

Universal Coding with Minimum Probability of Codeword Length Overflow

Neri Merhav

Abstract—Lossless block-to-variable length source coding is studied for finite-state, finite-alphabet sources. We aim to minimize the probability that the normalized length of the codeword will exceed a given threshold B , subject to the Kraft inequality. It is shown that the Lempel–Ziv (LZ) algorithm asymptotically attains the optimal performance in the sense just defined, independently of the source and the value of B . For the subclass of unifilar Markov sources, faster convergence to the asymptotic optimum performance can be accomplished by using the minimum description length (MDL) universal code for this subclass. It is demonstrated that these universal codes are also nearly optimal in the sense of minimizing buffer overflow probability, and asymptotically optimal in a competitive sense.

Index Terms—Universal noiseless coding, Lempel–Ziv algorithm, length overflow, buffer overflow, finite-state sources, large deviations, competitive optimality.

I. INTRODUCTION

LOSSLESS block-to-variable length source coding schemes are usually examined under the criterion of minimizing either the expected length $E(L)$ of a codeword [1], [2] or the cost function $\lambda^{-1} \log E(2^{\lambda L})$ [2], [3], where L denotes the codeword length, $\lambda > 0$ is a given constant, and $\log(\cdot) \triangleq \log_2(\cdot)$ throughout the sequel. It is well known that the Shannon entropy and the Rényi entropy, respectively, are achievable lower bounds to these cost functions.

We introduce a different criterion for the performance of a lossless code. Given a constant $B > 0$, we seek a uniquely decipherable block code, or a class of block codes, which asymptotically minimize the probability that $L > Bn$, where n is the input block size. This criterion might be useful in applications of fixed rate coding schemes, where a block of length n is first encoded to a binary string of variable length L , which if $L > Bn$, is truncated to Bn bits, resulting in a constant rate of B bits per input letter. We would like to minimize the probability that $L > Bn$, i.e., the probability that information is lost.

For the class of finite-state, finite alphabet sources, it is shown that the Lempel–Ziv (LZ) algorithm [4] is asymptotically optimal in the above sense, independently of the

source and the value of B . In other words, even if the source statistics are known, the best code, in that sense, is the LZ code, which does not use the knowledge of the source. The technique used to prove this result is similar to that of [5]. However, all the pertinent results from [5] are rederived here for the sake of completeness.

Next, it is demonstrated that for unifilar finite-state sources [6, p. 187], that is, sources for which the underlying state sequence is uniquely determined by the observations (e.g., Markov sources), the probability of length overflow associated with the LZ algorithm tends to zero exponentially fast as the size n of the block tends to infinity. If, in addition, the source is known *a priori* to be unifilar, then the convergence of $n^{-1} \log \Pr\{n^{-1}L > B\}$ to the asymptotic optimal exponent can be accelerated if one uses the minimum description length (MDL) universal code [7], [8] for the subclass of unifilar sources.

Two additional aspects of these universal codes are also discussed. First, it is shown that when buffer is used in variable length coding of fixed rate memoryless sources, the LZ algorithm, as well as the MDL universal code, nearly maximize the exponential decay rate of the buffer overflow probability. This is done by showing that these codes asymptotically attain the Rényi entropy (as well as the Shannon entropy) for every memoryless source. Second, it is shown that these universal codes are asymptotically optimal in a competitive sense [9], i.e., most of the time they provide a codeword shorter than that of any competing code, within a vanishingly small redundancy term.

In Section II we formulate the problem and state the main theorem. The case of unifilar sources is considered in Section III. In Section IV these results, along with additional aspects, are discussed. Finally, in Section V a proof of the main result is provided.

II. PROBLEM FORMULATION AND MAIN RESULT

Let $\mathbf{x} = x_1, x_2, \dots, x_i, \dots, x_n$, be a sequence of observable random variables taking values in a finite set X with cardinality $|X| = X$. Similarly, let $\mathbf{s} = s_1, s_2, \dots, s_i, \dots, s_n$, be another sequence of random variables, called states, which take values in another finite set S of size $|S| = S$. A probabilistic source P is called finite-state (with S states) if

$$P(\mathbf{x}, \mathbf{s}) = \prod_{i=1}^n p(x_i, s_i | s_{i-1}), \quad (1)$$

Manuscript received October 20, 1989; revised August 29, 1990. This work was presented at the 1990 AAAI Spring Symposium on the Theory and Applications of Minimum Length Coding, Stanford University, Stanford, CA, March 27–29, 1990.

The author is with the Department of Electrical Engineering, Technion—Israel Institute of Technology, Haifa 32000, Israel.

IEEE Log Number 9041980.

0018-9448/91/0500-0556\$01.00 ©1991 IEEE

where $p(x_i, s_i | s_{i-1})$ is the joint probability of a letter x_i and a state s_i given the previous state s_{i-1} , $P(x, s)$ is the joint probability of x and s , and $s_0 \in S$ is a fixed initial state. Often the state sequence s is unobservable. The class of finite-state sources with no more than S states will be denoted by \mathcal{P}_S .

A length function of a binary lossless code is a mapping L_n from the set X^n of all possible observation sequences x to the set N of positive integers, which satisfies the Kraft inequality,

$$\sum_{x \in X^n} 2^{-L_n(x)} \leq 1. \quad (2)$$

Given a finite-state source P and a constant $B > 0$, the problem is that of finding a length function of a uniquely decipherable code $L_n(\cdot)$ that minimizes the probability of length overflow, defined as

$$\Pr\{n^{-1}L_n(x) > B\} = \sum_{x: n^{-1}L_n(x) > B} P(x), \quad (3)$$

where $P(x) = \sum_{s \in S^n} P(x, s)$ and $P(x, s)$ is as in (1). It is assumed that

$$H < B < \log X, \quad (4)$$

where H is the Shannon entropy, which for a stationary source is given by

$$H \triangleq - \lim_{n \rightarrow \infty} n^{-1} \sum_{x \in X^n} P(x) \log P(x). \quad (5)$$

The assumption $B > H$ is necessary for the existence of a sequence of length functions $\{L_n(\cdot)\}_{n \geq 1}$ such that the probability of overflow, (3), will vanish as $n \rightarrow \infty$. On the other hand, the maximum value of B for which the problem is interesting is $\log X$, beyond which (3) can be made zero in a trivial manner.

Let $U_{LZ}(x)$ be the length function of the Lempel–Ziv [4] code, that is, $U_{LZ}(x)$ is the number of bits associated with the Lempel–Ziv codeword of a sequence x . The following theorem establishes the asymptotic optimality of the Lempel–Ziv algorithm in the sense of minimum probability of length overflow (3).

