

Minimax Universal Decoding With an Erasure Option

Neri Merhav, *Fellow, IEEE*, and Meir Feder, *Fellow, IEEE*

Abstract—Motivated by applications of rateless coding, decision feedback, and automatic repeat request (ARQ), we study the problem of universal decoding for unknown channels in the presence of an erasure option. Specifically, we harness the competitive minimax methodology developed in earlier studies, in order to derive a universal version of Forney's classical erasure/list decoder, which in the erasure case, optimally trades off between the probability of erasure and the probability of undetected error. The proposed universal erasure decoder guarantees universal achievability of a certain fraction ξ of the optimum error exponents of these probabilities (in a sense to be made precise in the sequel). A single-letter expression for ξ , which depends solely on the coding rate and the Neyman–Pearson threshold (to be defined), is provided. The example of the binary-symmetric channel is studied in full detail, and some conclusions are drawn.

Index Terms—Channel uncertainty, competitive minimax, erasure, error exponent, generalized likelihood ratio test (GLRT), rateless codes, universal decoding.

I. INTRODUCTION

WHEN communicating across an unknown channel, classical channel coding at any fixed rate, however small, is inherently problematic since this fixed rate might be larger than the unknown capacity of the underlying channel. It makes sense then to try to adapt the coding rate to the channel conditions, which can be learned online at the transmitter whenever a feedback link, from the receiver to the transmitter, is available.

One of the recent promising approaches to this end is rateless coding proposed in [17], [18] (see also [5]–[7], [14], [20], and references therein). Independently, rateless codes were also proposed in a networking scenario for the packet erasure channel [2], [3], [15], where they have been referred to as *fountain codes*. Fountain codes also have a low-density structure that allows computationally efficient decoding. In rateless coding, there is a fixed number M of messages, each one being represented by a codeword of unlimited length, in principle. A possible receiver for a rateless code examines, after each symbol has been received, whether it can decode the message, with “reasonably good confidence,” or alternatively, to request, via the feedback link, an additional symbol.¹ Upon receiving the new channel

output, again, the receiver either makes a decision, or requests another symbol from the transmitter, and so on. The coding rate is then defined by $\log M$ divided by the expected number of symbols transmitted before the decoder makes a decision. Clearly, at every time instant, the receiver of a rateless communication system operates just like an *erasure decoder* [10],² which partitions the space of channel output vectors into $(M + 1)$ regions, M for each one of the possible messages, and an additional region for “erasure,” which, in the rateless regime, is used for requesting an additional symbol. Keeping the erasure probability small is then motivated by the desire to keep the expected transmission time, for each message, small. Although these two criteria are not completely equivalent, they are strongly related.

When the channel is unknown at the decoder, it was suggested in some of the quoted references to use a universal decoder, which is inspired by the maximum mutual information (MMI) decoder [4]: by using a certain threshold, the receiver decides whether to make a decision or ask for another symbol. While this approach works fairly well, there is no evidence of optimality.

These observations, as well as techniques such as automatic repeat request (ARQ) and decision feedback, motivate us to study the problem in a more systematic manner. Specifically, we consider the problem of universal decoding with an erasure option, for the class of discrete memoryless channels (DMCs) indexed by an unknown parameter vector θ (e.g., the set of channel transition probabilities). We harness the competitive minimax methodology proposed in [9], in order to derive a universal version of Forney's classical erasure/list decoder. For a given DMC with parameter θ , a given coding rate R , and a given threshold parameter T (all to be formally defined later), Forney's erasure/list decoder optimally trades off between the exponent $E_1(R, T, \theta)$ of the probability of the erasure event, \mathcal{E}_1 , and the exponent, $E_2(R, T, \theta) = E_1(R, T, \theta) + T$, of the probability of undetected error event, \mathcal{E}_2 , in the random coding regime.

The universal erasure decoder, proposed in this paper, guarantees universal achievability of an erasure exponent $\hat{E}_1(R, T, \theta)$, which is at least as large as $\xi \cdot E_1(R, T, \theta)$ for all θ , for some constant $\xi \in [0, 1]$, that is independent of θ (but does depend on R and T), and at the same time, an undetected error exponent $\hat{E}_2(R, T, \theta) \geq \xi \cdot E_1(R, T, \theta) + T$ for all θ (in the random coding sense). At the very least, this guarantees that whenever the probabilities of \mathcal{E}_1 and \mathcal{E}_2 decay exponentially for a known channel, so they do even when the channel is unknown, using the proposed universal decoder. The question is, of course: what is the largest value of ξ for which the preceding statement holds? We partially answer this question by deriving a single-letter expression for a lower bound to the largest value of ξ , denoted henceforth by $\xi^*(R, T)$, that is guaranteed to be attainable by this decoder. While $\xi^*(R, T)$ is only a lower bound to the universally achievable fraction of the error exponent, for $T = 0$

Manuscript received April 15, 2006; revised February 19, 2007. This work was supported by the Israel Science Foundation under Grant 223/05. The material in this paper will be presented at the IEEE International Symposium on Information Theory, Nice, France, June 2007.

N. Merhav is with the Department of Electrical Engineering, Technion—Israel Institute of Technology, Technion City, Haifa 32000, Israel (e-mail: merhav@ee.technion.ac.il).

M. Feder is with the Department of Electrical Engineering—Systems, Tel-Aviv University, Ramat-Aviv 69978, Israel (e-mail: meir@eng.tau.ac.il).

Communicated by G. Kramer, Associate Editor for Shannon Theory.

Digital Object Identifier 10.1109/TIT.2007.894695

¹Alternatively, the receiver can use the feedback link only to notify the transmitter when it reached a decision regarding the current message (and keep silent at all other times). In network situations, this would not load the network much as it is done only once per each message.

²See also [21], [1], [13], [12] and references therein for later studies.

(i.e., essentially “no erasure”) and for the BSCs it equals unity, the optimal true value. But for $T > 0$, $\xi^*(R, T)$ may, in general, be less than unity (as we show in some examples). If we conjecture that the true universally achievable fraction of the error exponent is also less than unity in general, then it means that there is a major difference between ordinary universal decoding and universal erasure decoding: While for the former, it is well known that optimum³ random coding error exponents are fully universally achievable (at least for some classes of channels and certain random coding distributions [4], [22], [8]), in the latter, when the erasure option is available, this may no longer be the case, in general. Explicit results, including numerical values of $\xi^*(R, T)$, are derived for the example of the binary-symmetric channel (BSC), parameterized by the crossover probability θ , and some conclusions are drawn.

The outline of the paper is as follows. In Section II, we establish the notation conventions and we briefly review some known results about erasure decoding. In Section III, we formulate the problem of universal decoding with erasures. In Section IV, we present the proposed universal erasure decoder and prove its asymptotic optimality in the competitive minimax sense. In Section V, we present the main results concerning the performance of the proposed universal decoder. Section VI is devoted to the example of the BSC. Finally, in Section VII, we summarize our conclusions.

II. NOTATION AND PRELIMINARIES

Throughout this paper, scalar random variables (RVs) will be denoted by capital letters, their sample values will be denoted by the respective lower case letters, and their alphabets will be denoted by the respective calligraphic letters. A similar convention will apply to random vectors of dimension n and their sample values, which will be denoted with same symbols in the bold face font. The set of all n -vectors with components taking values in a certain alphabet, will be denoted as the same alphabet superscripted by n . Thus, for example, a random vector $\mathbf{X} = (X_1, \dots, X_n)$ may assume a specific vector value $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$ as each component takes values in \mathcal{X} . Channels will be denoted generically by the letter P , or P_θ , when we wish to emphasize that the channel is indexed or parametrized by a certain scalar or vector θ , taking on values in some set Θ . Information-theoretic quantities, such as entropies and conditional entropies, will be denoted following the usual conventions of the information-theory literature, e.g., $H(X)$, $H(X|Y)$, and so on. With a slight abuse of notation, when we wish to emphasize the dependence of the entropy on the underlying probability distribution P , we denote it by $H(P)$. The cardinality of a finite set \mathcal{A} will be denoted by $|\mathcal{A}|$.

Consider a DMC with a finite input alphabet \mathcal{X} , finite output alphabet \mathcal{Y} , and single-letter transition probabilities $\{P(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$. As the channel is fed by an input vector $\mathbf{x} \in \mathcal{X}^n$, it generates an output vector $\mathbf{y} \in \mathcal{Y}^n$ according to the sequence conditional probability distributions (cf. [16])

$$P(y_i|x_1, \dots, x_i, y_1, \dots, y_{i-1}) = P(y_i|x_i) \quad (1)$$

³Optimum exponents—corresponding to optimum maximum-likelihood (ML) decoding.

for $i = 1, 2, \dots, n$, where for $i = 1$, (y_1, \dots, y_{i-1}) is understood as the null string. A rate- R block code of length n consists of $M = e^{nR}$ n -vectors $\mathbf{x}_m \in \mathcal{X}^n$, $m = 1, 2, \dots, M$, which represent M different messages. We will assume that all possible messages are *a priori* equiprobable, i.e., $P(m) = 1/M$ for all $m = 1, 2, \dots, M$.

