# On Limited-Delay Lossy Coding and Filtering of Individual Sequences

Tsachy Weissman, *Student Member, IEEE,* and Neri Merhav, *Fellow, IEEE*

*Abstract*—We continue the study of adaptive schemes for the sequential lossy coding of individual sequences which was recently initiated by Linder and Lugosi. Specifically, we consider fixed-rate lossy coding systems of fixed (or zero) delay where the encoder (which is allowed to use randomization) and the decoder are connected via a noiseless channel of a given capacity. It is shown that for any finite set of such coding schemes of a given rate, there exists a source code (adhering to the same structural and delay limitations) with the same rate whose distortion is with high probability almost as small as that of the best scheme in that set, uniformly for all individual sequences. Applications of this result to reference classes of special interest are outlined. These include the class of scalar quantizers, trellis encoders with sliding block decoders, and differential pulse code modulator (DPCM)-based source codes. In particular, for the class of all scalar quantizers, a source code is obtained with (normalized) distortion redundancy relative to the best scheme in the reference class of order $n^{-1/3} \log n$ (where $n$ is the sequence length). This improves the $n^{-1/5} \log n$ rate achieved by Linder and Lugosi. More importantly, the decoder here is deterministic and, in particular, does not assume a common randomization sequence available at both encoder and decoder. Finally, we consider the case where the individual sequence is corrupted by noise prior to reaching the coding system, whose goal now is to reconstruct a sequence with small distortion relative to the clean individual sequence. It is shown that for the case of a finite alphabet and an invertible channel transition probability matrix, for any finite set of sliding-window schemes of a given rate, there exists a source code (allowed to use randomization yet adhering to the same delay constraints) whose performance is, with high probability, essentially as good as the best scheme in the class, for all individual sequences.

*Index Terms*—Individual sequences, limited-delay coding, lossy source coding, noisy source coding, sequential coding, universal coding.

## I. INTRODUCTION

**T**HE requirement for a limited decoding delay arises naturally in an increasing variety of coding applications. A standard model for a fixed-rate ($R$) source code with limited delay is the following. The source sequence $x_1, x_2, \ldots$ is transformed into channel symbols $y_1, y_2, \ldots$ taking values in $\{1, 2, \ldots, M\}$, $M = 2^R$, which are, in turn, transformed into a reconstruction sequence $\hat{x}_1, \hat{x}_2, \ldots$. The encoder–decoder is said to have overall delay of no more than $\delta$ if each channel symbol $y_n$ depends only on $x_1, \ldots, x_{n+d_1}$ and each recon-

struction symbol $\hat{x}_n$ depends only on $y_1, \ldots, y_{n+d_2}$ such that $d_1 + d_2 \leq \delta$.

In the first part of this work, we consider such lossy source codes of limited delay in an individual sequence setting. Namely, the source sequence $x_1, x_2, \ldots$ is considered deterministic and statements are made which hold uniformly for all possible sequences. This setting naturally models situations where delay is crucial (e.g., cellular telephony, teleconferencing), yet very little is known about the source. Given a family of source codes (a reference class) constrained to have delay no larger than $\delta$ ($\delta$ can be any positive integer or infinity), our goal is to construct a code, adhering to the same delay constraint, which does essentially as well (i.e., has the same rate and incurs essentially the same average distortion) as the best code in the family, uniformly over all individual sequences. It is shown that this goal is achievable for an arbitrary finite reference class of delay-constrained codes, of any given delay, provided that the decoders associated with these codes are of finite (and uniformly bounded) memory. It is also illustrated how this result can be specialized to deal with infinite, parametrizable reference classes of codes, by obtaining coding schemes which compete with reference classes of particular practical interest.

The subsequent part of this work is dedicated to the case where the individual sequence is corrupted by a noisy memoryless channel (independent and identically distributed (i.i.d.) noise) prior to reaching the coder. Specifically, the source sequence $x_1, x_2, \ldots$ is fed into a memoryless channel whose output is $Z_1, Z_2, \ldots$, and only the latter sequence is available to the encoder. This is a setting of practical interest which naturally models audio and imaging applications in which the underlying signal (which has no natural probabilistic model and, hence, is considered an individual sequence) is corrupted by noise, and delay limitations are inherent in the application. Typical schemes involving predictive coding of noisy images, for example, can be formulated to fall within this setting for zero delay.

It may be intuitively expected that the two goals associated with such a noisy scenario, namely, filtering and compression, will be in concord. This fact is well known and has been made precise and exploited in the probabilistic context (without delay limitations) cf. [8], [4], [18], [6], [13], [16]. This intuition is consolidated for the individual sequence setting as well by (constructively) establishing the existence of a code for this noisy scenario, which is guaranteed to do essentially as well as the best in a reference class. That is, the combined filtering-compression goal can be efficiently achieved in the individual sequence setting as well. A reference class which will be explicitly treated in

this context is one that consists of time-invariant sliding-window encoder–decoder pairs, with the implication that other classes, possibly of codes having somewhat different structural properties, can be handled using similar ideas and tools.

The idea behind the construction of such codes is, in essence, quite similar to that underlying efficient prediction schemes for individual sequences. "Tracking" or trying to "imitate" those schemes in the reference class which have been proven efficient on the past sequence by exponentially weighting the extent to which each scheme is followed, according to its past performance. The implementation of such an approach in the lossy coding situation of the present setting, however, is not as straightforward as in the prediction setting. The main reason for that, in the noise-free setting, is the fact that the decoder, which accesses the reconstructed sequence only (and does not know the source sequence), does not have a precise picture of the past performance of each of the schemes in the class. In the noisy case, an additional level of difficulty is due to the fact that even the encoder, which accesses the noisy sequence, does not know the losses associated with the schemes in the reference class, as these depend on the unseen underlying clean individual sequence. As will be shown, however, this difficulty can be alleviated by employing suitable estimators for the unobserved cumulative distortion associated with the schemes in the reference class.

The idea of harnessing the exponential weighting approach to the present setting of delay-limited coding was instigated by Linder and Lugosi [14], where the study of zero-delay lossy source coding in the individual sequence setting was initiated. The main result of that work was a construction of a delay-less sequential adaptive coding scheme, for bounded, real-valued individual sequences, which asymptotically achieves the average distortion of the best scalar quantizer matched to the sequence. The basic idea underlying the coding schemes that we construct here is similar to that presented in [14] in that both are based on the exponential weighting principle. There are, however, some essential differences which, among other things, eliminate the need for the availability of the common randomization sequence assumed in [14]. This point is further elaborated on in Section IV. The setting of the present work can be considered a generalization of that of [14] in several directions. Any finite delay (not necessarily zero) is allowed, richer and more general reference classes (than that of scalar quantizers), arbitrary alphabets and distortion measures (other than the squared-error distortion measure of [14]), and the case where the original data sequence is corrupted by noise.

The remainder of the paper is organized as follows. In Section II, a formal description of the problem for the noise-free setting is given. Section III is dedicated to a generic result which constructively establishes the existence of a delay-$\delta$ code having essentially the same distortion as the best in a given finite family of delay-$\delta$ codes. In Section IV, the result of Section III is applied to specific reference classes of practical interest. In particular, for the class of scalar quantizers (Section IV-A), a delayless source code is obtained which, as will be discussed, is superior to that of [14]. Finally, Section V is devoted to the setting where the source sequence is corrupted by noise: Section V-A formalizes the problem and Section V-B is dedicated to the construc-

tion and performance analysis of a delay-limited source code (or filtering scheme) which does essentially as well as the best sliding-window encoder–decoder pair in a given set.

## II. PROBLEM FORMULATION FOR THE NOISE-FREE SETTING

Throughout the paper, for any integers $m \le n$, we let $a_m^n$ denote the vector $(a_m, \ldots, a_n)$ and $a^n = a_1^n$. Equalities and inequalities between random variables, when not explicitly specified, should be understood in the almost-sure sense. For any set $\mathcal{S}$ we let $|\mathcal{S}|$ denote its cardinality. For any collection $\{Z_i\}_{i \in \mathcal{I}}$ of random variables defined on a common probability space we shall let $\sigma(\{Z_i\}_{i \in \mathcal{I}})$ denote the smallest sigma-algebra with respect to which all $Z_i$, $i \in \mathcal{I}$, are measurable.

A delay-$\delta$ ($\delta$ a nonnegative integer or $\infty$) sequential source code of fixed rate $R = \log M$ with a randomized encoder is given by a pair $(E, D)$. The randomized encoder $E$ is given by a sequence $\{E_i\}_{i=1}^{\infty}$, where

$$E_i \colon \mathcal{X}^{i+\delta} \times [0, 1]^i \to \{1, 2, \ldots, M\}$$

$\mathcal{X}$ being the source alphabet. The decoder $D$ is given by a sequence $\{D_i\}_{i=1}^{\infty}$, where

$$D_i \colon \{1, 2, \ldots, M\}^i \to \hat{\mathcal{X}}$$

$\hat{\mathcal{X}}$ being the reproduction alphabet. The source code operates as follows. The encoder produces the $i$th-channel symbol $y_i \in \{1, 2, \ldots, M\}$ based on $x^{i+\delta}$ and on the random sequence $U^i$ according to $y_i = E_i(x^{i+\delta}, U^i)$, where $\{U_i\}_{i=1}^{\infty}$ is a randomization sequence of i.i.d. random variables, uniformly distributed on $[0, 1]$. The decoder emits the reconstructed sequence $\hat{x}_1, \hat{x}_2, \ldots$ according to $\hat{x}_i = D_i(y^i)$. We let $\mathcal{F}^\delta(R)$ denote the class of all such source codes.[1] We shall, in the sequel, assume a given fixed rate $R$ and simply write $\mathcal{F}^\delta$.

The cumulative distortion of a source code $(E, D) \in \mathcal{F}^\delta$ is denoted by

$$d_{(E, D)}^n(\boldsymbol{x}) = \sum_{i=1}^n d(x_i, \hat{x}_i) \tag{1}$$

where $d \colon \mathcal{X} \times \hat{\mathcal{X}} \to [0, B]$ is a bounded distortion measure ($B < \infty$) and the $\hat{x}_i$'s on the right-hand side are generated by feeding the sequence $\boldsymbol{x} = x_1, x_2, \ldots$ into $(E, D)$ as described above. Note that, though the dependence is suppressed in the notation, $d_{(E, D)}^n(\boldsymbol{x})$ is a random variable which depends on the realization of the randomization sequence $U^n$. We similarly denote, for $n_1 \le n_2$

$$d_{(E, D)}^{n_1, n_2}(\boldsymbol{x}) = \sum_{i=n_1}^{n_2} d(x_i, \hat{x}_i).$$

A decoder $D$ is said to be of finite memory $s \ge 0$ if for all $i \ge s$ and all $y^i, z^i \in \hat{\mathcal{X}}^i$ such that $y_{i-s}^i = z_{i-s}^i$, $D_i(y^i) =$

[1] Note that there is no essential loss in restricting the decoder to "causality." This is true since any finite-delay source code with a noncausal decoder (of delay, say, $\delta_2$) can be represented by an almost-equivalent source code in $\mathcal{F}^\delta(R)$ for some $\delta$ by a time shift. The latter source code will be equivalent to the original one except, possibly, on the first $\delta_2$ symbols because the causal decoder must emit, e.g., $\hat{x}_1$ after receiving $y_1$ while the noncausal decoder has received $y_1, \ldots, y_{\delta_2}$ when producing $\hat{x}_1$.

