

# Discrete Denoising with Shifts \*

Taesup Moon      Tsachy Weissman

February 1, 2008

## Abstract

We introduce S-DUDE, a new algorithm for denoising DMC-corrupted data. The algorithm, which generalizes the recently introduced DUDE (Discrete Universal DENOISER) of Weissman et al., aims to compete with a genie that has access, in addition to the noisy data, also to the underlying clean data, and can choose to switch, up to  $m$  times, between sliding window denoisers in a way that minimizes the overall loss. When the underlying data form an individual sequence, we show that the S-DUDE performs essentially as well as this genie, provided that  $m$  is sub-linear in the size of the data. When the clean data is emitted by a piecewise stationary process, we show that the S-DUDE achieves the optimum distribution-dependent performance, provided that the same sub-linearity condition is imposed on the number of switches. To further substantiate the universal optimality of the S-DUDE, we show that when the number of switches is allowed to grow linearly with the size of the data, *any* (sequence of) scheme(s) fails to compete in the above senses. Using dynamic programming, we derive an efficient implementation of the S-DUDE, which has complexity (time and memory) growing only linearly with the data size and the number of switches  $m$ . Preliminary experimental results are presented, suggesting that S-DUDE has the capacity to significantly improve on the performance attained by the original DUDE in applications where the nature of the data abruptly changes in time (or space), as is often the case in practice.

*Index Terms*- Discrete denoising, competitive analysis, individual sequence, universal algorithms, piecewise stationary processes, dynamic programming, discrete memoryless channel (DMC), switching experts, forward-backward recursions.

## 1 Introduction

Discrete denoising is the problem of reconstructing the components of a finite-alphabet sequence based on the *entire* observation of its Discrete Memoryless Channel (DMC)-corrupted version. The quality of the reconstruction is evaluated via a user-specified (single-letter) loss function. *Universal* discrete denoising, in which no statistical or other properties are known a priori about the underlying clean data and the goal is to attain optimum performance, was considered and solved in [1]. The main problem setting there is the “semi-stochastic” one, in which the underlying signal is assumed to be an “individual sequence,” and the randomness is due solely to the channel noise. In this setting, it is unreasonable to expect to attain the best performance among all the denoisers in the world, since for every given sequence, there exists a denoiser that recovers all the sequence components perfectly. Thus, [1] limits the comparison class, a.k.a. expert class, and uses the competitive analysis approach. Specifically, it is shown that regardless of what the underlying individual sequence may be, the Discrete Universal DENOISER (DUDE) essentially attains the performance of the best sliding window denoiser that would be chosen by a genie with access to the

---

\* Authors are with the Department of Electrical Engineering, Stanford University, Stanford, CA 94305, USA. T. Weissman is also with the Department of Electrical Engineering, Technion, Haifa 32000, Israel. E-mails: {tsmoon, tsachy}@stanford.edu. This work was partially supported by NSF awards 0512140 and 0546535, and by a Samsung scholarship.

underlying clean sequence, in addition to the observed noisy sequence. This semi-stochastic setting result is shown in [1] to imply the stochastic setting result, i.e., that for any underlying stationary signal, the DUDE attains the optimal distribution-dependent performance.

The setting of an arbitrary individual sequence, combined with competitive analysis, has been very popular in many other research areas, especially for problems of sequential decision-making. Examples include universal compression [4], universal prediction [5], universal filtering [2], repeated game playing [6, 7, 8], universal portfolios [9], online learning [10, 11], zero-delay coding [12, 13], and much more. A comprehensive account of this line of research can be found in [14]. The beauty of this approach is the fact that it leads to the construction of schemes that perform, on every individual sequence, essentially as well as the best in a class of experts, which is the performance of a genie that had hindsight on the entire sequence before selecting his actions. Moreover, if the expert class is judiciously chosen, the relative sense of such a performance guarantees can, in many cases, imply optimum performance in absolute senses as well.

One extension to this approach is competition with an expert class and a genie that has the freedom to form a *compound* action, which breaks the sequence into a certain (limited) number of segments, applies different experts in each segment, and achieves an even better performance overall. Note that the optimal segmentation of the sequence and the choice of the best expert in each segment is also determined by hindsight. Clearly, competing with the best compound action is more challenging, since the number of possible compound actions is exponential in the sequence length  $n$ , and the brute-force vanilla implementation of the ordinary universal scheme requires prohibitive complexity. However, clever schemes with linear complexity that successfully track the best segments and experts have been devised in many different areas, such as online learning, universal prediction [15, 16], universal compression [17, 18], online linear regression [19], universal portfolios [20], and zero-delay lossy source coding [22].

In this paper, we expand the idea of compound actions and apply it to the discrete denoising problem. The motivation of this expansion is natural: the characteristics of the underlying data in the denoising problem often tend to be time- or space-varying. In this case, determining the best segmentation and the best expert for each segment requires complete knowledge of both clean and noisy sequences. Therefore, whereas the challenge in sequential decision-making problems is to track the shift of the best expert based on the past, true observation, the challenge in the denoising problem is to learn the shift based on the entire, but noisy, observation. We extend DUDE to meet this challenge and provide results that parallel and strengthen those of [1].

Specifically, we introduce the S-DUDE and show first that, for every underlying noiseless sequence, it attains the performance of the best compound finite-order sliding window denoiser (concretely defined later), both in expectation and in a high probability sense. We develop our scheme in the semi-stochastic setting as in [1]. The toolbox for the construction and analysis of our scheme draws on ideas developed in [2]. We circumvent the difficulty of not knowing the exact true loss by using an observable unbiased estimate of it. This kind of an estimate has proved to be very useful in [2] and [3] to devise schemes for filtering and for denoising with dynamic contexts. Building on this semi-stochastic setting result, we also establish a stochastic setting result, which can be thought of as a generalization and strengthening of the stochastic setting results of [1], from the world of stationary processes to that of piecewise stationary processes.

Our stochastic setting has connections to other areas, such as change-point detection problems in statistics [23, 24] and switching linear dynamical systems in machine learning and signal processing [25, 26]. Both of these lines of research share a common approach with S-DUDE, in that they try to learn the change of the underlying time-varying parameter or state of stochastic models, based on noisy observations of the parameter or state. One difference is that, whereas our goal is the noncausal estimation, i.e., denoising, of the general underlying piecewise stationary process, the change-point detection problems mainly focus on sequentially detecting the time point where the change of model happened. Another difference is in

that the switching linear dynamical systems focus on a special class of underlying processes, the linear dynamical system. In addition, they deal with continuous-valued signals, whereas our focus is the discrete case, with finite-alphabet signals.

As we explain in detail, the S-DUDE can be practically implemented using a two-pass algorithm with complexity (both space and time) linear in the sequence length and the number of switches. We also present initial experimental results that demonstrate the S-DUDE's potential to outperform the DUDE on both simulated and real data.

The remainder of the paper is organized as follows. Section 2 provides the notation, preliminaries and background for the paper; in Section 3 we present our scheme and establish its strong universality properties via an analysis of its performance in the semi-stochastic setting. Section 4 establishes the universality of our scheme in a fully stochastic setting, where the underlying noiseless sequence is emitted by a piecewise stationary process. Algorithmic aspects and complexity of the actual implementation of the scheme is considered in Section 5, and some experimental results are displayed in Section 6. In Section 7 we conclude with a summary of our findings and some possible future research directions.

## 2 Notation, Preliminaries, and Motivation

### 2-A Notation

We use a combination of notation of [1] and [2]. Let  $\mathcal{X}$ ,  $\mathcal{Z}$ ,  $\hat{\mathcal{X}}$  denote, respectively, the alphabet of the clean, noisy, and reconstructed sources, which are assumed to be finite. As in [1] and [2], the noisy sequence is a DMC-corrupted version of the clean one, where the channel matrix  $\mathbf{\Pi} = \{\Pi(x, z)\}_{x \in \mathcal{X}, z \in \mathcal{Z}}$ ,  $\Pi(x, z)$  denoting the probability of a noisy symbol  $z$  when the underlying clean symbol is  $x$ , is assumed to be *known and fixed* throughout the paper, and *of full row rank*. The  $z$ -th column of  $\mathbf{\Pi}$  will be denoted as  $\pi_z$ . Upper case letters will denote random variables as usual; lower case letters will denote either individual deterministic quantities or specific realizations of random variables.

Without loss of generality, the elements of any finite set  $\mathcal{V}$  will be identified with  $\{0, 1, \dots, |\mathcal{V}| - 1\}$ . We let  $\mathcal{V}^\infty$  denote the set of one-sided infinite sequences with  $\mathcal{V}$ -valued components, i.e.,  $\mathbf{v} \in \mathcal{V}^\infty$  is of the form  $\mathbf{v} = (v_1, v_2, \dots)$ ,  $v_i \in \mathcal{V}$ ,  $i \geq 1$ . For  $\mathbf{v} \in \mathcal{V}^\infty$ , let  $v^n = (v_1, \dots, v_n)$  and  $v_m^n = (v_m, \dots, v_n)$ . Furthermore, we let  $v^{n \setminus t}$  denote the sequence  $v^{t-1}v_{t+1}^n$ .  $\mathbb{R}^\mathcal{V}$  is a space of  $|\mathcal{V}|$ -dimensional column vectors with real-valued components indexed by the elements of  $\mathcal{V}$ . The  $a$ -th component of  $q \in \mathbb{R}^\mathcal{V}$  will be denoted by either  $q_a$  or  $q[a]$ . Subscripting a vector or a matrix by “max” will represent the difference between the maximum and minimum of all its components. Thus, for example, if  $\Gamma$  is a  $|\mathcal{Z}| \times |\mathcal{X}|$  matrix, then  $\Gamma_{\max}$  stands for  $\max_{x \in \mathcal{X}, z \in \mathcal{Z}} \Gamma(z, x) - \min_{x \in \mathcal{X}, z \in \mathcal{Z}} \Gamma(z, x)$  (in particular, if the components of  $\Gamma$  are nonnegative and  $\Gamma(z, x) = 0$  for some  $z$  and  $x$ , then  $\Gamma_{\max} = \max_{z \in \mathcal{Z}, z \in \mathcal{X}} \Gamma(z, x)$ .) In addition,  $\mathbf{1}_{\{\cdot\}}$  denotes an indicator of the event inside  $\{\cdot\}$ .

Generally, let the finite sets  $\mathcal{Y}$ ,  $\mathcal{A}$  be, respectively, a source alphabet and an action space. For a general loss function  $l : \mathcal{Y} \times \mathcal{A} \rightarrow \mathbb{R}$ , a *Bayes response* for  $\zeta \in \mathbb{R}^\mathcal{Y}$  under the loss function  $l$  is given as

$$b_l(\zeta) = \arg \min_{a \in \mathcal{A}} \zeta^T \cdot \mathcal{L}_a, \quad (1)$$

where  $\mathcal{L}_a$  denotes the column of the matrix of the loss function  $l$  corresponding to the  $a$ -th action, and ties are resolved lexicographically. The corresponding *Bayes envelope* is denoted as

$$U_l(\zeta) = \min_{a \in \mathcal{A}} \zeta^T \cdot \mathcal{L}_a. \quad (2)$$

Note that when  $\zeta$  is a probability, namely, it has non-negative components summing to one,  $U_l(\zeta)$  is the minimum achievable expected loss (as measured under the loss function  $l$ ) in guessing the value of  $Y \in \mathcal{Y}$  which is distributed according to  $\zeta$ . The associated optimal guess is  $b_l(\zeta)$ .

An  $n$ -block denoiser is a collection of  $n$  mappings  $\hat{\mathbf{X}}^n = \{\hat{X}_t\}_{1 \leq t \leq n}$ , where  $\hat{X}_t : \mathcal{Z}^n \rightarrow \hat{\mathcal{X}}$ . We assume a given loss function  $\Lambda : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow [0, \infty)$ , where the maximum single-letter loss is denoted by  $\Lambda_{\max}$ , and  $\lambda_{\hat{x}}$  denotes the  $\hat{x}$ -th column of the loss matrix. The normalized cumulative loss of the denoiser  $\hat{\mathbf{X}}^n$  on the individual sequence pair  $(x^n, z^n)$  is represented as

$$L_{\hat{\mathbf{X}}^n}(x^n, z^n) = \frac{1}{n} \sum_{t=1}^n \Lambda(x_t, \hat{X}_t(z^n)).$$

In words,  $L_{\hat{\mathbf{X}}^n}(x^n, z^n)$  is the normalized (per-symbol) loss, as measured under the loss function  $\Lambda$ , when using the denoiser  $\hat{\mathbf{X}}^n$  and when the observed noisy sequence is  $z^n$  while the underlying clean one is  $x^n$ . The notation  $L_{\hat{\mathbf{X}}^n}$  is extended for  $1 \leq i \leq j \leq n$ ,

$$L_{\hat{\mathbf{X}}^n}(x_i^j, z^n) = \frac{1}{j-i+1} \sum_{t=i}^j \Lambda(x_t, \hat{X}_t(z^n))$$

denoting the normalized (per-symbol) loss between (and including) locations  $i$  and  $j$ .