Theorem 1: For any length function $L_n(\cdot)$, every $B \in (H, \log X)$, every finite-state source P , and all large n ,

$$\Pr\{n^{-1}U_{LZ}(x) > B + \epsilon(n)\} \leq (1 + n^2 2^{-n/\sqrt{\log n}}) \Pr\{n^{-1}L_n(x) > B\}, \quad (6)$$

where $\epsilon(n) = O(1/\sqrt{\log n})$ is a positive sequence depending on X and S .

The proof appears in Section V.

Theorem 1 states that the Lempel–Ziv algorithm attains the best tail behavior of the distribution of codeword lengths, i.e., the best large deviations performance. In [4] Ziv and Lempel have shown their algorithm to yield the shortest length, uniformly for every sufficiently long sequence x , among all information lossless *finite-state* encoders. Here the result is dual to [4] in the sense that the *source* is limited to be finite-state, but *any* competing

lossless encoder is allowed, not necessarily a finite-state encoder.

III. UNIFILAR SOURCES

The convergence of $\epsilon(n)$ in (6) can be accelerated if P is known to belong to the subclass $\mathcal{P}_U \subset \mathcal{P}_S$ of *unifilar* sources. For these sources the state s_i , at time instant i , obeys the recursion

$$s_i = f(x_i, s_{i-1}), \quad (7)$$

where $f: X \times S \rightarrow S$ is a known deterministic mapping. Clearly, in this case, given s_0 , which is fixed, one can reconstruct the state sequence s recursively upon observing x . If, in addition, f is a 1–1 map from X to S given its second argument, then s uniquely determines x as well, by the inverse map. An important special case of this model is a k th order Markov source where $s_i = (x_i, x_{i-1}, \dots, x_{i-k+1})$.

We now demonstrate that by using the MDL universal code for the subclass \mathcal{P}_U , the term $\epsilon(n)$ can be reduced to $O(\log n/n)$. Let,

$$q_x(x, s) = \frac{1}{n} \sum_{i=1}^n \delta(x_i = x, s_{i-1} = s), \quad x \in X, s \in S, \quad (8)$$

where $\delta(x_i = x, s_{i-1} = s)$ is the indicator function for $x_i = x$ jointly with $s_{i-1} = s$. Also, let $q_x(s) = \sum_{x \in X} q_x(x, s)$ and

$$q_x(x|s) = \begin{cases} q_x(x, s)/q_x(s), & q_x(s) > 0 \\ 0, & q_x(s) = 0. \end{cases} \quad (9)$$

We denote by Q_x the empirical distribution $Q_x \triangleq \{q_x(x, s), x \in X, s \in S\}$, and define the empirical entropy as

$$H(Q_x) = - \sum_{s \in S} \sum_{x \in X} q_x(x, s) \log q_x(x|s), \quad (10)$$

and the Kullback–Leibler divergence between the empirical distribution Q_x and the source P is defined as

$$D(Q_x \| P) \triangleq \sum_{s \in S} \sum_{x \in X} q_x(x, s) \log \frac{q_x(x|s)}{p(x|s)}, \quad (11)$$

where $p(x|s) = \Pr\{x_i = x | s_{i-1} = s\}$. Note that in the unifilar case considered here,

$$P(x) = \exp_2\{-n[H(Q_x) + D(Q_x \| P)]\}. \quad (12)$$

It is not difficult to show that, for the subclass of unifilar finite-state sources, it is possible to attain an exponentially vanishing probability of codeword length overflow, that is,

$$\lim_{n \rightarrow \infty} n^{-1} \log \Pr\{n^{-1}L_n(x) > B\} < 0. \quad (13)$$

To do this, the first step is to show, by a technique similar to that in Section V, that for an arbitrary length function $L_n(x)$ there exists another length function $L'_n(x)$, depending on x only through the empirical distribution $Q_x =$

$\{q_x(x, s), x \in X, s \in S\}$, such that

$$n^{-1} \log \Pr \left\{ n^{-1} L'_n(x) > B + O\left(\frac{\log n}{n}\right) \right\} \\ \leq n^{-1} \log \Pr \{ n^{-1} L_n(x) > B \} + O\left(\frac{\log n}{n}\right). \quad (14)$$

Therefore, without loss in asymptotic performance, the optimal length function can be assumed to depend on x only through Q_x . Let T_x be the type of x , namely,

$$T_x = \{x' \in X^n: Q_{x'} = Q_x\}. \quad (15)$$

Bounds on the probability of T_x are given in the following Lemma.

Lemma 1: For every $P \in P_U$, where f is a 1-1 mapping given its second argument, and for every $x \in X^n$,

$$\left| \frac{1}{n} \log \Pr \{T_x\} + D(Q_x \| P) \right| \leq \epsilon(n), \quad (16)$$

where $\epsilon(n) = O(n^{-1} \log n)$ is independent of x .

Proof: The proof is a straightforward extension of [10, Lemma 1] where this argument was proved for the specific case of a k th order Markov source.

From (12), (16), and the fact that all sequences in T_x are equally likely, it follows that

$$\left| \frac{1}{n} \log |T_x| - H(Q_x) \right| = \left| \frac{1}{n} \log \frac{\Pr \{T_x\}}{P(x)} - H(Q_x) \right| \leq \epsilon(n). \quad (17)$$

Following (14), let all members of T_x have the same codeword length $L_n(x)$. Clearly, for a uniquely decipherable code, this length must be at least as large as the base 2 logarithm of the cardinality of T_x , namely,

$$L_n(x) \geq \log |T_x| \geq nH(Q_x) - n\epsilon(n), \quad (18)$$

where we have used (17). Note, that (18) holds even if $L_n(\cdot)$ does not satisfy the Kraft inequality (2). The only assumption is that $L_n(\cdot)$ is a length function associated with a 1-1 mapping from X^n to a set of codewords. From (18), we have

$$n^{-1} \log \Pr \{ n^{-1} L_n(x) \geq B \} \\ \geq n^{-1} \log \Pr \{ H(Q_x) \geq B + \epsilon(n) \} \\ = n^{-1} \log \sum_{T_x: H(Q_x) \geq B + \epsilon(n)} \Pr \{T_x\} \\ \geq n^{-1} \log \max_{\{Q_x: H(Q_x) \geq B + \epsilon(n)\}} \\ \cdot \exp_2 \{ -n [D(Q_x \| P) + \epsilon(n)] \} \\ = - \min_{\{Q_x: H(Q_x) \geq B + \epsilon(n)\}} D(Q_x \| P) - \epsilon(n), \quad (19)$$