$D_i(z^i)$. We let $\mathcal{F}_s^\delta(R)$ denote the class of all source codes in $\mathcal{F}^\delta(R)$ with a decoder of finite-memory $s$. Note that, similarly to a finite-memory decoder, one can define a finite-*state* decoder (cf. [21]). While, admittedly, not every finite-state decoder is a finite-memory decoder, one can show, using the techniques of [7] and [15], that finite-memory machines perform asymptotically as well as finite-state machines. One particularly interesting subset of $\mathcal{F}_s^\delta(R)$ consists of those source codes having a decoder which, in addition to being limited to a memory of at most $s$ channel symbols back, is also time-invariant. This subset is relevant for the modeling of a variety of coding schemes in applications which require a finite (or zero) delay. In many such situations, practically any randomized encoder (situated, e.g., in some base station where algorithmic resources are abundant) can be implemented, yet the decoder (situated, e.g., in some small and low-cost handset) has no access to a randomization sequence and is limited in memory or in algorithmic resources.

## III. GENERIC RESULT

We dedicate this section to the construction of a finite-delay coding scheme which competes with an arbitrary finite set of limited-delay schemes in the sense of operating at the same rate, and suffering a cumulative distortion which is, at most, negligibly higher than that of the best in the set, for all individual sequences. More precisely, we have the following:

*Theorem 1:* Let $\mathcal{A}$ be a finite subset of $\mathcal{F}_s^\delta(R)$ for some $s \geq 0$, $0 \leq \delta \leq \infty$ and $R = \log M$. For any $0 < \varepsilon \leq 1$ and $N$ sufficiently large such that

$$N > C_1[(\log |\mathcal{A}|)(\log(|\mathcal{A}|e^{sR}))]/\varepsilon^3$$

there exists a source code $(E, D) \in \mathcal{F}^\delta(R)$ such that for all $\boldsymbol{x} \in \mathcal{X}^\infty$ we have both

$$\mathbb{E}\left\{\frac{1}{N}\left[d_{(E,D)}^N(\boldsymbol{x}) - \min_{(E',D')\in\mathcal{A}} d_{(E',D')}^N(\boldsymbol{x})\right]\right\}$$
$$\leq C_2\left[(\log |\mathcal{A}|)(\log(|\mathcal{A}|e^{sR}))\right]^{1/3} \cdot N^{-1/3} \quad (2)$$

and

$$\Pr\left\{\frac{1}{N}\left[d_{(E,D)}^N(\boldsymbol{x}) - \min_{(E',D')\in\mathcal{A}} d_{(E',D')}^N(\boldsymbol{x})\right] \geq \varepsilon\right\}$$
$$\leq \exp\left\{-C_3\left(\frac{\log |\mathcal{A}|}{\log(|\mathcal{A}|e^{sR})}\right)^{1/3}\varepsilon^2 N^{2/3}\right\} \quad (3)$$

where $C_1$, $C_2$, $C_3$ are positive constants which depend only on $B$ and $R$.

*Discussion:* The explicit values of the constants $C_1$, $C_2$, $C_3$ appearing in the above theorem will be apparent in the proof. We note that the coding scheme suggested in the above theorem is a member of $\mathcal{F}^\delta(R)$. In particular, *its delay is no more than those of the schemes in the reference class and it has exactly the same rate $R$* (and not even a bit more as is sometimes the case in related universal coding scenarios). Theorem 1 holds with no assumptions on the source and reconstruction alphabets and the only assumption on the distortion measure is that it is bounded. It should also be emphasized that Theorem 1 is an individual-sequence result and the expectation and probability appearing in

(2) and (3) are with respect only to the randomization sequence, i.e., the encoders' local randomness.

Note that the reference class $\mathcal{A}$ allowed in Theorem 1 can be any finite subset of $\mathcal{F}_s^\delta(R)$. In particular, the encoders associated with the schemes in $\mathcal{A}$ may use randomization. In most conceivable applications of Theorem 1 (and, in particular, in the examples considered in the next section) to reference classes of practical interest, the associated encoder would not be randomized. In Theorem 1, however, we allow the competing schemes to use randomization at the encoder for two principal reasons. The first is to maximize the generality. The second is one of "fairness": since we allow the universal scheme that we construct to use randomization at the encoder, it is only fair that the schemes in the reference class be given the same courtesy.

It should also be observed that the delay $\delta$ in Theorem 1 can be any nonnegative integer, or infinite. The result is more interesting for relatively small delay since for large or infinite delay one can use the results from Ziv's work [21], or the later Yang and Kieffer codes [19], for the lossy compression of individual sequences. The codes in [21] and [19] offer little guidance regarding the problem of coding under delay limitations (cf. discussion in [14]). Indeed, theory regarding the delay-constrained lossy source coding of individual sequences was virtually nonexistent prior to [14].

We also remark that, while in most probabilistic contexts it is usually satisfactory and informative enough to make statements regarding the *expected* performance (distortion) of a source code, in the individual sequence setting considered here this is not the case. The whole point of the individual sequence setting is to have a complete picture of what is really happening (actual rather than expected distortion) for *every* possible sequence. This is why a point has been made to obtain the concentration inequality (3), which guarantees the actual performance of the source code, in addition to (2).

The source code $(E, D)$ in the above theorem depends on the length $N$ of the individual sequence to be encoded (this is referred to as "horizon-dependence" in the prediction literature). While in some applications (e.g., image coding) the length of the sequence is indeed known in advance, in others it may be desirable to have a source code guaranteed to be doing well at *all* points along the sequence (this is referred to as the "strong sequentiality" property). We merely remark here that it is a straightforward exercise to obtain a strongly sequential source code, for which (2) and (3) hold for *all* sufficiently large $N$. The Borel–Cantelli lemma can then be applied to obtain an almost sure performance guarantee for such a source code. This is done by employing the horizon-dependent source code of Theorem 1 on blocks of exponentially increasing lengths (cf., e.g., [1], [17] for typical examples of such constructions).

*Proof:* Fix an $l \ll N$ and divide the time-axis, $i = 1$, $2, \ldots, N$, into $n = N/l$ consecutive nonoverlapping blocks (assume $l$ divides $N$). We construct the source code

$$(E, D) \in \mathcal{F}^\delta(R)$$

as follows. At the beginning of the $k$th block $1 \leq k \leq n$, i.e., at the $i = (k-1)l+1$th channel use, when $x^{i+\delta}$ and $U^i$ are available (so that $d_{(E',D')}^{(k-1)l}(\boldsymbol{x})$ is known for all $(E', D') \in \mathcal{A}$), the encoder

uses $U_i$ to generate $(E^{(k)}, D^{(k)})$, a $\mathcal{A}$-valued random variable with distribution satisfying almost surely

$$\Pr\left\{\left(E^{(k)}, D^{(k)}\right) = (E', D')\Big| U^{(k-1)l}\right\}$$

$$= \frac{\exp\left\{-\eta\, d_{(E',D')}^{(k-1)l}(\boldsymbol{x})\right\}}{\sum\limits_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l}(\boldsymbol{x})\right\}}, \qquad \forall\,(E', D')\in\mathcal{A} \quad (4)$$

where $\eta > 0$ is a parameter to be chosen later. The conditioning on $U^{(k-1)l}$ on the left-hand side of (4) is necessary due to the fact that the cumulative distortion of each of the schemes in the reference class which, by the definition of $\mathcal{F}_s^\delta(R)$ are allowed to use the randomization sequence, will, in general, depend on the realization of the randomization sequence. The encoder, operating at the same rate $R$, now dedicates the first $\lceil\frac{1}{R}\log|\mathcal{A}|\rceil$ channel symbols at the beginning of the $k$th block, i.e., $y_i$ for

$$i = (k-1)l+1, \ldots, (k-1)l+\left\lceil\frac{1}{R}\log|\mathcal{A}|\right\rceil$$

to convey to the decoder the identity of $D^{(k)}$. At the remainder of the block, i.e., at times

$$i = (k-1)l + \left\lceil\frac{1}{R}\log|\mathcal{A}|\right\rceil + 1, \ldots, kl$$

the encoder produces the channel symbols

$$y_i = E_i^{(k)}(x^{i+\delta}, U^i).$$

At the same time, on the decoder's side, at the beginning of the block at times

$$i = (k-1)l+1, \ldots, (k-1)l+\left\lceil\frac{1}{R}\log|\mathcal{A}|\right\rceil + s - 1$$

the decoder outputs an arbitrary reproduction sequence of $\hat{x}_i$'s. From time

$$i = (k-1)l + \left\lceil\frac{1}{R}\log|\mathcal{A}|\right\rceil + s$$

up to the end of the block $i = kl$, the decoder, knowing the identity of $D^{(k)}$ and the output of $E^{(k)}$ at least $s-1$ channel uses back, outputs the reproduction sequence according to

$$\hat{x}_i = D_i^{(k)}(y^i) = D_i^{(k)}(y_{i-s}^i)$$

(the second equality is due to the fact that $(E^{(k)}, D^{(k)}) \in \mathcal{F}_s^\delta$). Note that the decoder in the source code $(E, D)$ we have constructed is indeed causal and that the encoder at each point $i$ along the sequence relies on knowledge of at most $x^{i+d}$ so that $(E, D)$ is a *bone fide* member of $\mathcal{F}^\delta(R)$. Furthermore, $(E, D)$ utilizes the randomization sequence in a rather economical way: it only uses randomization once at the beginning of each block (cf. Fig. 1 for a schematic description of the construction of $(E, D)$).

To establish (2) and (3), we use some standard ingredients from the theory of prediction of individual sequences (cf., e.g., [2, Theorem 1]). Define for each $k \geq 1$ the random variable

$$W_k = \sum_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l}(\boldsymbol{x})\right\} \qquad (5)$$
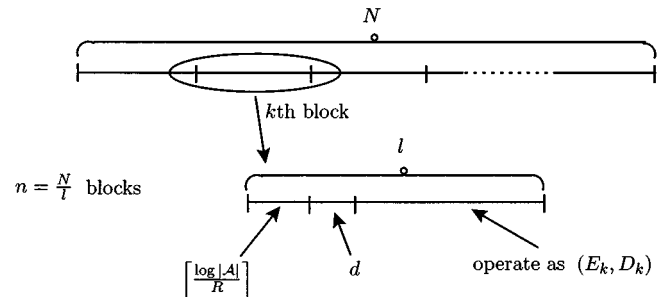


Fig. 1.    Structure of coding scheme of Theorem 1.