Now, consider the set  $\mathcal{S} = \{s : \mathcal{Z} \rightarrow \hat{\mathcal{X}}\}$ , which is the (finite) set of mappings that take  $\mathcal{Z}$  into  $\hat{\mathcal{X}}$ . We refer to elements of  $\mathcal{S}$  as “single-symbol denoisers”, since each  $s \in \mathcal{S}$  can be thought of as a rule for estimating  $X \in \mathcal{X}$  on the basis of  $Z \in \mathcal{Z}$ . Now, for any  $s \in \mathcal{S}$ , an unbiased estimator for  $\Lambda(x, s(Z))$  (based on  $Z$  only), where  $x$  is a deterministic symbol and  $Z$  is the output of the DMC when the input is  $x$ , can be obtained as in [2]. First, pick a function  $h : \mathcal{Z} \rightarrow \mathbb{R}^{\hat{\mathcal{X}}}$  with the property that, for  $a, b \in \mathcal{X}$ ,

$$\begin{aligned} E_a h_b(Z) &= \sum_{z \in \mathcal{Z}} h_b(z) \Pi(a, z) \\ &= \delta(a, b) \triangleq \begin{cases} 1, & \text{if } a = b \\ 0, & \text{otherwise} \end{cases}, \end{aligned} \quad (3)$$

where  $E_a$  denotes expectation over the channel output  $Z$  when the underlying channel input is  $a$ , and  $h_b(z)$  denotes the  $b$ -th component of  $h(z)$ . Let  $H$  denote the  $|\mathcal{Z}| \times |\mathcal{X}|$  matrix whose  $z$ -th row is  $h^T(z)$ , i.e.,  $H(z, b) = h_b(z)$ . To see that our assumption of a channel matrix with full row rank guarantees the existence of such an  $h$ , note that (3) can equivalently be stated in matrix form as

$$\Pi H = I, \quad (4)$$

where  $I$  is the  $|\mathcal{X}| \times |\mathcal{X}|$  identity matrix. Thus, e.g., any  $H$  of the form  $H = \Gamma^T (\Pi \Gamma^T)^{-1}$ , for any  $\Gamma$  such that  $\Pi \Gamma^T$  is invertible, satisfies (4). In particular,  $\Gamma = \Pi$  is a valid choice ( $\Pi \Pi^T$  is invertible, since  $\Pi$  is of full row rank) corresponding to the Moore-Penrose generalized inverse [27]. Now, for any  $s \in \mathcal{S}$ ,  $\rho(s) \in \mathbb{R}^{\mathcal{X}}$  denotes the column vector with  $x$ -th component

$$\begin{aligned} \rho_x(s) &= \sum_z \Lambda(x, s(z)) \Pi(x, z) \\ &= E_x \Lambda(x, s(Z)). \end{aligned} \quad (5)$$

In words,  $\rho_x(s)$  is the expected loss using the single-symbol denoiser  $s$ , while the underlying symbol is  $x$ . Considering  $\mathcal{S}$  as an action space alphabet, we define a loss function  $\ell : \mathcal{Z} \times \mathcal{S} \rightarrow \mathbb{R}$  as

$$\ell(z, s) = h(z)^T \cdot \rho(s). \quad (6)$$

We observe from (3) and (5) that  $\ell(Z, s)$  is an unbiased estimate of  $\Lambda(x, s(Z))$  since

$$E_x \ell(Z, s) = E_x h(Z)^T \cdot \rho(s) = \sum_{x'} E_x h_{x'}(Z) \rho_{x'}(s) = \sum_{x'} \delta(x, x') \rho_{x'}(s) = \rho_x(s) = E_x \Lambda(x, s(Z)) \quad \forall x \in \mathcal{X}. \quad (7)$$

For  $\xi \in \mathbb{R}^{\mathcal{Z}}$ , let  $B_H(\xi, \cdot) \in \mathcal{S}$  be defined by

$$B_H(\xi, z) = \arg \min_{\hat{x}} \xi^T \cdot H \cdot [\lambda_{\hat{x}} \odot \pi_z], \quad (8)$$

where, for vectors  $v_1$  and  $v_2$  of equal dimensions,  $v_1 \odot v_2$  denotes the vector obtained by component-wise multiplication. Note that, similarly as in [2, (88),(89) ],

$$\begin{aligned} B_H(\xi, \cdot) &= \arg \min_{s \in \mathcal{S}} \sum_z \xi^T \cdot H \cdot [\lambda_{s(z)} \odot \pi_z] \\ &= \arg \min_{s \in \mathcal{S}} \xi^T \cdot H \cdot \rho(s) \\ &= \arg \min_{s \in \mathcal{S}} \sum_z \xi_z \cdot [h^T(z) \cdot \rho(s)] \\ &= \arg \min_{s \in \mathcal{S}} \sum_z \xi_z \cdot \ell(z, s) = b_\ell(\xi). \end{aligned} \quad (9)$$

Thus,  $B_H(\xi, \cdot)$  is a Bayes response for  $\xi$  under the loss function  $\ell$  defined in (6).

## 2-B Preliminaries

In this section, we summarize the results from [1] and motivate the approach underlying the construction of our new class of denoisers. Analogously as in [2], the  $n$ -block denoiser  $\hat{\mathbf{X}}^n = \{\hat{X}_t\}_{1 \leq t \leq n}$  can be associated with  $\mathbf{F}^n = \{F_t\}_{1 \leq t \leq n}$ , where  $F_t : \mathcal{Z}^{n \setminus t} \rightarrow \mathcal{S}$  is defined as follows:  $F_t(z^{n \setminus t}, \cdot)$  is the single-symbol denoiser in  $\mathcal{S}$  satisfying

$$\hat{X}_t(z^n) = F_t(z^{n \setminus t}, z_t) \quad \forall z_t. \quad (10)$$

Therefore, we can adopt the view that at each time  $t$ , an  $n$ -block denoiser is choosing a single-symbol denoiser based on all the noisy sequence components but  $z_t$ , and applying that single-symbol denoiser on  $z_t$  to yield the  $t$ -th reconstruction  $\hat{x}_t$ . Conversely, any sequence of mappings into single-symbol denoisers  $\mathbf{F}^n$  defines a denoiser  $\hat{\mathbf{X}}^n$ , again via (10). We will adhere to this viewpoint in what follows.

One special class of widely used  $n$ -block denoisers is that of  $k$ -th order ‘‘sliding window’’ denoisers, which we denote by  $\hat{\mathbf{X}}^{n, \mathcal{S}_k}$ . Such denoisers are of the form

$$\hat{X}_t^{s_k}(z^n) = s_k(z_{t-k}^{t+k}), \quad t = k+1, \dots, n-k, \quad (11)$$

where  $s_k$  is an element of  $\mathcal{S}_k = \{s_k : \mathcal{Z}^{2k+1} \rightarrow \hat{\mathcal{X}}\}$ , the (finite) set of mappings from  $\mathcal{Z}^{2k+1}$  into  $\hat{\mathcal{X}}$ .<sup>1</sup> We also refer to  $s_k \in \mathcal{S}_k$  as a ‘‘ $k$ -th order denoiser’’. Note that  $\mathcal{S}_0 = \mathcal{S}$ . From the definition (11), it follows that

$$\hat{X}_i^{s_k}(z^n) = \hat{X}_j^{s_k}(z^n) \quad \text{whenever } z_{i-k}^{i+k} = z_{j-k}^{j+k}. \quad (12)$$

Following the association in (10), we can adopt an alternative view that the  $k$ -th order sliding window denoiser chooses a single-symbol denoiser  $s_k(z_{t-k}^{t-1}, z_{t+1}^{t+k}, \cdot) \in \mathcal{S}$  at time  $t$  on the basis of the context, and  $\hat{X}_t^{s_k}(z^n) = s_k(z_{t-k}^{t-1}, z_{t+1}^{t+k}, z_t)$ .

We denote  $\mathbf{c}_t \triangleq (z_{t-k}^{t-1}, z_{t+1}^{t+k})$  as a (two-sided) context for  $z_t$ , and define the set of all possible  $k$ -th order contexts,  $\mathbf{C}_k \triangleq \{(u_{-k}^{-1}, u_1^k) : (u_{-k}^{-1}, u_1^k) \in \mathcal{Z}^{2k}\}$ . Then, for given  $z^n$  and for each  $\mathbf{c} \in \mathbf{C}_k$ , we define

$$\mathcal{T}(\mathbf{c}) \triangleq \{t : \mathbf{c}_t = \mathbf{c}, \quad k+1 \leq t \leq n-k\} = \{t : (z_{t-k}^{t-1}, z_{t+1}^{t+k}) = \mathbf{c}, \quad k+1 \leq t \leq n-k\}, \quad (13)$$

---

<sup>1</sup>The value of  $\hat{X}_t^{s_k}(z^n)$  for  $t \leq k$  and  $t > n-k$  is defined, for concreteness and simplicity, as an arbitrary fixed symbol in  $\hat{\mathcal{X}}$ .

the set of indices where the context equals  $\mathbf{c}$ . Now, an equivalent interpretation for (12) is that for each  $\mathbf{c} \in \mathbf{C}_k$ , the  $k$ -th order sliding window denoiser employs a time-invariant single-symbol denoiser,  $s_k(\mathbf{c}, \cdot)$ , at all points  $t \in \mathcal{T}(\mathbf{c})$ . In other words, the sequence  $z^n$  is partitioned into the subsequences associated with the various contexts, and on each such subsequence a time-invariant single-symbol scheme is employed.

In [1], for integers  $k \geq 0$  and  $n > 2k$ , the  $k$ -th order minimum loss of  $(x^n, z^n)$  is defined by

$$\begin{aligned} D_k(x^n, z^n) &\triangleq \min_{\hat{\mathbf{X}}^n \in \hat{\mathbf{X}}^n, \mathcal{S}_k} L_{\hat{\mathbf{X}}^n}(x_{k+1}^{n-k}, z^n) \\ &= \min_{s_k \in \mathcal{S}_k} \frac{1}{n-2k} \sum_{t=k+1}^{n-k} \Lambda(x_t, s_k(\mathbf{c}_t, z_t)). \end{aligned} \quad (14)$$

The identity of the element  $s_k \in \mathcal{S}_k$  that achieves (14) depends not only on  $z^n$ , but also on  $x^n$ , since (14) can be expressed as

$$\frac{1}{n-2k} \sum_{\mathbf{c} \in \mathbf{C}_k} \left[ \min_{s \in \mathcal{S}} \sum_{\tau \in \mathcal{T}(\mathbf{c})} \Lambda(x_\tau, s(z_\tau)) \right],$$

and at each time  $t$ , the best  $k$ -th order sliding window denoiser that achieves (14) will employ the single-symbol denoiser

$$\arg \min_{s \in \mathcal{S}} \sum_{\tau \in \mathcal{T}(\mathbf{c}_t)} \Lambda(x_\tau, s(z_\tau)), \quad (15)$$

which is determined from the joint empirical distribution of pairs  $\{(x_\tau, z_\tau) : \tau \in \mathcal{T}(\mathbf{c}_t)\}$ .

It was shown in [1] that, despite the lack of knowledge of  $x^n$ ,  $D_k(x^n, Z^n)$  is achievable in a sense made precise below, in the limit of growing  $n$ , by a scheme that only has access to  $Z^n$ . This scheme is dubbed in [1] as the Discrete Universal DENOISER (DUDE),  $\hat{\mathbf{X}}_{\text{univ}}^{n,k}$ . The algorithm is defined by

$$\hat{\mathbf{X}}_{\text{univ},t}^k(z^n) = B_H(\mathbf{m}(z^n, z_{t-k}^{t-1}, z_{t+1}^{t+k}), z_t), \quad (16)$$

where  $\mathbf{m}(z^n, \mathbf{c})$  is the vector of counts of the appearances of the various symbols within the context  $\mathbf{c}$  along the sequence  $z^n$ . That is, for all  $\beta \in \mathcal{Z}$ ,  $\mathbf{m}(z^n, \tilde{z}_{-k}^{-1}, \tilde{z}_1^k)$  is the  $|\mathcal{Z}|$ -dimensional column vector whose  $\beta$ -th component is

$$\mathbf{m}(z^n, \tilde{z}_{-k}^{-1}, \tilde{z}_1^k)[\beta] = |\{t : k+1 \leq t \leq n-k, z_{t-k}^{t+k} = \tilde{z}_{-k}^{-1} \beta \tilde{z}_1^k\}|,$$

namely, the number of appearances of  $\tilde{z}_{-k}^{-1} \beta \tilde{z}_1^k$  along the sequence  $z^n$ .