where again we have used (16) for a lower bound on $\Pr \{T_x\}$. Since $D(\cdot \| P)$ and $H(\cdot)$ are continuous functions and since the set of rational empirical distributions $\{Q_x\}$ becomes dense in P_U as $n \rightarrow \infty$, the right-most side of (19) tends to a constant D (see also [14]) given by,

$$D \triangleq \min_{\{Q: H(Q) \geq B\}} D(Q \| P). \quad (20)$$

An alternative expression [11], [12] for D in the memory-

less case ($S = 1$) is

$$D = \max_{\lambda > 0} \lambda (B - H_\lambda), \quad (21)$$

where H_λ is the Rényi entropy [2] given by

$$H_\lambda \triangleq \frac{\lambda + 1}{\lambda} \log \sum_{x \in X} p(x)^{1/(1+\lambda)}, \quad (22)$$

$\{p(x)\}_{x \in X}$ being the letter probabilities. Clearly, to attain the optimal overflow exponent D , $L_n(x)$ should be as close as possible to $nH(Q_x)$, which is essentially the lower bound (18). In particular, consider a simple universal code which consists of $\log(n+1)$ bits (neglecting roundoff terms) to encode each one of the $S(X-1)$ source parameters $\{p(x|s)\}$ estimated from x , followed by $-\log Q_x(x)$ bits for Huffman coding of x with respect to the source estimate Q_x . The resulting length function is

$$L_n^*(x) = -\log Q_x(x) + S(X-1) \log(n+1) \\ = nH(Q_x) + S(X-1) \log(n+1). \quad (23)$$

This code is similar to the well known MDL universal code (see e.g., [8]), which was shown to be asymptotically optimal in the sense of uniformly minimizing the expected redundancy. A Bayesian approach [7], in which a Huffman code is designed with respect to a (uniform) mixture of all sources in the class, also yields a length function which is asymptotically equivalent to the MDL. Similarly to (19), by using (17) for an upper bound on the cardinality of T_x , and the fact that the number of distinct types is smaller than $(n+1)^{XS}$, one arrives at the upper bound

$$n^{-1} \log \Pr \{ n^{-1} L_n^*(x) \geq B \} \\ = n^{-1} \log \sum_{T_x: L_n^*(x) \geq Bn} \Pr \{T_x\} \\ \leq n^{-1} \log \sum_{Q_x: H(Q_x) \geq B - \epsilon(n)} \\ \cdot \exp_2 \{ -n [D(Q_x \| P) - \epsilon(n)] \} \\ \leq n^{-1} \log(n+1)^{XS} \\ \cdot \exp_2 \left\{ -n \left[\min_{\{Q_x: H(Q_x) \geq B - \epsilon(n)\}} D(Q_x \| P) - \epsilon(n) \right] \right\} \\ \leq - \min_{\{Q_x: H(Q_x) \geq B - \epsilon(n)\}} D(Q_x \| P) \\ + \epsilon(n) + \frac{1}{n} XS \log(n+1), \quad (24)$$

where $\epsilon'(n) = n^{-1} S(X-1) \log(n+1)$. Clearly, the right-hand side of (24) also tends to D as $n \rightarrow \infty$. This establishes the asymptotic optimality of $L_n^*(x)$ in the sense of (3).

IV. DISCUSSION

Note that no matter whether the true source $P \in P_U$ is given or not, the optimal length function $L_n^*(x)$ depends neither on P nor on B . In other words, knowledge of P or B cannot improve the asymptotic overflow expo-

nent D . Only a *second-order* improvement, associated with the rate of convergence to D , can be attained if it is known *a priori* that P belongs to a certain subset of P_U . For instance, if the source P is known to be memoryless, then the second term in (23) can be reduced to $(X-1)\log(n+1)$, resulting in a slightly faster convergence to D .

It should be pointed out that, rather than truncating variable length codewords $L_n(x)$ into a fixed length of Bn bits, one can consider a codebook for the set G_n of the 2^{Bn} most likely input n -tuples x , and use Bn bits to represent the codeword index. In this case, information is lost whenever $x \in G_n^c$. It is shown in [13, Theorem 2.15, p. 37]) that the best achievable exponential rate of the probability of information loss, $\lim_{n \rightarrow \infty} [-n^{-1} \log \Pr\{G_n^c\}]$, again equals D as defined in (20). (See also Marton [14, (4)] for the distortionless case.) This means that there is no loss in asymptotic optimality in that sense, when using the truncation approach.

Since Rényi's entropy is attained by the length function

$$L_n(x) = \left\lceil -\log \frac{P(x)^{1/(1+\lambda)}}{\sum_{x \in X^n} P(x)^{1/(1+\lambda)}} \right\rceil, \quad (25)$$

where $P(x)$ in turn depends on x only through Q_x , the same consideration as in (14) can be used to show that $L_n^*(x)$, defined in (23), asymptotically minimizes $(n\lambda)^{-1} \log E(2^{\lambda L_n(x)})$, uniformly for any unifilar finite state source, and any $\lambda > 0$ (see Appendix B). Following [15] (see also Appendix A), the chain of inequalities

$$\begin{aligned} U_{LZ}(x) &\leq -\log \max_{P \in P_S} P(x) + n\delta(n) \\ &\leq -\log \max_{P \in P_U} P(x) + n\delta(n) \\ &= nH(Q_x) + n\delta(n) \\ &\leq L_n^*(x) + n\delta(n), \end{aligned} \quad (26)$$

which holds for any $x \in X^n$ with $\delta(n) \rightarrow 0$ uniformly for every x , implies that the LZ algorithm is asymptotically optimal as well, in that sense. It has been shown in [15] that $\delta(n) = O(\log \log n / \log n)$, thus, the rate of convergence to Rényi's entropy, in the LZ case, might be slower than the $O(n^{-1} \log n)$ rate of the MDL universal code. These facts might have application to universal coding with minimum buffer overflow probability. Specifically, suppose that n is fixed and we use a buffer of size K to convert the variable rate coding to a fixed rate of R bits per input symbol, in the following manner. For each input block $x \in X^n$, a codeword of length $L_n(x)$ is stored sequentially in the buffer and at the same time, a previously stored block of size Rn is removed from the buffer and transmitted. For the memoryless case, it has been shown [16], [17] that if $E\{L_n(x)\} < Rn$, then the probability of buffer overflow P_{bo} can be made exponentially small as $K \rightarrow \infty$, and the maximum buffer overflow exponent $\lim_{K \rightarrow \infty} (-K^{-1} \log P_{bo})$ is attained by the length function (25), with $\lambda = \lambda_0$ being the unique positive solution of the

equation $H_\lambda = R$, where H_λ is the Rényi entropy given in (22). This solution λ_0 also equals the maximum attainable overflow exponent. Clearly, this length function depends on $P(\cdot)$ and R . However, it is easy to show (see Appendix B) that, if, one instead uses the universal code $L_n^*(x)$ (with $S=1$), which depends neither on $P(\cdot)$ nor on R , the overflow exponent degrades but only by a quantity which vanishes as fast as $n^{-1} \log n$.