A decoder with an erasure option is a partition of \mathcal{Y}^n into $(M + 1)$ regions, $\mathcal{R}_0, \mathcal{R}_1, \dots, \mathcal{R}_M$. Such a decoder works as follows: If \mathbf{y} falls into \mathcal{R}_m , $m = 1, 2, \dots, M$, then a decision is made in favor of message number m . If $\mathbf{y} \in \mathcal{R}_0$, no decision is made and an erasure is declared. We will refer to \mathcal{R}_0 as the *erasure event*.

Given a code $\mathcal{C} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ and a decoder $\mathcal{R} = (\mathcal{R}_0, \mathcal{R}_1, \dots, \mathcal{R}_M)$, let us now define two additional undesired events. The event \mathcal{E}_1 is the event of not making the right decision. This event is the disjoint union of the erasure event and the event \mathcal{E}_2 , which is the *undetected error* event, namely, the event of making the wrong decision. The probabilities of all three events are defined as follows:

$$\Pr\{\mathcal{E}_1\} = \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{R}_m^c} P(\mathbf{x}_m, \mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{R}_m^c} P(\mathbf{y}|\mathbf{x}_m) \quad (2)$$

$$\begin{aligned} \Pr\{\mathcal{E}_2\} &= \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} P(\mathbf{x}_{m'}, \mathbf{y}) \\ &= \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} P(\mathbf{y}|\mathbf{x}_{m'}) \end{aligned} \quad (3)$$

$$\Pr\{\mathcal{R}_0\} = \Pr\{\mathcal{E}_1\} - \Pr\{\mathcal{E}_2\}. \quad (4)$$

Forney [10] assumes that the DMC is known to the decoder, and shows, using the Neyman–Pearson methodology, that the best tradeoff between $\Pr\{\mathcal{E}_1\}$ and $\Pr\{\mathcal{E}_2\}$ (or, equivalently, between $\Pr\{\mathcal{R}_0\}$ and $\Pr\{\mathcal{E}_2\}$) is attained by the decoder $\mathcal{R}^* = (\mathcal{R}_0^*, \mathcal{R}_1^*, \dots, \mathcal{R}_M^*)$ defined by

$$\begin{aligned} \mathcal{R}_m^* &= \left\{ \mathbf{y} : \frac{P(\mathbf{y}|\mathbf{x}_m)}{\sum_{m' \neq m} P(\mathbf{y}|\mathbf{x}_{m'})} \geq e^{nT} \right\}, \quad m = 1, 2, \dots, M \\ \mathcal{R}_0^* &= \bigcap_{m=1}^M (\mathcal{R}_m^*)^c \end{aligned} \quad (5)$$

where $(\mathcal{R}_m^*)^c$ is the complement of \mathcal{R}_m^* , and where $T \geq 0$ is a parameter, henceforth referred to as the *threshold*, which controls the balance between the probabilities of \mathcal{E}_1 and \mathcal{E}_2 .

Forney devotes the remaining part of his paper [10] to derive lower bounds to the random coding exponents (associated with \mathcal{R}^*), $E_1(R, T)$ and $E_2(R, T)$, of $\Pr\{\mathcal{E}_1\}$ and $\Pr\{\mathcal{E}_2\}$, the average⁴ probabilities of \mathcal{E}_1 and \mathcal{E}_2 , respectively, and to investigate their properties. Specifically, Forney shows, among other things, that for the ensemble of randomly chosen codes, where each codeword is chosen independently under an independent and identically distributed (i.i.d.) distribution $Q_n(\mathbf{x}) = \prod_{i=1}^n Q(x_i)$

$$E_1(R, T) = \max_{0 \leq s \leq \rho \leq 1} \max_Q [E_0(s, \rho, Q) - \rho R - sT] \quad (6)$$

⁴Here, “average” means with respect to (w.r.t.) the ensemble of randomly selected codes.

where

$$E_0(s, \rho, Q) = -\ln \left[\sum_{y \in \mathcal{Y}} \left(\sum_{x \in \mathcal{X}} Q(x) P^{1-s}(y|x) \right) \times \left(\sum_{x' \in \mathcal{X}} Q(x') P^{s/\rho}(y|x') \right)^\rho \right] \quad (7)$$

and

$$E_2(R, T) = E_1(R, T) + T. \quad (8)$$

A simple observation that we will need, before passing to the case of an unknown channel, is that the same decision rule \mathcal{R}^* would be obtained if rather than adopting the Neyman–Pearson approach, one would consider a Lagrange function

$$\Gamma(\mathcal{C}, \mathcal{R}) \triangleq \Pr\{\mathcal{E}_2\} + e^{-nT} \Pr\{\mathcal{E}_1\} \quad (9)$$

for a given code $\mathcal{C} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ and a given threshold T , as the figure of merit, and seek a decoder \mathcal{R} that minimizes it. To see that this is equivalent, let us rewrite $\Gamma(\mathcal{C}, \mathcal{R})$ as follows:

$$\Gamma(\mathcal{C}, \mathcal{R}) = \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_{m,m'} \neq m} P(\mathbf{y}|\mathbf{x}_{m'}) + \sum_{\mathbf{y} \in \mathcal{R}_m^c} e^{-nT} P(\mathbf{y}|\mathbf{x}_m) \right] \quad (10)$$

and it is now clear that for each m , the bracketed expression (which has the form of weighted error of a binary hypothesis testing problem) is minimized by \mathcal{R}_m^* as defined above. Since this decision rule is identical to Forney's, it is easy to see that the resulting exponential decay of the ensemble average

$$\mathbf{E}\{\Gamma(\mathcal{C}, \mathcal{R}^*)\} = \overline{\Pr}\{\mathcal{E}_2\} + e^{-nT} \overline{\Pr}\{\mathcal{E}_1\}$$

is $E_2(R, T)$, as $\overline{\Pr}\{\mathcal{E}_1\}$ decays according to $e^{-nE_1(R, T)}$, $\overline{\Pr}\{\mathcal{E}_2\}$ decays according to $e^{-nE_2(R, T)}$, and $E_2(R, T) = E_1(R, T) + T$, as mentioned earlier. This Lagrangian approach will be more convenient to work with, when we next move on to the case of an unknown DMC, because it allows us to work with one figure of merit instead of a tradeoff between two.

III. UNKNOWN CHANNEL—PROBLEM DESCRIPTION

We now move on to the case of an unknown channel. While our techniques can be applied to quite general classes of channels, here, for the sake of concreteness and conceptual simplicity, and following in [10], we confine attention to DMCs. Consider then a family of DMCs

$$\{P_\theta(y|x), x \in \mathcal{X}, y \in \mathcal{Y}, \theta \in \Theta\}$$

where θ is the parameter, or the index of the channel in the class, taking values in some set Θ . For example, θ may be a positive integer, denoting the index of the channel within a finite or a countable index set. As another example, θ may simply represent the set of all $|\mathcal{X}| \cdot (|\mathcal{Y}| - 1)$ single-letter transition probabilities that define the DMC, and if there are some symmetries (like in the BSC), these reduce the dimensionality of θ . The basic questions are now the following.

1. How to devise a good erasure decoder when the underlying channel is known to belong to the class $\{P_\theta(y|x), x \in \mathcal{X}, y \in \mathcal{Y}, \theta \in \Theta\}$, but θ is unknown?
2. What are the resulting error exponents of \mathcal{E}_1 and \mathcal{E}_2 and how do they compare to Forney's exponents for known θ ?

In the quest for universal schemes for decoding with an erasure option, two difficulties⁵ are encountered in light of [10]. The first difficulty is that here we have two figures of merits, the probabilities of \mathcal{E}_1 and \mathcal{E}_2 . But this difficulty can be alleviated by adopting the Lagrangian approach, described at the end of the previous section. The second difficulty is somewhat deeper: Classical derivations of universal decoding rules for ordinary decoding (without erasures) over the class of DMCs, like the MMI decoder [4] and its variants, were based on ideas that are deeply rooted in considerations of joint typicality between the channel output \mathbf{y} and each hypothesized codeword \mathbf{x}_m . These considerations were easy to apply in ordinary decoding, where the score function (or, the “metric”) associated with the optimum maximum likelihood (ML) decoding, $\log P_\theta(\mathbf{y}|\mathbf{x}_m)$, involves only *one* codeword at a time, and that this function depends on \mathbf{x}_m and \mathbf{y} only via their joint empirical distribution, or, in other words, their joint type. Moreover, in the case of decoding without erasures, given the true transmitted codeword \mathbf{x}_m and the resulting channel output \mathbf{y} , the scores associated with all other randomly chosen codewords are independent of each other, a fact that facilitates the analysis to a great extent. This is very different from the situation in erasure decoding, where Forney's optimum score function for each codeword

$$\frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{\sum_{m' \neq m} P_\theta(\mathbf{y}|\mathbf{x}_{m'})}$$

depends on *all* codewords at the same time. Consequently, in a random coding analysis, it is rather complicated to apply joint typicality considerations, or to analyze the statistical behavior of this expression, let alone the statistical dependency between the score functions associated with the various codewords.