(so that, in particular, $W_1 = |\mathcal{A}|$). We then have for $n = N/l$

$$\log\frac{W_{n+1}}{W_1} = \log\sum_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{nl}(\boldsymbol{x})\right\} - \log|\mathcal{A}|$$

$$\geq \log\left(\max_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{nl}(\boldsymbol{x})\right\}\right) - \log|\mathcal{A}|$$

$$= -\eta\min_{(\tilde{E},\tilde{D})\in\mathcal{A}} d_{(\tilde{E},\tilde{D})}^{N}(\boldsymbol{x}) - \log|\mathcal{A}|. \qquad (6)$$

On the other hand, for each $1 \leq k \leq n$

$$\log\frac{W_{k+1}}{W_k}$$

$$= \log\frac{\sum\limits_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l+1,kl}(\boldsymbol{x})\right\}\cdot\exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l}(\boldsymbol{x})\right\}}{\sum\limits_{(\tilde{E},\tilde{D})\in\mathcal{A}} \exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l}(\boldsymbol{x})\right\}}$$

$$= \log\mathbb{E}_{Q_k}\exp\left\{-\eta\, d_{(\tilde{E},\tilde{D})}^{(k-1)l+1,\,kl}(\boldsymbol{x})\right\}$$

$$\leq -\eta\mathbb{E}_{Q_k} d_{(\tilde{E},\tilde{D})}^{(k-1)l+1,\,kl}(\boldsymbol{x}) + \frac{\eta^2 l^2 B^2}{8}$$

$$\leq -\eta\left\{\mathbb{E}\left[d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x})\Big| U^{(k-1)l}\right] - B\left(\frac{1}{R}\log|\mathcal{A}| + s\right)\right\}$$

$$+ \frac{\eta^2 l^2 B^2}{8} \qquad (7)$$

where $\mathbb{E}_{Q_k}$ denotes expectation with respect to the (random) distribution $Q_k$ on $\mathcal{A}$ which assigns a probability proportional to $\exp\{-\eta d_{(\tilde{E},\tilde{D})}^{(k-1)l}(\boldsymbol{x})\}$ to each $(\tilde{E},\tilde{D})\in\mathcal{A}$. The first inequality follows from an application of Hoeffding's bound (cf. [3, Lemma 8.1]). The second inequality follows from the construction of the source code $(E, D)$ described above (note that the cumulative distortion of $(E, D)$ from the beginning of the $k$th block up to time $i = (k-1)l + \lceil\frac{1}{R}\log|\mathcal{A}|\rceil + s - 1$ can be no more than $B(\frac{1}{R}\log|\mathcal{A}| + s)$ and from that time up to the end of the block it is exactly the loss of the pair $(E^{(k)}, D^{(k)})$ generated at the beginning of the block). Summing up over $k$ we obtain almost surely

$$\log\frac{W_{n+1}}{W_1} \leq -\eta\sum_{k=1}^{n}\mathbb{E}\left[d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x})\Big| U^{(k-1)l}\right]$$

$$+ \eta B n\left(\frac{1}{R}\log|\mathcal{A}| + s\right) + \frac{\eta^2 l^2 B^2 n}{8}. \quad (8)$$

Combining (8) and (6) gives almost surely

$$\sum_{k=1}^{n} \mathbb{E}\left[ d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) \Big| U^{(k-1)l} \right] - \min_{(\tilde{E},\tilde{D})\in\mathcal{A}} d_{(\tilde{E},\tilde{D})}^{N}(\boldsymbol{x})$$

$$\leq \frac{\log|\mathcal{A}|}{\eta} + \frac{\eta l^2 B^2 n}{8} + Bn\left(\frac{1}{R}\log|\mathcal{A}| + s\right)$$

$$= B\left(\sqrt{\log|\mathcal{A}|/2}\right) l n^{1/2} + Bn\left(\frac{1}{R}\log|\mathcal{A}| + s\right) \quad (9)$$

where the equality follows upon taking the minimizing value $\eta = \sqrt{8\log|\mathcal{A}|/(l^2 B^2 n)}$. For notational convenience, we now denote

$$\alpha = B\sqrt{\log|\mathcal{A}|/2} \quad \text{and} \quad \beta = B(\tfrac{1}{R}\log|\mathcal{A}| + s)$$

so that the right-hand side of (9) becomes

$$\alpha l n^{1/2} + \beta n = \alpha(N/n)n^{1/2} + \beta n = \alpha N n^{-1/2} + \beta n.$$

Minimizing with respect to $n$, we take $n = (\alpha/(2\beta))^{2/3} N^{2/3}$ and obtain an expression upper bounded by $3\alpha^{2/3}\beta^{1/3}N^{2/3}$. Plugging in the values of $\alpha, \beta$, we finally obtain

$$\sum_{k=1}^{n} \mathbb{E}\left[ d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) \Big| U^{(k-1)l} \right] - \min_{(\tilde{E},\tilde{D})\in\mathcal{A}} d_{(\tilde{E},\tilde{D})}^{N}(\boldsymbol{x})$$

$$\leq \frac{3B}{(2R)^{1/3}} \left[ (\log|\mathcal{A}|)(\log(|\mathcal{A}|e^{sR})) \right]^{1/3} N^{2/3}. \quad (10)$$

Note that our choice of the number of blocks

$$n = (\alpha/(2\beta))^{2/3} N^{2/3}$$

implies that the length of each block is

$$l = N/n = (2\beta/\alpha)^{2/3} N^{1/3}.$$

Since the above derivation assumes that the block length is greater than the overhead of $\left(\frac{1}{R}\log|\mathcal{A}| + s\right) = \beta/B$ channel uses at its beginning, we verify that this is indeed the case. We require $l = (2\beta/\alpha)^{2/3} N^{1/3} \geq \beta/B$ which is equivalent to

$$N \geq (\alpha^2\beta)/(4B^3) = \frac{1}{8R}(\log|\mathcal{A}|)(\log(|\mathcal{A}|e^{sR}))$$

an inequality which holds by hypothesis (for a suitable constant $C_1$). Hence, taking expectation in (10) establishes (2). Turning to establish (3), we denote

$$V_k = d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) - \mathbb{E}\left[ d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) \Big| U^{(k-1)l} \right]$$

and observe that $\{V_k, \sigma(U^{kl})\}_{k\geq 1}$ is a martingale difference sequence with the $V_k$'s almost-surely bounded in magnitude by $Bl$. Applying Hoeffding's bound for martingale difference sequences (cf., e.g., [3, Theorem 9.1]) gives for any $n$ and $\varepsilon > 0$

$$\Pr\left\{ \frac{1}{N}\left[ d_{(E,D)}^{N}(\boldsymbol{x}) - \sum_{k=1}^{n} \mathbb{E}\left[ d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) \Big| U^{(k-1)l} \right] \right] \geq \varepsilon \right\}$$

$$= \Pr\left\{ \frac{1}{nl}\sum_{k=1}^{n} V_k \geq \varepsilon \right\}$$

$$\leq \exp\left\{ -\frac{n\varepsilon^2}{2B^2} \right\}. \quad (11)$$

Finally, combining (10) and (11), we have for any $N$ large enough such that the right-hand side of (10) is less than $N\varepsilon/2$, namely, for $N \geq \frac{216B^3}{2R}\left[(\log|\mathcal{A}|)(\log(|\mathcal{A}|e^{sR}))\right]/(\varepsilon^3)$

$$\Pr\left\{ \frac{1}{N}\left[ d_{(E,D)}^{N}(\boldsymbol{x}) - \min_{(E',D')\in\mathcal{A}} d_{(E',D')}^{N}(\boldsymbol{x}) \right] \geq \varepsilon \right\}$$

$$\leq \Pr\left\{ \frac{1}{N}\left[ d_{(E,D)}^{N}(\boldsymbol{x}) - \sum_{k=1}^{n} \right.\right.$$

$$\left.\left. \cdot\, \mathbb{E}\left[ d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x}) \Big| U^{(k-1)l} \right] \right] \geq \varepsilon/2 \right\}$$

$$\leq \exp\left\{ -\frac{n(\varepsilon/2)^2}{2B^2} \right\}$$

$$= \exp\left\{ -\frac{(\alpha/(2\beta))^{2/3} N^{2/3} \varepsilon^2}{8B^2} \right\} \quad (12)$$

which, upon plugging in the values of $\alpha$ and $\beta$, establishes (3) and concludes the proof for the case where the above values of $n = (\alpha/(2\beta))^{2/3} N^{2/3}$ and $l = (2\beta/\alpha)^{2/3} N^{1/3}$ are integers. Otherwise, take $l = \lfloor (2\beta/\alpha)^{2/3} N^{1/3} \rfloor$ and let the coding scheme $(E, D)$ behave arbitrarily on the last (incomplete) block. It is straightforward to verify that the above derivation carries over for this case as well with, possibly, slightly modified constants $C_1, C_2, C_3$ to accommodate the "edge effects." □

Note that the source code constructed in the above proof must, effectively, run all the schemes in the reference class in parallel. This fact renders this scheme impractical for implementation when the reference class is excessively large (cf. discussion in Section VI).

It might be tempting to simplify the source code constructed in Theorem 1 in the following way. Rather than generate the member of the reference class to be followed on the $k$th block according to (4), simply decide, deterministically, to follow that member which has been proven most competent on the past sequence. This would eliminate the need for a randomization sequence altogether. Unfortunately, as is known from the theory of prediction of individual sequences, such a scheme does not have the necessary adaptivity properties. In fact, one can construct simple cases of reference classes containing as little as two source codes for which the normalized distortion redundancy of such a deterministic scheme relative to the reference class is lower-bounded by a nonvanishing term. As one simple example consider the following: Let $\mathcal{X} = \hat{\mathcal{X}} = \{0, 1, 2\}$, $M = 2$, and $d_H(\cdot, \cdot)$ be the Hamming distortion measure. A scalar quantizer of rate $R = \log M = \log 2$ for this case is a map $Q: \{0, 1, 2\} \to \mathcal{C}$, where $\mathcal{C} \subset \{0, 1, 2\}$ with $|\mathcal{C}| = 2$. For an individual sequence $\boldsymbol{x} \in \mathcal{X}^\infty$ let

$$d_Q^N(\boldsymbol{x}) = \sum_{i=1}^{N} d_H(x_i, Q(x_i)).$$

Each such $Q$ defines a time-invariant member of $\mathcal{F}_{s=0}^{\delta=0}$ in the obvious way. Define the quantizers

$$Q_0(x) = \begin{cases} 2, & x = 2 \\ 1, & x = 0, 1 \end{cases} \quad (13)$$

and

$$Q_1(x) = \begin{cases} 2, & x = 2 \\ 0, & x = 0, 1. \end{cases} \qquad (14)$$

Let now $(E, D) \in \mathcal{F}_{s=0}^{\delta=0}$ be any horizon-dependent scheme which partitions the data of length $N$ into any number of blocks and chooses, at the beginning of each block, to follow either the scheme associated with $Q_0$ or that associated with $Q_1$ according to some deterministic rule (e.g., that suggested above). It is clear that for any such scheme we can construct a sequence $\boldsymbol{x}$ which is constant on each block and such that $d_{(E, D)}^N(\boldsymbol{x}) = N$. On the $k$th block, if $Q_0$ is followed we let the corresponding subsequence of $\boldsymbol{x}$ consist solely of 0's and if $Q_1$ is followed it will consist of 1's. On the other hand, for *any* sequence we obviously have $\min\{d_{Q_0}^N(\boldsymbol{x}), d_{Q_1}^N(\boldsymbol{x})\} \leq N/2$. Evidently, we have constructed a sequence for which

$$d_{(E, D)}^N(\boldsymbol{x}) - \min\left\{d_{Q_0}^N(\boldsymbol{x}), d_{Q_1}^N(\boldsymbol{x})\right\} \geq N/2. \qquad (15)$$

In particular, the normalized distortion redundancy is lower bounded by $1/2$.