Another aspect of the MDL and the LZ algorithm is asymptotic competitive optimality. Cover [9] has shown that for a known memoryless source, the Huffman code $L_n^H(x)$ is within one bit optimal in a competitive sense, namely, for any competing code $L_n(x)$,

$$\Pr\{L_n^H(x) > L_n(x) + 1\} \leq \Pr\{L_n^H(x) < L_n(x) + 1\}. \quad (27)$$

This means that, within a normalized redundancy term of the order of $1/n$, the Huffman code provides most of the time a codeword shorter than that of any other competing code. We now show that this result extends to universal coding, again, by using $L_n^*(x)$ (or the MDL universal code), at the expense of increasing the redundancy term to $O(n^{-1} \log n)$. Specifically, we claim that for any unifilar finite-state source and any competing code $L_n(\cdot)$,

$$\Pr\{L_n^*(x) > L_n(x) + n\epsilon(n)\} \leq \Pr\{L_n^*(x) < L_n(x) + n\epsilon(n)\}, \quad (28)$$

where $\epsilon(n) = n^{-1}[L_n^*(x) - nH(Q_x)] = O(n^{-1} \log n)$. To show that (28) holds, we use a technique similar to [9]. First observe that,

$$\begin{aligned} &\Pr\{L_n^*(x) > L_n(x) + n\epsilon(n)\} \\ &\quad - \Pr\{L_n^*(x) < L_n(x) + n\epsilon(n)\} \\ &= E\{\text{sgn}[L_n^*(x) - L_n(x) - n\epsilon(n)]\}. \end{aligned} \quad (29)$$

Hence, it is sufficient to prove that the expectation on the right-hand side of (29) is nonpositive. To see this, we use the inequality $\text{sgn}(k) \leq \exp_2(k) - 1$, (k -integer), and obtain

$$\begin{aligned} &E\{\text{sgn}[L_n^*(x) - L_n(x) - n\epsilon(n)]\} \\ &\leq E\{\exp_2[L_n^*(x) - L_n(x) - n\epsilon(n)]\} - 1 \\ &= \sum_x P(x) \exp_2[L_n^*(x) - L_n(x) - n\epsilon(n)] - 1. \end{aligned} \quad (30)$$

Finally, since $P(x) \leq \exp_2[-nH(Q_x)]$ for every $P \in P_U$ by (12), we get

$$\begin{aligned} &E\{\text{sgn}[L_n^*(x) - L_n(x) - n\epsilon(n)]\} \\ &\leq \sum_x \exp_2[-nH(Q_x)] \\ &\quad \cdot \exp_2[L_n^*(x) - L_n(x) - n\epsilon(n)] - 1 \\ &\leq \sum_x \exp_2[-nH(Q_x)] \exp_2\{[nH(Q_x)] - L_n(x)\} - 1 \\ &\leq \sum_x \exp_2[-L_n(x)] - 1 \leq 0, \end{aligned} \quad (31)$$

where the last step follows from the assumption that $L_n(\cdot)$ satisfies the Kraft inequality. Note that this result extends easily to any parametric class of sources, under

fairly mild regularity conditions, where $L_n^*(x)$ should be interpreted as a Huffman code with respect to the maximum likelihood estimator of the source, followed by $O(\log n)$ bits to encode a quantized version of the parameter estimate. Again, inequality (26) implies that for the class of finite state sources, $L_n^*(x)$ can be replaced by $U_{LZ}(x)$ at the expense of increasing the normalized redundancy to $O(\log \log n / \log n)$ bits per source letter.

V. PROOF OF THEOREM 1

Assume that l divides n and parse x into n/l vectors of length l :

$$x = x_1^l, x_2^l, \dots, x_{n/l}^l, \quad x_i^l \in X^l. \quad (32)$$

Let s^l denote the sequence of initial states of the resulting phrases $\{x_i^l\}$,

$$s^l = s_0^l, s_1^l, \dots, s_{n/l}^l. \quad (33)$$

Fix s^l and let $K(x|s^l)$ denote the set of n -vectors x' obtained from x by permuting phrases x_i^l and x_j^l of x , for which $s_i^l = s_j^l$ and $s_{i+1}^l = s_{j+1}^l$. Clearly, for every $x' \in K(x|s^l)$ we have from (1) that $P(x, s^l) = P(x', s^l)$, and hence $K(x|s^l)$ can be thought of as a conditional type.

Given an arbitrary length function $L_n(\cdot)$ satisfying the Kraft inequality, we first find a modified length function $L'_n(\cdot|s^l)$, depending on a given s^l , which satisfies Kraft's inequality as well, and which has an additional property of mapping "large portions" of each conditional type $K(x|s^l)$ onto the same length. This property will be useful later in the derivation of the desired lower bound on the length overflow probability. We first define some notation.

Let $L_n(\cdot)$ be as above and $L_{\max} \triangleq \max_x L_n(x)$. Next, let $I \triangleq \{1, 2, \dots, \lfloor Bn \rfloor, L_{\max}\}$, where $\lfloor Bn \rfloor$ is the largest integer not exceeding Bn (floor function), and define a partition $\Omega = \{\Omega_j\}_{j \in I}$ of X^n by $\Omega_j = \{x: L_n(x) = j\}$ for $j = 1, \dots, \lfloor Bn \rfloor$ and $\Omega_{L_{\max}} = \{x: L_n(x) > Bn\}$. Let $\{\rho_n\}_{n \geq 1}$ be a positive sequence with the following properties:

$$\rho_n \leq 1, \quad (34a)$$

$$\lim_{n \rightarrow \infty} \rho_n = 0, \quad (34b)$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \rho_n = 0. \quad (34c)$$

The exact form of $\{\rho_n\}_{n \geq 1}$ will be chosen later. Define

$$j^* = \max \left\{ j \in I: |\Omega_j \cap K(x|s^l)| = \max_{i \in I} |\Omega_i \cap K(x|s^l)| \right\}. \quad (35)$$

As the number of classes $\{\Omega_j\}$ never exceeds $(Bn + 1)$, clearly,

$$|K(x|s^l) \cap \Omega_{j^*}| \geq \frac{1}{Bn + 1} |K(x|s^l)|. \quad (36)$$

It follows from (34a) and (36) that the set

$$V_{x^l} \triangleq \left\{ j \in I: |K(x|s^l) \cap \Omega_j| \geq \frac{\rho_n}{Bn + 1} |K(x|s^l)| \right\} \quad (37)$$

is nonempty. In particular, $j^* \in V_{x^l}$. Henceforth the sub-

script x^l will be omitted whenever clear from the context and we define $V_* = V - \{j^*\}$.