This difficulty is avoided if the competitive minimax methodology, proposed and developed in [9], is applied. Specifically, let $\Gamma_\theta(\mathcal{C}, \mathcal{R})$ denote the above defined Lagrangian, where we now emphasize the dependence on the index of the channel θ . Let us also define $\bar{\Gamma}_\theta^* = \mathbf{E}\{\min_{\mathcal{R}} \Gamma_\theta(\mathcal{C}, \mathcal{R})\}$, i.e., the ensemble average of the minimum of the above Lagrangian (achieved by Forney's optimum decision rule) w.r.t. the channel $\{P_\theta(y|x)\}$ for a given θ . Note that the exponential order of $\bar{\Gamma}_\theta^*$

⁵These difficulties may also be related to the observation discussed in the Introduction, that optimum error exponents may not be universally achievable in the erasure decoding setting.

is $e^{-n[E_1(R,T,\theta)+T]} = e^{-nE_2(R,T,\theta)}$, where $E_1(R,T,\theta)$ and $E_2(R,T,\theta)$ are the new notations for $E_1(R,T)$ and $E_2(R,T)$, respectively, with the dependence on the channel index θ , made explicit. In principle, we would have been interested in a decision rule \mathcal{R} that achieves

$$\min_{\mathcal{R}} \max_{\theta \in \Theta} \frac{\Gamma_{\theta}(\mathcal{C}, \mathcal{R})}{\bar{\Gamma}_{\theta}^*}, \quad (11)$$

or equivalently

$$\min_{\mathcal{R}} \max_{\theta \in \Theta} \frac{\Gamma_{\theta}(\mathcal{C}, \mathcal{R})}{e^{-n[E_1(R,T,\theta)+T]}}. \quad (12)$$

However, as is discussed in [9] (in the analogous context of ordinary decoding, without erasures), such an ambitious minimax criterion of competing with the optimum performance may be too optimistic: If $[E_1(R,T,\theta) + T]$ is not universally achievable, then the value of the above minimax may grow exponentially with n , and then there might be values of θ for which the numerator does not tend to zero at all, whereas the denominator still does. A better approach would be to compete with a similar expression of the exponential behavior, but where the term $E_1(R,T,\theta)$ is being multiplied by a constant $\xi \in (0, 1]$, which we would like to choose as large as possible. In other words, we are interested in the competitive minimax criterion

$$K_n(\mathcal{C}) \triangleq \min_{\mathcal{R}} \max_{\theta \in \Theta} \frac{\Gamma_{\theta}(\mathcal{C}, \mathcal{R})}{e^{-n[\xi E_1(R,T,\theta)+T]}}. \quad (13)$$

Similarly as in [9], we wish to find the largest value of ξ such that the ensemble average $\bar{K}_n \triangleq \mathbf{E}\{K_n(\mathcal{C})\}$ would *not* grow exponentially fast, i.e.,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \bar{K}_n \leq 0. \quad (14)$$

The rationale behind this is the following: If \bar{K}_n is subexponential in n , for some ξ , then this guarantees that there exists a code $\hat{\mathcal{C}}$ and a universal erasure decoder $\hat{\mathcal{R}}$, such that for every $\theta \in \Theta$, the exponential order of $\Gamma_{\theta}(\hat{\mathcal{C}}, \hat{\mathcal{R}})$ is no worse than $e^{-n[\xi E_1(R,T,\theta)+T]}$. This, in turn, implies that *both* terms of $\Gamma_{\theta}(\hat{\mathcal{C}}, \hat{\mathcal{R}})$ decay at least as $e^{-n[\xi E_1(R,T,\theta)+T]}$, which means that for the decoder $\hat{\mathcal{R}}$, the exponent of $\bar{\Pr}\{\mathcal{E}_1\}$ is at least $\xi \cdot E_1(R,T,\theta)$ and the exponent of $\bar{\Pr}\{\mathcal{E}_2\}$ is at least $\xi \cdot E_1(R,T,\theta) + T$, both for every $\theta \in \Theta$. Thus, the difference between the two (guaranteed) exponents remains T as before (as the weight of the term $\bar{\Pr}\{\mathcal{E}_1\}$ in $\Gamma(\mathcal{R}, \mathcal{C})$ is e^{-nT}), but the other term, $E_1(R,T,\theta)$, is now scaled by a factor of ξ .

The remaining parts of this paper focus on deriving a universal decoding rule that asymptotically achieves $\bar{K}_n(\mathcal{C})$ for a given ξ , and on analyzing its performance, i.e., finding the maximum value of ξ such that \bar{K}_n still grows subexponentially rapidly.

IV. DERIVATION OF A UNIVERSAL ERASURE DECODER

For a given $\xi \in (0, 1]$, let us define

$$f(\mathbf{x}_m, \mathbf{y}) \triangleq \max_{\theta \in \Theta} \left\{ e^{n[\xi E_1(R,T,\theta)+T]} P_{\theta}(\mathbf{y}|\mathbf{x}_m) \right\} \quad (15)$$

and consider the decoder $\hat{\mathcal{R}}$ whose decision regions are

$$\hat{\mathcal{R}}_m = \left\{ \mathbf{y} : \frac{f(\mathbf{x}_m, \mathbf{y})}{\sum_{m' \neq m} f(\mathbf{x}_{m'}, \mathbf{y})} \geq e^{nT} \right\}, \quad m = 1, 2, \dots, M$$

$$\hat{\mathcal{R}}_0 = \bigcap_{m=1}^M \hat{\mathcal{R}}_m^c. \quad (16)$$

Note that this can be thought of as an extension of a decoder in the spirit of the generalized-likelihood ratio test (GLRT), where the unknown parameter θ is estimated by the ML estimator for each term $P_{\theta}(\mathbf{y}|\mathbf{x}_i)$ individually. While this GLRT-like decoder is a special case of the above, corresponding to $\xi = 0$, the more general decoder, proposed here, assigns higher weights to good channels, as discussed in [9]. Denoting

$$K_n(\mathcal{C}, \mathcal{R}) = \max_{\theta \in \Theta} \frac{\Gamma_{\theta}(\mathcal{C}, \mathcal{R})}{e^{-n[\xi E_1(R,T,\theta)+T]}} \quad (17)$$

for a given encoder $\mathcal{C} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ and decoder \mathcal{R} , our first main result establishes the asymptotic optimality of $\hat{\mathcal{R}}$ in the competitive minimax sense, namely, that $K_n(\mathcal{C}, \hat{\mathcal{R}})$ is within a subexponential factor as small as $K_n(\mathcal{C}) = \min_{\mathcal{R}} K_n(\mathcal{C}, \mathcal{R})$, and therefore, $\mathbf{E}\{K_n(\mathcal{C}, \hat{\mathcal{R}})\}$ is within the same subexponential factor as small as $\bar{K}_n = \mathbf{E}\{K_n(\mathcal{C})\}$.

Theorem 1: For every code \mathcal{C}

$$K_n(\mathcal{C}, \hat{\mathcal{R}}) \leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} K_n(\mathcal{C}). \quad (18)$$

Proof: The result and the proof technique is similar to those of [9]. As \mathbf{x} and \mathbf{y} exhaust their spaces, \mathcal{X}^n and \mathcal{Y}^n , let Θ_n denote set of values of θ that achieve $\{f(\mathbf{x}, \mathbf{y}), \mathbf{x} \in \mathcal{X}^n, \mathbf{y} \in \mathcal{Y}^n\}$. Observe that for every θ , the expression $[e^{n[\xi E_1(R,T,\theta)+T]} P_{\theta}(\mathbf{y}|\mathbf{x})]$ depends on (\mathbf{x}, \mathbf{y}) only via their joint empirical distribution (or, the joint type). Consequently, the value of θ that achieves $f(\mathbf{x}, \mathbf{y})$ also depends on (\mathbf{x}, \mathbf{y}) only via their joint empirical distribution. Since the number of joint empirical distributions of (\mathbf{x}, \mathbf{y}) never exceeds $(n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1}$ (see [4]), then obviously

$$|\Theta_n| \leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \quad (19)$$

as well. Now, for every encoder \mathcal{C} and decoder \mathcal{R} , we get (20) at the top of the following page. Thus, we have defined $\hat{K}_n(\mathcal{C}, \mathcal{R})$ and sandwiched it between $K_n(\mathcal{C}, \mathcal{R})$ and $(n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \cdot K_n(\mathcal{R}, \mathcal{C})$ uniformly for every \mathcal{C} and \mathcal{R} . Now, obviously, $\hat{\mathcal{R}}$ minimizes $\hat{K}_n(\mathcal{C}, \mathcal{R})$, and so, for every \mathcal{R}

$$K_n(\mathcal{C}, \hat{\mathcal{R}}) \leq \hat{K}_n(\mathcal{C}, \hat{\mathcal{R}}) \leq \hat{K}_n(\mathcal{C}, \mathcal{R}) \leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \cdot K_n(\mathcal{C}, \mathcal{R}) \quad (21)$$

where the first and the third inequalities were just proved in the chain of inequalities (20), and the second inequality follows from the optimality of $\hat{\mathcal{R}}$ w.r.t. $\hat{K}_n(\mathcal{C}, \mathcal{R})$. Since we have shown that

$$K_n(\mathcal{C}, \hat{\mathcal{R}}) \leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \cdot K_n(\mathcal{C}, \mathcal{R})$$

for every \mathcal{R} , we can now minimize the right-hand side (RHS) w.r.t. \mathcal{R} and the assertion of Theorem 1 is obtained. This completes the proof of Theorem 1.

