathcal{A}_i\}}(y^t). \quad (95b)$$

Using (94) and (95) together with the data processing inequality for Kullback–Leibler divergence, we get

$$\begin{aligned} & \sum_{y^t \in \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \ln \frac{\mathbf{P}(y^t)}{\mathbf{P}_{\{\mathcal{A}_i\}}(y^t)} \\ & \geq \mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] \ln \frac{\mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]}{\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]} \\ & \quad + \mathbf{P}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] \ln \frac{\mathbf{P}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]}{\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]}. \end{aligned}$$

Since $0 \leq \mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] \leq 1$, we have

$$\begin{aligned} & \sum_{y^t \in \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \ln \frac{\mathbf{P}(y^t)}{\mathbf{P}_{\{\mathcal{A}_i\}}(y^t)} \geq -h\left(\mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]\right) \\ & \quad + \left(1 - \mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]\right) \ln \frac{1}{\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]}. \quad (96) \end{aligned}$$

Note that if $\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})$ and $M \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})$, then $\widehat{M} \neq M$. Consequently

$$\begin{aligned} \mathbf{P}\left[\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] &= \mathbf{P}\left[\left\{\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i}), M \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right\}\right] \\ & \quad + \mathbf{P}\left[\left\{\widehat{M} \in \mathcal{A}_i(\mathbf{Y}^{\tau_i}), M \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right\}\right] \\ & \leq P_e + \mathbf{P}\left[M \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]. \quad (97) \end{aligned}$$

Since the binary entropy function $h(\cdot)$ is increasing on the interval $[0, 1/2]$ if $P_e + \mathbf{P}\left[M \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right] \leq 1/2$, (96) and (97) imply

$$\begin{aligned} & \sum_{y^t \in \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \ln \frac{\mathbf{P}(y^t)}{\mathbf{P}_{\{\mathcal{A}_i\}}(y^t)} \geq -h(P_e + \mathbf{P}\left[M \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]) \\ & \quad + (1 - P_e - \mathbf{P}\left[M \in \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]) \ln \frac{1}{\mathbf{P}_{\{\mathcal{A}_i\}}\left[\widehat{M} \notin \mathcal{A}_i(\mathbf{Y}^{\tau_i})\right]}. \quad (98) \end{aligned}$$

Let \mathbf{B} , \mathbf{B}^* , and \mathbf{B}_τ be

$$\mathbf{B} \triangleq \ln \frac{\mathbf{P}(\mathbf{Y}^{\tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau})} \quad (99a)$$

$$\mathbf{B}^* \triangleq \ln \frac{\mathbf{P}(\mathbf{Y}^{\tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau})} \mathbb{1}_{\{\tau < \infty\}} \quad (99b)$$

$$\mathbf{B}_\tau \triangleq \ln \frac{\mathbf{P}(\mathbf{Y}^{\tau \wedge \tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau \wedge \tau})} \quad \forall \tau \in \{1, 2, \dots\} \quad (99c)$$

where $\tau \wedge \tau$ is the minimum of τ and τ .

Note that as τ goes to infinity, $\mathbf{B}_\tau \rightarrow \mathbf{B}$ and $\mathbf{B}_\tau \rightarrow \mathbf{B}^*$ with probability one. Since $|\mathbf{B}_\tau| \leq \tau \ln \frac{1}{\lambda}$ and $\mathbf{E}[\tau] < \infty$, we can apply the dominated convergence theorem [10, Th. 3, p. 187] to obtain

$$\mathbf{E}[\mathbf{B}] = \mathbf{E}[\mathbf{B}^*] = \lim_{\tau \rightarrow \infty} \mathbf{E}[\mathbf{B}_\tau]. \quad (100)$$

Finally for \mathbf{B} and \mathbf{B}^* defined in (99), we have

$$\mathbf{E}[\mathbf{B}] = \mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}^{\tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau})}\right] \quad (101a)$$

$$\mathbf{E}[\mathbf{B}^*] = \sum_{y^t \in \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \ln \frac{\mathbf{P}(y^t)}{\mathbf{P}_{\{\mathcal{A}_i\}}(y^t)}. \quad (101b)$$

Thus, as a result of (100) and (101), we have

$$\mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}^{\tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau})}\right] = \sum_{y^t \in \mathcal{Y}^{\tau^*}} \mathbf{P}(y^t) \ln \frac{\mathbf{P}(y^t)}{\mathbf{P}_{\{\mathcal{A}_i\}}(y^t)}. \quad (102)$$

Furthermore, using the definition of $\mathbf{P}_{\{\mathcal{A}_i\}}$ given in (39), we get

$$\begin{aligned} \mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}^{\tau})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}^{\tau})}\right] &= \mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_i})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_i})}\right] \\ &= \sum_{j=i}^k \xi_{i,j} \quad (103) \end{aligned}$$

where for all $i \geq 1$ and $j > i$

$$\xi_{i,j} \triangleq \begin{cases} 0 & \text{if } \mathbf{P}[T_{j+1} = T_j] = 1 \\ \mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_j})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_j})}\right] & \text{if } \mathbf{P}[T_{j+1} = T_j] < 1 \end{cases} \quad (104)$$

Assume for the moment that

$$\xi_{i,j} \leq \mathbf{E}[T_{j+1} - T_j] J(r_{j+1}) \quad (105)$$

where $T_{k+1} = T$ and r_j is defined in (47).

Then, Lemma (5) follows from (98), (102), (103), and (105).

Above, we have proved Lemma 5 by assuming that the inequality given in (105) holds for all i in $\{1, 2, \dots, k\}$ and j in $\{(i+1), \dots, (k+1)\}$; in the following, we prove that fact.

First note that if $\mathbf{P}[T_{j+1} = T_j] = 1$, then as result of (47) and (103), (105) is equivalent to $0 \leq 0J(0)$ which holds trivially. Thus, we assume hence forth that $\mathbf{P}[T_{j+1} = T_j] < 1$, which implies $\mathbf{E}[T_{j+1} - T_j] > 0$.

Let us consider the stochastic sequence

$$\begin{aligned} \mathbf{U}_\tau = & \left[-\ln \frac{\mathbf{P}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_j})}{\mathbf{P}_{\{\mathcal{A}_i\}}(\mathbf{Y}_{T_{j+1}}^{\tau} | \mathbf{Y}^{\tau_j})} + \sum_{k=T_{j+1}}^{\tau} J(\mathbf{I}(M; \mathbf{Y}_k | \mathbf{Y}^{k-1})) \right] \mathbb{1}_{\{\tau > T_j\}} \quad (106) \end{aligned}$$

where $\mathbf{I}(M; \mathbf{Y}_k | \mathbf{Y}^{k-1})$ is the conditional mutual information between M and \mathbf{Y}_k given \mathbf{Y}^{k-1} , defined as

$$\mathbf{I}(M; \mathbf{Y}_k | \mathbf{Y}^{k-1}) \triangleq \mathbf{E}\left[\ln \frac{\mathbf{P}(\mathbf{Y}_k | M, \mathbf{Y}^{k-1})}{\mathbf{P}(\mathbf{Y}_k | \mathbf{Y}^{k-1})} \middle| \mathbf{Y}^{k-1}\right].$$

Note that as it was the case for conditional entropy, while defining the conditional mutual information, we do not take the average over the conditioned random variable. Thus, $\mathbf{I}(M; \mathbf{Y}_k | \mathbf{Y}^{k-1})$ is itself a random variable.

For U_τ defined in (106), we have

$$U_{\tau+1} - U_\tau = \left(-\ln \frac{\mathbb{P}(Y_{\tau+1}|Y^\tau)}{\mathbb{P}_{\{\mathcal{A}_j\}}(Y_{\tau+1}|Y^\tau)} + J(\mathbf{I}(M; Y_{\tau+1} | Y^\tau)) \right) \mathbb{1}_{\{\tau \geq T_j\}}. \quad (107)$$

Conditioned on Y^τ random variables $M - X_{\tau+1} - Y_{\tau+1}$ form a Markov chain: thus, as a result of the data processing inequality for the mutual information, we have $\mathbf{I}(X_{\tau+1}; Y_{\tau+1} | Y^\tau) > \mathbf{I}(M; Y_{\tau+1} | Y^\tau)$. Since $J(\cdot)$ is a decreasing function, this implies that

$$J(\mathbf{I}(M; Y_{\tau+1} | Y^\tau)) \geq J(\mathbf{I}(X_{\tau+1}; Y_{\tau+1} | Y^\tau)). \quad (108)$$

Furthermore, because of the definitions of $J(\cdot)$, \mathbb{P} , and $\mathbb{P}_{\{\mathcal{A}_j\}}$ given in (17), (36), and (39), the convexity of Kullback–Leibler divergence and Jensen’s inequality, we have

$$J(\mathbf{I}(X_{\tau+1}; Y_{\tau+1} | Y^\tau)) \geq \mathbf{E} \left[\ln \frac{\mathbb{P}(Y_{\tau+1}|Y^\tau)}{\mathbb{P}_{\{\mathcal{A}_j\}}(Y_{\tau+1}|Y^\tau)} \middle| Y^\tau \right]. \quad (109)$$

Using (107)–(109), we get

$$\mathbf{E}[U_{\tau+1} | Y^\tau] \geq U_\tau. \quad (110)$$

Recall that $\min_{x,y} W_x(y) = \lambda$ and $|J(\cdot)| \leq D$. Thus, as a result of (107), we have

$$\mathbf{E}[|U_{\tau+1} - U_\tau| | Y^\tau] \leq \ln \frac{1}{\lambda} + D. \quad (111)$$

As a result of (110) and (111) and the fact that $U_0 = 0$, U_τ is a submartingale.