Next, define a modification $\Omega' = \{\Omega'_j\}_{j \in I}$ of Ω as follows. If $x \in \Omega_j$ with $j \in V$ then also $x \in \Omega'_{j^*}$. Otherwise, if $x \in \Omega_j$ with $j \in V^c$ then $x \in \Omega'_j$. Finally, define the modified length function as follows. $L'_n(x|s^l) = j$ for all $x \in \Omega'_j$, $j \neq j^*$. For $x \in \Omega'_{j^*}$, $L'_n(x|s^l) = j^* + \lceil \log(Bn + 1) \rceil$, where $\lceil a \rceil$ is the least integer not less than a (ceiling function). Thus, sequences in Ω_j , for any $j \in V_*$, are unaffected by the modification and all other sequences will have modified codewords of length $j^* + \lceil \log(Bn + 1) \rceil$ bits.

Comment: Note, that even for a fixed s^l , not all sequences in Ω'_j have modified codewords of the same length $L'_n(x|s^l)$. The reason is that j^* depends not only on s^l , but also on the conditional type $K(x|s^l)$, and therefore, within a given Ω'_j , for some of the types, j might coincide with j^* , in which case, $L'_n(x|s^l) = j + \lceil \log(Bn + 1) \rceil$, and for other types, $L'_n(x|s^l) = j$. However, all sequences from the same type within given Ω'_j and s^l , are assigned to modified codewords of the same length.

We now show that $L'_n(\cdot|s^l)$ satisfies Kraft's inequality for every $s^l \in S^{n/l+1}$. Since the original length function $L_n(\cdot)$ satisfies Kraft's inequality by assumption and since for every s^l , the set of all conditional types $\{K(x|s^l)\}$ form a partition of X^n , we obtain

$$\begin{aligned} 1 &\geq \sum_{x \in X^n} 2^{-L_n(x)} = \sum_{\{K(x|s^l)\}_{s^l} \subseteq X^n} \sum_{x' \in K(x|s^l)} 2^{-L_n(x')} \\ &\geq \sum_{\{K(x|s^l)\}} \left[\sum_{x' \in \Omega_{j^*} \cap K(x|s^l)} 2^{-j^*} + \sum_{j \in V_*, x' \in \Omega_j \cap K(x|s^l)} \sum 2^{-j} \right] \\ &= \sum_{\{K(x|s^l)\}} \left[|\Omega_{j^*} \cap K(x|s^l)| 2^{-j^*} + \sum_{j \in V_*, x' \in \Omega_j \cap K(x|s^l)} \sum 2^{-j} \right] \\ &\geq \sum_{\{K(x|s^l)\}} \left[\frac{1}{Bn + 1} |K(x|s^l)| 2^{-j^*} \right. \\ &\quad \left. + \sum_{j \in V_*, x' \in \Omega_j \cap K(x|s^l)} \sum 2^{-j} \right] \\ &\geq \sum_{\{K(x|s^l)\}} \left[\frac{1}{Bn + 1} \sum_{j \in V_*^c} |\Omega_j \cap K(x|s^l)| 2^{-j^*} \right. \\ &\quad \left. + \sum_{j \in V_*, x' \in K(x|s^l) \cap \Omega_j} \sum 2^{-j} \right] \\ &\geq \sum_{\{K(x|s^l)\}} \left[\sum_{j \in V_*^c} \sum_{x' \in \Omega_j \cap K(x|s^l)} 2^{-(j^* + \lceil \log(Bn + 1) \rceil)} \right. \\ &\quad \left. + \sum_{j \in V_*, x' \in K(x|s^l) \cap \Omega_j} \sum 2^{-j} \right] \\ &= \sum_{x \in X^n} 2^{-L'_n(x|s^l)}. \quad (38) \end{aligned}$$

Let $\Pr\{L_n(x) > Bn, s^l\}$ denote the joint probability under P (1) of the event $\{L_n(x) > Bn\}$ and the occurrence of a state sequence s^l . For every fixed s^l , we now lower bound

$\Pr\{L_n(\mathbf{x}) > Bn, s^l\}$ by means of the auxiliary length function $L'_n(\mathbf{x}|s^l)$. Let

$$U = \{j: j > Bn\} \quad (39)$$

$$U' = \{j: j > Bn + \lceil \log(Bn + 1) \rceil\} \quad (40)$$

$$\delta_U(j) = \begin{cases} 1, & j \in U \\ 0, & j \in U^c \end{cases} \quad (41)$$

and $\delta_{U'}(j)$ be defined similarly to $\delta_U(j)$. Clearly, $\delta_U(j) = \delta_{U'}(j + \lceil \log(Bn + 1) \rceil)$. Hence,

$$\begin{aligned} & \Pr\{L_n(\mathbf{x}) > Bn, s^l\} \\ &= \sum_{j \in U} \sum_{\mathbf{x} \in \Omega_j} P(\mathbf{x}, s^l) \\ &= \sum_{\{K(\mathbf{x}|s^l)\}} \sum_{j \in U} \sum_{\mathbf{x}' \in \Omega_j \cap K(\mathbf{x}|s^l)} P(\mathbf{x}', s^l) \\ &\geq \sum_{\{K(\mathbf{x}|s^l)\}} \left[\delta_U(j^*) |\Omega_{j^*} \cap K(\mathbf{x}|s^l)| P(\mathbf{x}, s^l) \right. \\ &\quad \left. + \sum_{j \in U \cap V_*} |\Omega_j \cap K(\mathbf{x}|s^l)| P(\mathbf{x}, s^l) \right]. \quad (42) \end{aligned}$$