## IV. APPLICATIONS

We dedicate this section to an application of the generic result of Section III to a few representative cases.

### A. Scalar Quantizers

Though the reference class in the above theorem is finite, it can be straightforwardly applied to obtain delay-$d$ source codes which compete with any reference class having a well-behaved effective covering number. One example for such an application is for the reference class of all scalar quantizers considered in [14]. Specifically, consider the case $\mathcal{X} = \hat{\mathcal{X}} = [0, 1]$ with the squared-error distortion $d(x_i, \hat{x}_i) = (x_i - \hat{x}_i)^2$. An $M$-level scalar quantizer $Q$ is a measurable mapping $[0, 1] \to \mathcal{C}$ where $\mathcal{C}$ is any finite subset of $[0, 1]$ with cardinality $|\mathcal{C}| = M$. Following [14], we let $\mathcal{Q}$ denote the class of all $M$-level $(M = 2^R)$ scalar quantizers and let

$$d_Q^N(\boldsymbol{x}) = \sum_{i=1}^{N} (x_i - Q(x_i))^2.$$

We then have the following.

*Corollary 2:* Let $\mathcal{X} = \hat{\mathcal{X}} = [0, 1]$ and $d(\cdot, \cdot)$ be the squared-error distortion measure. For any $0 < \varepsilon \leq 1$ and $N$ sufficiently large such that $N/(\log N)^2 > K_1/\varepsilon^3$ there exists a zero-delay source code $(E, D) \in \mathcal{F}^0$ such that for all $\boldsymbol{x} \in \mathcal{X}^\infty$ we have both

$$\mathbb{E} \frac{1}{N} d_{(E, D)}^N(\boldsymbol{x}) - \inf_{Q \in \mathcal{Q}} \frac{1}{N} d_Q^N(\boldsymbol{x}) \leq K_2(\log N)N^{-1/3} \quad (16)$$

and

$$\Pr\left\{\frac{1}{N}\left[d_{(E, D)}^N(\boldsymbol{x}) - \inf_{Q \in \mathcal{Q}} d_Q^N(\boldsymbol{x})\right] \geq \varepsilon\right\}$$
$$\leq \exp\left\{-K_3\varepsilon^2 N^{2/3}\right\} \quad (17)$$

where $K_1, K_2, K_3$ are positive constants which depend only on $R$.

As can be seen by examining the constants in the proofs (of both Corollary 2 and Theorem 1), the dependence of the constants on the rate $R$ or, equivalently, on $M$ is according to

$K_1 \propto 1/(\log M)$, $K_3 \propto (\log M)^{2/3}$, and the right side of (16) can be replaced by

$$C_2[-M\log(MC_2N^{-1/3})]N^{-1/3} + \frac{1}{(MC_2N^{-1/3})^{-1} - 1}$$

where $C_2$ is the constant from Theorem 1 which behaves as $C_2 \propto (\log M)^{-1/3}$. Thus, we see that better constants are attained for higher rates. The intuition behind this fact is that the higher the instantaneous rate, the less channel uses must be dedicated at the beginning of each block to convey to the decoder the identity of the chosen scheme and, hence, the lesser portion of the time is the decoder idle.

Note that Corollary 2 (inequality (16) in particular) improves the main result of [14] in two directions. The first is in establishing the existence of a source code with a redundancy term upper bounded by $C(\log N)N^{-1/3}$, where that associated with the source code of [14] was shown to be upper-bounded by $C(\log N)N^{-1/5}$. The second, perhaps more essential improvement, is in the fact that the source code employed here belongs to $\mathcal{F}^{d=0}$. In particular, it does not use randomization at the decoder and, *a fortiori*, does not need to know the realization of the randomization sequence used by the encoder. This is in contrast to the source code constructed in [14] which not only required the availability of a randomization sequence at the decoder's side as well but required that both the encoder and the decoder dispose of the *same* randomizing sequence (as in models where a "subtractive dither" is assumed, cf. [22], [20]).

*Proof:* It is straightforward to show (cf. [14, Lemma 2]) that for any $A \geq 2$ there exists a set $\mathcal{Q}_A = \{Q_1, \ldots, Q_A\} \subset \mathcal{Q}$ such that for all $N$ and $\boldsymbol{x} \in \mathcal{X}^\infty$

$$\min_{Q \in \mathcal{Q}_A} d_Q^N(\boldsymbol{x}) \leq \inf_{Q \in \mathcal{Q}} d_Q^N(\boldsymbol{x}) + N\frac{1}{A^{1/M} - 1}. \quad (18)$$

Therefore, letting $(E, D) \in \mathcal{F}^0$ be the source code of Theorem 1, tailored for the set $\mathcal{Q}_A \subset \mathcal{F}_{s=0}^{\delta=0}$, we obtain

$$\mathbb{E}\frac{1}{N}d_{(E, D)}^N(\boldsymbol{x}) - \inf_{Q \in \mathcal{Q}}\frac{1}{N}d_Q^N(\boldsymbol{x})$$
$$\leq \mathbb{E}\frac{1}{N}d_{(E, D)}^N(\boldsymbol{x}) - \min_{Q \in \mathcal{Q}_A}\frac{1}{N}d_Q^N(\boldsymbol{x}) + \frac{1}{A^{1/M} - 1}$$
$$\leq C_2(\log A)^{2/3}N^{-1/3} + \frac{1}{A^{1/M} - 1} \quad (19)$$

where the second inequality follows from Theorem 1 (assuming $N$ is sufficiently large, a fact we shortly verify to follow from the hypothesis of the corollary). Taking $A = \lceil(MC_2N^{-1/3})^{-M}\rceil$, we obtain

$$\mathbb{E}\frac{1}{N}d_{(E, D)}^N(\boldsymbol{x}) - \inf_{Q \in \mathcal{Q}}\frac{1}{N}d_Q^N(\boldsymbol{x})$$
$$\leq C_2\left[-M\log(MC_2N^{-1/3})\right]N^{-1/3}$$
$$+ \frac{1}{(MC_2N^{-1/3})^{-1} - 1}$$
$$\leq K_2(\log N)N^{-1/3}. \quad (20)$$

We note that for (19), and hence for (20), to hold we need (by Theorem 1) to require $N > (C_1(\log A)^2)/(\varepsilon^3)$. But, by our choice of $A$, $\log A$ is, up to a multiplicative constant (depending only on $M$), equivalent to $\log N$ and, therefore, it is enough

to require $N > (K_1 (\log N)^2)/(\varepsilon^3)$ for some $K_1$ which establishes (16). To arrive at (17), it is easy to verify that the hypothesis $N > (K_1 (\log N)^2)/(\varepsilon^3)$ implies also $1/(A^{1/M} - 1) \leq \varepsilon/2$ (for a suitably chosen $K_1$). Hence

$$\Pr \left\{ \frac{1}{N} d_{(E, D)}^N (\boldsymbol{x}) - \inf_{Q \in \mathcal{Q}} \frac{1}{N} d_Q^N (\boldsymbol{x}) \geq \varepsilon \right\}$$

$$\leq \Pr \left\{ \frac{1}{N} d_{(E, D)}^N (\boldsymbol{x}) - \min_{Q \in \mathcal{Q}_A} \frac{1}{N} d_Q^N (\boldsymbol{x}) \geq \varepsilon/2 \right\}$$

$$\leq \exp \left\{ -K_3 \varepsilon^2 N^{2/3} \right\} \tag{21}$$

where the first inequality follows from (18) and the second one follows from Theorem 1 (inequality (3)). □

### B. Sliding Block and Trellis Source Codes

A class to which Theorem 1 can be directly applied is that of sliding-block codes (cf., e.g., [9]). Assume here an arbitrary distortion measure. A finite-constraint length, time-invariant encoder with constraint-length $l$, memory $l_M$, and delay $l_d = l - l_M - 1$ is a mapping $f : \mathcal{X}^l \to \{1, 2, \ldots, M\}$ yielding the channel symbols $\{y_i\}$ defined by $y_i = f(x_{i - l_M}^{i + l_d})$. Similarly, a sliding-block decoder with constraint length $k$ is a mapping $g : \{1, 2, \ldots, M\}^k \to \hat{\mathcal{X}}$ yielding a reproduction process $\hat{x}_i = g(y_{i - k + 1}^i)$. When the source and reproduction alphabets are finite, it is easy to see that the cardinality of the class of all such sliding-block source codes with respective constraint lengths of $l$ and $k$ is $M^{|\mathcal{X}|^l} \cdot |\hat{\mathcal{X}}|^{M^k}$. Therefore, Theorem 1 guarantees the existence of a source code with delay $l_d$ which achieves the lowest distortion attainable by any sliding-block source code with constraint lengths $l$, $k$ and delay $l_d$, uniformly for all individual sequences in the sense of (2) and (3) (with $|\mathcal{A}| = M^{|\mathcal{X}|^l} \cdot |\hat{\mathcal{X}}|^{M^k}$). For the case where the alphabets are not of finite cardinality, competing with *all* such sliding-block source codes of a given order is clearly an overambitious task. However, any sufficiently smoothly parametrizable subset of this class can be dealt with using appropriate grids, similarly as was done in the proof of Corollary 2.

Another important related family of source codes that are covered by Theorem 1 consists of the block trellis codes (cf., e.g., [9], [11]). In particular, note that a block trellis coding scheme of constraint length $K$ and search depth $L$, which consists of a constraint-length $K$ sliding-block decoder $g$ and a depth-$L$ trellis search encoder matched to $g$, is a member of $\mathcal{F}_K^{L-1}$.

### C. Adaptive Differential Pulse Code Modulation (DPCM)

We dedicate this subsection to specializing the approach demonstrated in Theorem 1 for obtaining an adaptive scheme which "tracks" the differential pulse code modulator (DPCM), out of a given family, which does best on the data. DPCM-based schemes are widely used in speech and image applications and are, therefore, of considerable practical interest. In particular, the study of such coding schemes within the individual-sequence framework which we consider here is especially relevant for many image compression situations in which there is no natural statistical model for the data.

The reader is referred, e.g., to [10], [12] and the references therein for a comprehensive account of the theory and practice of DPCM-based coding schemes. Following is a laconic description intended primarily to introduce notation for later use. A DPCM delayless source code of order $s$ is an element of $\mathcal{F}^0$, which is fully characterized by a predictor $P : \hat{\mathcal{X}}^s \to \mathbb{R}$ and an $M$-level quantizer $Q \in \mathcal{Q}$. The encoder produces the $i$th channel symbol representing the value of $Q(e_i)$, where $e_i = x_i - P(\hat{x}_{i-s}^{i-1})$. The decoder gives the $i$th reconstruction according to $\hat{x}_i = Q(e_i) + P(\hat{x}_{i-s}^{i-1})$. For concreteness in initializing the predictor take, say, $\hat{x}_{-s+1}^0 = (a, \ldots, a)$ for some $a \in \hat{\mathcal{X}}$. An important subset of this family is that where the source and reconstruction alphabets are (possibly subsets of) the real line and the predictor $P$ is a linear time-invariant predictor, i.e., $P(\hat{x}_{i-s}^{i-1}) = \sum_{j=1}^s h_j \hat{x}_{i-j}$.[2] We let $\mathcal{F}_{\mathrm{DPCM}}(s, \eta, R)$ denote the class of all DPCM delayless source codes of order up to $s$ having a linear time-invariant predictor with impulse response $h$ satisfying $\|h\|_1 = \sum_{j=1}^s |h_j| \leq 1 - \eta$ and operating at rate $R$. The reason we are interested in such impulse responses is the strong stability of the decoder (linear feedback system) that they induce. Specifically, this property ensures that the respective outputs of a DPCM (having such an $h$) fed with similar inputs will also be similar. Assume that the source and reconstruction alphabets in what follows consist of a bounded subset of the real line and that $d$ is the squared-error distortion. In particular, we can assume that the distortion measure and the magnitude of the components of the source and reconstruction sequences are bounded by $B > 0$.