Recall that we have assumed that $\mathbf{P}[T_{j+1} \leq T] = 1$ and $\mathbf{E}[T] < \infty$; consequently

$$\mathbf{E}[T_{j+1}] < \infty. \quad (112)$$

Because of (111) and (112), we can apply a version of Doob’s optional stopping theorem [10, Th. 2, p. 487] to the submartingale U_τ and the stopping time T_{j+1} to obtain $\mathbf{E}[U_{T_{j+1}}] \geq \mathbf{E}[U_0] = 0$. Consequently

$$\mathbf{E} \left[\ln \frac{\mathbb{P}(Y_{T_{j+1}}^{T_{j+1}} | Y^{T_j})}{\mathbb{P}_{\{\mathcal{A}_j\}}(Y_{T_{j+1}}^{T_{j+1}} | Y^{T_j})} \right] \leq \mathbf{E} \left[\sum_{\tau=T_j+1}^{T_{j+1}} J(\mathbf{I}(M; Y_\tau | Y^{\tau-1})) \right]. \quad (113)$$

Note that as a result of the concavity of $J(\cdot)$ and Jensen’s inequality, we have

$$\begin{aligned} & \mathbf{E} \left[\sum_{\tau=T_j+1}^{T_{j+1}} J(\mathbf{I}(M; Y_\tau | Y^{\tau-1})) \right] \\ &= \mathbf{E}[T_{j+1} - T_j] \mathbf{E} \left[\sum_{\tau \geq 1} \frac{\mathbb{1}_{\{T_{j+1} \geq \tau > T_j\}} J(\mathbf{I}(M; Y_\tau | Y^{\tau-1}))}{\mathbf{E}[T_{j+1} - T_j]} \right] \\ &\leq \mathbf{E}[T_{j+1} - T_j] J \left(\frac{\mathbf{E} \left[\sum_{\tau \geq 1} \mathbb{1}_{\{T_{j+1} \geq \tau > T_j\}} \mathbf{I}(M; Y_\tau | Y^{\tau-1}) \right]}{\mathbf{E}[T_{j+1} - T_j]} \right). \end{aligned} \quad (114)$$

In order to calculate the argument of $J(\cdot)$ in (114), consider the stochastic sequence

$$V_\tau = H(M|Y^\tau) + \sum_{j=1}^{\tau} \mathbf{I}(M; Y_j | Y^{j-1}). \quad (115)$$

Clearly, $\mathbf{E}[V_{\tau+1} | Y^\tau] = V_\tau$ and $\mathbf{E}[|V_\tau|] \leq \ln |\mathcal{M}| + C\tau < \infty$. Hence, V_τ is a martingale.

Furthermore

$$\mathbf{E}[|V_{\tau+1} - V_\tau| | Y^\tau] \leq \ln |\mathcal{M}| + C. \quad (116)$$

Recall that we have assumed that $\mathbf{P}[T_j \leq T_{j+1} \leq T] = 1$ and $\mathbf{E}[T] < \infty$; consequently

$$\mathbf{E}[T_j] \leq \mathbf{E}[T_{j+1}] < \infty. \quad (117)$$

As a result of (116) and (117), we can apply Doob’s optimal stopping theorem, [10, Th. 2, p. 487] to V_τ both at stopping time T_j and at stopping time T_{j+1} , i.e., $\mathbf{E}[V_{T_{j+1}}] = \mathbf{E}[V_0]$ and $\mathbf{E}[V_{T_j}] = \mathbf{E}[V_0]$. Consequently

$$\begin{aligned} & \mathbf{E} \left[\sum_{\tau \geq 1} \mathbb{1}_{\{T_{j+1} \geq \tau > T_j\}} \mathbf{I}(M; Y_\tau | Y^{\tau-1}) \right] \\ &= \mathbf{E}[H(M|Y^{T_j}) - H(M|Y^{T_{j+1}})]. \end{aligned} \quad (118)$$

Using (113), (114), and (118)

$$\begin{aligned} & \mathbf{E} \left[\ln \frac{\mathbb{P}(Y_{T_{j+1}}^{T_{j+1}} | Y^{T_j})}{\mathbb{P}_{\{\mathcal{A}_j\}}(Y_{T_{j+1}}^{T_{j+1}} | Y^{T_j})} \right] \\ &\leq \mathbf{E}[T_{j+1} - T_j] J \left(\frac{\mathbf{E}[H(M|Y^{T_j}) - H(M|Y^{T_{j+1}})]}{\mathbf{E}[T_{j+1} - T_j]} \right). \end{aligned} \quad (119)$$

Hence, inequality given in (105) is not only when $\mathbf{P}[T_{j+1} = T_j] = 1$ but also when $\mathbf{P}[T_{j+1} = T_j] < 1$. ■

F. Proof of Lemma 6 for the Case $\mathbf{E}[T] < \infty$

Proof: In order to bound $P_{e|m}$ from below, we apply Lemma 5 for (T_1, \mathcal{A}_1) and (T_2, \mathcal{A}_2) given in (48)–(51) and use the fact that $J(\cdot) \leq D$ to get

$$\ln P_{e|m} \geq \frac{-h(P_e + |\mathcal{M}|^{-1}) - \mathbf{E}[T_2] J \left(\frac{\mathbf{E}[H(M) - H(M|Y^{T_2})]}{\mathbf{E}[T_2]} \right) - \mathbf{E}[T - T_2] D}{1 - P_e - |\mathcal{M}|^{-1}} \quad (120a)$$

$$\ln \mathbf{P}_{\{\mathcal{A}_2\}} \left[\widehat{M} \notin \mathcal{A}_2(Y^{T_2}) \right] \geq \frac{-h(P_e + \mathbf{P}[M \in \mathcal{A}_2(Y^{T_2})]) - \mathbf{E}[T - T_2] D}{1 - P_e - \mathbf{P}[M \in \mathcal{A}_2(Y^{T_2})]} \quad (120b)$$

provided that $|\mathcal{M}|^{-1} + P_e \leq 1/2$ and $\mathbf{P}[M \in \mathcal{A}_2(Y^{T_2})] + P_e \leq 1/2$.

We start with bounding $\mathbf{P}_{\{\mathcal{A}_2\}} \left[\widehat{M} \notin \mathcal{A}_2(Y^{T_2}) \right]$ from above and $\mathbf{P}[M \notin \mathcal{A}_2(Y^{T_2})]$ from below.

1) Since $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) = \lambda$, the posterior probability of a message at time $\tau + 1$ cannot be smaller than λ times the posterior probability of the same message at time τ . Hence,

for the stopping time T_2 defined in (50), random³³ set \mathcal{A}_2 defined in (51), and $\delta < \frac{1}{2}$, we have

$$\mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathbf{Y}^{T_2} = \mathbf{y}^{t_2}] > \lambda\delta \quad \forall \mathbf{y}^{t_2} \in \mathcal{Y}^{T_2*}. \quad (121)$$

As a result of the definition of $\mathbf{P}_{\{\mathcal{A}_2\}}(m, \mathbf{y}^t)$ given in (39), we have

$$\mathbf{P}_{\{\mathcal{A}_2\}}(m, \mathbf{y}^t) < \mathbf{P}(m, \mathbf{y}^t) \frac{\mathbb{1}_{\{m \in \mathcal{A}_2(\mathbf{y}^{t_2})\}}}{\lambda\delta} \quad \forall m \in \mathcal{M}, \mathbf{y}^t \in \mathcal{Y}^{T*}. \quad (122)$$