$$\begin{aligned}
K_n(\mathcal{C}, \mathcal{R}) &= \max_{\theta \in \Theta} \frac{\Gamma_\theta(\mathcal{C}, \mathcal{R})}{e^{-n[\xi E_1(R, T, \theta) + T]}} \\
&= \max_{\theta \in \Theta} \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&\leq \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \max_{\theta \in \Theta} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \max_{\theta \in \Theta} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&= \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} f(\mathbf{x}_{m'}, \mathbf{y}) + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} f(\mathbf{x}_m, \mathbf{y}) \right] \\
&\triangleq \hat{K}_n(\mathcal{C}, \mathcal{R}) \\
&= \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \max_{\theta \in \Theta_n} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \max_{\theta \in \Theta_n} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&\leq \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \left(\sum_{\theta \in \Theta_n} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right) + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \left(\sum_{\theta \in \Theta_n} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right) \right] \\
&= \sum_{\theta \in \Theta_n} \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&\leq |\Theta_n| \cdot \max_{\theta \in \Theta_n} \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&\leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \cdot \max_{\theta \in \Theta} \frac{1}{M} \sum_{m=1}^M \left[\sum_{\mathbf{y} \in \mathcal{R}_m} \sum_{m' \neq m} \frac{P_\theta(\mathbf{y}|\mathbf{x}_{m'})}{e^{-n[\xi E_1(R, T, \theta) + T]}} + e^{-nT} \sum_{\mathbf{y} \in \mathcal{R}_m^c} \frac{P_\theta(\mathbf{y}|\mathbf{x}_m)}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \\
&\leq (n+1)^{|\mathcal{X}| \cdot |\mathcal{Y}| - 1} \cdot K_n(\mathcal{C}, \mathcal{R}). \tag{20}
\end{aligned}$$

V. PERFORMANCE

In this section, we present an upper bound to \bar{K}_n from which we derive a lower bound to ξ^* , the largest value of ξ for which \bar{K}_n is subexponential in n .

We begin with a few definitions. The empirical distribution $\hat{P}_{\mathbf{x}}$ of \mathbf{x} is the vector of relative frequencies $\{\hat{P}_{\mathbf{x}}(a) = n_{\mathbf{x}}(a)/n, a \in \mathcal{X}\}$, $n_{\mathbf{x}}(a)$ being the number of occurrences of $a \in \mathcal{X}$ within $\mathbf{x} \in \mathcal{X}^n$. The type class $T_{\mathbf{x}}$ of $\mathbf{x} \in \mathcal{X}^n$ is the set of all $\mathbf{x}' \in \mathcal{X}^n$ such that $\hat{P}_{\mathbf{x}'} = \hat{P}_{\mathbf{x}}$. We next define the class \mathcal{Q} of the sequences of random coding distributions $\{Q_n\}$ that we assume. For every positive integer n , let Q_n be a random coding distribution of the following form:

$$Q_n(\mathbf{x}) = \frac{Q_n(T_{\mathbf{x}})}{|T_{\mathbf{x}}|} \tag{22}$$

where, of course, $\sum_{T_{\mathbf{x}}} Q_n(T_{\mathbf{x}}) = 1$. Let

$$\Delta_n(P_{\mathbf{x}}) = -\frac{1}{n} \ln Q_n(T_{\mathbf{x}})$$

and let $\Delta_n^*(P)$ be an extension of the function $\Delta_n^*(P_{\mathbf{x}})$ that is defined over the continuum of probability distributions over \mathcal{X} (rather than just the set of rational probability distributions with denominator n). A sequence of random coding distributions $\{Q_n\}_{n \geq 1}$ is said to belong to the class \mathcal{Q} if there exists

such an extension $\Delta_n^*(P)$ that converges, as $n \rightarrow \infty$, to a certain nonnegative functional $\Delta^*(P)$, uniformly over all probability distributions $\{P\}$ over \mathcal{X} .

It is easy to see that the class \mathcal{Q} essentially covers all random coding distributions that are customarily used (and much more). In particular, to approximate a random coding distribution which is uniform within a small neighborhood of one type class—corresponding to a probability distribution P_0 , and which vanishes elsewhere, we set $\Delta^*(P) = 0$ for every P in that neighborhood of P_0 , and $\Delta^*(P) = \infty$ elsewhere. For the case where Q is i.i.d.,

$$\Delta^*(P) = D(P||Q) = \sum_{a \in \mathcal{X}} P(a) \ln [P(a)/Q(a)].$$

the Kullback–Leibler divergence between Q and P . In particular, if $Q(\mathbf{x}) = 1/|\mathcal{X}|^n$ for all $\mathbf{x} \in \mathcal{X}^n$, then $\Delta^*(P) = \ln |\mathcal{X}| - H(P)$, $H(P)$ being the entropy associated with the distribution P .

Given a distribution P_y on \mathcal{Y} , a positive real λ , and a value of θ , let

$$F(P_y, \lambda, \theta) \triangleq \min_{P_{x|y}} \left[I(X; Y) + \Delta^* \left(\sum_{b \in \mathcal{Y}} P_y(b) P_{x|y}(\cdot|b) \right) - \lambda \mathbf{E} \ln P_\theta(Y|X) \right] \tag{23}$$

where $\mathbf{E}\{\cdot\}$ is the expectation and $I(X;Y)$ is the mutual information w.r.t. a generic joint distribution $P_{xy}(x,y) = P_y(y)P_{x|y}(x|y)$ of the RVs (X,Y) . Next, for a pair $(\theta, \tilde{\theta}) \in \Theta^2$, and for two real numbers s and ρ , $0 \leq s \leq \rho \leq 1$, define

$$E(\theta, \tilde{\theta}, \rho, s) = \min_{P_y} [F(P_y, 1-s, \theta) + \rho F(P_y, s/\rho, \tilde{\theta}) - H(Y)] \quad (24)$$

where $H(Y)$ is the entropy of Y induced by P_y . Finally, let

$$\xi^*(R, T) \triangleq \min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} \frac{E(\theta, \tilde{\theta}, s, \rho) - \rho R - sT}{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})} \quad (25)$$

with the convention that if the denominator vanishes, then $\xi^*(R, T) \triangleq 1$. Our main result, in this section is the following theorem.

Theorem 2: Consider a sequence of ensemble of codes where each codeword is drawn independently, under a distribution Q_n , where the sequence $\{Q_n\}_{n \geq 1}$ is a member of the class \mathcal{Q} . Then

1. for every $\xi \leq \xi^*(R, T)$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \tilde{K}_n \leq 0;$$

2. there exists a sequence of encoders and decoders such that for every $\theta \in \Theta$

$$\liminf_{n \rightarrow \infty} \left[-\frac{1}{n} \log \Gamma_\theta(\mathcal{C}, R) \right] \geq \xi^*(R, T) \cdot E_1(R, T, \theta) + T.$$

The proof of the first part of Theorem 2 appears in the Appendix. The second part follows immediately as discussed after (14).

We now pause to discuss Theorem 2 and some of its aspects.

Theorem 2 suggests a conceptually simple strategy: Given R and T , first compute $\xi^*(R, T)$ using (25). This may require some nontrivial optimization procedures, but it has to be done only once, and since this is a single-letter expression, it can be carried at least numerically, if closed-form analytic expressions are not apparent to be available (see the example of the BSC below). Once $\xi^*(R, T)$ has been computed, apply the decoding rule $\hat{\mathcal{R}}$ with $\xi = \xi^*(R, T)$, and the theorem guarantees that the resulting random coding error exponents of \mathcal{E}_1 and \mathcal{E}_2 are at least $\xi^*(R, T) \cdot E_1(R, T, \theta)$ and $\xi^*(R, T) \cdot E_1(R, T, \theta) + T$, respectively.

The theorem is interesting only when $\xi^*(R, T) > 0$, which is the case iff

$$R < R_0(T) \triangleq \min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} \left[\frac{E(\theta, \tilde{\theta}, s, \rho) - sT}{\rho} \right]$$

or equivalently iff

$$T < T_0(R) \triangleq \min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} \left[\frac{E(\theta, \tilde{\theta}, s, \rho) - \rho R}{s} \right].$$

When $\xi^*(R, T) > 0$, the proposed universal decoder with $\xi = \xi^*(R, T)$ has the important property that whenever Forney's optimum decoder yields an exponential decay of $\Pr\{\mathcal{E}_1\}$ ($E_1(R, T, \theta) > 0$), then so does the corresponding exponent of the proposed decoder $\hat{\mathcal{R}}$. It should be pointed out that the exponential rates $\xi^*(R, T) \cdot E_1(R, T, \theta)$ and

$\xi^*(R, T) \cdot E_1(R, T, \theta) + T$, guaranteed by Theorem 2, are only lower bounds to the real exponential rates, and that true exponential rate, at some points in Θ , might be larger.