*Theorem 3:* Let $\mathcal{A}$ be a finite subset of $\mathcal{F}_{\mathrm{DPCM}}(s, \eta, R)$ for some $s \geq 0$, $R$, and $0 < \eta \leq 1$. Then for any $0 \leq \varepsilon \leq B$ and

$$N > C_1 [(\log |\mathcal{A}|)(\log |\mathcal{A}| + C_4 (s, \eta) R)] / (\varepsilon^3)$$

there exists a delayless source code $(E, D) \in \mathcal{F}^0 (R)$ such that for all $\boldsymbol{x} \in \mathcal{X}^\infty$

$$\mathbb{E} \left\{ \frac{1}{N} \left[ d_{(E, D)}^N (\boldsymbol{x}) - \min_{\mathcal{A}} d_{(E', D')}^N (\boldsymbol{x}) \right] \right\}$$

$$\leq C_2 [(\log |\mathcal{A}|)(\log |\mathcal{A}| + C_4 (s, \eta) R)]^{1/3} \cdot N^{-1/3} \tag{22}$$

and

$$\Pr \left\{ \frac{1}{N} \left[ d_{(E, D)}^N (\boldsymbol{x}) - \min_{\mathcal{A}} d_{(E', D')}^N (\boldsymbol{x}) \right] \geq \varepsilon \right\}$$

$$\leq \exp \left\{ -C_3 \left( \frac{\log |\mathcal{A}|}{\log |\mathcal{A}| + C_4 (s, \eta) R} \right)^{1/3} \varepsilon^2 N^{2/3} \right\} \tag{23}$$

where $C_1$, $C_2$, $C_3$ are positive constants which depend only on $B$ and $R$ and $C_4 (s, \eta) > 0$ for all $s, \eta > 0$.

Note that Theorem 3 does not follow directly from Theorem 1 since a source code belonging to $\mathcal{F}_{\mathrm{DPCM}}(s, \eta)$ will, in general (for $s \geq 1$), not have a finite-memory decoder (note that the decoder's output depends on the present channel symbol and on past *reconstruction* symbols). The proof of Theorem 3, however, which we outline below, is very similar to that of Theorem 1. In particular, the construction of the source code $(E, D)$ satisfying (22) and (23) is essentially the same as that in Theorem 1.

---

[2]Note that in this case the decoder is a simple linear time-invariant feedback system.

*Sketch of Proof:* The construction of $(E, D)$ follows that in the proof of Theorem 1 essentially verbatim. The only difference is that now, at the decoder's side, from time

$$i = (k-1)l + \left\lceil \frac{1}{R} \log |\mathcal{A}| \right\rceil + 1$$

up to the end of the block $i = kl$, the decoder, knowing the identity of $D^{(k)}$ (equivalently, of the linear predictor $h^{(k)}$ and of the quantizer $Q^{(k)}$), outputs a reproduction sequence by feeding the incoming $Q^{(k)^{-1}}(y_i)$, where $Q^{(k)^{-1}}(\cdot)$: $, \{1, \ldots, M\} \to [0, 1]$ denotes the inverse transformation of the quantizer $Q^{(k)}$, into the linear feedback system characterizing the decoder $D^{(k)}$, where the decoder $D$ assumes that up to this point in time the input to the system was constant at say zero. The only thing left to verify is that the performance of this scheme on the $k$th block is not too heavily deteriorated, relative to the performance of $(E^{(k)}, D^{(k)})$, by the fact that the input to the decoder $D^{(k)}$ up to time $i$ was in fact $\{Q^{(k)^{-1}}(y_j)\}_{j<i}$ and not constant at zero as the decoder $D$ is assuming. This can be done via the following crude calculation: if $h$ is the linear filter associated with any scheme in $\mathcal{F}_{\mathrm{DPCM}}(s, \eta, R)$ and $H$ is its Fourier transform then the Fourier transform of the impulse response $f$ of the decoder (the linear feedback system) is given by $F = 1/(1 - H)$ or, more explicitly, by

$$F(e^{j\omega}) = \sum_{l=0}^{\infty} \left[ \sum_{k=1}^{s} h_k e^{-jk\omega} \right]^l. \tag{24}$$

The fact that $\|h\|_1 \leq 1 - \eta$ clearly implies that for arbitrary $\boldsymbol{y}, \tilde{\boldsymbol{y}} \in \hat{\mathcal{X}}$ which coincide between any time $t - r$ and $t$, i.e., $y_{t'} = \tilde{y}_{t'}$ for $t - r \leq t' \leq t$ then

$$|(\boldsymbol{y} * f)_t - (\tilde{\boldsymbol{y}} * f)_t| \leq B \cdot \sum_{m=\lfloor r/s \rfloor}^{\infty} (1 - \eta)^m$$

$$\leq B \frac{(1-\eta)^{r/s}}{\eta(1-\eta)} = (\tilde{B}/\eta)(1-\eta)^{r/s}.$$

Now, if $a, b, c \in [-B, B]$ and $|b - c|$ is small then

$$(a - b)^2 \leq (a - c)^2 + 3B|b - c|.$$

It thus follows that the differences between the (unnormalized) cumulative distortions of $(E, D)$ and $(E^{(k)}, D^{(k)})$ on all blocks (of all lengths) are uniformly upper-bounded by

$$\sum_{r=0}^{\infty} (\tilde{B}/\eta)(1-\eta)^{r/s} = C(s, \eta) < \infty.$$

Hence, an inequality is obtained which is analogous to (7), with $s$ replaced by $s + C(s, \eta)$. The remainder of the proof carries through verbatim. □

## V. THE NOISY SETTING

In this section, we consider the case where the individual sequence is corrupted by noise. Limited-delay coding schemes are sought, for this case, which operate on the noisy sequence and produce a reconstruction sequence which is judged with respect to the clean individual sequence. Contrary to the noise-free case of the previous sections, we shall consider the case where the reference class consists of time-invariant sliding-window schemes. It is currently unknown whether the approach we will present can be similarly applied to handle reference classes of other types.

### A. Problem Formulation

We now formalize the notion of a limited-delay sequential source code for the case where the individual sequence to encode is corrupted by noise. Specifically, we assume now, as in the noise-free case, that there is an individual sequence $x_1, x_2, \ldots, x_i \in \mathcal{X}$ to encode. The encoder–decoder pair, however, accesses the sequence $Z_1, Z_2, \ldots, Z_i \in \mathcal{Z}$, which is the output of the fixed memoryless channel whose input is the individual sequence of interest $x_1, x_2, \ldots$. For simplicity of the exposition, we assume that $\mathcal{X} = \mathcal{Z}$ is finite and we let $M$ denote the $\mathcal{X} \times \mathcal{X}$ channel transition probability matrix, which we assume invertible. The approach we will present can be applied to more general cases.

For this setting, we define a delay $\delta = d_e + d_d$ $(d_e, d_d \geq 0)$ sequential scheme for combined filtering and compression of fixed rate $R = \log M$ with a randomized encoder by a pair $(E, D)$. The randomized encoder $E$ is given by a sequence $\{E_i\}_{i=1}^{\infty}$, where $E_i$: $\mathcal{X}^{i+d_e} \times [0, 1]^i \to \{1, 2, \ldots, M\}$. The decoder $D$ is given, as in the noise-free setting, by a sequence $\{D_i\}_{i=1}^{\infty}$, where $D_i$: $\{1, 2, \ldots, M\}^{i+d_d} \to \hat{\mathcal{X}}$. This scheme operates as follows. The encoder produces the $i$th channel symbol $Y_i \in \{1, 2, \ldots, M\}$ based on $Z^{i+d_e}$ and on the random sequence $U^i$ according to $Y_i = E_i(Z^{i+d_e}, U^i)$, where $\{U_i\}_{i=1}^{\infty}$ is a randomization sequence of i.i.d. random variables, uniformly distributed on $[0, 1]$ and independent of $\{Z_i\}$. The decoder emits the reconstructed sequence $\hat{x}_1, \hat{x}_2, \ldots$ according to $\hat{x}_i = D_i(Y^{i+d_d})$. We continue to let $\mathcal{F}^{\delta}(R)$ $(\mathcal{F}^{\delta})$, where $R = \log M$, denote the class of all such combined filtering and compression schemes operating at the fixed rate $R$ since they have exactly the same structure as in the noise-free case.

The definition of the cumulative distortion of a scheme $(E, D) \in \mathcal{F}^{\delta}$ remains as in the noise-free setting

$$d_{(E, D)}^n(\boldsymbol{x}) = \sum_{i=1}^{n} d(x_i, \hat{x}_i) \tag{25}$$

where $d$: $\mathcal{X} \times \hat{\mathcal{X}} \to [0, B]$ is a bounded distortion measure $(B < \infty)$ and here the $\hat{x}_i$s on the right-hand side are generated by feeding the noisy sequence $\boldsymbol{Z} = Z_1, Z_2, \ldots$ into $(E, D)$ as described above. Note that $d_{(E, D)}^n(\boldsymbol{x})$ is a random variable which depends on the realization of the channel noise and of the randomization sequence $U^n$. As before, we similarly denote, for $n_1 \leq n_2$

$$d_{(E, D)}^{n_1, n_2}(\boldsymbol{x}) = \sum_{i=n_1}^{n_2} d(x_i, \hat{x}_i).$$

### B. Time-Invariant Sliding-Window Schemes

A member $(E, D)$ of $\mathcal{F}^{\delta}$ is a time-invariant sliding-window scheme if there exist $d_1, d_2, d_3, d_4 \geq 0$, $f$: $\mathcal{Z}^{d_1+d_2+1} \to \{1, 2, \ldots, M\}$ and $g$: $\{1, 2, \ldots, M\}^{d_3+d_4+1} \to \hat{\mathcal{X}}$ such that for all $i$, the $i$th-channel symbol generated by the encoder is given by $Y_i = f(Z_{i-d_1}^{i+d_2})$ and the $i$th reconstruction symbol by

$\hat{x}_i = g(Y_{i-d_3}^{i+d_4})$. We refer to $d_1$ ($d_3$) as the encoder- (decoder-) memory and to $d_2$ ($d_4$) as the encoder- (decoder-) horizon. We let $\mathcal{F}^{d_1,d_2,d_3,d_4} \subset \mathcal{F}^{d_2+d_4}$ denote the family of all such time-invariant sliding-window schemes.