If the decoded message $\widehat{\mathbf{M}}(\mathbf{y}^t)$ is not in $\mathcal{A}_2(\mathbf{y}^{t_2})$ and message m is in $\mathcal{A}_2(\mathbf{y}^{t_2})$, then $\widehat{\mathbf{M}}(\mathbf{y}^t) \neq m$

$$\mathbb{1}_{\{\widehat{\mathbf{M}}(\mathbf{y}^t) \notin \mathcal{A}_2(\mathbf{y}^{t_2})\}} \mathbb{1}_{\{m \in \mathcal{A}_2(\mathbf{y}^{t_2})\}} \leq \mathbb{1}_{\{\widehat{\mathbf{M}}(\mathbf{y}^t) \neq m\}} \quad \forall m \in \mathcal{M}, \mathbf{y}^t \in \mathcal{Y}^T. \quad (123)$$

Using (122) and (123), we get

$$\mathbf{P}_{\{\mathcal{A}_2\}}(m, \mathbf{y}^t) \mathbb{1}_{\{\widehat{\mathbf{M}}(\mathbf{y}^t) \notin \mathcal{A}_2(\mathbf{y}^{t_2})\}} < \mathbf{P}(m, \mathbf{y}^t) \frac{\mathbb{1}_{\{\widehat{\mathbf{M}}(\mathbf{y}^t) \neq m\}}}{\lambda\delta} \quad \forall m \in \mathcal{M}, \mathbf{y}^t \in \mathcal{Y}^{T*}. \quad (124)$$

If we sum over all (m, \mathbf{y}^t) 's in $\mathcal{M} \times \mathcal{Y}^{T*}$ and use (37) and (41), we get

$$\mathbf{P}_{\{\mathcal{A}_2\}}[\widehat{\mathbf{M}} \notin \mathcal{A}_2(\mathbf{Y}^{T_2})] < \frac{\mathbf{P}[\widehat{\mathbf{M}} \neq \mathbf{M}]}{\lambda\delta} = \frac{P_e}{\lambda\delta}. \quad (125)$$

2) The probability of an event Γ_1 is lower bounded by the probability of its intersection with any event Γ_2 , i.e., $\mathbf{P}[\Gamma_1] \geq \mathbf{P}[\{\Gamma_1, \Gamma_2\}]$:

$$\begin{aligned} P_e &= \mathbf{P}[\widehat{\mathbf{M}} \neq \mathbf{M}] \\ &\geq \mathbf{P}[\{\widehat{\mathbf{M}} \neq \mathbf{M}, \mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}\}] \\ &= \mathbf{P}[\widehat{\mathbf{M}} \neq \mathbf{M} | \mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] \mathbf{P}[\mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}]. \end{aligned} \quad (126)$$

Note that if $\mathcal{A}_2(\mathbf{y}^{t_2}) = \mathcal{M}$, then T is reached before any of the messages reach a posterior probability of $1 - \delta$. Thus

$$\mathbf{P}[\widehat{\mathbf{M}} \neq \mathbf{M} | \mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] > \delta. \quad (127)$$

Thus, as a result of (126) and (127), we have

$$\mathbf{P}[\mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] < \frac{P_e}{\delta}. \quad (128)$$

On the other hand, if $\mathcal{A}_2(\mathbf{y}^{t_2}) \neq \mathcal{M}$, then the most likely message with a probability at least $(1 - \delta)$ is excluded from $\mathcal{A}_2(\mathbf{y}^{t_2})$. Thus

$$\mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathcal{A}_2(\mathbf{Y}^{T_2}) \neq \mathcal{M}] \leq \delta. \quad (129)$$

³³The set \mathcal{A}_2 is random in the sense that it depends on previous channel outputs.

Using (128) and (129) together with total probability formula, we get

$$\begin{aligned} \mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})] &= \mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] \mathbf{P}[\mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] \\ &\quad + \mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathcal{A}_2(\mathbf{Y}^{T_2}) \neq \mathcal{M}] \mathbf{P}[\mathcal{A}_2(\mathbf{Y}^{T_2}) \neq \mathcal{M}] \\ &\leq \mathbf{P}[\mathcal{A}_2(\mathbf{Y}^{T_2}) = \mathcal{M}] + \mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathcal{A}_1(\mathbf{Y}^{T_1}) \neq \mathcal{M}] \\ &< \frac{P_e}{\delta} + \delta. \end{aligned} \quad (130)$$

We plug the bounds on $\mathbf{P}_{\{\mathcal{A}_2\}}[\widehat{\mathbf{M}} \notin \mathcal{A}_2(\mathbf{Y}^{T_2})]$ and $\mathbf{P}[\mathbf{M} \notin \mathcal{A}_2(\mathbf{Y}^{T_2})]$ given in (125) and (130) in (120) to get

$$\ln P_{e|m} \geq \frac{-h(\tilde{\epsilon}_1) - \mathbf{E}[T_2]J \left(\frac{\mathbf{E}[H(\mathbf{M}) - H(\mathbf{M} | \mathbf{Y}^{T_2})]}{\mathbf{E}[T_2]} \right) - \mathbf{E}[T - T_2]D}{1 - \tilde{\epsilon}_1} \quad (131a)$$

$$\ln \frac{P_e}{\lambda\delta} \geq \frac{-h(\tilde{\epsilon}_1) - \mathbf{E}[T - T_2]D}{1 - \tilde{\epsilon}_1} \quad (131b)$$

provided that $\tilde{\epsilon}_1 \leq 1/2$ where $\tilde{\epsilon}_1 = P_e + \delta + \frac{P_e}{\delta} + |\mathcal{M}|^{-1}$.

Now, we bound $\mathbf{E}[H(\mathbf{M} | \mathbf{Y}^{T_2})]$ from below. Note that $\mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}}$ is a discrete random variable that is either zero or one; its conditional entropy given \mathbf{Y}^{T_2} is given by

$$H(\mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}} | \mathbf{Y}^{T_2}) = h(\mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathbf{Y}^{T_2}]). \quad (132)$$

Furthermore, since $\mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}}$ is a function of \mathbf{Y}^{T_2} and \mathbf{M} , chain rule entropy implies that

$$\begin{aligned} H(\mathbf{M} | \mathbf{Y}^{T_2}) &= H(\mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}} | \mathbf{Y}^{T_2}) \\ &\quad + \mathbf{E}[H(\mathbf{M} | \mathbf{Y}^{T_2}, \mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}}) | \mathbf{Y}^{T_2}]. \end{aligned} \quad (133)$$

Since $\mathcal{A}_2(\mathbf{Y}^{T_2})$ has at most $|\mathcal{M}|$ elements and its complement, $\mathcal{M} \setminus \mathcal{A}_2(\mathbf{Y}^{T_2})$, has at most one element, we can bound the conditional entropy of the messages as follows:

$$H(\mathbf{M} | \mathbf{Y}^{T_2}, \mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}}) \leq \mathbb{1}_{\{\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2})\}} \ln |\mathcal{M}|. \quad (134)$$

Thus, using (132)–(134), we get

$$\begin{aligned} H(\mathbf{M} | \mathbf{Y}^{T_2}) &\leq h(\mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathbf{Y}^{T_2}]) \\ &\quad + \mathbf{P}[\mathbf{M} \in \mathcal{A}_2(\mathbf{Y}^{T_2}) | \mathbf{Y}^{T_2}] \ln |\mathcal{M}|. \end{aligned} \quad (135)$$

Then, using concavity of the binary entropy function $h(\cdot)$ together with (130) and (135), we get

$$\mathbf{E}[H(\mathbf{M} | \mathbf{Y}^{T_2})] < h(\delta + \frac{P_e}{\delta}) + (\delta + \frac{P_e}{\delta}) \ln |\mathcal{M}| \quad (136)$$

provided that $\delta + \frac{P_e}{\delta} \leq 1/2$.

If we plug in (136) and the identity $H(\mathbf{M}) = \ln |\mathcal{M}|$ in (131), we get

$$(1 - \tilde{\epsilon}_1) \frac{\ln P_{e|m}}{\mathbf{E}[T]} \geq -\frac{h(\tilde{\epsilon}_1)}{\mathbf{E}[T]} - \eta J \left(\frac{(1 - \tilde{\epsilon}_1)R - h(\tilde{\epsilon}_1) / \mathbf{E}[T]}{\eta} \right) - (1 - \eta)D \quad (137a)$$

$$-(1 - \tilde{\epsilon}_1)E \geq \frac{-h(\tilde{\epsilon}_1) + \ln \lambda\delta}{\mathbf{E}[T]} - (1 - \eta)D \quad (137b)$$

provided that $\tilde{\epsilon}_1 \leq 1/2$ where $\eta = \frac{\mathbf{E}[T_2]}{\mathbf{E}[T]}$, $\tilde{\epsilon}_1 = P_e + \delta + \frac{P_e}{\delta} + |\mathcal{M}|^{-1}$, $R = \frac{|\mathcal{M}|}{\mathbf{E}[T]}$, and $E = \frac{-\ln P_e}{\mathbf{E}[T]}$.