The main result of [1] is the following theorem, pertaining to the semi-stochastic setting of an individual sequence  $\mathbf{x} = (x_1, x_2, \dots)$  corrupted by a DMC that yields the stochastic noisy sequence  $\mathbf{Z} = (Z_1, Z_2, \dots)$ .

**Theorem 1** ([1, Theorem 1]) *Take  $k = k_n$  satisfying  $k_n |\mathcal{Z}|^{2k_n} = o(n/\log n)$ . Then, for all  $\mathbf{x} \in \mathcal{X}^\infty$ , the sequence of denoisers  $\{\hat{\mathbf{X}}_{\text{univ}}^{n,k_n}\}$  defined in (16) satisfies:*

a)

$$\lim_{n \rightarrow \infty} \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k_n}}(x^n, Z^n) - D_{k_n}(x^n, Z^n) \right] = 0 \quad a.s.$$

b)

$$E \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k_n}}(x^n, Z^n) - D_{k_n}(x^n, Z^n) \right] = O \left( \sqrt{\frac{k_n |\mathcal{Z}|^{2k_n}}{n}} \right).$$

Theorem 1 was further shown in [1] to imply the universality of the DUDE in the fully stochastic setting where the underlying sequence is emitted by a stationary source (and the goal is to attain the performance of the optimal distribution-dependent denoiser).

From (16), it is apparent that the DUDE ends up employing a  $k$ -th order sliding window denoiser (where the sliding window scheme the DUDE chooses depends on  $z^n$ ). Moreover, (9) implies that, at each time  $t$ , DUDE is merely employing the single-symbol denoiser  $B_H(\mathbf{m}(z^n, z_{t-k}^{t-1}, z_{t+1}^{t+k}), \cdot) \in \mathcal{S}$ , which can be obtained by finding the Bayes response  $b_\ell(\mathbf{m}(z^n, z_{t-k}^{t-1}, z_{t+1}^{t+k}))$  or, equivalently, the mapping in  $\mathcal{S}$  given by

$$\arg \min_{s \in \mathcal{S}} \sum_{\tau \in \mathcal{T}(\mathbf{c}_t)} \ell(z_\tau, s), \quad (17)$$

where  $\ell(z, s)$  is the loss function defined in (6). By comparing (15) with (17), and from Theorem 1, we observe that working with the estimated loss  $\ell(z_\tau, s)$  in lieu of the genie-aided  $\Lambda(x_\tau, s(z_\tau))$  allows us to essentially achieve the genie-aided performance in (14).

## 2-C Motivation

Our motivation for this paper is based on the observation that the  $k$ -th order sliding window denoisers ignore the time-varying nature of the underlying sequence  $x^n$ . That is, as discussed above, for time instances with the same contexts, the single-symbol denoiser employed along the associated subsequence is time-invariant. In other words, for each  $t$ , only the empirical distribution of the sequence  $\{(x_\tau, z_\tau) : \tau \in \mathcal{T}(\mathbf{c}_t)\}$  matters, but its order of composition, i.e., its time-varying nature, is not considered. It is clear, however, that when the characteristics of the underlying clean sequence  $x^n$  are changing, the (normalized) cumulative loss that is achieved by sliding window denoisers that can shift from one rule to another along the sequence may be strictly lower (better) than (14). We now devise and analyze our new scheme that achieves this more ambitious target performance.

## 3 The Shifting Denoiser (S-DUDE)

In this section, we derive our new class of denoisers and analyze their performance. In Subsection 3-A, we begin with the simplest case, competing with shifting symbol-by-symbol denoisers, or, in other words, shifting 0-th order denoisers. The argument is generalized to shifting  $k$ -th order denoisers in Subsection 3-B, and the framework and results include Subsection 3-A as a special case. We will use the notation  $\mathcal{S}_0$ , instead of  $\mathcal{S}$ , for consistency in denoting the class of single-symbol denoisers. Throughout this section, we assume the semi-stochastic setting.

### 3-A Switching between symbol-by-symbol (0-th order) denoisers

Consider an  $n$ -tuple of single-symbol denoisers  $\mathbf{S} = \{s_1, \dots, s_n\} \in \mathcal{S}_0^n$ . Then, as mentioned in Section 2-B, for such  $\mathbf{S}$ , we can define the associated  $n$ -block denoiser  $\hat{\mathbf{X}}^{n, \mathbf{S}}$  as

$$\hat{X}_t^{\mathbf{S}}(z^n) = s_t(z_t). \quad (18)$$

Note that in this case, the single-symbol denoiser applied at each time may depend on the time  $t$  (but not on  $z^{n \setminus t}$ , as would be the case for a general denoiser). We also denote the estimated normalized cumulative loss as

$$\tilde{L}_{\mathbf{S}}(z^n) \triangleq \frac{1}{n} \sum_{t=1}^n \ell(z_t, s_t), \quad (19)$$

whose property is given in the following lemma, which parallels [2, Theorem 4].

**Lemma 1** Fix  $\epsilon > 0$ . For fixed  $\mathbf{S} \in \mathcal{S}_0^n$ , and all  $x^n \in \mathcal{X}^n$ ,

$$P\left(L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n) > \epsilon\right) \leq \exp\left(-n \frac{2\epsilon^2}{L_{\max}^2}\right) \quad \text{and} \quad (20)$$

$$P\left(\tilde{L}_{\mathbf{S}}(Z^n) - L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, Z^n) > \epsilon\right) \leq \exp\left(-n \frac{2\epsilon^2}{L_{\max}^2}\right), \quad (21)$$

where  $L_{\max} = \Lambda_{\max} + \ell_{\max}$ .

In words, the lemma shows that for every  $\mathbf{S} \in \mathcal{S}_0^n$ , the estimated loss  $\tilde{L}_{\mathbf{S}}(Z^n)$  is concentrated around the true loss  $L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, Z^n)$  with high probability, as  $n$  becomes large, regardless of the underlying sequence  $x^n$ .

*Proof of Lemma 1:* See Appendix 8-A. ■

Now, let the integer  $0 \leq m \leq \lfloor \frac{n}{2} \rfloor$  denote the maximum number of shifts allowed along the sequence. Then, define a set  $\mathcal{S}_{0,m}^n \subseteq \mathcal{S}_0^n$  as

$$\mathcal{S}_{0,m}^n = \left\{ \mathbf{S} \in \mathcal{S}_0^n : \sum_{t=2}^n \mathbf{1}_{\{s_{t-1} \neq s_t\}} \leq m \right\}, \quad (22)$$

namely,  $\mathcal{S}_{0,m}^n$  is the set of  $n$ -tuples of single-symbol denoisers with at most  $m$  shifts from one mapping to another.<sup>2</sup> Analogously to (14), for the class of  $n$ -block denoisers  $\hat{\mathbf{X}}^{n,\mathbf{S}}$  with  $\mathbf{S} \in \mathcal{S}_{0,m}^n$ , we define

$$\begin{aligned} D_{0,m}(x^n, z^n) &\triangleq \min_{\mathbf{S} \in \mathcal{S}_{0,m}^n} L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, z^n) \\ &= \min_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \frac{1}{n} \sum_{t=1}^n \Lambda(x_t, s_t(z_t)), \end{aligned} \quad (23)$$

which is the minimum normalized cumulative loss that can be achieved for  $(x^n, z^n)$  by the sequence of  $n$  single-symbol denoisers that allow at most  $m$  shifts. Our goal in this section is to build a universal scheme that only has access to  $Z^n$ , but still essentially achieves  $D_{0,m}(x^n, Z^n)$ .

As hinted by the DUDE, we build our universal scheme by working with the estimated loss. That is, define

$$\hat{\mathbf{S}} = \hat{\mathbf{S}}(z^n) \triangleq \arg \min_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \tilde{L}_{\mathbf{S}}(z^n), \quad (24)$$

and our  $(0, m)$ -Shifting Discrete Universal DEnoiser (S-DUDE),  $\hat{\mathbf{X}}_{\text{univ}}^{n,0,m}$ , is defined as  $\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}}$ . It is clear that, by definition,  $L_{\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}}}(x^n, z^n) \geq D_{0,m}(x^n, z^n)$  for all  $x^n$  and  $z^n$ , but we can also show that, with high probability,  $L_{\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}}}(x^n, Z^n)$  does not exceed  $D_{0,m}(x^n, Z^n)$  by much, as stated in the following theorem.

**Theorem 2** Let  $\hat{\mathbf{X}}_{\text{univ}}^{n,0,m}$  be defined as  $\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}}$ , where  $\hat{\mathbf{S}}$  is given in (24). Then, for all  $\epsilon > 0$  and  $x^n \in \mathcal{X}^n$ ,

$$P\left(L_{\hat{\mathbf{X}}_{\text{univ}}^{n,0,m}}(x^n, Z^n) - D_{0,m}(x^n, Z^n) > \epsilon\right) \leq 2 \exp\left(-n \left[ \frac{\epsilon^2}{2L_{\max}^2} - 2 \left\{ h\left(\frac{m}{n}\right) + \frac{(m+1) \ln N}{n} \right\} \right]\right),$$

where  $h(x) = -x \ln x - (1-x) \ln(1-x)$  for  $0 \leq x \leq 1$ , and  $N = |\mathcal{S}| = |\mathcal{Z}|^{|\mathcal{X}|}$ . In particular, the right-hand side of the inequality is exponentially small, provided  $m = o(n)$ .

<sup>2</sup>Note that, when  $m = 0$ ,  $\mathcal{S}_{0,0}^n$  is the set of constant  $n$ -tuples consisting of the same single-symbol denoiser.

*Remark:* It is reasonable to expect this theorem to hold, given Lemma 1. That is, since, for fixed  $\mathbf{S} \in \mathcal{S}_{0,m}^n$ ,  $\tilde{L}_{\mathbf{S}}(Z^n)$  is concentrated on  $L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, Z^n)$ , it is plausible that  $\hat{\mathbf{S}}$  that achieves  $\min_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \tilde{L}_{\mathbf{S}}(Z^n)$  will have a loss  $L_{\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}}}(x^n, Z^n)$  close to  $\min_{\mathbf{S} \in \mathcal{S}_{0,m}^n} L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x^n, Z^n)$ , i.e.,  $D_{0,m}(x^n, Z^n)$ .

*Proof of Theorem 2:* See Appendix 8-B. ■

### 3-B Switching between $k$ -th order denoisers

Now, we extend the result from Subsection 3-A to the case of shifting between  $k$ -th order denoisers. The argument parallels that of Subsection 3-A. Let  $\{s_{k,t}\}_{t=k+1}^{n-k}$  be an arbitrary sequence of the  $k$ -th order denoiser mappings, i.e.,  $s_{k,t} \in \mathcal{S}_k$  for  $k+1 \leq t \leq n-k$ . Now, for given  $z^n$ , define an  $(n-2k)$ -tuple of ( $k$ -th order denoiser induced) single-symbol denoisers

$$\mathbf{S}_k(z^n) \triangleq \{s_{k,t}(\mathbf{c}_t, \cdot)\}_{t=k+1}^{n-k} \in \mathcal{S}_0^{n-2k}, \quad (25)$$

where, to recall,  $\mathbf{c}_t = (z_{t-k}^{t-1}, z_{t+1}^{t+k})$ , and  $s_{k,t}(\mathbf{c}_t, \cdot)$  is the single-symbol denoiser induced from  $s_{k,t} \in \mathcal{S}_k$  and  $\mathbf{c}_t$ . For brevity of notation, we will suppress the dependence on  $z^n$  in  $\mathbf{S}_k(z^n)$  and denote it as  $\mathbf{S}_k$ . Then, as in (18), we define the associated  $n$ -block denoiser  $\hat{\mathbf{X}}^{n,\mathbf{S}_k}$  as <sup>3</sup>

$$\hat{\mathbf{X}}_t^{\mathbf{S}_k}(z^n) = s_{k,t}(\mathbf{c}_t, z_t). \quad (26)$$

In addition, extending (19), the estimated normalized cumulative loss is given as

$$\tilde{L}_{\mathbf{S}_k}(z^n) = \frac{1}{n-2k} \sum_{t=k+1}^{n-k} \ell(z_t, s_{k,t}(\mathbf{c}_t, \cdot)). \quad (27)$$

Then, we have the following lemma, which parallels Lemma 1.