Our last comment concerns the choice of the threshold T . Thus far, we assumed that T is a constant, independent of θ . However, in some situations, it makes sense to let T depend on the quality of the channel, and hence on the parameter θ . Intuitively, for fixed T , if the signal-to-noise ratio (SNR) becomes very high, the erasure option will be used so rarely that it will effectively be nonexistent. This means that we are actually no longer "enjoying" the benefits of the erasure option, and hence not the gain in the undetected error exponent that is associated with it. An alternative approach is to let $T = T_\theta$ depend on θ in a certain way. In this case, $K_n(C)$ would be redefined as follows:

$$K_n(C) = \min_{\mathcal{R}} \max_{\theta \in \Theta} \frac{e^{-nT_\theta} \Pr\{\mathcal{E}_1\} + \Pr\{\mathcal{E}_2\}}{e^{-n[\xi E_1(R, T_\theta, \theta) + T_\theta]}}. \quad (26)$$

The corresponding generalized version of the competitive minimax decision rule $\hat{\mathcal{R}}$ would now be

$$\begin{aligned} \hat{\mathcal{R}}_m &= \left\{ \mathbf{y} : g(\mathbf{x}_m, \mathbf{y}) \geq \sum_{m' \neq m} h(\mathbf{x}_{m'}, \mathbf{y}) \right\}, \quad m = 1, \dots, M \\ \hat{\mathcal{R}}_0 &= \bigcap_{m=1}^M \hat{\mathcal{R}}_m^c \end{aligned} \quad (27)$$

where

$$g(\mathbf{x}_m, \mathbf{y}) \triangleq \max_{\theta} [P_\theta(\mathbf{y}|\mathbf{x}_m) \cdot e^{n\xi E_1(R, T_\theta, \theta)}] \quad (28)$$

and

$$h(\mathbf{x}_m, \mathbf{y}) \triangleq \max_{\theta} [P_\theta(\mathbf{y}|\mathbf{x}_m) \cdot e^{n[\xi E_1(R, T_\theta, \theta) + T_\theta]}]. \quad (29)$$

By extending the performance analysis carried out in the Appendix, the resulting expression of ξ^* now becomes

$$\xi^*(R) \triangleq \min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} \frac{E(\theta, \tilde{\theta}, s, \rho) - \rho R - sT_\theta}{(1-s)E_1(R, T_\theta, \theta) + sE_1(R, T_\theta, \tilde{\theta})}. \quad (30)$$

The main question that naturally arises, in this case, is: which function T_θ would be reasonable to choose? A plausible guideline could be based on the typical behavior of

$$\tau_\theta = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{E} \ln \frac{P_\theta(\mathbf{Y}|\mathbf{x}_m)}{\sum_{m' \neq m} P_\theta(\mathbf{Y}|\mathbf{x}_{m'})}$$

which can be assessed, using standard bounding techniques, under the hypothesis that \mathbf{x}_m is the correct message. For example, T_θ may be given by $\alpha\tau_\theta$ with some constant $\alpha \in [0, 1]$, or $\tau_\theta - \beta$ for some $\beta > 0$. This will make the probability of erasure (exponentially) small, but not *too* small, so that there would be some gain in the undetected error exponent for every θ .

VI. EXAMPLE—THE BINARY-SYMMETRIC CHANNEL

Consider the BSC, where $\mathcal{X} = \mathcal{Y} = \{0, 1\}$, and where θ designates the crossover probability, and let the sequence of random coding distributions be uniform, i.e., $Q_n(\mathbf{x}) = 1/|\mathcal{X}|^n$ for all $\mathbf{x} \in \mathcal{X}^n$, which as mentioned earlier, belongs to the class \mathcal{Q}

with $\Delta^*(P) = \ln|\mathcal{X}| - H(P) = \ln 2 - H(P)$. We would like to examine, more closely, the expression of $\xi^*(R, T)$ and its behavior in this case. Let $h_2(u)$ denote the binary entropy function $-u \ln u - (1-u) \ln(1-u)$, $u \in [0, 1]$. Denoting the modulo 2 sum of X and Y by $X \oplus Y$, we then have

$$\begin{aligned}
F(P_y, \lambda, \theta) &= \min_{P_{x|y}} [I(X;Y) + (\ln 2 - H(X)) - \lambda \mathbf{E} \ln P(Y|X)] \\
&= \ln 2 - \max_{P_{x|y}} [H(X|Y) + \lambda \mathbf{E} \ln P(Y|X)] \\
&= \ln 2 \\
&\quad - \max_{P_{x|y}} \left\{ H(X|Y) + \lambda \mathbf{E} \ln \left[(1-\theta) \left(\frac{\theta}{1-\theta} \right)^{X \oplus Y} \right] \right\} \\
&= \ln 2 - \lambda \ln(1-\theta) \\
&\quad - \max_{P_{x|y}} \left[H(X|Y) + \left(\lambda \ln \frac{\theta}{1-\theta} \right) \cdot \mathbf{E}(X \oplus Y) \right] \\
&= \ln 2 - \lambda \ln(1-\theta) \\
&\quad - \max_{P_{x|y}} \left[H(X \oplus Y|Y) + \left(\lambda \ln \frac{\theta}{1-\theta} \right) \cdot \mathbf{E}(X \oplus Y) \right] \\
&\geq \ln 2 - \lambda \ln(1-\theta) \\
&\quad - \max_{P_{x|y}} \left[H(X \oplus Y) + \left(\lambda \ln \frac{\theta}{1-\theta} \right) \cdot \mathbf{E}(X \oplus Y) \right] \\
&= \ln 2 - \lambda \ln(1-\theta) - \max_u \left[h_2(u) + \left(\lambda \ln \frac{\theta}{1-\theta} \right) \cdot u \right] \\
&= \ln 2 - \lambda \ln(1-\theta) - \ln \left[1 + \left(\frac{\theta}{1-\theta} \right)^\lambda \right] \\
&= \ln 2 - \ln[\theta^\lambda + (1-\theta)^\lambda] \quad (31)
\end{aligned}$$

where the inequality is, in fact, an equality achieved by a backward $P_{x|y}$ where $X \oplus Y$ is independent of Y . Since $F(P_y, \lambda, \theta)$ is independent of P_y , this easily yields

$$\begin{aligned}
E(\theta, \tilde{\theta}, \rho, s) &= \rho \ln 2 - \ln[\theta^{1-s} + (1-\theta)^{1-s}] \\
&\quad - \rho \ln[\tilde{\theta}^{s/\rho} + (1-\tilde{\theta})^{s/\rho}] \quad (32)
\end{aligned}$$

and so, we get (33) at the bottom of the page with (34), also at the bottom of the page. This expression, although still involves nontrivial optimizations, is much more explicit than the general one. We next offer a few observations regarding the function $\xi^*(R, T)$ for the example of the BSC.

First, observe that if Θ is a singleton, i.e., we are back to the case of a known channel, then $\theta = \tilde{\theta}$, and the numerator, after maximization over ρ and s , becomes $E_1(R, T, \theta)$, and so does the denominator, thus $\xi^*(R, T) = 1$, as expected.

We next demonstrate that $\xi^*(R, 0) = 1$. This result is expected, as the case $T = 0$ is asymptotically equivalent (cf. [10])

to the case without erasures in the sense that $E_1(R, 0, \theta) = E_2(R, 0, \theta)$ coincide with Gallager's random coding exponent [11] (although erasures are still possible). This is in agreement with the aforementioned full universality result for ordinary universal decoding.

Referring to the definition of the Gallager function $E(\theta, \rho)$ for the BSC

$$E(\theta, \rho) = \rho \ln 2 - (1+\rho) \ln[\theta^{1/(1+\rho)} + (1-\theta)^{1/(1+\rho)}] - \rho R \quad (35)$$

let us define $\rho' = 1/(1-s) - 1$ and $\rho'' = \rho/s - 1$, and rewrite the numerator of the expression for $\xi^*(R, 0)$ as follows:

$$\begin{aligned}
&\rho \ln 2 - \ln[\theta^{1-s} + (1-\theta)^{1-s}] - \rho \ln[\tilde{\theta}^{s/\rho} + (1-\tilde{\theta})^{s/\rho}] - \rho R \\
&= \rho \ln 2 - \ln[\theta^{1/(1+\rho')} + (1-\theta)^{1/(1+\rho')}] \\
&\quad - \rho \ln[\tilde{\theta}^{1/(1+\rho'')} + (1-\tilde{\theta})^{1/(1+\rho'')}] - \rho R \\
&= \frac{1}{1+\rho'} \{ \rho' \ln 2 - (1+\rho') \ln[\theta^{1/(1+\rho')} + (1-\theta)^{1/(1+\rho')}] - \rho' R \} + \\
&\quad + \frac{\rho}{1+\rho''} \{ \rho'' \ln 2 - (1+\rho'') \ln[\tilde{\theta}^{1/(1+\rho'')} + (1-\tilde{\theta})^{1/(1+\rho'')}] - \rho'' R \} \\
&= (1-s)E(\theta, \rho') + sE(\tilde{\theta}, \rho'') \\
&= (1-s)E\left(\theta, \frac{1}{1-s} - 1\right) + sE\left(\tilde{\theta}, \frac{\rho}{s} - 1\right). \quad (36)
\end{aligned}$$

Now, let us choose $s = \rho/(1+\tilde{\rho})$, where $\tilde{\rho}$ is the achiever of $E^*(\tilde{\theta}) = \max_{0 \leq \rho \leq 1} E(\tilde{\theta}, \rho)$, and $\rho = \rho^*(1+\tilde{\rho})/(1+\rho^*)$, where ρ^* is the achiever of $E^*(\theta) = \max_{0 \leq \rho \leq 1} E(\theta, \rho)$ (observing that $\rho^*(1+\tilde{\rho})/(1+\rho^*) \leq 1$, therefore, this is choice is feasible). With this choice, the numerator of $\xi^*(R, 0)$ becomes equal to the denominator, and so, $\xi^*(R, 0) = 1$.