Note that the inequality given in (137b) bounds the value of η from above

$$\eta \leq 1 - \frac{(1-\tilde{\epsilon}_1)E-\tilde{\epsilon}_2}{D} \quad (138)$$

where $\tilde{\epsilon}_2 = \frac{h(\tilde{\epsilon}_1)-\ln \lambda \delta}{\mathbf{E}[T]}$.

Furthermore, for any $\eta_1 \leq \eta_2 \leq \frac{\tilde{R}}{C}$, as a result of concavity of $J(\cdot)$, we have

$$\begin{aligned} \eta_1 J\left(\frac{\tilde{R}}{\eta_1}\right) + (1-\eta_1)D &= \eta_1 J\left(\frac{\tilde{R}}{\eta_1}\right) + (\eta_2 - \eta_1)J(0) \\ &\quad + (1-\eta_2)D \\ &\leq \eta_2 J\left(\frac{\tilde{R}}{\eta_2}\right) + (1-\eta_2)D. \end{aligned} \quad (139)$$

Using (138) and (139) we see that the bound in (137a) is lower bounded by its value at $\eta = 1 - \frac{(1-\tilde{\epsilon}_1)E-\tilde{\epsilon}_2}{D}$ if $E \geq \frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1}$ and by its value at $\eta = 1$ otherwise, i.e.,

$$\frac{\ln P_{e|m}}{\mathbf{E}[T]} \geq \begin{cases} -E - \left(1 - \frac{E-\tilde{\epsilon}}{D}\right) J\left(\frac{R-\frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1}}{1-\frac{E-\tilde{\epsilon}}{D}}\right) & \text{if } E \geq \frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1} \\ -\frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1} - \frac{1}{1-\tilde{\epsilon}_1} J((1-\tilde{\epsilon}_1)R - \tilde{\epsilon}_2) & \text{if } E < \frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1} \end{cases}$$

where $\tilde{\epsilon} = \frac{\tilde{\epsilon}_1 D + \tilde{\epsilon}_2}{1-\tilde{\epsilon}_1}$.

Then, for the case $E \geq \frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1}$, Lemma 6 follows from the fact that $J(\cdot)$ is a nonnegative decreasing function. For the case $E < \frac{\tilde{\epsilon}_2}{1-\tilde{\epsilon}_1}$, Lemma 6 follows from the fact that $J(\cdot)$ is a concave nonnegative decreasing function. ■

G. Proof of Lemma 7 for the Case $\mathbf{E}[T] < \infty$

Proof: We start with proving the bounds given in (57) and (58).

1) Let us start with the bound on $\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})]$ given in (57). Since $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) = \lambda$, the posterior probability of a $m^i \in \mathcal{M}^i$ at time $\tau + 1$ cannot be smaller than λ times its value at time τ . Hence, for $\delta < 1/2$, as a result of definitions of T_i and $\mathcal{A}_i(Y^{T_i})$ given in (55) and (56), we have

$$\begin{aligned} \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | Y^{T_i} = y^{t_i}] &> \lambda \delta \\ &\forall y^{t_i} \in \mathcal{Y}^{T_i^*}, i \in \{1, 2, \dots, \ell\}. \end{aligned}$$

Then, as a result of the definition of $\mathbf{P}_{\{\mathcal{A}_i\}}(m, y^t)$ given in (39), we have

$$\begin{aligned} \mathbf{P}_{\{\mathcal{A}_i\}}(m, y^t) &< \mathbf{P}(m, y^t) \frac{\mathbb{1}_{\{m \in \mathcal{A}_i(y^{t_i})\}}}{\lambda \delta} \\ &\forall m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*}, i \in \{1, 2, \dots, \ell\}. \end{aligned} \quad (140)$$

For $\mathcal{A}_i(y^{t_i})$ given in (56), if the decoded message $\widehat{M}(y^t)$ is not in $\mathcal{A}_i(y^{t_i})$ but m is in $\mathcal{A}_i(y^{t_i})$, then $\widehat{M}(y^t) \neq m^i$

$$\begin{aligned} \mathbb{1}_{\{\widehat{M}(y^t) \notin \mathcal{A}_i(y^{t_i})\}} \mathbb{1}_{\{m \in \mathcal{A}_i(y^{t_i})\}} &\leq \mathbb{1}_{\{\widehat{M}^i \neq m^i\}} \\ &\forall m \in \mathcal{M}, y^t \in \mathcal{Y}^T, i \in \{1, 2, \dots, \ell\}. \end{aligned} \quad (141)$$

Using (140) and (141), we get

$$\begin{aligned} \mathbf{P}_{\{\mathcal{A}_i\}}(m, y^t) \mathbb{1}_{\{\widehat{M}(y^t) \notin \mathcal{A}_i(y^{t_i})\}} &< \mathbf{P}(m, y^t) \frac{\mathbb{1}_{\{\widehat{M}^i \neq m^i\}}}{\lambda \delta} \\ &\forall m \in \mathcal{M}, y^t \in \mathcal{Y}^{T^*}, i \in \{1, 2, \dots, \ell\}. \end{aligned}$$

If we sum over all (m, y^t) 's in $\mathcal{M} \times \mathcal{Y}^{T^*}$ and use (37) and (41), we get

$$\begin{aligned} \mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})] &< \frac{\mathbf{P}[\widehat{M}^i \neq M^i]}{\lambda \delta} \\ &= \frac{P_e(i)}{\lambda \delta} \quad \forall i \in \{1, 2, \dots, \ell\}. \end{aligned} \quad (142)$$

2) Let us now prove the bound on $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$ given in (58).

a) If $\mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}$, then at T_i , there is a m^i with posterior probability $(1-\delta)$ and all the messages m of the form $m = (m^i, m_{i+1}, \dots, m_k)$ are excluded from \mathcal{A}_i . Consequently, we have

$$\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | \mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}] < \delta. \quad (143)$$

b) If $\mathcal{A}_i(Y^{T_i}) = \mathcal{M}$, then at T_i , there is no m^i with posterior probability $(1-\delta)$ and $T_i = T$. Since $\widehat{M}^i \neq M^i$ implies that $\widehat{M} \neq M$, we have

$$\mathbf{P}[\widehat{M} \neq M | \mathcal{A}_i(Y^{T_i}) = \mathcal{M}] \geq \delta. \quad (144)$$

As a result of total probability formula for $\mathbf{P}[\widehat{M} \neq M]$, we have

$$\begin{aligned} P_e &= \mathbf{P}[\widehat{M} \neq M | \mathcal{A}_i(Y^{T_i}) = \mathcal{M}] \mathbf{P}[\mathcal{A}_i(Y^{T_i}) = \mathcal{M}] \\ &\quad + \mathbf{P}[\widehat{M} \neq M | \mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}] \mathbf{P}[\mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}] \\ &\geq \delta \mathbf{P}[\mathcal{A}_i(Y^{T_i}) = \mathcal{M}]. \end{aligned} \quad (145)$$

If use the total probability formula for $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$ together with (143) and (145), we get

$$\begin{aligned} \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})] &= \mathbf{P}[\{M \in \mathcal{A}_i(Y^{T_i}), \mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}\}] \\ &\quad + \mathbf{P}[\{M \in \mathcal{A}_i(Y^{T_i}), \mathcal{A}_i(Y^{T_i}) = \mathcal{M}\}] \\ &\leq \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | \mathcal{A}_i(Y^{T_i}) \neq \mathcal{M}] \\ &\quad + \mathbf{P}[\mathcal{A}_i(Y^{T_i}) = \mathcal{M}] \\ &\leq \delta + \frac{P_e}{\delta}. \end{aligned}$$

We apply Lemma 5 for $(T_1, \mathcal{A}_1), \dots, (T_k, \mathcal{A}_k)$ defined in (55) and (56); use the bounds on $\mathbf{P}_{\{\mathcal{A}_i\}}[\widehat{M} \notin \mathcal{A}_i(Y^{T_i})]$ and $\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]$ given in (57) and (58). Then, we can conclude that if $P_e + \delta + P_e/\delta \leq 1/2$, then

$$(1-\tilde{\epsilon}_3)E_i \leq \tilde{\epsilon}_5 + \sum_{j=i+1}^{\ell+1} \nu_j J(r_j) \quad i = 1, 2, \dots, \ell \quad (146)$$

where $R_i, E_i, \tilde{\epsilon}_3$, and $\tilde{\epsilon}_5$ are defined in Lemma 7, r_j 's are defined in (47) of Lemma 5, and ν_j 's are defined as follows³⁴:

$$\nu_j \triangleq \frac{\mathbf{E}[T_j] - \mathbf{E}[T_{j-1}]}{\mathbf{E}[T]} \quad \forall j \in \{1, 2, \dots, \ell + 1\}. \quad (147)$$

Depending on the values of ν_j and r_j , the bound in (146) takes different values. However, ν_j and r_j are not changing freely. As a result of (118) and the fact that $\mathbf{I}(M; Y_{t+1} | Y^t) \leq C$, we have

$$r_j \leq C \quad j \in \{1, 2, \dots, (\ell + 1)\}. \quad (148)$$

In addition, ν_j 's and r_j 's are constrained by the definitions of T_j and $\mathcal{A}_j(Y^{T_j})$ given in (55) and (56). At T_j with high probability, one element of \mathcal{M}^j has a posterior probability $(1 - \delta)$. In the following, we use this fact to bound $\mathbf{E}[H(M|Y^{T_j})]$ from above. Then, we turn this bound into a constraint on the values of ν_j 's and r_j 's and use that constraint together with (146) and (148) to bound E_i 's from above.