By the construction of $L'_n(\cdot|s^l)$,

$$\begin{aligned} & |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| \\ &= |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| + \sum_{j \in V^c} |K(\mathbf{x}|s^l) \cap \Omega_j| \\ &< |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| + (Bn + 1) \cdot \frac{\rho_n}{Bn + 1} |K(\mathbf{x}|s^l)| \\ &\leq |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| + \rho_n (Bn + 1) |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| \\ &= |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| \cdot [1 + \rho_n (Bn + 1)], \quad (43) \end{aligned}$$

or, equivalently,

$$|K(\mathbf{x}|s^l) \cap \Omega_{j^*}| \geq \frac{|K(\mathbf{x}|s^l) \cap \Omega_{j^*}|}{1 + \rho_n (Bn + 1)}. \quad (44)$$

By plugging (44) into (42) and using the fact that $\Omega_j \cap K(\mathbf{x}|s^l) = \Omega'_j \cap K(\mathbf{x}|s^l)$ for every $j \in V_*$, we can further lower bound $\Pr\{L_n(\mathbf{x}) > Bn, s^l\}$ as follows.

$$\begin{aligned} & \Pr\{L_n(\mathbf{x}) > Bn, s^l\} \\ &\geq \sum_{\{K(\mathbf{x}|s^l)\}} \left[[1 + \rho_n (Bn + 1)]^{-1} |K(\mathbf{x}|s^l) \cap \Omega_{j^*}| \right. \\ &\quad \cdot P(\mathbf{x}, s^l) \delta_{U'}(j^* + \lceil \log(Bn + 1) \rceil) \\ &\quad \left. + \sum_{j \in U' \cap V_*} |\Omega'_j \cap K(\mathbf{x}|s^l)| P(\mathbf{x}, s^l) \right] \\ &\geq [1 + \rho_n (Bn + 1)]^{-1} \sum_{\{K(\mathbf{x}|s^l)\}} \left[|K(\mathbf{x}|s^l) \cap \Omega_{j^*}| P(\mathbf{x}, s^l) \right. \\ &\quad \cdot \delta_{U'}(j^* + \lceil \log(Bn + 1) \rceil) \\ &\quad \left. + \sum_{j \in U' \cap V_*} |\Omega'_j \cap K(\mathbf{x}|s^l)| P(\mathbf{x}, s^l) \right] \\ &= [1 + \rho_n (Bn + 1)]^{-1} \Pr\{L'_n(\mathbf{x}|s^l) > Bn \\ &\quad + \lceil \log(Bn + 1) \rceil, s^l\}. \quad (45) \end{aligned}$$

Now, by construction of Ω'_j , for any nonempty Ω'_j , we have $|K(\mathbf{x}|s^l) \cap \Omega'_j| \geq (Bn + 1)^{-1} \rho_n |K(\mathbf{x}|s^l)|$. Since $L'_n(\cdot|s^l)$ is a uniquely decipherable code for any fixed s^l and since all n -tuples in $K(\mathbf{x}|s^l) \cap \Omega'_j$ have modified codewords of the same length, then, similarly to (18), this length must be at least as large as the base 2 logarithm of the cardinality of this set. Thus,

$$\begin{aligned} L'_n(\mathbf{x}|s^l) &\geq \log |K(\mathbf{x}|s^l) \cap \Omega'_j| \\ &\geq \log |K(\mathbf{x}|s^l)| + \log \frac{\rho_n}{Bn + 1}. \quad (46) \end{aligned}$$

It is shown in [5, Lemma 1], [15] (see also Appendix A) that one can choose a sequence $l_n \rightarrow \infty$ such that

$$\log |K(\mathbf{x}|s^{l_n})| \geq U_{LZ}(\mathbf{x}) - n\delta(n), \quad (47)$$

where $\delta(n) = O(1/\sqrt{\log n})$ uniformly for any \mathbf{x} . By substituting (47) into (46) with $l = l_n$, we get

$$L'_n(\mathbf{x}|s^{l_n}) \geq U_{LZ}(\mathbf{x}) - n\delta(n) - \log \left(\frac{Bn + 1}{\rho_n} \right). \quad (48)$$

From (45) and (48) we obtain the lower bound

$$\begin{aligned} & \Pr\{L_n(\mathbf{x}) > Bn, s^{l_n}\} \\ &\geq [1 + \rho_n (Bn + 1)]^{-1} \Pr\{L'_n(\mathbf{x}|s^{l_n}) > Bn \\ &\quad + \lceil \log(Bn + 1) \rceil, s^{l_n}\} \\ &\geq [1 + \rho_n (Bn + 1)]^{-1} \\ &\quad \cdot \Pr\left\{ U_{LZ}(\mathbf{x}) > Bn + \lceil \log(Bn + 1) \rceil \right. \\ &\quad \left. + n \left[\delta(n) + \frac{1}{n} \log \left(\frac{Bn + 1}{\rho_n} \right) \right], s^{l_n} \right\}. \quad (49) \end{aligned}$$

By choosing $\rho_n = \min\{1, (Bn + 1)^{-1} n^2 2^{-n/\sqrt{\log n}}\}$, which satisfies (34), the term $n^{-1} \log[(Bn + 1)/\rho_n]$ on the right-most side of (49) becomes $O(1/\sqrt{\log n})$, which is the same order as $\delta(n)$. Finally, by summing the left-most and the right-most sides of (49), over all possible state sequences s^{l_n} , we complete the proof of Theorem 1.

ACKNOWLEDGMENT

The author wishes to thank Jacob Ziv, Toby Berger, Aaron D. Wyner, Amir Dembo, and Esther Levin for helpful discussions and suggestions. Useful comments made by the anonymous referees are greatly appreciated.