Finally, in Table I, we provide some numerical results pertaining to the function $\xi^*(R, T)$, where all minimizations and maximizations were carried out by an exhaustive search with a step-size of 0.01 in each dimension. As can be seen, at the left-most column, corresponding to $T = 0$, we indeed obtain $\xi^*(R, 0) = 1$. As can also be seen, $\xi^*(R, T)$ is always strictly less than unity for $T > 0$, and it in general decreases as T grows.

VII. CONCLUSION

We have addressed the problem of universal decoding with erasures, using the competitive minimax methodology proposed in [9], which proved useful. This is in contrast to earlier approaches for deriving universal decoders, based on joint typicality considerations, for which we found no apparent extensions to accommodate Forney's erasure decoder. In order

$$\xi^*(R, T) = \min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} \frac{\rho \ln 2 - \ln[\theta^{1-s} + (1-\theta)^{1-s}] - \rho \ln[\tilde{\theta}^{s/\rho} + (1-\tilde{\theta})^{s/\rho}] - \rho R - sT}{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})} \quad (33)$$

$$E_1(R, T, \theta) = \max_{0 \leq s \leq \rho \leq 1} \{ \rho \ln 2 - \ln[\theta^{1-s} + (1-\theta)^{1-s}] - \rho \ln[\theta^{s/\rho} + (1-\theta)^{s/\rho}] - \rho R - sT \}. \quad (34)$$

TABLE I
NUMERICAL VALUES OF $\xi^*(R, T)$ FOR VARIOUS VALUES OF R AND T

	$T = 0.000$	$T = 0.025$	$T = 0.050$	$T = 0.075$	$T = 0.100$	$T = 0.125$	$T = 0.150$
$R = 0.00$	1.000	0.364	0.523	0.418	0.396	0.422	0.298
$R = 0.05$	1.000	0.756	0.713	0.656	0.535	0.562	0.495
$R = 0.10$	1.000	0.858	0.774	0.648	0.655	0.585	0.518
$R = 0.15$	1.000	0.877	0.809	0.720	0.713	0.662	0.622
$R = 0.20$	1.000	0.905	0.815	0.729	0.729	0.684	0.647
$R = 0.25$	1.000	0.912	0.832	0.763	0.706	0.661	0.627
$R = 0.30$	1.000	0.896	0.850	0.788	0.738	0.644	0.613

to guarantee the uniform achievability of a certain fraction of the exponent, the competitive minimax approach was applied to the Lagrangian, pertaining to a weighted sum of the two error probabilities.

The analysis of the minimax ratio, \bar{K}_n , resulted in a single-letter lower bound to the largest universally achievable fraction $\xi^*(R, T)$ of Forney's exponent. An interesting problem for future work would be to derive a (hopefully compatible) lower bound. This requires the derivation of an exponentially tight lower bound to K_n , which is a challenge.

Our results cover performance analysis of competitive-minimax universal decoders with various types of random coding distributions in a considerable wide class \mathcal{Q} . This is in contrast to earlier works (see, e.g., [4], [22]), which were firmly based on the assumption that the random coding distribution is uniform within a set. A similar analysis technique can be applied also to universal decoding without erasures.

Finally, we analyzed the example of the BSC in full detail and demonstrated that $\xi^*(R, 0) = 1$. We have also provided some numerical results for this case.

APPENDIX PROOF OF THEOREM 2

For a given subset $\mathcal{E} \subseteq \mathcal{Y}^n$, let $1\{\mathbf{y}|\mathcal{E}\}$ denote the indicator function of \mathcal{E} , i.e., $1\{\mathbf{y}|\mathcal{E}\} = 1$ if $\mathbf{y} \in \mathcal{E}$ and $1\{\mathbf{y}|\mathcal{E}\} = 0$ otherwise. First, observe that

$$\begin{aligned} 1\{\mathbf{y}|\hat{\mathcal{R}}_m\} &= 1\left\{\mathbf{y} \mid f(\mathbf{x}_m, \mathbf{y}) \geq e^{nT} \sum_{m' \neq m} f(\mathbf{x}_{m'}, \mathbf{y})\right\} \\ &\leq \min_{0 \leq s \leq 1} \left[\frac{f(\mathbf{x}_m, \mathbf{y})}{e^{nT} \sum_{m' \neq m} f(\mathbf{x}_{m'}, \mathbf{y})} \right]^{1-s} \end{aligned} \quad (\text{A1})$$

and similarly

$$1\{\mathbf{y} \in \hat{\mathcal{R}}_m^c\} \leq \min_{0 \leq s \leq 1} \left[\frac{e^{nT} \sum_{m' \neq m} f(\mathbf{x}_{m'}, \mathbf{y})}{f(\mathbf{x}_m, \mathbf{y})} \right]^s. \quad (\text{A2})$$

Then, we have

$$\begin{aligned} \bar{K}_n &\leq \mathbf{E}\{K_n(\hat{\mathcal{R}}, \mathcal{C})\} \\ &= \mathbf{E}\left\{\max_{\theta} \left[\frac{e^{-nT} \Pr\{\mathcal{E}_1\} + \Pr\{\mathcal{E}_2\}}{e^{-n[\xi E_1(R, T, \theta) + T]}} \right] \right\} \\ &= \mathbf{E}\left\{\max_{\theta} \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \left[e^{-nT} \left(P_{\theta}(\mathbf{y}|\mathbf{X}_m) \cdot e^{n[\xi E_1(R, T, \theta) + T]} \right) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m^c\} + \right. \right. \\ &\quad \left. \left(\sum_{m' \neq m} P_{\theta}(\mathbf{y}|\mathbf{X}_{m'}) \cdot e^{n[\xi E_1(R, T, \theta) + T]} \right) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m\} \right] \right\} \\ &\stackrel{(a)}{\leq} \mathbf{E}\left\{\frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \left[e^{-nT} \max_{\theta} \left(P_{\theta}(\mathbf{y}|\mathbf{X}_m) \cdot e^{n[\xi E_1(R, T, \theta) + T]} \right) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m^c\} + \right. \right. \\ &\quad \left. \left(\sum_{m' \neq m} \max_{\theta} [P_{\theta}(\mathbf{y}|\mathbf{X}_{m'}) \cdot e^{n[\xi E_1(R, T, \theta) + T]}] \right) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m\} \right] \right\} \end{aligned}$$

$$\begin{aligned}
&= \mathbf{E} \left\{ \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \left[e^{-nT} f(\mathbf{X}_m, \mathbf{y}) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m^c\} + \left(\sum_{m' \neq m} f(\mathbf{X}_{m'}, \mathbf{y}) \right) \cdot 1\{\mathbf{y}|\hat{\mathcal{R}}_m\} \right] \right\} \\
&\stackrel{(b)}{\leq} \mathbf{E} \left\{ \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \left[e^{-nT} f(\mathbf{X}_m, \mathbf{y}) \cdot \min_{0 \leq s \leq 1} \left(\frac{e^{nT} \sum_{m' \neq m} f(\mathbf{X}_{m'}, \mathbf{y})}{f(\mathbf{X}_m, \mathbf{y})} \right)^s + \right. \right. \\
&\quad \left. \left(\sum_{m' \neq m} f(\mathbf{X}_{m'}, \mathbf{y}) \right) \cdot \min_{0 \leq s \leq 1} \left(\frac{f(\mathbf{X}_m, \mathbf{y})}{e^{nT} \sum_{m' \neq m} f(\mathbf{X}_{m'}, \mathbf{y})} \right)^{1-s} \right] \Big\} \\
&= \mathbf{E} \left\{ \frac{2}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} f^{1-s}(\mathbf{X}_m, \mathbf{y}) \left(\sum_{m' \neq m} f(\mathbf{X}_{m'}, \mathbf{y}) \right)^s \right\} \right\} \\
&= \mathbf{E} \left\{ \frac{2}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left(\max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right)^{1-s} \cdot \right. \right. \\
&\quad \left. \left(\sum_{m' \neq m} \max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]} \right)^s \right\} \Big\} \\
&\stackrel{(c)}{\leq} \mathbf{E} \left\{ \frac{2}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[\max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \right. \\
&\quad \left. \left(\sum_{m' \neq m} \sum_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]} \right)^s \right\} \Big\} \\
&= \mathbf{E} \left\{ \frac{2}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[\max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \right. \\
&\quad \left. \left(\sum_{\theta \in \Theta_n} \sum_{m' \neq m} P_\theta(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]} \right)^s \right\} \Big\} \\
&\stackrel{(d)}{\leq} \mathbf{E} \left\{ \frac{2}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[\max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \right. \\
&\quad \left. \left(|\Theta_n| \cdot \max_{\theta \in \Theta_n} \sum_{m' \neq m} P_\theta(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]} \right)^s \right\} \Big\} \\
&\leq \mathbf{E} \left\{ \frac{2|\Theta_n|}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[\max_{\theta \in \Theta_n} P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \right. \\
&\quad \left. \left(\max_{\theta \in \Theta_n} \sum_{m' \neq m} P_\theta(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]} \right)^s \right\} \Big\} \tag{A3}
\end{aligned}$$

where (a) follows from the fact that the maximum (over θ) of a summation is upper-bounded by the summation of the maxima, (b) follows from (A1) and (A2), and (c) and (d) follow from the fact that if $g(\theta)$ is nonnegative then

$$\max_{\theta \in \Theta_n} g(\theta) \leq \sum_{\theta \in \Theta_n} g(\theta) \leq |\Theta_n| \cdot \max_{\theta \in \Theta_n} g(\theta). \tag{A4}$$