For all j in $\{1, 2, \dots, \ell\}$, $\mathbb{1}_{\{M \in \mathcal{A}_j(Y^{T_j})\}}$ is a discrete random variable that is either zero or one; its conditional entropy is given by

$$H(\mathbb{1}_{\{M \in \mathcal{A}_j(Y^{T_j})\}} | Y^{T_j}) = h(\mathbf{P}[M \in \mathcal{A}_j(Y^{T_j}) | Y^{T_j}]). \quad (149)$$

Furthermore, since $\mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}}$ is a function of Y^{T_i} and M , the chain rule entropy implies that

$$H(M|Y^{T_i}) = H(\mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}} | Y^{T_i}) + \mathbf{E}\left[H(M|Y^{T_i}, \mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}}) | Y^{T_i}\right]. \quad (150)$$

Note that $\mathcal{A}_i(Y^{T_i})$ has at most $|\mathcal{M}|$ elements and its complement, $\mathcal{M} \setminus \mathcal{A}_i(Y^{T_i})$, has at most $\frac{|\mathcal{M}|}{|\mathcal{M}^i|}$ elements. We can bound the conditional entropy of the messages $H(M|Y^{T_i}, \mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}})$ as follows:

$$\begin{aligned} H(M|Y^{T_i}, \mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}}) &\leq \mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}} \ln |\mathcal{M}| \\ &\quad + \mathbb{1}_{\{M \notin \mathcal{A}_i(Y^{T_i})\}} \ln \frac{|\mathcal{M}|}{|\mathcal{M}^i|} \\ &= \ln \frac{|\mathcal{M}|}{|\mathcal{M}^i|} + \mathbb{1}_{\{M \in \mathcal{A}_i(Y^{T_i})\}} \ln |\mathcal{M}^i|. \end{aligned} \quad (151)$$

Thus, using (149)–(151), we get

$$\begin{aligned} H(M|Y^{T_i}) &\leq h(\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | Y^{T_i}]) + \ln \frac{|\mathcal{M}|}{|\mathcal{M}^i|} \\ &\quad + \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i}) | Y^{T_i}] \ln |\mathcal{M}^i|. \end{aligned} \quad (152)$$

If we take the expectation of both sides of the inequality (152) and use the concavity of the binary entropy function, we get

$$\begin{aligned} \mathbf{E}[H(M|Y^{T_i})] &\leq h(\mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})]) + \ln \frac{|\mathcal{M}|}{|\mathcal{M}^i|} \\ &\quad + \mathbf{P}[M \in \mathcal{A}_i(Y^{T_i})] \ln |\mathcal{M}^i|. \end{aligned}$$

Using the inequality given (58) and the fact that binary entropy function is an increasing function on the interval $[0, 1/2]$, we see that

$$\begin{aligned} \mathbf{E}[H(M|Y^{T_i})] &< h(P_e + \delta + \frac{P_e}{\delta}) + \ln \frac{|\mathcal{M}|}{|\mathcal{M}^i|} \\ &\quad + (P_e + \delta + \frac{P_e}{\delta}) \ln |\mathcal{M}^i| \end{aligned} \quad (153)$$

provided that $P_e + \delta + \frac{P_e}{\delta} \leq 1/2$.

Note that as a result of Fano's inequality for $\mathbf{E}[H(M|Y^T)]$, we have

$$\mathbf{E}[H(M|Y^T)] < h(P_e) + P_e \ln |\mathcal{M}|. \quad (154)$$

If we divide both sides of the inequalities (153) and (154) by $\mathbf{E}[T]$, we see that following bounds holds:

$$\frac{\mathbf{E}[H(M|Y^{T_i})]}{\mathbf{E}[T]} \leq \tilde{\epsilon}_4 + R - \sum_{j=1}^i R_j + \tilde{\epsilon}_3 \sum_{j=1}^i R_j \quad i = 1, 2, \dots, \ell \quad (155a)$$

$$\mathbf{E}[H(M|Y^T)] \leq h(P_e) + P_e R. \quad (155b)$$

Note that

$$\frac{\mathbf{E}[H(M|Y^{T_i})]}{\mathbf{E}[T]} = R - \sum_{j=1}^i \nu_j r_j \quad i = 1, 2, \dots, (\ell + 1). \quad (156)$$

Using (155) and (156), we get

$$\sum_{j=1}^i \nu_j r_j \geq (1 - \tilde{\epsilon}_3) \sum_{j=1}^i R_j - \tilde{\epsilon}_4 \quad i = 1, 2, \dots, \ell \quad (157a)$$

$$\sum_{j=1}^{\ell+1} \nu_j r_j \geq (1 - P_e)R - \frac{h(P_e)}{\mathbf{E}[T]} \quad (157b)$$

where r_j 's and ν_j 's are given in (47) and (147) respectively.

Thus, using (47), (146)–(148), and (157), we reach the following conclusion. For any variable-length block code satisfying the hypothesis of Lemma 7 and for any positive δ such that $P_e + \delta + \frac{P_e}{\delta} \leq \frac{1}{2}$

$$(1 - \tilde{\epsilon}_3)E_i - \tilde{\epsilon}_5 \leq \sum_{j=i+1}^{\ell+1} \nu_j J(r_j) \quad i = 1, 2, \dots, \ell \quad (158a)$$

$$(1 - \tilde{\epsilon}_3) \sum_{j=1}^i R_j - \tilde{\epsilon}_4 \leq \sum_{j=1}^i \nu_j r_j \quad i = 1, 2, \dots, \ell \quad (158b)$$

$$(1 - P_e)R - \frac{h(P_e)}{\mathbf{E}[T]} \leq \sum_{j=1}^{\ell+1} \nu_j r_j \quad (158c)$$

for some $(\nu_1^{\ell+1}, r_1^{\ell+1})$ such that

$$r_i \in [0, C] \quad i = 1, 2, \dots, (\ell + 1) \quad (159a)$$

$$\nu_i \geq 0 \quad i = 1, 2, \dots, (\ell + 1) \quad (159b)$$

$$\sum_{i=1}^{\ell+1} \nu_i = 1. \quad (159c)$$

³⁴We use the convention $T_0 = 0$ and $T_{\ell+1} = T$.

We show in the following if the constraints given in (158) is satisfied for some $(\nu_1^{\ell+1}, r_1^{\ell+1})$ satisfying (159), and constraints given in (60) is satisfied for some (η_1^ℓ) satisfying (61).

One can confirm numerically that

$$(1 - \tilde{\epsilon}_3) \ln 2 > h(\tilde{\epsilon}_3) \quad \forall \tilde{\epsilon}_3 \in [0, \frac{1}{5}].$$

Recall that we have assumed that $P_e + \delta + \frac{P_e}{\delta} \leq \frac{1}{5}$, i.e., $\tilde{\epsilon}_3 \leq \frac{1}{5}$. Thus

$$(1 - \tilde{\epsilon}_3)R_1 - \tilde{\epsilon}_4 > 0. \quad (160)$$

Let $\eta_1, \tilde{r}_1, \tilde{\nu}_2$, and \tilde{r}_2 be

$$\begin{aligned} \eta_1 &= \frac{(1 - \tilde{\epsilon}_3)R_1 - \tilde{\epsilon}_4}{r_1} \\ \tilde{r}_1 &= r_1 \\ \tilde{\nu}_2 &= \nu_2 + \nu_1 - \eta_1 \\ \tilde{r}_2 &= \frac{r_2 \nu_2 + (\nu_1 - \eta_1)r_1}{\nu_2}. \end{aligned}$$

Note that $(\eta_1, \tilde{\nu}_2, \nu_3^{\ell+1}, \tilde{r}_1, \tilde{r}_2, r_3^{\ell+1})$ satisfies (158b), (158c), and (159) by construction. Furthermore, as a result of concavity of $J(\cdot)$, we have

$$\nu_1 J(r_1) + \nu_2 J(r_2) \leq \eta_1 J(\tilde{r}_1) + \tilde{\nu}_2 J(\tilde{r}_2).$$

Thus, $(\eta_1, \tilde{\nu}_2, \nu_3^{\ell+1}, \tilde{r}_1, \tilde{r}_2, r_3^{\ell+1})$ also satisfies (158a).