APPENDIX A

Proof of (47) and (26): We first prove (47) and then, based on this proof, we derive (26). Let $K(\mathbf{x}|s^l)$ be defined as in the first paragraph of Section V. Define the following empirical distributions.

$$q_{\mathbf{x}^l s^l}^l(u, s, s') = \frac{1}{n} \sum_{i=1}^{n/l} \delta(x_i^l = u, s_{i-1}^l = s, s_i^l = s'), \quad u \in \mathcal{X}^l, s, s' \in \mathcal{S}, \quad (\text{A.1})$$

where $s_{n/l+1}^l \triangleq s_0$, $\delta(x_i^l = u, s_{i-1}^l = s, s_i^l = s')$ is the indica-

tor function for $x_i^l = u$, $s_{i-1}^l = s$ and $s_i^l = s'$. Let

$$q_x^l(u) = \sum_{s, s' \in S} q_{xs'}^l(u, s, s') \quad (\text{A.2})$$

$$q_{s'}^l(s, s') = \sum_{u \in X^l} q_{xs'}^l(u, s, s') \quad (\text{A.3})$$

$$q_{x|s'}^l(u|s, s') = \begin{cases} q_{xs'}^l(u, s, s')/q_{s'}^l(s, s'), & q_{s'}^l(s, s') > 0 \\ 0, & q_{s'}^l(s, s') = 0 \end{cases} \quad (\text{A.4})$$

$$H(q_{xs'}^l) = - \sum_{s, s' \in S} \sum_{u \in X^l} q_{xs'}^l(u, s, s') \log q_{xs'}^l(u, s, s') \quad (\text{A.5})$$

$$H(q_{s'}^l) = - \sum_{s, s' \in S} q_{s'}^l(s, s') \log q_{s'}^l(s, s') \quad (\text{A.6})$$

$$H(q_x^l) = - \sum_{u \in X^l} q_x^l(u) \log q_x^l(u). \quad (\text{A.7})$$

$$H(q_{x|s'}^l) = - \sum_{s, s' \in S} \sum_{u \in X^l} q_{xs'}^l(u, s, s') \log q_{x|s'}^l(u|s, s'). \quad (\text{A.8})$$

It is easy to show that

$$\begin{aligned} H(q_{x|s'}^l) &= H(q_{xs'}^l) - H(q_{s'}^l) \geq H(q_x^l) - H(q_{s'}^l) \\ &\geq H(q_x^l) - 2 \log S. \end{aligned} \quad (\text{A.9})$$

Let $n(u, s, s') = l^{-1} n q_{xs'}^l(u, s, s')$ and $n(s, s') = \sum_{u \in X^l} n(u, s, s')$. By the definition of $K(x|s')$, we have

$$|K(x|s')| = \prod_{s, s' \in S} \frac{n(s, s')!}{\prod_{u \in X^l} n(u, s, s')!}. \quad (\text{A.10})$$

By using the inequality $n \log n - n \log e \leq \log n! \leq n \log n$, we obtain

$$\begin{aligned} &\frac{1}{n} \log |K(x|s')| \\ &= \frac{1}{n} \sum_{s, s' \in S} \left[\log n(s, s')! - \sum_{u \in X^l} \log n(u, s, s')! \right] \\ &\geq \frac{1}{n} \sum_{s, s' \in S} \left[n(s, s') \log n(s, s') - n(s, s') \log e \right. \\ &\quad \left. - \sum_{u \in X^l} n(u, s, s') \cdot \log n(u, s, s') \right] \\ &= \frac{1}{l} H(q_{x|s'}^l) - \frac{1}{l} \log e \\ &\geq \frac{1}{l} [H(q_x^l) - 2 \log S - \log e] \\ &\triangleq \frac{1}{l} [H(q_x^l) - C]. \end{aligned} \quad (\text{A.11})$$

It is known [8] that there exists a uniquely decipherable

code with a length function

$$L_n(x) \leq \frac{n}{l} [H(q_x^l) + 1]. \quad (\text{A.12})$$

Furthermore, this coding algorithm can be implemented by an encoder with X^l states. Hence, by [4, Theorem 1],

$$L_n(x) \geq (c(x) + X^{2l}) \log \frac{c(x) + X^{2l}}{4X^{2l}}, \quad (\text{A.13})$$

where $c(x)$ is the maximum number of distinct strings (phrases) in x [4]. On the other hand, by [4, Theorem 2],

$$U_{LZ}(x) \leq (c(x) + 1) \log [2X(c(x) + 1)]. \quad (\text{A.14})$$

Hence, by (A.11)–(A.14),

$$\begin{aligned} &\frac{1}{n} \log |K(x|s')| \\ &\geq \frac{1}{l} H(q_x^l) - \frac{C}{l} \\ &\geq \frac{1}{n} (c(x) + X^{2l}) \log \frac{c(x) + X^{2l}}{4X^{2l}} - \frac{C+1}{l} \\ &\geq \frac{1}{n} U_{LZ}(x) - \frac{c(x) + X^l}{n} \log 8X^{2l+1} - \frac{C+1}{l} \\ &\geq \frac{1}{n} U_{LZ}(x) - \frac{c(x) + X^l}{n} \log X^{4l+2} - \frac{C+1}{l} \\ &\triangleq \frac{1}{n} U_{LZ}(x) - \delta_0(n, l), \end{aligned} \quad (\text{A.15})$$

where in the last inequality we have used the fact that $X \geq 2$ and $l \geq 1$. By [4, (6)], $c(x) \leq n \log X / [(1 - \epsilon_n) \log n]$ for some $\epsilon_n \rightarrow 0$. Thus, for n and l sufficiently large, there exist constants C_1 , C_2 , and C_3 such that

$$\delta_0(n, l) \leq \frac{C_1}{l} + \frac{C_2 l}{\log n} + \frac{C_3 l X^l}{n}. \quad (\text{A.16})$$

By choosing $l = l_n = \lceil \sqrt{\log n} \rceil$, the right-hand side of (A.16) becomes $O(1/\sqrt{\log n})$, completing the proof of (47).

Finally, to prove (26), we return to (A.15) and observe that, for any s' ,

$$1 \geq \sum_{x' \in K(x|s')} P(x', s') = |K(x|s')| P(x, s'), \quad (\text{A.17})$$

and hence,

$$\begin{aligned} P(x) &= \sum_{s'} P(x, s') \\ &\leq \sum_{s'} |K(x|s')|^{-1} \\ &\leq S^{n/l+1} \exp_2 \{ -U_{LZ}(x) + n\delta_0(n, l) \} \\ &\leq \exp_2 \left\{ -U_{LZ}(x) + n \left[\delta_0(n, l) + \left(\frac{1}{n} + \frac{1}{l} \right) \log S \right] \right\} \\ &\triangleq \exp_2 \{ -U_{LZ}(x) + n\delta_1(n, l) \}. \end{aligned} \quad (\text{A.18})$$

Hence, by choosing again $l_n = \lceil \sqrt{\log n} \rceil$, we get

$$U_{LZ}(x) \leq -\log P(x) + n\delta'(n), \quad (\text{A.19})$$

where $\delta'(n) = \delta_1(n, l_n) = O(1/\sqrt{\log n})$. Tighter bounds

where $\delta'(n) = O(\log \log n / \log n)$ can be found in [15]. Since (A.19) holds for any finite-state source P , the right-hand side can be minimized with respect to P , resulting in

$$U_{LZ}(x) \leq -\log \max_{P \in \mathcal{P}_S} P(x) + n\delta'(n). \quad (\text{A.20})$$