*1) Estimating the Distortion of $(E, D) \in \mathcal{F}^{d_1,d_2,d_3,d_4}$:* Unfortunately, the cumulative distortion of a scheme $(E, D) \in \mathcal{F}^\delta$, which was defined in (25), depends on the individual sequence $\boldsymbol{x}$ and, therefore, is not available when only its noisy version $\boldsymbol{Z}$ is accessible. Motivated by the approach which guided the construction of the coding schemes in the noise-free setting, our first goal in the noisy setting is to obtain an efficient estimator for the distortion suffered by a scheme in $\mathcal{F}^{d_1,d_2,d_3,d_4}$, which is only based on the observed noisy sequence $\boldsymbol{Z}$. To this end, let $(E, D) \in \mathcal{F}^{d_1,d_2,d_3,d_4}$ be given and denote

$$\hat{x}_i = g(Y_{i-d_3}^{i+d_4}) = g(\{f(Z_{i'-d_1}^{i'+d_2})\}_{i'=i-d_3}^{i+d_4}) \equiv h(Z_{i-(d_1+d_3)}^{i+(d_2+d_4)}).$$

We can now write

$$d_{(E,D)}^n(\boldsymbol{x}) = \sum_{i=1}^n d\left(x_i, h\left(Z_{i-(d_1+d_3)}^{i+(d_2+d_4)}\right)\right) \tag{26}$$

$$= \sum_{i=1}^n d\left(x_i, h(Z_{i-j}^{i+k})\right) \tag{27}$$

$$= \sum_{(\boldsymbol{a},\boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{Z}^m} d(a_{j+1}, h(\boldsymbol{b})) N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b}) \tag{28}$$

where we denote $d_1 + d_3 = j$, $d_2 + d_4 = k$, $m = j + k + 1$, the $j + 1$th component of $\boldsymbol{a}$ by $a_{j+1}$, and

$$N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b}) \triangleq |\{1 \le i \le n : x_{i-j}^{i+k} = \boldsymbol{a}, Z_{i-j}^{i+k} = \boldsymbol{b}\}|.$$

Evidently, efficient estimators for the unobserved

$$\{N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})\}_{(\boldsymbol{a},\boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{Z}^m}$$

would lead to an efficient estimator for $d_{(E,D)}^n(\boldsymbol{x})$ by plugging these into (28).

Let $M(m)$ denote the transition matrix characterizing the memoryless noisy channel for vector inputs of length $m$, induced by the channel matrix $M$. Specifically, $M(m)$ is a $|\mathcal{X}|^m \times |\mathcal{X}|^m$ matrix whose entry at $(\boldsymbol{a}, \boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m$ (where the rows and columns of $M(m)$ are arranged according to some lexicographic ordering of $|\mathcal{X}|^m$) is given by the probability for noisy $Z_{i-j}^{i+k} = \boldsymbol{b}$ when $x_{i-j}^{i+k} = \boldsymbol{a}$. Since the channel is fixed and memoryless, this means

$$M(m)_{\boldsymbol{a},\boldsymbol{b}} = \prod_{k=1}^m M_{a_k, b_k}.$$

We will further let $M(m)_{\boldsymbol{a},\boldsymbol{b}}^{-1}$ denote the entry corresponding to $(\boldsymbol{a}, \boldsymbol{b})$ in the inverse matrix of $M(m)$, which exists by the existence of $M^{-1}$ which is assumed throughout (recall Section V-A). The first observation we make toward constructing an estimator for $N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})$ is the following.

*Lemma 4:* For all $i \ge 1$ and $x_{i-j}^{i+k} \in \mathcal{X}^m$, the observable $M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot M(m)_{Z_{i-j}^{i+k},\boldsymbol{a}}^{-1}$ is an unbiased estimator for (the unobserved) $1_{\{x_{i-j}^{i+k}=\boldsymbol{a}, Z_{i-j}^{i+k}=\boldsymbol{b}\}}$. That is,

$$E\left\{M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot M(m)_{Z_{i-j}^{i+k},\boldsymbol{a}}^{-1}\right\} = E1_{\{x_{i-j}^{i+k}=\boldsymbol{a}, Z_{i-j}^{i+k}=\boldsymbol{b}\}}. \tag{29}$$

*Proof:*

$$E\left\{M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot M(m)_{Z_{i-j}^{i+k},\boldsymbol{a}}^{-1}\right\}$$

$$= E\left\{M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot \sum_{\boldsymbol{b}' \in |\mathcal{X}|^m} 1_{\{Z_{i-j}^{i+k}=\boldsymbol{b}'\}} M(m)_{\boldsymbol{b}',\boldsymbol{a}}^{-1}\right\}$$

$$= E\left\{M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot \sum_{\boldsymbol{b}' \in |\mathcal{X}|^m} \cdot \left(\sum_{\boldsymbol{a}' \in |\mathcal{X}|^m} 1_{\{x_{i-j}^{i+k}=\boldsymbol{a}', Z_{i-j}^{i+k}=\boldsymbol{b}'\}}\right) M(m)_{\boldsymbol{b}',\boldsymbol{a}}^{-1}\right\}$$

$$= M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot \sum_{\boldsymbol{b}' \in |\mathcal{X}|^m} \cdot \left(\sum_{\boldsymbol{a}' \in |\mathcal{X}|^m} 1_{\{x_{i-j}^{i+k}=\boldsymbol{a}'\}} \cdot M(m)_{\boldsymbol{a}',\boldsymbol{b}'}\right) M(m)_{\boldsymbol{b}',\boldsymbol{a}}^{-1}$$

$$= M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot \sum_{\boldsymbol{b}' \in |\mathcal{X}|^m} M(m)_{x_{i-j}^{i+k},\boldsymbol{b}'} M(m)_{\boldsymbol{b}',\boldsymbol{a}}^{-1}$$

$$= M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot I_{x_{i-j}^{i+k},\boldsymbol{a}}$$

$$= M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot 1_{\{x_{i-j}^{i+k}=\boldsymbol{a}\}}$$

$$= E1_{\{x_{i-j}^{i+k}=\boldsymbol{a}, Z_{i-j}^{i+k}=\boldsymbol{b}\}} \tag{30}$$

where the third equality follows from the fact that

$$E1_{\{x_{i-j}^{i+k}=\boldsymbol{a}', Z_{i-j}^{i+k}=\boldsymbol{b}'\}} = 1_{\{x_{i-j}^{i+k}=\boldsymbol{a}'\}} \cdot M(m)_{\boldsymbol{a}',\boldsymbol{b}'}$$

and where $I$ denotes the identity matrix. $\square$

We now let

$$\hat{N}_{\boldsymbol{a},\boldsymbol{b}}^n(Z^{n+k}) \triangleq \sum_{i=1}^n M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot M(m)_{Z_{i-j}^{i+k},\boldsymbol{a}}^{-1} \tag{31}$$

be our estimator for the unobserved $N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})$. Note that the estimation error

$$N^n(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b}) - \hat{N}_{\boldsymbol{a},\boldsymbol{b}}^n(Z^{n+k})$$

$$= \sum_{i=1}^n \left(1_{\{x_{i-j}^{i+k}=\boldsymbol{a}, Z_{i-j}^{i+k}=\boldsymbol{b}\}} - M(m)_{\boldsymbol{a},\boldsymbol{b}} \cdot M(m)_{Z_{i-j}^{i+k},\boldsymbol{a}}^{-1}\right)$$

$$\triangleq \sum_{i=1}^n \delta_i, \tag{32}$$

is, by Lemma 4, a sum of zero-mean, bounded random variables. Furthermore, these random variables have a well-behaved dependence structure. In particular, for all $i \ge 1$, $\delta_i$ is independent of $\{\delta_j\}_{j>i+m}$. Hence, by standard techniques (cf., e.g., [5]), it can be shown that the estimation error will typically be $O(\sqrt{n})$, that the probability of its magnitude exceeding $n\varepsilon$, for any $\varepsilon > 0$, decays exponentially rapidly, pointwise asymptotic bounds on its magnitude in the order of $\sqrt{n \log \log n}$, etc. Motivated by this, we finally let our estimator for the unavailable cumulative distortion $d_{(E,D)}^n(\boldsymbol{x})$ be

$$\hat{d}_{(E,D)}^n(\boldsymbol{Z}) = \sum_{(\boldsymbol{a},\boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m} d(a_{j+1}, h(\boldsymbol{b})) \cdot \hat{N}_{\boldsymbol{a},\boldsymbol{b}}^n(Z^{n+k}). \tag{33}$$

*2) The Coding-Filtering Scheme for the Noisy Case:* Equipped with an estimator for the unobserved cumulative distortion

of a sliding-window scheme in the noisy setting, we now construct an adaptive scheme which is similar, in principle, to that of Section III. The essential difference is in the fact that the randomization is performed with respect to a distribution assigning weights which are exponentially proportional to the *estimated* past distortion. This approach has recently lead to efficient schemes in the context of prediction of individual sequences corrupted by noise [17].

*Theorem 5:* Let $\mathcal{A}$ be a finite subset of $\mathcal{F}^{d_1, d_2, d_3, d_4}(R)$ for some $d_1, d_2, d_3, d_4 \geq 0$ and a given $R = \log M$. Then there exists a source code $(E, D) \in \mathcal{F}^{\delta}(R)$ $(\delta = d_2 + d_4)$ such that for any $\varepsilon > 0$, $N$ sufficiently large so that

$$B\sqrt{\log |\mathcal{A}|/2} \cdot N^{-1/4} + BN^{-1/2}\left(\frac{1}{R}\log |\mathcal{A}| + \delta\right) < \varepsilon/3$$

and all $\boldsymbol{x} \in \mathcal{X}^{\infty}$

$$\Pr\left\{\frac{1}{N}\left[d_{(E, D)}^{N}(\boldsymbol{x}) - \min_{(E', D') \in \mathcal{A}} d_{(E', D')}^{N}(\boldsymbol{x})\right] \geq \varepsilon\right\}$$
$$\leq 5\sqrt{N}|\mathcal{A}||\mathcal{X}|^{2m} \exp\left\{-\sqrt{N}\varepsilon^2 C\right\} \quad (34)$$

where $C$ is a positive constant (which depends only on $B$, $R$, $\mathcal{A}$, $m$, and $\max_{(\boldsymbol{a}, \boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m} M(m)_{\boldsymbol{a}, \boldsymbol{b}}^{-1}$, the explicit value of which will be apparent in the proof).

Note that here, as in the noise-free setting, the scheme which we construct, in addition to complying with the delay limitations of the sliding-window schemes with which it competes, operates at the exact same rate.

We further remark, as we did in the noise-free case, that while the source code $(E, D)$ depends on the length $N$ of the sequence, it is straightforward to use it to obtain a strongly sequential (horizon-independent) scheme for which a concentration inequality in the form of (34) holds for all sufficiently large $N$. One can then apply the Borel–Cantelli lemma to obtain almost-sure asymptotic performance guarantees. Finally, we remark that in the above theorem, $\mathcal{A}$ can be taken as the set of all sliding-window schemes of up to a certain order (memory and horizon), which is a finite set in the finite alphabet case. By gradually increasing the order and employing the scheme of Theorem 5 tailored for this order on consecutive blocks of lengths which increase at the proper rate, a scheme can be obtained which asymptotically almost surely achieves the performance of any sliding-block scheme, no matter what the underlying individual sequence may be.