For $j \geq 2$, we use $\tilde{\nu}_j$ and \tilde{r}_j to define $\eta_j, \tilde{\nu}_{j+1}$ and \tilde{r}_{j+1} as follows:

$$\eta_j = \frac{(1 - \tilde{\epsilon}_3)R_j}{r_j} \quad (161a)$$

$$\tilde{\nu}_{j+1} = \nu_{j+1} + \tilde{\nu}_j - \eta_j \quad (161b)$$

$$\tilde{r}_{j+1} = \frac{r_{j+1}\nu_{j+1} + (\tilde{\nu}_j - \eta_j)\tilde{r}_j}{\nu_{j+1}}. \quad (161c)$$

Using the fact that $(\eta_1^{j-1}, \tilde{\nu}_j, \nu_{j+1}^{\ell+1}, \tilde{r}_1^j, r_{j+1}^{\ell+1})$ satisfies (158) and (159) and the concavity of $J(\cdot)$, we can show that $(\eta_1^j, \tilde{\nu}_{j+1}, \nu_{j+2}^{\ell+1}, \tilde{r}_1^{j+1}, r_{j+2}^{\ell+1})$ also satisfies (158) and (159). We repeat the iteration given in (161) until we reach $\tilde{\nu}_{\ell+1}$ and $\tilde{r}_{\ell+1}$ and we let $\eta_{\ell+1} = \tilde{\nu}_{\ell+1}$.

Then, we conclude that for any variable-length block code satisfying the hypothesis of the Lemma 7 and for any positive δ such that $P_e + \delta + \frac{P_e}{\delta} \leq \frac{1}{5}$

$$(1 - \tilde{\epsilon}_3)E_i - \tilde{\epsilon}_5 \leq \sum_{j=i+1}^{\ell+1} \eta_j J(\tilde{r}_j) \quad i=1, 2, \dots, \ell \quad (162a)$$

$$(1 - \tilde{\epsilon}_3)R_i - \tilde{\epsilon}_4 \mathbb{1}_{\{i=1\}} = \tilde{r}_i \eta_i \quad i=1, 2, \dots, \ell \quad (162b)$$

$$(\tilde{\epsilon}_3 - P_e)R + \frac{h(\tilde{\epsilon}_3) - h(P_e)}{E[T]} \leq \tilde{r}_{\ell+1} \eta_{\ell+1} \quad (162c)$$

for some $(\eta_1, \dots, \eta_{\ell+1}, \tilde{r}_1, \dots, \tilde{r}_{\ell+1})$ such that³⁵

$$\tilde{r}_i \in [0, C] \quad i=1, 2, \dots, (\ell+1) \quad (163a)$$

$$\eta_i \geq 0 \quad i=1, 2, \dots, (\ell+1) \quad (163b)$$

$$\sum_{i=1}^{\ell+1} \eta_i = 1. \quad (163c)$$

Lemma 7 follows from the fact that $J(\cdot) \leq D$. \blacksquare

³⁵One can replace the inequality in (162c) by equality because $J(\cdot)$ is a decreasing function.

H. Codes with Infinite Expected Decoding Time on Channels with Positive Transition Probabilities

In this section, we consider variable-length block codes on DMCs with positive transition probabilities, i.e., $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) > 0$, and derive lower bounds to the probabilities of various error events. These bounds, i.e., (166), (172), and (175), enable us to argue that Lemmas 6 and 7 hold for variable-length block codes with infinite expected decoding time, i.e., $\mathbf{E}[T] = \infty$.

1) $P_e > 0$: On DMC such that $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) = \lambda$, the posterior probability of any message $m \in \mathcal{M}$ at time τ is lower bounded as

$$\mathbf{P}[M = m | Y^\tau] \geq \left(\frac{\lambda}{1-\lambda}\right)^\tau \frac{1}{|\mathcal{M}|}.$$

Then, conditioned on the event $\{T = \tau\}$, the probability of erroneous decoding is lower bounded as

$$\mathbf{P}[\hat{M} \neq M | T = \tau] \geq \frac{|\mathcal{M}|-1}{|\mathcal{M}|} \left(\frac{\lambda}{1-\lambda}\right)^\tau. \quad (164)$$

Note that since $\mathbf{P}[T < \infty] = 1$, the error probability of any variable-length code satisfies

$$P_e = \sum_{\tau=1}^{\infty} \mathbf{P}[M \neq \hat{M} | T = \tau] \mathbf{P}[T = \tau]. \quad (165)$$

Using (164) and (165), we get

$$P_e \geq \frac{|\mathcal{M}|-1}{|\mathcal{M}|} \mathbf{E} \left[\left(\frac{\lambda}{1-\lambda}\right)^T \right]. \quad (166)$$

Note that (166) implies that for a variable-length code with infinite expected decoding time, not only the rate R but also the error exponent E is zero.

2) If $P_e + \frac{1}{|\mathcal{M}|} < 1$, then $\min_m P_{e|m} > 0$: Note that since $\mathbf{P}[T < \infty] = 1$ and $|\mathcal{M}| < \infty$

$$\mathbf{P}[T < \infty | M = m] = 1 \quad \forall m \in \mathcal{M}.$$

For any variable-length block code such that $P_e + \frac{1}{|\mathcal{M}|} < 1$, let τ^* be

$$\tau^* = \min \left\{ \tau : \max_{m \in \mathcal{M}} \mathbf{P}[T > \tau | M = m] \leq \frac{|\mathcal{M}|-1}{|\mathcal{M}|} - P_e \right\}. \quad (167)$$

Since $\mathbf{P}[T < \infty | M = m] = 1$ for all m in \mathcal{M} and \mathcal{M} is finite, τ^* is finite.

Note that for any τ, m , and \tilde{m} we have

$$\mathbf{P}[Y^\tau = y^\tau | M = m] \geq \left(\frac{\lambda}{1-\lambda}\right)^\tau \mathbf{P}[Y^\tau = y^\tau | M = \tilde{m}]. \quad (168)$$

Then, using (168), we get

$$\begin{aligned} P_{e|m} &\geq \sum_{\tilde{m} \neq m} \mathbf{P} \left[\left\{ \hat{M} = \tilde{m}, T \leq \tau^* \right\} \middle| M = m \right] \\ &= \left(\frac{\lambda}{1-\lambda}\right)^{\tau^*} \sum_{\tilde{m} \neq m} \mathbf{P} \left[\left\{ \hat{M} = \tilde{m}, T \leq \tau^* \right\} \middle| M = \tilde{m} \right] \\ &\geq \left(\frac{\lambda}{1-\lambda}\right)^{\tau^*} \sum_{\tilde{m} \neq m} \left(\mathbf{P} \left[\hat{M} = \tilde{m} \middle| M = \tilde{m} \right] - \mathbf{P}[T > \tau^* | M = \tilde{m}] \right). \end{aligned} \quad (169)$$

Note that as a result of (167), we have

$$\mathbf{P}[\mathsf{T} > \tau^* | \mathsf{M} = \tilde{m}] \leq \left(\frac{|\mathcal{M}|-1}{|\mathcal{M}|} - P_e \right) \quad \forall \tilde{m} \in \mathcal{M}. \quad (170)$$

Furthermore

$$\sum_{\tilde{m} \neq m} \mathbf{P}[\hat{\mathsf{M}} = \tilde{m} | \mathsf{M} = \tilde{m}] \geq |\mathcal{M}|(1 - P_e) - 1. \quad (171)$$

Thus, using (169)–(171), we get

$$\min_{m \in \mathcal{M}} P_{e|m} \geq \left(\frac{\lambda}{1-\lambda} \right)^{\tau^*} \left(1 - \frac{1}{|\mathcal{M}|} - P_e \right) \quad (172)$$

where τ^* is a finite integer defined in (167).