**Lemma 2** Fix  $\epsilon > 0$ . For any fixed sequence  $\{s_{k,t}\}_{t=k+1}^{n-k}$ , and all  $x^n \in \mathcal{X}^n$ ,

$$\Pr\left(L_{\hat{\mathbf{X}}^{n,\mathbf{S}_k}}(x_{k+1}^{n-k}, Z^n) - \tilde{L}_{\mathbf{S}_k}(Z^n) > \epsilon\right) \leq (k+1) \exp\left(-\frac{2(n-2k)\epsilon^2}{(k+1)L_{\max}^2}\right) \quad \text{and} \quad (28)$$

$$\Pr\left(\tilde{L}_{\mathbf{S}_k}(Z^n) - L_{\hat{\mathbf{X}}^{n,\mathbf{S}_k}}(x_{k+1}^{n-k}, Z^n) > \epsilon\right) \leq (k+1) \exp\left(-\frac{2(n-2k)\epsilon^2}{(k+1)L_{\max}^2}\right), \quad (29)$$

where  $L_{\max} = \Lambda_{\max} + \ell_{\max}$ .

*Remark:* Note that when  $k = 0$ , this lemma coincides with Lemma 1. The proof of this lemma combines Lemma 1 and the de-interleaving argument in the proof of [1, Theorem 2]. Namely, we de-interleave  $Z^n$  into  $(k+1)$  subsequences consisting of symbols separated by blocks of  $k$  symbols, and exploit the conditional independence of symbols in each subsequence, given all symbols not in that subsequence, to use Lemma 1.

*Proof of Lemma 2:* See Appendix 8-C. ■

Now, for an integer  $0 \leq m \leq \lfloor \frac{n-2k}{2} \rfloor$  and given  $z^n$ , let  $n(\mathbf{c}) \triangleq |\mathcal{T}(\mathbf{c})|$ , and  $m(\mathbf{c}) \triangleq \min\{n(\mathbf{c}), m\}$  for  $\mathbf{c} \in \mathbf{C}_k$ . Then, analogously as in (22), we define

$$\mathcal{S}_{k,m}^n(z^n) = \left\{ \mathbf{S}_k(z^n) \in \mathcal{S}_0^{n-2k} : \{s_{k,\tau}(\mathbf{c}, \cdot)\}_{\tau \in \mathcal{T}(\mathbf{c})} \in \mathcal{S}_{0,m(\mathbf{c})}^{n(\mathbf{c})} \quad \text{for all } \mathbf{c} \in \mathbf{C}_k \right\}. \quad (30)$$

<sup>3</sup>Again, the value of  $\hat{\mathbf{X}}_t^{\mathbf{S}_k}(z^n)$  for  $t \leq k$  and  $t > n-k$  can be defined as an arbitrary fixed symbol, since it will be inconsequential in subsequent development.

In words,  $\mathcal{S}_{k,m}^n(z^n)$  is the set of  $(n - 2k)$ -tuples of ( $k$ -th order denoiser induced) single-symbol denoisers that allow at most  $m(\mathbf{c})$  shifts within the subsequence  $\{t : t \in \mathcal{T}(\mathbf{c})\}$  for each context  $\mathbf{c} \in \mathbf{C}_k$ .<sup>4</sup> Again, for brevity, the dependence on  $z^n$  in  $\mathcal{S}_{k,m}^n(z^n)$  is suppressed, and we write simply  $\mathcal{S}_{k,m}^n$ . It is worth noting that  $\mathcal{S}_{k,m}^n$  is a larger class than the class of  $k$ -th order ‘sliding window’ denoisers that are allowed to shift at most  $m$  times. The reason is that in  $\mathcal{S}_{k,m}^n$ , the shift within each subsequence associated with each context can occur at any time, regardless of the shifts in other subsequences, whereas in the latter class, the shifts in each subsequence occur together with other shifts in other subsequences.

For integers  $k \geq 0$  and  $n > 2k$ , we now define, for the class of  $n$ -block denoisers  $\hat{\mathbf{X}}^{n,\mathbf{S}}$  with  $\mathbf{S} \in \mathcal{S}_{k,m}^n$ ,

$$\begin{aligned} D_{k,m}(x^n, z^n) &\triangleq \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x_{k+1}^{n-k}, z^n) \\ &= \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \frac{1}{n - 2k} \sum_{t=k+1}^{n-k} \Lambda(x_t, s_{k,t}(\mathbf{c}_t, z_t)), \end{aligned} \quad (31)$$

the minimum normalized cumulative loss of  $(x^n, z^n)$  that can be achieved by the sequence of  $k$ -th order denoisers that allow at most  $m$  shifts within each context. Now, to build a legitimate (non genie-aided) universal scheme achieving (31) on the basis of  $Z^n$  only, we define

$$\hat{\mathbf{S}}_{k,m} = \arg \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \tilde{L}_{\mathbf{S}}(z^n), \quad (32)$$

and the  $(k, m)$ -S-DUDE,  $\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}$ , is defined as  $\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}_{k,m}}$ . Note that when  $m = 0$ ,  $\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}_{k,m}}$  coincides with the DUDE in [1]. The following theorem generalizes Theorem 2 to the case of general  $k \geq 0$ .

**Theorem 3** *Let  $\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}$  be given by  $\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}_{k,m}}$ , where  $\hat{\mathbf{S}}_{k,m}$  is defined in (32). Then, for all  $\epsilon > 0$  and  $x^n \in \mathcal{X}^n$ ,*

$$\begin{aligned} &\Pr\left(L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x_{k+1}^{n-k}, Z^n) - D_{k,m}(x^n, Z^n) > \epsilon\right) \\ &\leq 2(k+1) \exp\left(- (n - 2k) \cdot \left[ \frac{\epsilon^2}{2(k+1)L_{\max}^2} - 2|\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\ln N}{n-2k} \right\} \right]\right), \end{aligned} \quad (33)$$

$$\leq 2(k+1) \exp\left(- (n - 2k) \cdot \left[ \frac{\epsilon^2}{2(k+1)L_{\max}^2} - 2|\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\ln N}{n-2k} \right\} \right]\right), \quad (34)$$

where  $h(x) = -x \ln x - (1-x) \ln(1-x)$  for  $0 \leq x \leq 1$ , and  $N = |\mathcal{S}| = |\mathcal{Z}|^{|\hat{\mathcal{X}}|}$ .

*Remark:* Note that when  $k = 0$ , this theorem coincides with Theorem 2. Similarly to the way Theorem 2 was plausible given Lemma 1, Theorem 3 can be expected given Lemma 2, since  $\hat{\mathbf{S}}_{k,m}$  achieves  $\min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \tilde{L}_{\mathbf{S}}(Z^n)$ , and we expect  $L_{\hat{\mathbf{X}}^{n,\hat{\mathbf{S}}_{k,m}}}(x_{k+1}^{n-k}, Z^n)$  to be close to  $D_{k,m}(x^n, Z^n)$  from the concentration of  $\tilde{L}_{\mathbf{S}}(Z^n)$  to  $L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(x_{k+1}^{n-k}, Z^n)$  for all  $\mathbf{S} \in \mathcal{S}_{k,m}^n$ .

*Proof of Theorem 3:* See Appendix 8-D. ■

From Theorem 3, we now easily obtain one of the main results of the paper, which extends Theorem 1 from the case  $m = 0$  to the case of general  $0 \leq m \leq \lfloor \frac{n-2k}{2} \rfloor$ . That is, the following theorem asserts that, for every underlying sequence  $\mathbf{x} \in \mathcal{X}^\infty$ , our  $(k, m)$ -S-DUDE performs essentially as well as the best shifting  $k$ -th order denoiser that allows at most  $m$  shifts within each context, both in high probability and expectation sense, provided a growth condition on  $k$  and  $m$  is satisfied.

**Theorem 4** *Suppose  $k = k_n$  and  $m = m_n$  are such that the right-hand side of (34) is summable in  $n$ . Then, for all  $\mathbf{x} \in \mathcal{X}^\infty$ , the sequence of denoisers  $\{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}\}$  satisfies*

<sup>4</sup>When  $m = 0$ ,  $\mathcal{S}_{k,0}^n(z^n)$  becomes the set of  $n$ -block  $k$ -th order ‘sliding window’ denoisers.

a)

$$\lim_{n \rightarrow \infty} \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x^n, Z^n) - D_{k,m}(x^n, Z^n) \right] = 0 \quad \text{a.s.} \quad (35)$$

b) For any  $\delta > 0$ ,

$$E \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x^n, Z^n) - D_{k,m}(x^n, Z^n) \right] = O \left( \sqrt{k_n |\mathcal{Z}|^{2k_n} \left( \frac{m_n}{n} \right)^{1-\delta}} \right). \quad (36)$$

*Remark:* It will be seen in Claim 1 below that the stipulation in the theorem implies  $\lim_{n \rightarrow \infty} k_n |\mathcal{Z}|^{2k_n} \left( \frac{m_n}{n} \right)^{1-\delta} = 0$ , which, when combined with (36), implies that the expected difference on the left hand side of (36) vanishes with increasing  $n$ . That in itself, however, can easily be deduced from (35) and bounded convergence. The more significant value of (36) is in providing a rate of convergence result for the ‘redundancy’ in the S-DUDE’s performance, as a function of both  $k$  and  $m$ . In particular, note that for any  $\eta > 0$ ,  $O(n^{-1/2+\eta})$  is achievable provided  $k_n = c \log n$  and  $m_n = n^\xi$ , for small enough positive constants  $c, \xi$ .

In what follows, we specify the maximal growth rates for  $k = k_n$  and  $m = m_n$  under which the summability condition stipulated in Theorem 4 holds.

**Claim 1** a) *Maximal growth rate for  $k$ : The summability condition in Theorem 4 is satisfied provided  $k_n = c_1 \log n$  with  $c_1 < \frac{1}{2 \log |\mathcal{Z}|}$  and  $m_n$  grows at any sub-polynomial rate. On the other hand, the condition is not satisfied for  $k_n = c_1 \log n$  with any  $c_1 \geq \frac{1}{2 \log |\mathcal{Z}|}$ , even when  $m$  is fixed (not growing with  $n$ ).*

b) *Maximal growth rate for  $m$ : The summability condition in Theorem 4 is satisfied for any sub-linear growth rate of  $m_n$ , provided  $k_n$  is taken to increase sufficiently slowly that  $k_n |\mathcal{Z}|^{2k_n} = o((n/m_n)^{1-\delta})$  for some  $\delta > 0$ . On the other hand, the condition is not satisfied whenever  $m_n$  grows linearly with  $n$ , even when  $k$  is fixed.*

*Proof of Claim 1:* See Appendix 8-E. ■

*Proof of Theorem 4:* See Appendix 8-F. ■

### 3-C A “strong converse”

In Claim 1, we have shown the necessity of  $m = o(n)$  for the condition required in Theorem 4 to hold. However, we can prove the necessity of  $m = o(n)$  in a much stronger sense, described in the following theorem.

**Theorem 5** *Suppose that  $\mathcal{X} = \hat{\mathcal{X}}$ , that  $\Lambda(x, \hat{x}) \geq 0$  for all  $x, \hat{x}$  with equality if and only if  $x = \hat{x}$ , and that  $\Pi(x, z) > 0$  for all  $x, z$ . If  $m = \Theta(n)$ , then for any sequence of denoisers  $\{\hat{\mathbf{X}}^n\}$ , there exists  $\mathbf{x}^\infty \in \mathcal{X}^\infty$  such that*

$$\limsup_{n \rightarrow \infty} E \left[ L_{\hat{\mathbf{X}}^n}(x^n, Z^n) - D_{0,m}(x^n, Z^n) \right] > 0. \quad (37)$$

*Remark:* The theorem establishes the fact that when  $m = o(n)$  does not hold, namely, when  $m = \Theta(n)$ , not only does the almost sure convergence in Theorem 4 not hold but, in fact, even the much weaker convergence in expectation would fail. Further, it shows that this would be the case for *any* sequence of denoisers, not necessarily the S-DUDE. Furthermore, (37) features  $D_{0,m}(x^n, Z^n)$ , pertaining to competition

with a genie that shifts among single-symbol denoisers so, *a fortiori*, it implies that for any fixed  $k > 0$  or  $k$  that grows with  $n$ ,

$$\limsup_{n \rightarrow \infty} E [L_{\hat{\mathbf{X}}^n}(x^n, Z^n) - D_{k,m}(x^n, Z^n)] > 0 \quad (38)$$

also holds since, by definition,  $D_{0,m}(x^n, z^n) \geq D_{k,m}(x^n, z^n)$  for all  $x^n, z^n$  and  $k \geq 0$ . Therefore, the theorem asserts that for *any* sequence of denoisers to compete with  $D_{k,m}(x^n, Z^n)$ , even in expectation sense,  $m = o(n)$  is necessary. Finally, we mention that the conditions stipulated in the statement of the theorem regarding the loss function and the channel can be considerably relaxed without compromising the validity of the theorem. These conditions are made to allow for the simple proof that we give in Appendix 8-G.