*Proof of Theorem 5:* The construction of $(E, D) \in \mathcal{F}^{d_2 + d_4}$ which follows is similar to that in the proof of Theorem 1, with the essential difference that the past cumulative distortions of the schemes in $\mathcal{A}$, which are now unavailable to the encoder, are replaced by their respective estimates for the purpose of randomization.

Fix an $l \ll N$ and divide the time axis $i = 1, 2, \ldots, N$ into $n = N/l$ consecutive nonoverlapping blocks (assume $l$ divides $N$). At the beginning of the $k$th block $1 \leq k \leq n$, i.e., at the $i = (k-1)l + 1$th channel use, after seeing $Z^{i+\delta}$ and when $U^i$ are available, the encoder uses $U_i$ to generate $(E^{(k)}, D^{(k)})$,

a $\mathcal{A}$-valued random variable with distribution satisfying almost surely

$$\Pr\left\{\left(E^{(k)}, D^{(k)}\right) = (E', D')|\mathcal{G}_Z\right\}$$
$$= \frac{\exp\left\{-\eta \hat{d}_{(E', D')}^{(k-1)l}(\boldsymbol{Z})\right\}}{\sum_{(\tilde{E}, \tilde{D}) \in \mathcal{A}} \exp\left\{-\eta \hat{d}_{(\tilde{E}, \tilde{D})}^{(k-1)l}(\boldsymbol{Z})\right\}}, \qquad \forall (E', D') \in \mathcal{A} \quad (35)$$

where $\mathcal{G}_Z$ is the smallest sigma-algebra with respect to (w.r.t.) which all $Z_i$'s are measurable and $\eta > 0$ is taken as in Theorem 1. For convenience, we let $f^{(k)}$ and $g^{(k)}$ denote the sliding-window encoder and decoder, respectively, characterizing the source code $(E^{(k)}, D^{(k)})$. The encoder now dedicates the first $\lceil\frac{1}{R}\log |\mathcal{A}|\rceil$ channel symbols at the beginning of the $k$th block, i.e., $Y_i$ for

$$i = (k-1)l + 1, \ldots (k-1)l + \left\lceil\frac{1}{R}\log |\mathcal{A}|\right\rceil$$

to convey to the decoder the identity of $g^{(k)}$. At the remainder of the block, i.e., at channel uses

$$i = (k-1)l + \left\lceil\frac{1}{R}\log |\mathcal{A}|\right\rceil + 1, \ldots, kl$$

the encoder produces the channel symbols

$$Y_i = f^{(k)}(Z_{i-d_1}^{i+d_2}).$$

On the decoder's side, at the beginning of the block, at channel uses

$$i = (k-1)l + 1, \ldots, (k-1)l + \left\lceil\frac{1}{R}\log |\mathcal{A}|\right\rceil + d_3$$

the decoder outputs an arbitrary reproduction sequence of $\hat{x}_i$'s. From the $i = (k-1)l + \lceil\frac{1}{R}\log |\mathcal{A}|\rceil + d_3 + 1$th channel use up to the end of the block $i = kl$, the decoder, knowing the identity of $g^{(k)}$ and the output of $f^{(k)}$ at least $d_3$ channel uses back, outputs the reproduction sequence according to $\hat{x}_i = g^{(k)}(Y_{i-d_3}^{i+d_4})$.

To analyze the performance of the above scheme and to establish (34), we will consider a genie-aided scheme. Specifically, let $(E_G, D_G)$ be a source code which is allowed to access the clean sequence $\boldsymbol{x}$ as well as its noisy version $\boldsymbol{Z}$. This scheme operates exactly as $(E, D)$ described above, with the only difference that the pair $(E^{(k)}, D^{(k)})$ is generated (using the randomization sequence) according almost surely to

$$\Pr\left\{\left(E^{(k)}, D^{(k)}\right) = (E', D')|\mathcal{G}_Z\right\}$$
$$= \frac{\exp\left\{-\eta d_{(E', D')}^{(k-1)l}(\boldsymbol{x})\right\}}{\sum_{(\tilde{E}, \tilde{D}) \in \mathcal{A}} \exp\left\{-\eta d_{(\tilde{E}, \tilde{D})}^{(k-1)l}(\boldsymbol{x})\right\}}, \qquad \forall (E', D') \in \mathcal{A} \quad (36)$$

rather than according to (35) (this is where the "genie" comes in since the generation of $(E^{(k)}, D^{(k)})$ from such a distribution requires knowledge of $\{d_{(\tilde{E}, \tilde{D})}^{(k-1)l}(\boldsymbol{x})\}_{(\tilde{E}, \tilde{D}) \in \mathcal{A}}$ which depends on the clean sequence $\boldsymbol{x}$). All the stages of the proof of Theorem 1 carry over essentially verbatim for $(E_G, D_G)$ (with $s$ replaced by $\delta$). In particular, we have, as in (9), almost surely

$$\sum_{k=1}^{n} \mathbb{E}\left[d_{(E_G, D_G)}^{(k-1)l+1, kl}(\boldsymbol{x})\Big|\mathcal{G}_Z\right] - \min_{(\tilde{E}, \tilde{D}) \in \mathcal{A}} d_{(\tilde{E}, \tilde{D})}^{N}(\boldsymbol{x})$$
$$\leq B\sqrt{\log |\mathcal{A}|/2} \cdot ln^{1/2} + Bn\left(\frac{1}{R}\log |\mathcal{A}| + d\right). \quad (37)$$

According to [14, Lemma 3] for arbitrary

$$(a_1, \ldots, a_N), (b_1, \ldots, b_N) \in \mathbb{R}^N$$

and $C > 0$

$$\sup_{(\alpha_1, \ldots, \alpha_N) \in [0, C]^N} \left| \frac{\sum_{j=1}^{N} e^{-a_j} \alpha_j}{\sum_{j=1}^{N} e^{-a_j}} - \frac{\sum_{j=1}^{N} e^{-b_j} \alpha_j}{\sum_{j=1}^{N} e^{-b_j}} \right| \leq 2C \max_{1 \leq j \leq N} |a_j - b_j|. \quad (38)$$

From the construction of the genie-aided $(E_G, D_G)$ and our legitimate $(E, D)$, it follows by an application of (38) that for all $1 \leq k \leq n$

$$\left| \mathbb{E}\left[ d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right] - \mathbb{E}\left[ d_{(E_G, D_G)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right] \right|$$
$$\leq 2Bl\eta \max_{(E', D') \in \mathcal{A}} \left| \hat{d}_{(E', D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E', D')}^{(k-1)l}(\boldsymbol{x}) \right| \text{ a.s.} \quad (39)$$

This leads to

$$\sum_{k=1}^{n} \mathbb{E}\left[ d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right] - \min_{(\check{E}, \check{D}) \in \mathcal{A}} d_{(\check{E}, \check{D})}^{N}(\boldsymbol{x})$$
$$\leq B\sqrt{\log|\mathcal{A}|/2} \cdot ln^{1/2} + Bn\left( \frac{1}{R}\log|\mathcal{A}| + d \right)$$
$$+ 2Bl\eta \sum_{k=1}^{n} \max_{(E', D') \in \mathcal{A}} \left| \hat{d}_{(E', D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E', D')}^{(k-1)l}(\boldsymbol{x}) \right| \text{ a.s.}$$
$$(40)$$

where we recall that $\eta$ is taken to be $\sqrt{8\log|\mathcal{A}|/(l^2 B^2 n)}$. To conclude, we show that the first term on the left-hand side of (40) is, with high probability, close to $d_{(E, D)}^{N}(\boldsymbol{x})$ and that the last term on the right-hand side of (40) is small with high probability. Specifically, it is shown in the Appendix that for all $\varepsilon > 0$

$$\Pr\left\{ \frac{1}{N}\left[ d_{(E, D)}^{N}(\boldsymbol{x}) - \sum_{k=1}^{n} \mathbb{E}\left[ d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right] \right] > \varepsilon \right\}$$
$$\leq \exp\left\{ -\frac{2N^2\varepsilon^2}{n(2lB)^2} \right\} \quad (41)$$

and

$$\Pr\left\{ \frac{1}{N} 2Bl\eta \sum_{k=1}^{n} \max_{(E', D') \in \mathcal{A}} \left| \hat{d}_{(E', D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E', D')}^{(k-1)l}(\boldsymbol{x}) \right| > \varepsilon \right\}$$
$$\leq 4n|\mathcal{A}| |\mathcal{X}|^{2m} \exp\left\{ -\frac{N}{n} C_1 \varepsilon^2 \right\} \quad (42)$$

where

$$C_1 = (32\log|\mathcal{A}|mB^2\tilde{B}^2|\mathcal{X}|^{4m})^{-1}$$

and

$$\tilde{B} = \max_{(\boldsymbol{a}, \boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m} M(m)_{\boldsymbol{a}, \boldsymbol{b}}^{-1}.$$

Choosing this time $n = l = \sqrt{N}$ (as opposed to the $n \propto N^{2/3}$, $l \propto N^{1/3}$ regime of the noise-free case), the above two upper bounds give the same exponential speed of decay ($\sqrt{N}$). Com-

bining inequalities (40)–(42), taking $n = l = \sqrt{N}$, applying the union bound, and noting that $C_1 \leq 1/(2B^2)$ gives

$$\Pr\left\{ \frac{1}{N}\left[ d_{(E, D)}^{N}(\boldsymbol{x}) - \min_{(\check{E}, \check{D}) \in \mathcal{A}} d_{(\check{E}, \check{D})}^{N}(\boldsymbol{x}) \right] \geq \varepsilon \right\}$$
$$\leq \exp\left\{ -\sqrt{N}(\varepsilon/3)^2 \frac{1}{2B^2} \right\}$$
$$+ 4\sqrt{N} |\mathcal{A}| |\mathcal{X}|^{2m} \exp\left\{ -\sqrt{N}(\varepsilon/3)^2 C_1 \right\}$$
$$\leq 5\sqrt{N} |\mathcal{A}| |\mathcal{X}|^{2m} \exp\left\{ -\sqrt{N}(\varepsilon/3)^2 C_1 \right\} \quad (43)$$

for all $N$ sufficiently large so that

$$B\sqrt{\log|\mathcal{A}|/2} \cdot N^{-1/4} + BN^{-1/2}(\tfrac{1}{R}\log|\mathcal{A}| + \delta) < \varepsilon/3. \quad \square$$

Note that the main idea in the above proof is similar to that underlying proof of Theorem 1 with the added ingredient that this time, the randomization at the beginning of each block is performed according to (35), where estimators for the unobserved distortions associated with the schemes in the reference class are employed. Evidently, Theorem 5 can be extended to account for reference classes other than time-invariant sliding-window schemes, provided a way is found to efficiently (uniformly for all individual sequences) estimate distortions associated with more general schemes. The construction of such estimators for general schemes has defied our efforts thus far.

## VI. FUTURE DIRECTIONS

Directions for related future research include.