This completes the proof of the first inequality in (26). The remaining inequalities in (26) are straightforward. \square

APPENDIX B

UNIVERSAL CODING WITH MINIMUM PROBABILITY OF BUFFER OVERFLOW

Let $\delta_n \triangleq n^{-1}X \log(n+1)$. We first upper bound $(n\lambda)^{-1} \log E(2^{\lambda L_n^*(x)})$.

$$\begin{aligned} & \frac{1}{n\lambda} \log E(2^{\lambda L_n^*(x)}) \\ & \leq \frac{1}{n\lambda} \log E(2^{n\lambda H(Q_x)}) + \delta_n \\ & = \frac{1}{n\lambda} \log \left[\sum_{x \in X^n} P(x) 2^{\lambda n H(Q_x)} \right] + \delta_n \\ & = \frac{1}{n\lambda} \log \left[\sum_{T_x \subset X^n} |T_x| 2^{-n[H(Q_x) + D(Q_x \| P)]} 2^{\lambda n H(Q_x)} \right] + \delta_n \\ & \leq \frac{1}{n\lambda} \log \left[\max_{x \in X^n} (n+1)^X 2^{nH(Q_x)} \right. \\ & \quad \left. \cdot 2^{-n[H(Q_x) + D(Q_x \| P)]} 2^{\lambda n H(Q_x)} \right] + \delta_n \\ & = \frac{1}{n\lambda} \log \left[\max_{x \in X^n} 2^{\lambda n H(Q_x)} 2^{-nD(Q_x \| P)} \right] + \left(1 + \frac{1}{\lambda}\right) \delta_n \\ & = \max_{x \in X^n} \left[H(Q_x) - \frac{1}{\lambda} D(Q_x \| P) \right] + \left(1 + \frac{1}{\lambda}\right) \delta_n. \quad (\text{B.1}) \end{aligned}$$

On the other hand, Rényi's entropy can be lower bounded as follows:

$$\begin{aligned} H_\lambda & = \frac{\lambda+1}{n\lambda} \log \sum_{x \in X^n} P(x)^{1/(1+\lambda)} \\ & = \frac{\lambda+1}{n\lambda} \log \sum_{T_x \subset X^n} |T_x| P(x)^{1/(1+\lambda)} \\ & \geq \frac{\lambda+1}{n\lambda} \log \max_{T_x \subset X^n} \left[(n+1)^{-X} 2^{nH(Q_x)} \right. \\ & \quad \left. \cdot 2^{-n/(1+\lambda)[H(Q_x) + D(Q_x \| P)]} \right] \\ & = \max_{x \in X^n} \left[H(Q_x) - \frac{1}{\lambda} D(Q_x \| P) \right] - \left(1 + \frac{1}{\lambda}\right) \delta_n. \quad (\text{B.2}) \end{aligned}$$

It now follows by (B.1) and (B.2) that

$$\frac{1}{n\lambda} \log E(2^{\lambda L_n^*(x)}) \leq H_\lambda + 2 \left(1 + \frac{1}{\lambda}\right) \delta_n. \quad (\text{B.3})$$

Thus, if $L_n^*(x)$ is used instead of the optimal code, then the buffer overflow exponent is lower bounded by the solution to the equation $H_\lambda = R - 2(1 + \lambda^{-1})\delta_n$, while the optimal exponent is given by λ_0 , the solution to $H_\lambda = R$. Since the solution $\lambda_0(R)$ is a continuous differentiable function of R in the range $H < R < \log X$, the resulting degradation in overflow exponent is also $O(\delta_n) = O(n^{-1} \log n)$.

REFERENCES

- [1] R. J. McEliece, *Theory of Information and Coding*. New York: Addison-Wesley, 1977.
- [2] J. Aczel and Z. Daroczy, *Measures of Information and Their Characterizations*. New York: Academic Press, 1975.
- [3] P. A. Humblet, "Generalization of Huffman coding to minimize the probability of buffer overflow," *IEEE Trans. Inform. Theory*, vol. IT-27, no. 2, pp. 230–237, Mar. 1981.
- [4] J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE Trans. Inform. Theory*, vol. IT-24, no. 5, pp. 530–536, Sept. 1978.
- [5] J. Ziv and N. Merhav, "Estimating the number of states of a finite-state source," submitted to *IEEE Trans. Inform. Theory*. Also summarized in *Proc. AAAI Spring Symp. Theory Applications of Minimal-Length Encoding*, Mar. 1990, pp. 80–84.
- [6] R. B. Ash, *Information Theory*. New York: Wiley, 1965.
- [7] L. D. Davission, "Universal noiseless coding," *IEEE Trans. Inform. Theory*, vol. IT-19, no. 6, pp. 783–793, Nov. 1973.
- [8] J. Rissanen, "Universal coding, information, prediction, and estimation," *IEEE Trans. Inform. Theory*, vol. IT-30, no. 4, pp. 629–636, July 1984.
- [9] T. M. Cover, "On the competitive optimality of Huffman codes," preprint.
- [10] M. Gutman, "Asymptotically optimal classification for multiple tests with empirically observed statistics," *IEEE Trans. Inform. Theory*, vol. 35, no. 2, pp. 401–408, Mar. 1989.
- [11] F. Jelinek and K. S. Schneider, "On variable-length-to-block coding," *IEEE Trans. Inform. Theory*, vol. IT-18, no. 6, p. 762, Nov. 1972.
- [12] F. Jelinek, *Probabilistic Information Theory*. New York, McGraw-Hill, 1968.
- [13] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. New York: Academic Press, 1981.
- [14] K. Marton, "Error exponent for source coding with a fidelity criterion," *IEEE Trans. Inform. Theory*, vol. IT-20, no. 2, pp. 197–199, Mar. 1974.
- [15] E. Plotnik, M. J. Weinberger, and J. Ziv, "Upper bounds on the probability of sequences emitted by finite-state sources and on the redundancy of the Lempel–Ziv Algorithm," submitted to *IEEE Trans. Inform. Theory*.
- [16] F. Jelinek, "Buffer overflow in variable length coding of fixed rate sources," *IEEE Trans. Inform. Theory*, vol. IT-14, no. 3, pp. 490–501, May 1968.
- [17] A. D. Wyner, "On the probability of buffer overflow under and arbitrary bounded input–output distribution," *SIAM J. Appl. Math.*, vol. 27, pp. 544–570, 1974.