3) For all $i \in \{1, 2, \dots, \ell\}$, $P_e(i) > 0$: For a variable-length block code with message set \mathcal{M} of the form $\mathcal{M} = \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_k$ on a DMC such that $\min_{x \in \mathcal{X}, y \in \mathcal{Y}} W_x(y) = \lambda$, the posterior probability of any element of \mathcal{M}_i at time τ is lower bounded as

$$\mathbf{P}[\mathsf{M}^i = m^i | \mathsf{Y}^\tau] \geq \left(\frac{\lambda}{1-\lambda} \right)^\tau \frac{1}{|\mathcal{M}_i|} \quad \forall m^i \in \mathcal{M}^i, \forall i \in \{1, 2, \dots, \ell\}.$$

Then, conditioned on the event $\{\mathsf{T} = \tau\}$, the probability of decoding any one of the first i sub-messages erroneously is lower bounded as

$$\mathbf{P}[\hat{\mathsf{M}}^i \neq \mathsf{M}^i | \mathsf{T} = \tau] \geq \frac{|\mathcal{M}^i|-1}{|\mathcal{M}^i|} \left(\frac{\lambda}{1-\lambda} \right)^\tau. \quad (173)$$

Since $\mathbf{P}[\mathsf{T} < \infty] = 1$, $P_e(i)$ satisfies

$$P_e = \sum_{\tau=1}^{\infty} \mathbf{P}[\mathsf{M}^i \neq \hat{\mathsf{M}}^i | \mathsf{T} = \tau] \mathbf{P}[\mathsf{T} = \tau]. \quad (174)$$

Using (173) and (174), we get

$$P_e(i) \geq \frac{|\mathcal{M}^i|-1}{|\mathcal{M}^i|} \mathbf{E} \left[\left(\frac{\lambda}{1-\lambda} \right)^\tau \right] \quad \forall i \in \{1, 2, \dots, \ell\}. \quad (175)$$

Equation (175) implies that for a variable-length code with infinite expected decoding time, not only the rates but also the error exponents of submessages are zero.

I. Proof of Theorem 1

Proof: In Section IV-C, it is shown that for any rate $R \in [0, C]$, error exponent $E \in [0, (1 - \frac{R}{C})D]$, there exists a reliable sequence \mathbb{Q} such that $R_{\mathbb{Q}} = R$, $E_{\mathbb{Q}} = E$, $E_{\text{md}, \mathbb{Q}} = E + (1 - \frac{E}{D})J\left(\frac{R}{1-E/D}\right)$. Thus, as a result of the definition of $E_{\text{md}}(R, E)$ given in (13), we have

$$E_{\text{md}}(R, E) \geq E + \left(1 - \frac{E}{D}\right) J\left(\frac{R}{1-E/D}\right). \quad (176)$$

In Section V-C, we have shown that any reliable sequence of codes \mathbb{Q} with rate $R_{\mathbb{Q}}$ and error exponent $E_{\mathbb{Q}}$ satisfies

$$E_{\text{md}, \mathbb{Q}} \leq E_{\mathbb{Q}} + \left(1 - \frac{E_{\mathbb{Q}}}{D}\right) J\left(\frac{R_{\mathbb{Q}}}{1-E_{\mathbb{Q}}/D}\right).$$

Thus, using the fact that $J(\cdot)$ is a decreasing concave function, we can conclude that

$$\max_{\substack{R_{\mathbb{Q}} \geq R \\ E_{\mathbb{Q}} \geq E}} E_{\text{md}, \mathbb{Q}} \leq E + \left(1 - \frac{E}{D}\right) J\left(\frac{R}{1-E/D}\right).$$

Consequently, as a result of the definition of $E_{\text{md}}(R, E)$ given in (13), we have

$$E_{\text{md}}(R, E) \leq E + \left(1 - \frac{E}{D}\right) J\left(\frac{R}{1-E/D}\right). \quad (177)$$

Thus, using (176) and (177), we can conclude that

$$E_{\text{md}}(R, E) = E + \left(1 - \frac{E}{D}\right) J\left(\frac{R}{1-E/D}\right). \quad (178)$$

In order to prove the concavity of $E_{\text{md}}(R, E)$ in (R, E) pair, let (R_a, E_a) and (R_b, E_b) be two pairs such that

$$R_a \in [0, C] \quad E_a \leq \left(1 - \frac{R_a}{C}\right)D \quad (179a)$$

$$R_b \in [0, C] \quad E_b \leq \left(1 - \frac{R_b}{C}\right)D. \quad (179b)$$

Then, for any $\alpha \in [0, 1]$, let R_α and E_α be

$$R_\alpha = \alpha R_a + (1 - \alpha)R_b \quad (180a)$$

$$E_\alpha = \alpha E_a + (1 - \alpha)E_b. \quad (180b)$$

From (179) and (180), we have

$$R_\alpha \in [0, C] \quad E_\alpha \leq \left(1 - \frac{R_\alpha}{C}\right)D. \quad (181)$$

Furthermore, using the concavity of $J(\cdot)$, we get

$$\begin{aligned} & \alpha E_{\text{md}}(R_a, E_a) + (1 - \alpha)E_{\text{md}}(R_b, E_b) \\ &= \alpha \left(E_a + \left(1 - \frac{E_a}{D}\right) J\left(\frac{R_a}{1-E_a/D}\right) \right) \\ & \quad + (1 - \alpha) \left(E_b + \left(1 - \frac{E_b}{D}\right) J\left(\frac{R_b}{1-E_b/D}\right) \right) \\ &= E_\alpha + \alpha \left(1 - \frac{E_a}{D}\right) J\left(\frac{R_a}{1-E_a/D}\right) \\ & \quad + (1 - \alpha) \left(1 - \frac{E_b}{D}\right) J\left(\frac{R_b}{1-E_b/D}\right) \\ & \leq E_\alpha + \left(1 - \frac{E_\alpha}{D}\right) J\left(\frac{\alpha R_a + (1-\alpha)R_b}{1-E_\alpha/D}\right) \\ &= E_{\text{md}}(R_\alpha, E_\alpha). \end{aligned} \quad (182)$$

Thus, $E_{\text{md}}(R, E)$ is jointly concave in rate exponent pairs. ■

J. Proof of Theorem 2

Proof: In Section IV-D, it is shown that for any positive integer ℓ , a rate-exponent vector (\vec{R}, \vec{E}) is achievable if there exists a time-sharing vector $\vec{\eta}$ such that

$$\begin{aligned} E_i & \leq \left(1 - \sum_{j=1}^{\ell} \eta_j\right)D \\ & \quad + \sum_{j=i+1}^{\ell} \eta_j J\left(\frac{R_j}{\eta_j}\right) \quad \forall i \in \{1, 2, \dots, \ell\} \end{aligned} \quad (183a)$$

$$R_i \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (183b)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (183c)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1. \quad (183d)$$

Thus, the existence of a time-sharing vector $\vec{\eta}$ satisfying (183) is a sufficient condition for the achievability of a rate-exponent vector (\vec{R}, \vec{E}) .

For any reliable code sequence \mathbb{Q} whose message sets are of the form $\mathcal{M}^{(\kappa)} = \mathcal{M}_1^{(\kappa)} \times \mathcal{M}_2^{(\kappa)} \times \dots \times \mathcal{M}_\ell^{(\kappa)}$, Lemma 7 with $\delta = \frac{-1}{\ln P_e}$ implies that there exists a sequence $\vec{\eta}_\kappa$ such that

$$(1 - \tilde{\epsilon}_{3,\kappa})E_{i,\kappa} - \tilde{\epsilon}_{5,\kappa} \leq (1 - \sum_{j=1}^{\ell} \eta_{j,\kappa})D + \sum_{j=i+1}^{\ell} \eta_{j,\kappa} J\left(\frac{(1 - \tilde{\epsilon}_{3,\kappa})R_{j,\kappa}}{\eta_{j,\kappa}}\right) \quad i = 1, 2, \dots, \ell \quad (184a)$$

$$(1 - \tilde{\epsilon}_{3,\kappa})R_{i,\kappa} - \tilde{\epsilon}_{4,\kappa} \mathbf{1}_{\{i=1\}} \leq C\eta_{i,\kappa} \quad i = 1, 2, \dots, \ell \quad (184b)$$

$$\eta_{i,\kappa} \geq 0 \quad i = 1, 2, \dots, \ell \quad (184c)$$

$$\sum_{j=1}^{\ell} \eta_{j,\kappa} \leq 1 \quad (184d)$$

where $R_{i,\kappa} = \frac{|\ln \mathcal{M}_i^{(\kappa)}|}{\mathbf{E}^{(\kappa)}[\mathbb{T}^{(\kappa)}]}$, $E_{i,\kappa} = \frac{-\ln P_e(i)^{(\kappa)}}{\mathbf{E}^{(\kappa)}[\mathbb{T}^{(\kappa)}]}$,
 $\tilde{\epsilon}_{3,\kappa} = \frac{P_e^{(\kappa)} + 1 - P_e^{(\kappa)} \ln P_e^{(\kappa)}}{-\ln P_e^{(\kappa)}}$, $\tilde{\epsilon}_{4,\kappa} = \frac{h(\tilde{\epsilon}_{3,\kappa})}{\mathbf{E}^{(\kappa)}[\mathbb{T}^{(\kappa)}]}$,
and $\tilde{\epsilon}_{5,\kappa} = \frac{h(\tilde{\epsilon}_{3,\kappa}) - \ln \lambda \delta}{\mathbf{E}^{(\kappa)}[\mathbb{T}^{(\kappa)}]}$.