1) As discussed in Section III, the implementation of the source code constructed therein is impractical for a large reference class as the algorithm is required to effectively run all the schemes of the reference class in parallel. It would be of interest to find a more practical scheme which is universal in the sense dealt with in this work. We mention, however, that for reference classes of a relatively small size (a few dozens at the most), the algorithm of Section III has been implemented and shown to be quite effective for various types of data. Details can be found at http://visl.technion.ac.il/projects/2001s16/.

2) Extension of the main result of Section V to reference classes other than those consisting of time-invariant sliding-window schemes.

## APPENDIX
### TECHNICAL STAGES IN THE PROOF OF THEOREM 5

*Proof of (41):* Since the randomization sequence $\boldsymbol{U}$ is independent of the channel's noisy output, for almost every realization of the noisy sequence $\boldsymbol{Z}$, the difference

$$d_{(E, D)}^{N}(\boldsymbol{x}) - \sum_{k=1}^{n} \mathbb{E}\left[ d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right]$$
$$= \sum_{k=1}^{n} d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) - \sum_{k=1}^{n} \mathbb{E}\left[ d_{(E, D)}^{(k-1)l+1, kl}(\boldsymbol{x}) \Big| \mathcal{G}_Z \right]$$

under $\Pr\{\cdot|\mathcal{G}_Z\}$ is distributed as a sum of $n$ zero-mean independent random variables bounded in magnitude by $l \cdot B$. Applying Hoeffding's inequality [3, Theorem 8.1] gives

$$\Pr\left\{\frac{1}{N}\left(d_{(E,D)}^N(\boldsymbol{x}) - \sum_{k=1}^{n}\mathbb{E}\left[d_{(E,D)}^{(k-1)l+1,\,kl}(\boldsymbol{x})\Big|\mathcal{G}_Z\right]\right) > \varepsilon\;\Big|\;\mathcal{G}_Z\right\}$$
$$\leq \exp\left\{-\frac{2N^2\varepsilon^2}{n(2lB)^2}\right\}\;\text{a.s.}\quad(A1)$$

Taking expectation over the above inequality completes the proof. $\square$

*Proof of (42):* Using the union bound generously, we have for all $\varepsilon > 0$

$$\Pr\left\{\sum_{k=1}^{n}\max_{(E',D')\in\mathcal{A}}\left|\hat{d}_{(E',D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E',D')}^{(k-1)l}(\boldsymbol{x})\right| > \varepsilon\right\}$$

$$\leq \sum_{k=1}^{n}\sum_{(E',D')\in\mathcal{A}}\Pr\left\{\left|\hat{d}_{(E',D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E',D')}^{(k-1)l}(\boldsymbol{x})\right| > \varepsilon/n\right\}$$

$$= \sum_{k=1}^{n}\sum_{(E',D')\in\mathcal{A}}\Pr\left\{\left|\sum_{(\boldsymbol{a},\boldsymbol{b})\in\mathcal{X}^m\times\mathcal{X}^m}d(a_{j+1}, h(\boldsymbol{b}))\right.\right.$$
$$\left.\left.\cdot\left(\hat{N}_{\boldsymbol{a},\boldsymbol{b}}^{(k-1)l}(Z^{n+k}) - N^{(k-1)l}(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})\right)\right| > \varepsilon/n\right\}$$

$$\leq \sum_{k=1}^{n}\sum_{(E',D')\in\mathcal{A}}\sum_{(\boldsymbol{a},\boldsymbol{b})\in\mathcal{X}^m\times\mathcal{X}^m}$$
$$\cdot\Pr\left\{\left|\hat{N}_{\boldsymbol{a},\boldsymbol{b}}^{(k-1)l}\left(Z^{(k-1)l+d}\right) - N^{(k-1)l}(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})\right|\right.$$
$$\left.> \frac{\varepsilon}{Bn|\mathcal{X}|^{2m}}\right\}.\quad(A2)$$

To bound each of the summands in (A2), recall that it was established in Section V-B1 that for each $(\boldsymbol{a}, \boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m$ and all $t$, $\hat{N}_{\boldsymbol{a},\boldsymbol{b}}^t(Z^{t+d}) - N^t(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})$ is a sum of $t$ random variables $\sum_{i=1}^{t}\delta_i$, where the $\delta_i$'s are almost surely bounded by

$$\tilde{B} = \max_{(\boldsymbol{a},\boldsymbol{b})\in\mathcal{X}^m\times\mathcal{X}^m}M(m)_{\boldsymbol{a},\boldsymbol{b}}^{-1}$$

and where $\{\delta_j\}_{j\leq i}$ is independent of $\{\delta_j\}_{j>i+m}$ for all $i$. For a given $t$, we can divide the set of indexes $\{1, \ldots, t\}$ into $\lceil t/(2m)\rceil$ "intervals" of length $2m$ each (except, possibly, for the last one which may be shorter). The sum from $1$ to $t$ can be decomposed into a sum over those indexes that belong to the first $m$ indexes of some interval, and a sum over the complementary set of indexes. By the aforementioned independence property of $\{\delta_i\}$, each of these two sums is a sum of (at the most) $\lceil t/(2m)\rceil$ independent, zero-mean random variables bounded in magnitude by $m \cdot \tilde{B}$. Therefore, applying Hoeffding's inequality to each of these two sums and employing the union bound gives for all $\varepsilon > 0$, $(\boldsymbol{a}, \boldsymbol{b}) \in \mathcal{X}^m \times \mathcal{X}^m$, and all $t$

$$\Pr\left\{\left|\hat{N}_{\boldsymbol{a},\boldsymbol{b}}^t\left(Z^{(k-1)l+d}\right) - N^t(\boldsymbol{x}, \boldsymbol{Z}, \boldsymbol{a}, \boldsymbol{b})\right| > \varepsilon\right\}$$
$$\leq 4\exp\left\{-\frac{2(\varepsilon/2)^2}{\lceil t/(2m)\rceil(m\tilde{B})^2}\right\}$$
$$\leq 4\exp\left\{-\frac{\varepsilon^2}{2tm\tilde{B}^2}\right\}.\quad(A3)$$

Consequently, each of the summands in the right-hand side of (A2) is upper-bounded by

$$4\exp\left\{-\frac{\left(\frac{\varepsilon}{Bn|\mathcal{X}|^{2m}}\right)^2}{2(k-1)lm\tilde{B}^2}\right\} \leq 4\exp\left\{-\frac{\left(\frac{\varepsilon}{Bn|\mathcal{X}|^{2m}}\right)^2}{2Nm\tilde{B}^2}\right\}\quad(A4)$$

which finally gives

$$\Pr\left\{\sum_{k=1}^{n}\max_{(E',D')\in\mathcal{A}}\left|\hat{d}_{(E',D')}^{(k-1)l}(\boldsymbol{Z}) - d_{(E',D')}^{(k-1)l}(\boldsymbol{x})\right| > \varepsilon\right\}$$
$$\leq 4n|\mathcal{A}||\mathcal{X}|^{2m}\exp\left\{-\frac{\varepsilon^2}{Nn^2mB^2\tilde{B}^2|\mathcal{X}|^{4m}}\right\}.\quad(A5)$$

Inequality (42) is exactly (A5) with the assignment $\varepsilon \to \varepsilon N/(2Bl\eta)$. $\square$

## ACKNOWLEDGMENT

## REFERENCES

[1] N. Cesa-Bianchi, Y. Freund, D. P. Helmbold, D. Haussler, R. Schapire, and M. K. Warmuth, "How to use expert advice," *J. Assoc. Comput. Mach.*, vol. 44, no. 3, pp. 427–485, 1997.

[2] N. Cesa-Bianchi and G. Lugosi, "On sequential prediction of individual sequences relative to a set of experts," *Ann. Statist.*, vol. 27, no. 6, pp. 1865–1895, 1999.

[3] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York: Springer-Verlag, 1996.

[4] R. L. Dobrushin and B. S. Tsybakov, "Information transmission with additional noise," *IEEE Trans. Inform. Theory*, vol. IT-8, pp. 293–304, Sept. 1962.

[5] E. Eberlein and M. S. Taqqu, Eds., *Dependence in Probability and Statistics*. Boston, MA: Birkhauser, 1986.

[6] Y. Ephraim and R. M. Gray, "A unified approach for encoding clean and noisy sources by means of waveform and autoregressive model vector quantization," *IEEE Trans. Inform. Theory*, vol. 34, pp. 826–834, July 1988.

[7] M. Feder, N. Merhav, and M. Gutman, "Universal prediction of individual sequences," *IEEE Trans. Inform. Theory*, vol. 38, pp. 1258–1270, July 1992.

[8] T. L. Fine, "Optimum mean-square quantization of a noisy input," *IEEE Trans. Inform. Theory*, vol. IT-11, pp. 293–294, Apr. 1965.

[9] R. M. Gray, "Time-invariant trellis encoding of ergodic discrete-time sources with a fidelity criterion," *IEEE Trans. Inform. Theory*, vol. IT-23, pp. 71–83, Jan. 1977.

[10] R. M. Gray and D. L. Neuhoff, "Quantization," *IEEE Trans. Inform. Theory*, vol. 44, pp. 2325–2383, Oct. 1998.

[11] M. E. Hellman, "Convolutional source encoding," *IEEE Trans. Inform. Theory*, vol. IT-21, pp. 651–656, Nov. 1975.

[12] N. S. Jayant and P. Noll, *Digital Coding of Waveforms: Principles and Applications to Speech and Video*. Englewood Cliffs, NJ: Prentice-Hall, 1984.

[13] I. Lev, "Universal signal enhancement by coding," Master's thesis, Technion–Israel Inst. Technol., Haifa, Israel, Apr. 1996.

[14] T. Linder and G. Lugosi, "A zero-delay sequential scheme for lossy coding of individual sequences," *IEEE Trans. Inform. Theory*, vol. 47, pp. 2533–2538, Sept. 2001.

[15] N. Merhav and M. Feder, "Universal schemes for sequential decision from individual data sequences," *IEEE Trans. Inform. Theory*, vol. 39, pp. 1280–1292, July 1993.

[16] B. Natarajan, "Filtering random noise via data compression," in *Proc. Data Compression Conf., DCC'93*, 1993, pp. 60–69.

[17] T. Weissman and N. Merhav, "Universal prediction of individual binary sequences in the presence of noise," *IEEE Trans. Inform. Theory*, vol. 47, pp. 2151–2173, Sept. 2001.

[18] J. K. Wolf and J. Ziv, "Transmission of noisy information to a noisy receiver with minimum distortion," *IEEE Trans. Inform. Theory*, vol. IT-16, no. 4, pp. 406–411, July 1970.

[19] E. H. Yang and J. Kieffer, "Simple universal lossy data compression schemes derived from the Lempel–Ziv algorithm," *IEEE Trans. Inform. Theory*, vol. 42, pp. 239–245, Jan. 1996.

[20] R. Zamir and M. Feder, "On universal quantization by randomized uniform/lattice quantizers," *IEEE Trans. Inform. Theory*, vol. IT-38, pp. 428–436, Mar. 1992.

[21] J. Ziv, "Distortion-rate theory for individual sequences," *IEEE Trans. Inform. Theory*, vol. IT-26, pp. 137–143, Mar. 1980.

[22] ——, "On universal quantization," *IEEE Trans. Inform. Theory*, vol. IT-31, pp. 344–347, May 1985.