## 4 The Stochastic Setting

In [1], the semi-stochastic setting result, [1, Theorem 1], was shown to imply the result for the stochastic setting as well. That is, when the underlying data form a stationary process, [1, Section VI] shows that the DUDE attains optimum distribution-dependent performance. Analogously, we can now use the results from the semi-stochastic setting of the previous section to generalize the results of [1, Section VI] and show that our S-DUDE attains optimum distribution-dependent performance when the underlying data form a piecewise stationary process. We first define the precise notion of the class of piecewise stationary processes in Subsection 4-A, and discuss the richness of this class in Subsection 4-B. Subsection 4-C gives the main result of this section: the stochastic setting optimality of the S-DUDE.

### 4-A Definition of the class of processes $\mathcal{P}\{m_n\}$

Let  $P_{\mathbf{X}}^{(1)}, \dots, P_{\mathbf{X}}^{(M)}$  be a finite collection of  $M$  probability distributions of stationary processes, with components taking the values in  $\mathcal{X}$ . Let  $\mathbf{A}$  be a process with components taking the values in  $\{1, \dots, M\}$ . Then, a piecewise stationary process  $\mathbf{X}$  is generated by shifting between the  $M$  processes in a way specified by the “switching process”  $\mathbf{A}$ , as we now describe.

First, denote  $r(A^n)$  as the number of shifts that have occurred along the  $n$ -tuple  $A^n$ , i.e.,

$$r(A^n) \triangleq \sum_{j=1}^{n-1} \mathbf{1}_{\{A_j \neq A_{j+1}\}}.$$

Thus, there are  $r(A^n) + 1$  “blocks” in  $A^n$ , where each block is a tuple of constant values that are different from the values of adjacent blocks. Now, for each  $1 \leq i \leq r(A^n) + 1$ , we define

$$\tau_i(A^n) \triangleq \begin{cases} \inf\{t : \sum_{j=1}^t \mathbf{1}_{\{A_j \neq A_{j+1}\}} = i\} & \text{if } 1 \leq i \leq r(A^n) \\ n & \text{if } i = r(A^n) + 1 \end{cases}$$

as the last time instance of the  $i$ -th block in  $A^n$ . In addition, define  $\tau_0(A^n) \triangleq 0$ . Clearly,  $r(A^n)$  and  $\tau_i(A^n)$  depend on  $A^n$  and, thus, are random variables. However, for brevity, we suppress the dependence on  $A^n$  when there is no confusion, and write simply  $r$  and  $\tau_i$ , respectively.

Using these definitions, and by denoting  $P_{A^n}$  as the  $n$ -th order marginal distribution of  $\mathbf{A}$ , we define a piecewise stationary process  $\mathbf{X}$  by characterizing its  $n$ -th order marginal distribution  $P_{X^n}$  as

$$\begin{aligned} P_{X^n}(X^n = x^n) &= \sum_{a^n} P_{A^n}(a^n) P(X^n = x^n | A^n = a^n) \\ &= \sum_{a^n} P_{A^n}(a^n) \prod_{i=1}^{r+1} P_{\mathbf{X}}^{(a_{\tau_i})}(x_{\tau_{i-1}+1}^{\tau_i}), \end{aligned} \quad (39)$$

for each  $n$ . The corresponding distribution of the process  $\mathbf{X}$  is denoted as  $P_{\mathbf{X}}$ .<sup>5</sup> In words,  $\mathbf{X}$  is constructed by following one of the  $M$  probability distributions in each block, switching from one to another depending on  $\mathbf{A}$ . Furthermore, conditioned on the realization of  $\mathbf{A}$ , each stationary block is independent of other blocks, even if the distribution of distinct blocks is the same. This property of conditional independence is reasonable for modeling many types of data arising in practice, since we can think of the  $M$  distributions as different ‘modes’; if the process returns to the same mode, it is reasonable to model the new block as a new independent realization of that same distribution. In other words, the ‘mode’ may represent the kind of ‘texture’ in a certain region of the data, but two different regions with the same ‘texture’ should have independent realizations from the texture-generating source. Our notion of a piecewise stationary process almost coincides with that developed in [21]. The main difference is that we allow an arbitrary distribution for the process  $\mathbf{A}$ .

Now, we define  $\mathcal{P}\{m_n\}$  to be the class of all process distributions that can be constructed as in (39) for some  $M$ , some collection  $P_{\mathbf{X}}^{(1)}, \dots, P_{\mathbf{X}}^{(M)}$  of stationary processes, and some switching process  $\mathbf{A}$  whose number of shifts satisfies

$$r(A^n) \leq m_n \quad a.s. \quad \forall n. \quad (40)$$

In words, a process  $\mathbf{X}$  belongs to<sup>6</sup>  $\mathcal{P}\{m_n\}$  if and only if it can be formed by switching between a finite collection of independent processes in which the number of switches by time  $n$  does not exceed  $m_n$ .

#### 4-B Richness of $\mathcal{P}\{m_n\}$

In this subsection, we examine how rich the class  $\mathcal{P}\{m_n\}$  is, in terms of the growth rate  $m_n$  and the existence of denoising schemes that are universal with respect to  $\mathcal{P}\{m_n\}$ . First, given any distribution on a noiseless  $n$ -tuple,  $P_{X^n}$ , we define

$$\mathbb{D}(P_{X^n}, \mathbf{\Pi}) \triangleq \min_{\hat{\mathbf{X}}^n \in \mathcal{D}_n} EL_{\hat{\mathbf{X}}^n}(X^n, Z^n), \quad (41)$$

where  $\mathcal{D}_n$  is the class of *all*  $n$ -block denoisers. The expectation on the right-hand side of (41) assumes that  $X^n$  is generated from  $P_{X^n}$  and that  $Z^n$  is the output of the DMC,  $\mathbf{\Pi}$ , whose input is  $X^n$ . Thus,  $\mathbb{D}(P_{X^n}, \mathbf{\Pi})$  is the optimum denoising performance (in the sense of expected per-symbol loss) attainable when the source distribution  $P_{X^n}$  is known.

What happens when the source distribution is unknown? Theorem 3 of [1] established the fact that<sup>7</sup>

$$\lim_{n \rightarrow \infty} \left[ EL_{\hat{\mathbf{X}}_{\text{DUDE}}^n}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \right] = 0 \quad \text{for all stationary } P_{\mathbf{X}}. \quad (42)$$

Note that our newly-defined class of processes,  $\mathcal{P}\{m_n\}$ , is simply the class of all stationary processes if one takes the sequence  $m_n$  to be  $m_n \equiv 0$  for all  $n$ . Thus, assuming  $m_n \equiv 0$ , (42) is equivalent to

$$\lim_{n \rightarrow \infty} \left[ EL_{\hat{\mathbf{X}}_{\text{DUDE}}^n}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \right] = 0 \quad \text{for all } P_{\mathbf{X}} \in \mathcal{P}\{m_n\}. \quad (43)$$

At the other extreme, when  $m_n = n$ ,  $\mathcal{P}\{m_n\}$  consists of all possible (not necessarily stationary) processes. We can observe this equivalence by having  $M = |\mathcal{X}|$  processes each be a constant process at a different symbol in  $\mathcal{X}$ , and creating any process by switching to the appropriate symbol. In this case, not only does

<sup>5</sup> $\{P_{X^n}\}_{n \geq 1}$  is readily verified to be a consistent family of distributions and, thus, by Kolmogorov’s extension theorem, uniquely defines the distribution of the process  $\mathbf{X}$ .

<sup>6</sup>The phrase “the process  $\mathbf{X}$  belongs to  $\mathcal{P}\{m_n\}$ ” is shorthand for “the distribution of the process  $\mathbf{X}$ ,  $P_{\mathbf{X}}$ , belongs to  $\mathcal{P}\{m_n\}$ ”.

<sup>7</sup>When  $P_{\mathbf{X}}$  is stationary, the limit  $\lim_{n \rightarrow \infty} \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \triangleq \mathbb{D}(P_{\mathbf{X}}, \mathbf{\Pi})$  was shown to exist in [1]. Thus, (42) was equivalently stated as  $\lim_{n \rightarrow \infty} EL_{\hat{\mathbf{X}}_{\text{DUDE}}^n} = \mathbb{D}(P_{\mathbf{X}}, \mathbf{\Pi})$  in [1, Theorem 3].

(43) not hold for the DUDE, but clearly (43) cannot hold under *any* sequence of denoisers. In other words,  $\mathcal{P}\{m_n\}$  is far too rich to allow for the existence of schemes that are universal with respect to it.

It is obvious then that  $\mathcal{P}\{m_n\}$  is significantly richer than the family of stationary processes whenever  $m_n$  grows with  $n$ . It is of interest then to identify the maximal growth rate of  $m_n$  that allows for the existence of schemes that are universal with respect to  $\mathcal{P}\{m_n\}$ , and to find such a universal scheme. In what follows, we offer a complete answer to these questions. Specifically, we show that if the growth rate of  $m_n$  allows for the existence of *any* scheme which is universal with respect to  $\mathcal{P}\{m_n\}$ , the S-DUDE is universal, too.

## 4-C Universality of S-DUDE

Here, we state our stochastic setting result, which establishes the universality of  $(k, m)$ -S-DUDE with respect to the class  $\mathcal{P}\{m_n\}$ .

**Theorem 6** *Let  $k = k_n$  and  $m = m_n$  satisfy the growth rate condition stipulated in Theorem 4, in addition to  $\lim_{n \rightarrow \infty} k_n = \infty$ . Then, the sequence of denoisers  $\{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}\}$  defined in Section 3 satisfy*

$$\lim_{n \rightarrow \infty} \left[ EL_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \right] = 0 \quad \text{for all } P_{\mathbf{X}} \in \mathcal{P}\{m_n\}. \quad (44)$$

*Remark 1:* Recall that, as noted in Claim 1,  $m_n = o(n)$  together with appropriately slowly growing  $k = k_n$  is sufficient to guarantee the growth rate condition stipulated in Theorem 4. Hence, by Theorem 6,  $m = o(n)$  and the sufficiently slowly growing  $k = k_n$  suffices for (44) to hold. Therefore, Theorem 6 implies the existence of schemes that are universal with respect to  $\mathcal{P}\{m_n\}$  whenever  $m_n$  increases sublinearly in  $n$ . Since, as discussed in Subsection 4-B, no universal scheme exists for  $\mathcal{P}\{m_n\}$  when  $m_n$  is linear in  $n$ , we conclude that the sub-linearity of  $m_n$  is the necessary and sufficient condition for a universal scheme to exist with respect to  $\mathcal{P}\{m_n\}$ . Moreover, Theorem 6 establishes the strong sense of optimality of the S-DUDE, as it shows that whenever  $\mathcal{P}\{m_n\}$  is universally “competable”, the S-DUDE does the job. This fact is somewhat analogous to the situation in [21], where the optimality of the universal lossless coding scheme presented therein for piecewise stationary sources was established under the condition that  $m = o(n)$ .

*Remark 2:* A pointwise result

$$\lim_{n \rightarrow \infty} \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \right] = 0 \quad \text{a.s.}$$

for all  $P_{\mathbf{X}} \in \mathcal{P}\{m_n\}$ , which is analogous to [1, Theorem 4], can also be derived. However, we omit such a result here since the details required for stating it rigorously would be convoluted, and its added value over the strong point-wise result we have already established in the semi-stochastic setting would be little.

*Proof of Theorem 6:* See Appendix 8-H. ■

## 5 Algorithm and Complexity

### 5-A An Efficient Implementation of S-DUDE

In the preceding two sections, we gave strong asymptotic performance guarantees for the new class of schemes, the S-DUDE. However, the question regarding the practical implementation of (32), i.e., obtaining

$$\hat{\mathbf{S}}_{k,m} = \arg \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \tilde{L}_{\mathbf{S}}(z^n),$$

for fixed  $k$ ,  $m$  and  $n$  remains and, at first glance, may seem to be a difficult combinatorial optimization problem. In this section, we devise an efficient two-pass algorithm, which yields (32) and performs denoising

with linear complexity in the sequence length  $n$ . A recursion similar to that in the first pass of the algorithm we present appears also in the study of tracking the best expert in on-line learning [15, 16].