Now, for every given \mathbf{y} and $\{\mathbf{x}_m\}$, let $\theta^* \in \Theta_n$ be the achiever of

$$\max_{\theta \in \Theta_n} (P_\theta(\mathbf{y}|\mathbf{x}_m) e^{n[\xi E_1(R, T, \theta) + T]})$$

and let $\theta^{**} \in \Theta_n$ be the achiever of

$$\max_{\theta \in \Theta_n} \sum_{m' \neq m} P_\theta(\mathbf{y}|\mathbf{x}_{m'}) e^{n[\xi E_1(R, T, \theta) + T]}.$$

Note that θ^* and θ^{**} depend on $\mathbf{x}_1, \dots, \mathbf{x}_M$ and \mathbf{y} , but not on the parameter s . Let us denote

$$W(\mathbf{x}_1, \dots, \mathbf{x}_M, \mathbf{y}, \theta, \tilde{\theta}) = \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[P_\theta(\mathbf{y}|\mathbf{x}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{x}_{m'}) e^{n[\xi E_1(R, T, \tilde{\theta}) + T]} \right)^s \right\}. \quad (\text{A5})$$

Now, obviously, we get

$$\begin{aligned} \bar{K}_n &\leq \mathbf{E} \left\{ \frac{2|\Theta_n|}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} W(\mathbf{X}_1, \dots, \mathbf{X}_M, \mathbf{y}, \theta^*, \theta^{**}) \right\} \\ &\leq \mathbf{E} \left\{ \frac{2|\Theta_n|}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \sum_{\theta \in \Theta_n} \sum_{\tilde{\theta} \in \Theta_n} W(\mathbf{X}_1, \dots, \mathbf{X}_M, \mathbf{y}, \theta, \tilde{\theta}) \right\} \\ &= \mathbf{E} \left\{ \frac{2|\Theta_n|}{M} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \sum_{\theta \in \Theta_n} \sum_{\tilde{\theta} \in \Theta_n} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \right. \\ &\quad \left. \left. \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \tilde{\theta}) + T]} \right)^s \right\} \right\} \\ &= \frac{2|\Theta_n|}{M} \sum_{\theta \in \Theta_n} \sum_{\tilde{\theta} \in \Theta_n} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \\ &\quad \left. \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \tilde{\theta}) + T]} \right)^s \right\} \\ &\stackrel{(a)}{\leq} \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta_n} \max_{\tilde{\theta} \in \Theta_n} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \min_{0 \leq s \leq 1} \left\{ e^{-nT(1-s)} \left[P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \\ &\quad \left. \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \tilde{\theta}) + T]} \right)^s \right\} \\ &\leq \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq 1} \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \left\{ e^{-nT(1-s)} \left[P_\theta(\mathbf{y}|\mathbf{X}_m) e^{n[\xi E_1(R, T, \theta) + T]} \right]^{1-s} \times \right. \\ &\quad \left. \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) e^{n[\xi E_1(R, T, \tilde{\theta}) + T]} \right)^s \right\} \\ &= \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq 1} e^{n[\xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} + sT]} \times \\ &\quad \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \left\{ P_\theta^{1-s}(\mathbf{y}|\mathbf{X}_m) \cdot \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) \right)^s \right\} \end{aligned} \quad (\text{A6})$$

where in (a) we used again (A4). Assuming that the codewords are drawn independently, we then have

$$\begin{aligned} \bar{K}_n &\leq \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq 1} e^{n[\xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} + sT]} \times \\ &\quad \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \{ P_\theta^{1-s}(\mathbf{y}|\mathbf{X}_m) \} \cdot \mathbf{E} \left\{ \left(\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) \right)^s \right\} \\ &= \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq 1} e^{n[\xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} + sT]} \times \\ &\quad \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E} \{ P_\theta^{1-s}(\mathbf{y}|\mathbf{X}_m) \} \cdot \min_{s \leq \rho \leq 1} \mathbf{E} \left\{ \left(\left[\sum_{m' \neq m} P_{\tilde{\theta}}(\mathbf{y}|\mathbf{X}_{m'}) \right]^{s/\rho} \right)^\rho \right\} \\ &\leq \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq 1} e^{n[\xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} + sT]} \times \end{aligned}$$

$$\begin{aligned}
& \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E}\{P_\theta^{1-s}(\mathbf{y}|\mathbf{X}_m)\} \cdot \min_{0 \leq s \leq \rho \leq 1} \mathbf{E} \left[\left(\sum_{m' \neq m} P_{\tilde{\theta}}^{s/\rho}(\mathbf{y}|\mathbf{X}_{m'}) \right)^\rho \right] \\
& \leq \frac{2|\Theta_n|^3}{M} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \times \\
& \quad \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{E}\{P_\theta^{1-s}(\mathbf{y}|\mathbf{X}_m)\} \cdot \left(\sum_{m' \neq m} \mathbf{E}\{P_{\tilde{\theta}}^{s/\rho}(\mathbf{y}|\mathbf{X}_{m'})\} \right)^\rho
\end{aligned} \tag{A7}$$

where in the last step we have used Jensen's inequality. Now, observe that the summands do not depend on m , therefore, the effects of the summation over m and the factor of $1/M$ cancel each other. Also, the sum of $M - 1$ contributions of identical expectations $\mathbf{E}\{P_{\tilde{\theta}}^{s/\rho}(\mathbf{y}|\mathbf{X}_{m'})\}$ create a factor of $M - 1$ (upper-bounded by M) raised to the power of ρ . Denoting

$$U(\mathbf{y}, \lambda, \theta) = \mathbf{E}\{P_\theta^\lambda(\mathbf{y}|\mathbf{X})\}$$

we have

$$\bar{K}_n \leq 2|\Theta_n|^3 \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} M^\rho \cdot e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \sum_{\mathbf{y} \in \mathcal{Y}^n} U(\mathbf{y}, 1-s, \theta) \cdot U^\rho(\mathbf{y}, s/\rho, \tilde{\theta}). \tag{A8}$$

To assess the exponential order of $U(\mathbf{y}, \lambda, \theta)$, we use the method of types [4] as well as the assumption that the sequence of random coding distributions belongs to the class \mathcal{Q} . Let $\hat{P}_{\mathbf{y}}$ denote the empirical distribution of \mathbf{y} , and let $T_{\mathbf{y}}$ denote the type class of \mathbf{y} , i.e., the set of \mathbf{y}' with $\hat{P}_{\mathbf{y}'} = \hat{P}_{\mathbf{y}}$. Let $\hat{H}_{\mathbf{y}}(Y)$ denote the corresponding empirical entropy of Y . Similarly, let $\hat{P}_{\mathbf{xy}}$ denote the empirical joint distribution of (\mathbf{x}, \mathbf{y}) and let $\hat{E}_{\mathbf{xy}}\{\cdot\}$ denote the corresponding empirical expectation, i.e., the expectation w.r.t. $\hat{P}_{\mathbf{xy}}$. Also, let $T_{\mathbf{x}|\mathbf{y}}$ denote the conditional type class of \mathbf{x} given \mathbf{y} , i.e., the set of \mathbf{x}' with $\hat{P}_{\mathbf{x}'|\mathbf{y}} = \hat{P}_{\mathbf{xy}}$ and let $\hat{I}_{\mathbf{xy}}(X; Y)$ denote the corresponding empirical mutual information between X and Y . Then, we get

$$\begin{aligned}
U(\mathbf{y}, \lambda, \theta) &= \sum_{\mathbf{x} \in \mathcal{X}^n} Q_n(\mathbf{x}) P_\theta^\lambda(\mathbf{y}|\mathbf{x}) \\
&= \sum_{T_{\mathbf{x}|\mathbf{y}} \subset \mathcal{X}^n} |T_{\mathbf{x}|\mathbf{y}}| \cdot \frac{e^{-n\Delta_n^*(\hat{P}_{\mathbf{x}})}}{|T_{\mathbf{x}}|} \cdot e^{\lambda n \hat{E}_{\mathbf{xy}} \ln P_\theta(Y|X)} \\
&\leq (n+1)^{|\mathcal{X}|-1} \sum_{T_{\mathbf{x}|\mathbf{y}} \subset \mathcal{X}^n} e^{-n\hat{I}_{\mathbf{xy}}(X;Y)} \cdot e^{-n[\Delta_n^*(\hat{P}_{\mathbf{x}}) - \epsilon_n]} \cdot e^{\lambda n \hat{E}_{\mathbf{xy}} \ln P_\theta(Y|X)} \\
&\leq (n+1)^{(|\mathcal{Y}|+1) \cdot (|\mathcal{X}|-1)} \cdot e^{-n[F(\hat{P}_{\mathbf{y}}, \lambda, \theta) - \epsilon_n]}
\end{aligned} \tag{A9}$$