Note that as a result of (184), all members of the sequence $\vec{\eta}_\kappa$ are from a compact metric space.³⁶ Thus, there exists a convergent subsequence, converging to a $\vec{\eta}$. Using (184), with definitions of $R_{\mathbb{Q},i}$ and $E_{\mathbb{Q},i}$ given in Definition 11, we can conclude that $\vec{\eta}$ satisfies

$$E_{\mathbb{Q},i} \leq (1 - \sum_{j=1}^{\ell} \eta_j) D + \sum_{j=i+1}^{\ell} \eta_j J\left(\frac{R_{\mathbb{Q},j}}{\eta_j}\right) \quad \forall i \in \{1, 2, \dots, \ell\} \quad (185a)$$

$$R_{\mathbb{Q},i} \leq C\eta_i \quad \forall i \in \{1, 2, \dots, \ell\} \quad (185b)$$

$$\eta_i \geq 0 \quad \forall i \in \{1, 2, \dots, \ell\} \quad (185c)$$

$$\sum_{j=1}^{\ell} \eta_j \leq 1. \quad (185d)$$

According to Definition 11 describing the *bit-wise* UEP problem, a rate-exponent vector (\vec{R}, \vec{E}) is achievable only if there exists a reliable code sequence \mathbb{Q} such that $(\vec{R}_{\mathbb{Q}}, \vec{E}_{\mathbb{Q}}) = (\vec{R}, \vec{E})$. Consequently, the existence of a time-sharing vector satisfying (183) is also a necessary condition for the achievability of a rate-exponent vector (\vec{R}, \vec{E}) .

Thus, we can conclude that a rate-exponent vector (\vec{R}, \vec{E}) is achievable if and only if there exists a $\vec{\eta}$ satisfying (183).

In order to prove the convexity of region of achievable rate-exponent vectors, let (\vec{R}_a, \vec{E}_a) and (\vec{R}_b, \vec{E}_b) be two achievable rate-exponent vectors. Then, there exist triples $(\vec{R}_a, \vec{E}_a, \vec{\eta}_a)$ and $(\vec{R}_b, \vec{E}_b, \vec{\eta}_b)$ satisfying (183).

For any $\alpha \in [0, 1]$, let \vec{R}_α , \vec{E}_α , and $\vec{\eta}_\alpha$ be

$$\vec{R}_\alpha = \alpha \vec{R}_a + (1 - \alpha) \vec{R}_b$$

$$\vec{E}_\alpha = \alpha \vec{E}_a + (1 - \alpha) \vec{E}_b$$

$$\vec{\eta}_\alpha = \alpha \vec{\eta}_a + (1 - \alpha) \vec{\eta}_b.$$

As $J(\cdot)$ is concave and the triples $(\vec{R}_a, \vec{E}_a, \vec{\eta}_a)$ and $(\vec{R}_b, \vec{E}_b, \vec{\eta}_b)$ satisfy the constraints given in (183), the triple $(\vec{R}_\alpha, \vec{E}_\alpha, \vec{\eta}_\alpha)$ also satisfies the constraints given in (183). Consequently, the rate-exponent vector $(\vec{R}_\alpha, \vec{E}_\alpha)$ is achievable and the region of achievable rate-exponent vectors is convex. ■

ACKNOWLEDGMENT

Authors would like to thank Gerhard Kramer and Reviewer II for pointing out the importance of the case of varying the number of groups of bits with expected block length. Reviewer II has also coined the term rates-exponents vector which replaced authors' previous attempt: rate vector error exponent vector pair. Authors would also like to thank all of the reviewers for their meticulous reviews and penetrating questions.

REFERENCES

- [1] P. Berlin, B. Nakiboğlu, B. Rimoldi, and E. Telatar, "A simple converse of Burnashev's reliability function," *IEEE Trans. Inf. Theory*, vol. 55, no. 7, pp. 3074–3080, Jul. 2009.
- [2] S. Borade, B. Nakiboğlu, and L. Zheng, "Unequal error protection: An information-theoretic perspective," *IEEE Trans. Inf. Theory*, vol. 55, no. 12, pp. 5511–5539, Dec. 2009.
- [3] M. V. Burnashev, "Data transmission over a discrete channel with feedback, random transmission time," *Problemy Perdachi Informatsii*, vol. 12, no. 4, pp. 10–30, 1976.
- [4] F. Chung and L. Lu, "Concentration inequalities and martingale inequalities: A survey," *Internet Math.*, vol. 3, no. 1, pp. 79–127, 2006.
- [5] I. Csiszár, "Joint source-channel error exponent," *Probl. Control Inf. Theory*, vol. 9, no. 5, pp. 315–328, 1980.
- [6] S. K. Gorantla, B. Nakiboğlu, T. P. Coleman, and L. Zheng, "Bit-wise unequal error protection for variable length blockcodes with feedback," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2010, pp. 241–245.
- [7] B. D. Kudryashov, "On message transmission over a discrete channel with noiseless feedback," *Problemy Perdachi Informatsii*, vol. 15, no. 1, pp. 3–13, 1973.
- [8] B. Nakiboğlu and L. Zheng, "Errors-and-erasures decoding for block codes with feedback," *IEEE Trans. Inf. Theory*, vol. 58, no. 1, pp. 24–49, Jan. 2012.
- [9] B. Nazer, Y. Shkel, and S. C. Draper, "The AWGN red alert problem 2012," arXiv:1102.4411[cs.IT]. DOI:10.1109/TIT.2012.2235120.
- [10] A. N. Shiriaev, *Probability*. New York: Springer-Verlag, 1996.
- [11] D. Wang, V. Chandar, S. Y. Chung, and G. W. Wornell, "On reliability functions for single-message unequal error protection," in *Proc. IEEE Int. Symp. Inf. Theory*, Jul. 2012, pp. 2934–2938.
- [12] H. Yamamoto and K. Itoh, "Asymptotic performance of a modified Schalkwijk-Barron scheme for channels with noiseless feedback," *IEEE Trans. Inf. Theory*, vol. 25, no. 6, pp. 729–733, Nov. 1979.

Barış Nakiboğlu received the B.S. degrees in electrical and electronics engineering and in physics from Middle East Technical University (METU), Ankara, Turkey, in 2002 and the M.S. degree in electrical engineering and computer science and the Ph.D. degree in electrical engineering from the Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, in 2005 and 2011, respectively. He is now a postdoctoral researcher at University of California Berkeley, working with Anant Sahai. He is also a research affiliate at Laboratory for Information and Decision Systems (LIDS) at MIT working with Sanjoy Mitter.

Siva K. Gorantla received the B.Tech. degree in electrical engineering from the Indian Institute of Technology, Madras, in 2007. He received the M.S. and Ph.D. degrees in electrical and computer engineering and M.S. in Statistics degree from the University of Illinois at Urbana-Champaign, Urbana, in 2009, 2012 and 2012 respectively. Since June 2012, he has been working as a Research Scientist at Adchemy Inc, CA. His research interests include stochastic control, information theory, statistical learning and prediction.

Mr. Gorantla has been awarded the James Henderson fellowship at UIUC and S Subramanian Award (Institute Merit Prize) from IIT Madras in the years 2008 and 2004.

³⁶Let the metric be $\|\vec{\eta} - \vec{\nu}\| = \max_j |\eta_j - \nu_j|$.

Lizhong Zheng received the B.S. and M.S. degrees, in 1994 and 1997 respectively, from the Department of Electronic Engineering, Tsinghua University, China, and the Ph.D. degree, in 2002, from the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley. Since 2002, he has been working in the Department of Electrical Engineering and Computer Sciences, where he is currently an associate professor. His research interests include information theory, wireless communications and wireless networks. He received Eli Jury award from UC Berkeley in 2002, IEEE Information Theory Society Paper Award in 2003, and NSF CAREER award in 2004, and the AFOSR Young Investigator Award in 2007.

Todd P. Coleman (S'01–M'05–SM'11) received the B.S. degrees in electrical engineering (*summa cum laude*) as well as computer engineering (*summa cum laude*) from the University of Michigan, Ann Arbor, in 2000. He received the M.S. and Ph.D. degrees in electrical engineering from the Massachusetts Institute of Technology (MIT), Cambridge, in 2002, and 2005. He was a postdoctoral scholar in neuroscience at MIT during the 2005–2006 academic year. He was an Assistant Professor in ECE and Neuroscience at the University of Illinois from 2006–2011.

Dr. Coleman is currently an Associate Professor in Bioengineering and director of the Neural Interaction Laboratory at UCSD, where his group builds flexible bio-electronics for neurological monitoring applications. His research is highly inter-disciplinary and lies at the intersection of bio-electronics, neuroscience, medicine, and applied mathematics. Dr. Coleman is a science advisor for the Science & Entertainment Exchange (National Academy of Sciences).