From the definition of  $\mathcal{S}_{k,m}^n$ , (30), we can see that obtaining (32) is equivalent to obtaining the best combination of single-symbol denoisers with at most  $m(\mathbf{c})$  shifts that minimizes the cumulative estimated loss along  $\{t : t \in \mathcal{T}(\mathbf{c})\}$ , for each  $\mathbf{c} \in \mathbf{C}_k$ . Thus, our problem breaks down to  $|\mathbf{C}_k|$  independent problems, each being a problem of competing with the best combination of single-symbol schemes allowing  $m$  switches.

To describe an algorithm that implements this parallelization efficiently, we first define variables. For  $(k, m)$ -S-DUDE, let  $I = m + 1, J = N + 1$ , where  $N = |\mathcal{S}| = |\mathcal{Z}|^{|\hat{\mathcal{X}}|}$ . Then, a matrix  $M_t \in \mathbb{R}^{I \times J}$  is defined for  $k + 1 \leq t \leq n - 2k$ , where  $M_t(i, j)$  for  $1 \leq i \leq I$  and  $1 \leq j \leq J - 1$  represents the minimum (un-normalized) cumulative estimated loss of the sequence of single-symbol denoisers along the time index  $\{\tau : \tau \leq t, \mathbf{c}_\tau = \mathbf{c}_t\}$ , allowing at most  $(i - 1)$  shifts between single-symbol denoisers and applying  $s_t = j$ . Moreover,  $M_t(i, J)$ , for  $1 \leq i \leq I$ , is the symbol-by-symbol denoiser that attains the minimum value of the  $i$ -th row of  $M_t$ , i.e.,  $\arg \min_{1 \leq j \leq J-1} M_t(i, j)$ . A time pointer  $T \in \mathbb{R}^D$ , where  $D = |\mathbf{C}_k| = |\mathcal{Z}|^{2k}$ , is defined to store the closest time index that has the same context as current time, during the first and second pass. That is,

$$T(\mathbf{c}_t) \triangleq \begin{cases} \max\{\tau : \tau < t, \mathbf{c}_\tau = \mathbf{c}_t\}, & \text{when first pass} \\ \min\{\tau : \tau > t, \mathbf{c}_\tau = \mathbf{c}_t\}, & \text{when second pass} \end{cases} \quad (45)$$

We also define  $r \in \mathbb{R}^D$  and  $q \in \mathbb{R}^D$  as variables for storing the pointer enabling our scheme to follow the best combination of single-symbol denoisers during the second pass. Thus, the total memory size required is  $O(mNn + |\mathcal{Z}|^{2k}) = O(mn)$  (assuming that  $k$  satisfies the growth rate stipulated in the previous sections, which implies  $|\mathcal{Z}|^{2k} = o(n)$ ).

Our two-pass algorithm has ingredients from both the DUDE and from the forward-backward recursions of hidden Markov models [28] and, in fact, the algorithm becomes equivalent to DUDE when  $m = 0$ . The first pass of the algorithm runs forward from  $t = k + 1$  to  $t = n - k$ , and updates the elements of  $M_t$  recursively. The recursions have a natural dynamic programming structure. For  $2 \leq i \leq I, 1 \leq j \leq J - 1$ ,  $M_t(i, j)$  is determined by

$$M_t(i, j) = \ell(z_t, j) + \min \left\{ M_{T(\mathbf{c}_t)}(i, j), M_{T(\mathbf{c}_t)}(i - 1, M_{T(\mathbf{c}_t)}(i - 1, J)) \right\}, \quad (46)$$

that is, adding the current loss to the best cumulative loss up to  $T(\mathbf{c}_t)$  along  $\{\tau : \tau < t, \mathbf{c}_\tau = \mathbf{c}_t\}$ . When  $i = 1$ , the second term in the minimum of (46) is not defined, and  $M_t(i, j)$  just becomes  $\ell(z_t, j) + M_{T(\mathbf{c}_t)}(i, j)$ . The validity of (46) can be verified by observing that there are two possible cases in achieving  $M_t(i, j)$ : either the  $(i - 1)$ -th shift to the single-symbol denoiser  $j$  occurred before  $t$ , or it occurred at time  $t$ . We can see that the first term in the minimum of (46) corresponds to the former case; the second term corresponds to the latter. Obviously, the minimum of these two (where ties may be resolved arbitrarily), leads to the value of  $M_t(i, j)$  as in (46). After updating all  $M_t$ 's during the first pass, the second pass runs backwards from  $t = n - k$  to  $t = k + 1$ , and extracts  $\hat{\mathbf{S}}_{k,m}$  from  $\{M_t\}_{t=k+1}^{n-2k}$  by following the best shifting between single-symbol denoisers. The actual denoising (i.e., assembling the reconstruction sequence  $\hat{X}^n$ ) is also performed in that pass. The pointers  $r(\mathbf{c}_t)$  and  $q(\mathbf{c}_t)$  are updated recursively, and they track the best shifting point and combination of single-symbol denoisers, respectively, for each of the subsequences associated with the various contexts. A succinct description of the algorithm is provided in Algorithm 1. The time complexity of the algorithm is readily seen to be  $O(mn)$  as well.

---

**Algorithm 1** The  $(k, m)$ -Shifting Discrete Denoising Algorithm

---

**Require:**  $M_t(i, j) \in \mathbb{R}^{I \times J}$ ,  $k + 1 \leq t \leq n - 2k$ ,  $1 \leq i \leq I$ ,  $1 \leq j \leq J$ ,  $T \in \mathbb{R}^D$ ,  $r \in \mathbb{R}^D$ ,  $q \in \mathbb{R}^D$ ,  $L \in \mathbb{R}$   
**Ensure:**  $\hat{\mathbf{S}}_k = \{s_{k,t}(\mathbf{c}_t, \cdot)\}_{t=k+1}^{n-2k}$  in (32) and the denoised output  $\{\hat{x}_t\}_{t=k+1}^{n-k}$   
 $\tau(\mathbf{c}) \leftarrow \phi$  for all  $\mathbf{c} \in \mathbf{C}_k$   
**for**  $t = k + 1$  to  $n - 2k$  **do**  
  **if**  $T(\mathbf{c}_t) = \phi$  **then**  
     $M_t(i, j) \leftarrow \ell(z_t, j)$  for  $1 \leq i \leq I$ ,  $1 \leq j \leq J - 1$   
     $M_t(i, J) \leftarrow \arg \min_{1 \leq j \leq J-1} M_t(i, j)$  for  $1 \leq i \leq I$   
  **else**  
     $M_{T(\mathbf{c}_t)}^*(i, j) \leftarrow \left\{ \begin{array}{ll} M_{T(\mathbf{c}_t)}(i, j) & \text{for } i = 1, \quad 1 \leq j \leq J - 1 \\ \min \{M_{T(\mathbf{c}_t)}(i, j), M_{T(\mathbf{c}_t)}(i - 1, M_{T(\mathbf{c}_t)}(i - 1, J))\} & \text{for } 2 \leq i \leq I, \quad 1 \leq j \leq J - 1 \end{array} \right\}$   
     $M_t(i, j) \leftarrow M_{T(\mathbf{c}_t)}^*(i, j) + \ell(z_t, j)$  for  $1 \leq i \leq I$ ,  $1 \leq j \leq J - 1$   
     $M_t(i, J) \leftarrow \arg \min_{1 \leq j \leq J-1} M_t(i, j)$  for  $1 \leq i \leq I$   
  **end if**  
   $T(\mathbf{c}_t) \leftarrow t$   
**end for**  
 $T(\mathbf{c}) \leftarrow \phi$  for all  $\mathbf{c} \in \mathbf{C}_k$   
**for**  $t = n - 2k$  to  $k + 1$  **do**  
  **if**  $T(\mathbf{c}_t) = \phi$  **then**  
     $r(\mathbf{c}_t) \leftarrow I$ ,  $q(\mathbf{c}_t) \leftarrow M_t(r(\mathbf{c}_t), J)$   
  **else**  
     $L \leftarrow M_{T(\mathbf{c}_t)}(r(\mathbf{c}_t), q(\mathbf{c}_t)) - \ell(z_t, q(\mathbf{c}_t))$   
    **if**  $L < M_t(r(\mathbf{c}_t), q(\mathbf{c}_t))$  **then**  
       $r(\mathbf{c}_t) \leftarrow r(\mathbf{c}_t) - 1$ ,  $q(\mathbf{c}_t) \leftarrow M_t(r(\mathbf{c}_t), J)$   
    **end if**  
  **end if**  
   $T(\mathbf{c}_t) \leftarrow t$ ,  $s_{k,t}(\mathbf{c}_t, \cdot) \leftarrow q(\mathbf{c}_t)$   
   $\hat{x}_t \leftarrow s_{k,t}(\mathbf{c}_t, z_t)$   
**end for**

---

## 5-B Extending the S-DUDE to Multi-Dimensional Data

As noted, our algorithm is essentially separately employing the same algorithm to compete with the best shifting single-symbol denoisers, on each subsequence associated with each context. The overall algorithm is the result of parallelizing the operations of the schemes for the different subsequences, which allows for a more efficient implementation than if these schemes were to be run completely independently of one another. This characteristic of running the same algorithm in parallel along each subsequence enables us to extend S-DUDE to the case of multi-dimensional data: run the same algorithm along each subsequence associated with each (this time multi-dimensional) context. It should be noted, however, that the extension of the S-DUDE to the multidimensional case is not as straightforward as the extension of the DUDE was, since, whereas the DUDE's output is independent of the ordering of the data within each context, this ordering may be very significant in its effect on the output and, hence, the performance of S-DUDE. Therefore, the choice of a scheme for scanning the data and capturing its local spatial stationarity, e.g., Peano-Hilbert scanning [29], is an important ingredient in extending S-DUDE to the denoising of multi-dimensional data. Findings from the recent study on universal scanning reported in [30, 31] can be brought to bear on such an extension.

## 6 Experimentation

In this section, we report some preliminary experimental results obtained by applying S-DUDE to several kinds of noise-corrupted data.

### 6-A Image denoising

In this subsection, we report some experimental results of denoising a binary image under the Hamming loss function. The first and most simplistic experiment is with the  $400 \times 400$  black-and-white binary image shown in Figure 1. The first figure is the clean underlying image. The image is passed through a binary symmetric channel (BSC) with crossover probability  $\delta = 0.1$ , to obtain the noisy image (second image in Figure 1). Note that in this case, there are only four symbol-by-symbol denoisers, namely,  $\mathcal{S} = \{0, 1, z, \bar{z}\}$ , representing always-say-0, always-say-1, say-what-you-see, and flip-what-you-see, respectively. The third image in Figure 1 is the DUDE output with  $k = 0$ , and the last image is the output of our S-DUDE with  $k = 0, m = 1$ . The DUDE with  $k = 0$  competes with the best time-invariant symbol-by-symbol denoiser

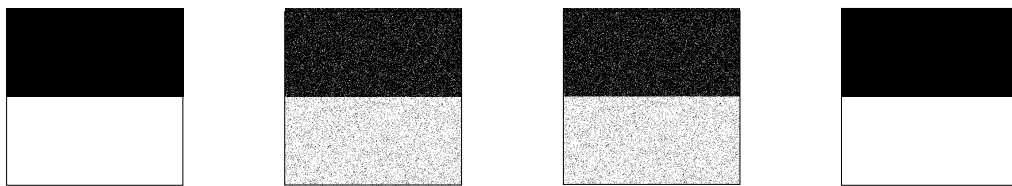


Figure 1:  $400 \times 400$  binary images.

which, in this case, is the say-what-you-see denoiser, since the empirical distribution of the clean image is  $(0.5, 0.5)$  and  $\delta = 0.1$ . Thus, the DUDE output is the same as the noisy image; hence, no denoising is performed. However, it is clear that, for this image, the best compound action of the symbol-by-symbol denoisers is always-say-0 for the first half and then a shift to always-say-1 for the remainder. We can see that our  $(0, 1)$ -S-DUDE successfully captures this shift from the noisy observations, and results in *perfect* denoising with zero bit errors.