where $\epsilon_n \rightarrow 0$ independently of $\hat{P}_{\mathbf{x}}$ by the uniform convergence assumption that defines the class \mathcal{Q} , and where $F(P_{\mathbf{y}}, \lambda, \theta)$ is defined as in (23). On substituting this bound into the upper bound on $K_n(\hat{\mathcal{R}})$, we get

$$\begin{aligned}
\bar{K}_n &\leq 2|\Theta_n|^3 (n+1)^{2|\mathcal{Y}| \cdot (|\mathcal{X}|-1)} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} \\
&\quad M^\rho \cdot e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \cdot \sum_{\mathbf{y} \in \mathcal{Y}^n} e^{-n[F(P_{\mathbf{y}}, 1-s, \theta) + \rho F(P_{\mathbf{y}}, s/\rho, \tilde{\theta})]} \\
&\leq 2|\Theta_n|^3 (n+1)^{2|\mathcal{Y}| \cdot (|\mathcal{X}|-1)} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} \\
&\quad M^\rho \cdot e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \cdot \sum_{T_{\mathbf{y}} \subset \mathcal{Y}^n} e^{n\hat{H}_{\mathbf{y}}(Y)} \cdot e^{-n[F(P_{\mathbf{y}}, 1-s, \theta) + \rho F(P_{\mathbf{y}}, s/\rho, \tilde{\theta})]} \\
&\leq 2|\Theta_n|^3 (n+1)^{3|\mathcal{Y}| \cdot (|\mathcal{X}|-1)} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} \\
&\quad M^\rho \cdot e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \cdot e^{-n \min_{P_{\mathbf{y}}} [F(P_{\mathbf{y}}, 1-s, \theta) + \rho F(P_{\mathbf{y}}, s/\rho, \tilde{\theta}) - H(Y)]} \\
&\leq 2|\Theta_n|^3 (n+1)^{3|\mathcal{Y}| \cdot (|\mathcal{X}|-1)} \max_{\theta \in \Theta} \max_{\tilde{\theta} \in \Theta} \min_{0 \leq s \leq \rho \leq 1} \\
&\quad M^\rho \cdot e^{n[\xi\{(1-s)E_1(R,T,\theta)+sE_1(R,T,\tilde{\theta})\}+sT]} \cdot e^{-nE(\theta, \tilde{\theta}, s, \rho)}.
\end{aligned} \tag{A10}$$

We would like to find the maximum value of ξ such that \tilde{K}_n would be guaranteed not to grow exponentially. To this end, we can now ignore the factor $2|\Theta_n|^3(n+1)^{3|\mathcal{Y}| \cdot (|\mathcal{X}|-1)}$, which is polynomial in n (cf. (19)). Thus, the latter upper bound will be subexponential in n as long as

$$\min_{\theta, \tilde{\theta}} \max_{0 \leq s \leq \rho \leq 1} [E(\theta, \tilde{\theta}, s, \rho) - \xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} - \rho R - sT] \geq 0 \quad (\text{A11})$$

holds or, equivalently, for every $(\theta, \tilde{\theta})$, there exist (ρ, s) , $0 \leq s \leq \rho \leq 1$, such that

$$E(\theta, \tilde{\theta}, s, \rho) \geq \xi \{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})\} + \rho R + sT \quad (\text{A12})$$

i.e.,

$$\xi \leq \frac{E(\theta, \tilde{\theta}, s, \rho) - \rho R - sT}{(1-s)E_1(R, T, \theta) + sE_1(R, T, \tilde{\theta})}. \quad (\text{A13})$$

In other words, for every $\xi \leq \xi^*(R, T)$, where $\xi^*(R, T)$ is defined as in (25)), $K_n(\hat{\mathcal{R}})$ is guaranteed not to grow exponentially with n . This completes the proof of Theorem 2.

ACKNOWLEDGMENT

The authors would like to thank the Associate Editor, Gerhard Kramer, and the three anonymous reviewers for their very useful comments, which certainly helped to improve the presentation of our results.

REFERENCES

- [1] R. Ahlswede, N. Cai, and Z. Zhang, "Erasure, list, and detection zero-error capacities for low noise and a relation to identification," *IEEE Trans. Inf. Theory*, vol. 42, no. 1, pp. 55–62, Jan. 1996.
- [2] J. Byers, M. Luby, and M. Mitzenmacher, "A digital fountain approach to asynchronous reliable multicast," *IEEE J. Sel. Areas Commun.*, vol. 20, no. 8, pp. 1528–1540, Oct. 2002.
- [3] J. Byers, M. Luby, M. Mitzenmacher, and A. Rege, "A digital fountain approach to reliable distribution of bulk data," in *Proc. ACM SIGCOMM '98*, Vancouver, BC, Canada, Sep. 1998, pp. 56–67.
- [4] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. New York: Academic, 1981.
- [5] S. Draper, B. J. Frey, and F. R. Kschischang, "Rateless coding for non-ergodic channels with decoder channel state information," *IEEE Trans. Inf. Theory*, submitted for publication.
- [6] S. C. Draper, B. J. Frey, and F. R. Kschischang, "Efficient variable length channel coding for unknown DMCs," in *Proc. IEEE Int. Symp. Information Theory*, Chicago, IL, Jun./Jul. 2004, p. 377.
- [7] U. Erez, G. W. Wornell, and M. D. Trott, "Rateless space-time coding," in *Proc. IEEE Int. Symp. Information Theory*, Adelaide, Australia, Sep. 2005, pp. 1937–1941.
- [8] M. Feder and A. Lapidoth, "Universal decoders for channels with memory," *IEEE Trans. Inf. Theory*, vol. 44, no. 5, pp. 1726–1745, Sep. 1998.
- [9] M. Feder and N. Merhav, "Universal composite hypothesis testing: A competitive minimax approach," *IEEE Trans. Inf. Theory, Special Issue in Memory of Aaron D. Wyner*, vol. 48, no. 6, pp. 1504–1517, Jun. 2002.
- [10] G. D. Forney, Jr., "Exponential error bounds for erasure, list, and decision feedback," *IEEE Trans. Inf. Theory*, vol. IT-14, no. 2, pp. 206–220, Mar. 1968.
- [11] R. G. Gallager, *Information Theory and Reliable Communication*. New York: Wiley, 1968.
- [12] T. Hashimoto, "Composite scheme LT+Th for decoding with erasures and its effective equivalence to Forney's rule," *IEEE Trans. Inf. Theory*, vol. 45, no. 1, pp. 78–93, Jan. 1999.
- [13] T. Hashimoto and M. Taguchi, "Performance and explicit error detection and threshold decision in decoding with erasures," *IEEE Trans. Inf. Theory*, vol. 43, no. 5, pp. 1650–1655, Sep. 1997.
- [14] J. Jiang and K. R. Narayanan, "Multilevel Coding for Channels with Non-Uniform Inputs and Rateless Transmission Over the BSC," Jan. 18, 2006 [Online]. Available: arXiv:cs.IT/0601083
- [15] M. Luby, M. Mitzenmacher, A. Shokrollahi, and D. Spielman, "Efficient erasure correction codes," *IEEE Trans. Inf. Theory*, vol. 47, no. 2, pp. 569–584, Feb. 2001.
- [16] J. L. Massey, "Causality, feedback and directed information," in *Proc. IEEE Int. Symp. Information Theory and Its Applications (ISITA '90)*, HI, 1990, pp. 303–305.
- [17] N. Shulman, "Communication over an unknown channel via common broadcasting," Ph.D. dissertation, Tel-Aviv Univ., Tel-Aviv, Israel, 2003.
- [18] N. Shulman and M. Feder, "Static broadcasting," in *Proc. IEEE Int. Symp. Information Theory*, Sorrento, Italy, Jun. 2000, p. 23.
- [19] N. Shulman and M. Feder, "The uniform distribution as a universal prior," *IEEE Trans. Inf. Theory*, vol. 50, no. 6, pp. 1356–1362, Jun. 2004.
- [20] A. Tchamkerten and I. E. Telatar, "Variable length coding over unknown channels," *IEEE Trans. Inf. Theory*, vol. 52, no. 6, pp. 2126–2145, May 2006.
- [21] A. J. Viterbi, "Error bounds for the white Gaussian and other very noisy memoryless channels with generalized decision regions," *IEEE Trans. Inf. Theory*, vol. IT-15, no. 2, pp. 279–287, Mar. 1969.
- [22] J. Ziv, "Universal decoding for finite-state channels," *IEEE Trans. Inf. Theory*, vol. IT-31, no. 4, pp. 453–460, Jul. 1985.