Now, we move on to a more realistic example. The first image in Figure 2, a concatenation of a half-toned Einstein image ( $300 \times 300$ ) and scanned Shannon’s 1948 paper ( $300 \times 300$ ), is the clean image. We pass the image through a binary symmetric channel (BSC) with crossover probability  $\delta = 0.1$ , to obtain the second noisy image, which we raster scan and employ the S-DUDE on the resulting one-dimensional sequence. Since the two concatenated images are of a very different nature, we expect our S-DUDE to perform better than the DUDE, because it is designed to adapt to the possibility of employing different schemes in different regions of the data. The plot shows the performance of our  $(k, m)$ -S-DUDE with various values of  $k$  and  $m$ . The horizontal axis reflects  $k$ , and the vertical axis represents the ratio of bit error per symbol (BER) to  $\delta = 0.1$ . Each curve represents the BER of schemes with different  $m = 0, 1, 2, 3$ . Note that  $m = 0$  corresponds to the DUDE. We can see that S-DUDE with  $m > 0$  mostly dominates the DUDE, with an additional BER reduction of  $\sim 11\%$ , including when  $k = 6$ , the best  $k$  value for the DUDE. The bottom three figures show the denoised images with  $(k, m) = (4, 0), (4, 2), (6, 1)$ , achieving BERs of  $\delta \times (0.744, 0.6630, 0.4991)$ , respectively. Thus, in this example,  $(4, 2)$ -S-DUDE achieves an additional BER reduction of  $11\%$  over the DUDE with  $k = 4$ , and the overall best performance is achieved by  $(6, 1)$ -S-DUDE. Given the nature of the image, which is a concatenation of two completely different types of images, each reasonably uniform in texture, it is not surprising to find that the S-DUDE with  $m = 1$  performs the best.

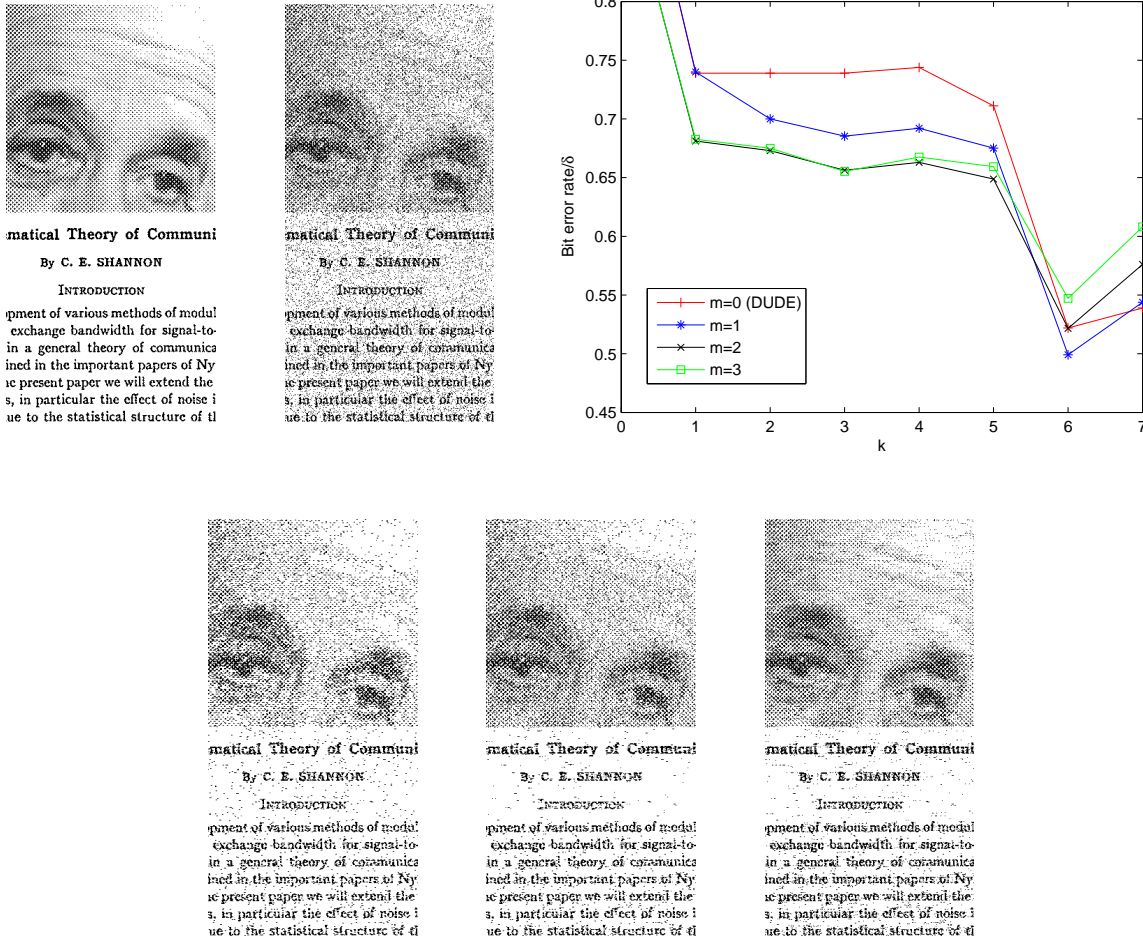


Figure 2: Clean and noisy images, the bit error rate plot for  $(k, m)$ -S-DUDE, and three denoised outputs for  $(k, m) = (4, 0), (4, 2), (6, 1)$ , respectively.

## 6-B State estimation for a switching binary hidden Markov process

Here, we give a stochastic setting experiment. A switching binary hidden Markov process in this example is defined as a binary symmetric Markov chain observed through a BSC, where the transition probabilities of the Markov chain switches over time. The goal of a denoiser here is to estimate the underlying Markov chain based on the noisy output.

In our example, we construct a simple switching binary hidden Markov process of length  $n = 10^6$ , in which the transition probability of the underlying binary symmetric Markov source switches from  $p = 0.01$  to  $p = 0.2$  at the midpoint of the sequence, and the crossover probability of BSC is  $\delta = 0.1$ . Then, we estimate the state of the underlying Markov chain based on the BSC output. The goodness of the estimation is again measured by the Hamming loss, i.e., the fraction of errors made. Slightly better than the optimal Bayesian distribution-dependent performance for this case can be obtained by employing the forward-backward recursion scheme, incorporating the varying transition probabilities with the help of a genie that knows the exact location of the change in the process distribution. Figure 3 plots the BER of  $(k, m)$ -S-DUDE with various  $k$  and  $m$ , compared to the genie-aided Bayes optimal BER. The horizontal

axis represents  $k$ , and the two curves refer to  $m = 0$  (DUDE) and  $m = 1$ . The vertical axis is the ratio of BER to  $\delta = 0.1$ .

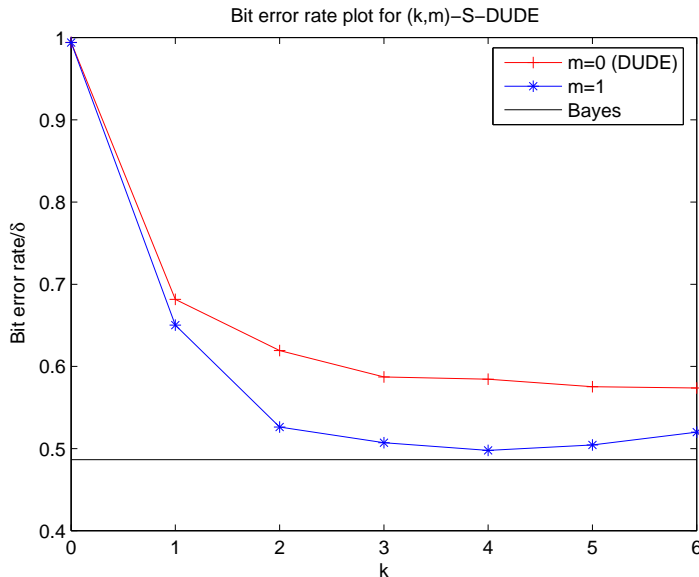


Figure 3: BER for switching binary hidden Markov process ( $\delta = 0.1, n = 10^6$ ). The switch of the underlying binary Markov chain occurs when  $t = 5 \times 10^5$ , from the transition probability  $p = 0.01$  to  $p = 0.2$ .

We can observe that the optimal Bayesian BER is (lower bounded by)  $0.4865 \times \delta$ . The best performance of the DUDE was achieved when  $k = 6$  with a BER of  $0.5738 \times \delta$ , which is far above (18% more than) the optimal BER. It is clear that, despite the size of the data, the DUDE fails to converge to the optimum, as it is confined to be employing the same sliding-window scheme throughout the whole data. However, we can see that the  $(4, 1)$ -S-DUDE achieves a BER of  $.4979 \times \delta$ , which is within 2.3% of the optimal BER. This example shows that our S-DUDE is competent in attaining the optimum performance for a class richer than that of the stationary processes. Specifically, it attains the optimum performance for piecewise stationary processes, on which the DUDE generally fails.

## 7 Conclusion and Some Future Directions

Inspired by the DUDE algorithm, we have developed a generalization that accommodates switching between sliding window rules. We have shown a strong semi-stochastic setting result for our new scheme in competing with shifting  $k$ -th order denoisers. This result implies a stochastic setting result as well, asserting that the S-DUDE asymptotically attains the optimal distribution-dependent performance for the case in which the underlying data is piecewise stationary. We also described an efficient low-complexity implementation of the algorithm, and presented some simple experiments that demonstrate the potential benefits of employing S-DUDE in practice.

There are several future research directions related to this work. The S-DUDE can be thought of as a generalization of the DUDE, with the introduction of a new component captured by the non-negative integer parameter  $m$ . Many previous extensions of the DUDE, such as the settings of channel with memory[34], channel uncertainty [33], applications to channel decoding[37], discrete-input, continuous-output data[35], denoising of analog data[32], and decoding in the Wyner-Ziv problem[36], may stand to benefit from a revision that would incorporate the viewpoint of switching between time-invariant schemes. Particularly,

extending S-DUDE to the case where the the data are analog as in [32] will be non-trivial and interesting from both a theoretical and a practical viewpoint. In addition, as mentioned in Section 5, an extension of the S-DUDE to the case of multi-dimensional data is not as straightforward as the extension of the DUDE was. Such an extension should prove interesting and practically important. Finally, it would be useful to devise guidelines, in the spirit of those in [38, 3], for the choice of  $k$  and  $m$  based on  $n$  and the noisy observation sequence  $z^n$ .

## Acknowledgments

The first author is grateful to Professor Manfred Warmuth for introducing him to a substantial amount of related work on the expert tracking problems in online learning.

## 8 Appendix

### 8-A Proof of Lemma 1

We first establish the fact that for all  $x^n \in \mathcal{X}^n$ , and for fixed  $\mathbf{S} \in \mathcal{S}_0^n$ ,

$$\left\{ n[L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n)] \right\}_{n \geq 1}$$

is a  $\{Z_n\}$ -martingale. This is not hard to see by following:

$$\begin{aligned} & E\left(n[L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n)] \middle| Z^{n-1}\right) \\ &= E\left(\sum_{t=1}^n \Lambda(x_t, s_t(Z_t)) - \sum_{t=1}^n \ell(Z_t, s_t) \middle| Z^{n-1}\right) \\ &= (n-1)[L_{\hat{\mathbf{X}}^{n-1}, \mathbf{S}}(x^{n-1}, Z^{n-1}) - \tilde{L}_{\mathbf{S}}(Z^{n-1})] + E\left(\Lambda(x_n, s_n(Z_n)) - \ell(Z_n, s_n) \middle| Z^{n-1}\right) \\ &= (n-1)[L_{\hat{\mathbf{X}}^{n-1}, \mathbf{S}}(x^{n-1}, Z^{n-1}) - \tilde{L}_{\mathbf{S}}(Z^{n-1})], \end{aligned} \quad (47)$$

where (47) follows from the fact that  $Z_n$  is independent of  $Z^{n-1}$ , and  $E\Lambda(x_n, s_n(Z_n)) = E\ell(Z_n, s_n)$ . Therefore,  $L_{\mathbf{S}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n)$  is a normalized sum of bounded martingale differences; therefore the inequalities (20) and (21) follow directly from the Hoeffding-Azuma inequality [14, Lemma A.7]. ■

### 8-B Proof of Theorem 2

Consider following chain of inequalities:

$$\begin{aligned} & P\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}}^n(x^n, Z^n) - D_{0,m}(x^n, Z^n) > \epsilon\right) \\ &= P\left(\max_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \{L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}}^n(x^n, Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}^n(x^n, Z^n)\} > \epsilon\right) \\ &\leq \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} P\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}}^n(x^n, Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}^n(x^n, Z^n) > \epsilon\right) \quad (48) \\ &\leq \underbrace{\sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} P\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}}^n(x^n, Z^n) - \tilde{L}_{\hat{\mathbf{S}}}^n(Z^n) > \epsilon/2\right)}_{(i)} + \underbrace{\sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} P\left(\tilde{L}_{\hat{\mathbf{S}}}^n(Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}^n(x^n, Z^n) > \epsilon/2\right)}_{(ii)}, \quad (49) \end{aligned}$$

where (48) follows from the union bound, and (49) follows from adding and subtracting  $\tilde{L}_{\hat{\mathbf{S}}}(Z^n)$ , and the union bound. For term (i) in (49),

$$(i) \leq \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} P\left(\max_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \{L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n)\} > \epsilon/2\right) \quad (50)$$

$$\leq \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \exp\left(-n \frac{\epsilon^2}{2L_{\max}^2}\right), \quad (51)$$

where (50) follows from  $L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}}(x^n, Z^n) - \tilde{L}_{\hat{\mathbf{S}}}(Z^n) \leq \max_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \{L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x^n, Z^n) - \tilde{L}_{\mathbf{S}}(Z^n)\}$ , and (51) follows from the union bound and (20). Similarly, for term (ii) in (49),

$$(ii) \leq \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} P\left(\tilde{L}_{\mathbf{S}}(Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x^n, Z^n) > \epsilon/2\right) \quad (52)$$

$$\leq \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \exp\left(-n \frac{\epsilon^2}{2L_{\max}^2}\right), \quad (53)$$

where (52) follows from  $\tilde{L}_{\hat{\mathbf{S}}}(Z^n) \leq \tilde{L}_{\mathbf{S}}(Z^n)$  a.s., and (53) follows from (21). Therefore, continuing (49), we obtain

$$(49) \leq 2 \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \sum_{\mathbf{S} \in \mathcal{S}_{0,m}^n} \exp\left(-n \frac{\epsilon^2}{2L_{\max}^2}\right) \\ = 2 \left[ \sum_{k=0}^m \binom{n-1}{k} N(N-1)^k \right]^2 \exp\left(-n \frac{\epsilon^2}{2L_{\max}^2}\right) \quad (54)$$

$$\leq 2 \exp\left(-n \left[ \frac{\epsilon^2}{2L_{\max}^2} - 2h\left(\frac{m}{n}\right) - \frac{2(m+1) \ln N}{n} \right]\right), \quad (55)$$

where (54) follows from  $|\mathcal{S}_{0,m}^n| = \sum_{k=0}^m \binom{n-1}{k} N(N-1)^k$ , and (55) follows from  $|\mathcal{S}_{0,m}^n| \leq N^{m+1} \exp(nh(\frac{m}{n}))$ . Hence, the theorem is proved. ■

## 8-C Proof of Lemma 2

We will prove (28) since the proof of (29) is essentially identical. As in [1], define

$$\mathcal{I}_d \triangleq \{t : k+1 \leq t \leq n-k, t \equiv d \pmod{k+1}\},$$

whose cardinality is denoted  $n_d = \lfloor (n-d-k)/(k+1) \rfloor$ . Then, by denoting  $\mathbf{C}_t = (Z_{t-k}^{t-1}, Z_{t+1}^{t+k})$ , we start the chain of inequalities,

$$\Pr\left(L_{\hat{\mathbf{X}}^n, \mathbf{S}_k}(x_{k+1}^{n-k}, Z^n) - \tilde{L}_{\mathbf{S}_k}(Z^n) > \epsilon\right) \\ \leq \Pr\left(\sum_{d=0}^k \sum_{\tau \in \mathcal{I}_d} \left\{ \Lambda(x_\tau, s_{k,\tau}(\mathbf{C}_\tau, Z_\tau)) - \ell(Z_\tau, s_{k,\tau}(\mathbf{C}_\tau, \cdot)) \right\} > (n-2k)\epsilon\right) \quad (56)$$

$$\leq \sum_{d=0}^k \Pr\left(\sum_{\tau \in \mathcal{I}_d} \left\{ \Lambda(x_\tau, s_{k,\tau}(\mathbf{C}_\tau, Z_\tau)) - \ell(Z_\tau, s_{k,\tau}(\mathbf{C}_\tau, \cdot)) \right\} > (n-2k)\gamma_d \epsilon\right), \quad (57)$$

where (56) follows from the triangle inequality, (57) follows from the union bound, and  $\{\gamma_d\}$  is a set of nonnegative constants (to be specified later) satisfying  $\sum_d \gamma_d = 1$ . In the sequel, for simplicity, we will denote  $\Lambda(x_\tau, s_{k,\tau}(\mathbf{C}_\tau, Z\tau))$  and  $\ell(Z_\tau, s_{k,\tau}(\mathbf{C}_\tau, \cdot))$  in (48) as  $\Lambda_\tau$  and  $\ell_\tau$ , respectively. Now, the collection of random variables  $Z(d)$  is defined to be

$$Z(d) \triangleq \{Z_t : 1 \leq t \leq n, t \notin \mathcal{I}_d\},$$

and  $z(d) \in \mathcal{Z}^{n-n_d}$  denotes a particular realization of  $Z(d)$ . Then, by conditioning, we have

$$(57) \leq \sum_{d=0}^k \sum_{z(d) \in \mathcal{Z}^{n-n_d}} \Pr(Z(d) = z(d)) \Pr\left(\sum_{\tau \in \mathcal{I}_d} \{\Lambda_\tau - \ell_\tau\} > (n-2k)\gamma_d \epsilon \middle| Z(d) = z(d)\right), \quad (58)$$

and let  $P_d$  denote the conditional probability of (58). Now, conditioned on  $Z(d) = z(d)$ ,  $\{Z_\tau\}_{\tau \in \mathcal{I}_d}$  are all independent, and the summation in  $P_d$  beomes

$$\sum_{\tau \in \mathcal{I}_d} \left\{ \Lambda(x_\tau, s_{k,\tau}(\mathbf{c}_\tau, Z\tau)) - \ell(Z_\tau, s_{k,\tau}(\mathbf{c}_\tau, \cdot)) \right\},$$

which is the sum of the absolute differences of the true and estimated losses of the symbol-by-symbol denoisers  $s_{k,\tau}(\mathbf{c}_\tau, \cdot)$  over  $\tau \in \mathcal{I}_d$ . Thus, we can apply (20), and obtain

$$\begin{aligned} P_d &= \Pr\left(\sum_{\tau \in \mathcal{I}_d} \{\Lambda_\tau - \ell_\tau\} > n_d \cdot \frac{(n-2k)\gamma_d \epsilon}{n_d} \middle| Z(d) = z(d)\right) \\ &\leq \exp\left(-\frac{2(n-2k)^2 \gamma_d^2 \epsilon^2}{L_{\max}^2 n_d}\right). \end{aligned} \quad (59)$$

Following [1], we choose  $\gamma_d = \frac{\sqrt{n_d}}{\sum_j \sqrt{n_j}}$ , and from the Cauchy-Schwartz inequality and  $\sum_d n_d = n - 2k$ , we arrive at

$$\frac{n_d}{\gamma_d^2} \leq (k+1) \sum_{d=0}^k n_d = (k+1)(n-2k),$$

and, hence,

$$P_d \leq \exp\left(-\frac{2(n-2k)\epsilon^2}{(k+1)L_{\max}^2}\right). \quad (60)$$

Therefore, plugging (60) into (58), we finally have

$$(58) \leq (k+1) \exp\left(-\frac{2(n-2k)\epsilon^2}{(k+1)L_{\max}^2}\right),$$

which proves the lemma.  $\blacksquare$

## 8-D Proof of Theorem 3

The proof resembles that of Theorem 2. Consider

$$\begin{aligned} & \Pr\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}_{k,m}}(x_{k+1}^{n-k}, Z^n) - D_{k,m}(x^n, Z^n) > \epsilon\right) \\ &= P\left(\max_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \{L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}_{k,m}}(x_{k+1}^{n-k}, Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x_{k+1}^{n-k}, Z^n)\} > \epsilon\right) \\ &\leq \sum_{\mathbf{S} \in \mathcal{S}_{k,m}^n} P\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}_{k,m}}(x_{k+1}^{n-k}, Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x_{k+1}^{n-k}, Z^n) > \epsilon\right) \end{aligned} \quad (61)$$

$$\leq \sum_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \left\{ P\left(L_{\hat{\mathbf{X}}^n, \hat{\mathbf{S}}_{k,m}}(x_{k+1}^{n-k}, Z^n) - \tilde{L}_{\hat{\mathbf{S}}_{k,m}}(Z^n) > \frac{\epsilon}{2}\right) + P\left(\tilde{L}_{\hat{\mathbf{S}}_{k,m}}(Z^n) - L_{\hat{\mathbf{X}}^n, \mathbf{S}}(x_{k+1}^{n-k}, Z^n) > \frac{\epsilon}{2}\right) \right\} \quad (62)$$

$$\leq 2(k+1) \sum_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \sum_{\mathbf{S} \in \mathcal{S}_{k,m}^n} \exp\left(-\frac{(n-2k)\epsilon^2}{2(k+1)L_{\max}^2}\right) \quad (63)$$

$$= 2(k+1) \left[ \sum_{k=0}^{m(\mathbf{c})} \binom{n(\mathbf{c})-1}{k} N(N-1)^k \right]^{2|\mathbf{C}_k|} \exp\left(-\frac{(n-2k)\epsilon^2}{2(k+1)L_{\max}^2}\right), \quad (64)$$

where (61)-(62) follow similarly as in (48)-(49); (63) follows from arguments similar to (50), (52), and Lemma 2 (which plays the role that Lemma 1 played there); and (64) follows from  $|\mathcal{S}_{k,m}^n| = \left[ \left( \sum_{k=0}^{m(\mathbf{c})} \binom{n(\mathbf{c})-1}{k} \right) N(N-1)^k \right]^{|\mathbf{C}_k|}$ . Now, for all  $\mathbf{c} \in \mathbf{C}_k$ ,

$$\begin{aligned} \sum_{k=0}^{m(\mathbf{c})} \binom{n(\mathbf{c})-1}{k} N(N-1)^k &\leq N^{m+1} \exp\left(n(\mathbf{c})h\left(\frac{m(\mathbf{c})}{n(\mathbf{c})}\right)\right) \\ &\leq N^{m+1} \exp\left((n-2k)h\left(\frac{m(\mathbf{c})}{n-2k}\right)\right) \end{aligned} \quad (65)$$

$$\leq N^{m+1} \exp\left((n-2k)h\left(\frac{m}{n-2k}\right)\right), \quad (66)$$

where (65) is based on the fact that  $\exp(nh(\frac{m}{n}))$  is an increasing function in  $n$ , and (66) follows from  $m \leq \lfloor \frac{n-2k}{2} \rfloor$ . Therefore, together with  $|\mathbf{C}_k| = |\mathcal{Z}|^{2k}$ , we have

$$(64) \leq 2(k+1) \exp\left(- (n-2k) \cdot \left[ \frac{\epsilon^2}{2(k+1)L_{\max}^2} - 2|\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\ln N}{n-2k} \right\} \right] \right), \quad (67)$$

which proves the theorem.  $\blacksquare$

## 8-E Proof of Claim 1

For part a), to show the necessity first, suppose  $c_1 \geq \frac{1}{2\log|\mathcal{Z}|}$ . Then, from  $|\mathcal{Z}|^{2k} = n^{\frac{2k \log|\mathcal{Z}|}{\log n}}$ , we have  $2|\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\ln N}{n-2k} \right\} = \Omega\left(n^{\frac{2k \log|\mathcal{Z}|}{\log n}} \left(\frac{m}{n}\right)^{1-\delta}\right)$ , which will grow to infinity as  $n$  grows, even when  $m$  is fixed. Therefore, the right-hand side of (34) is not summable. On the other hand,  $k = c_1 \log n$  with  $c_1 < \frac{1}{2\log|\mathcal{Z}|}$  is readily verified to suffice for the summability, provided that  $m = m_n$  grows at any sub-polynomial rate, i.e., grows more slowly than  $n^\alpha$  for any  $\alpha > 0$  (e.g.,  $c_2 \log n$ ).

For part b), to show the necessity, suppose  $m = \Theta(n)$ . Then,  $h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\ln N}{n-2k} = \Theta(1)$ , and, thus, for sufficiently small  $\epsilon$ ,  $\frac{\epsilon^2}{2(k+1)L_{\max}^2} - |\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\log N}{n-2k} \right\} < 0$  even for  $k$  fixed. Therefore, the

right-hand side of (34) is not summable. Hence,  $m = o(n)$  is necessary for the summability. For sufficiency, suppose  $m = m_n$  is any rate, such that  $\lim_{n \rightarrow \infty} \frac{m_n}{n} = 0$ . Then,

$$\begin{aligned} & \frac{\epsilon^2}{2(k+1)L_{\max}^2} - 2|\mathcal{Z}|^{2k} \cdot \left\{ h\left(\frac{m}{n-2k}\right) + \frac{(m+1)\log N}{n-2k} \right\} \\ &= \frac{1}{k} \left\{ \frac{\epsilon^2}{2\left(1 + \frac{1}{k}\right)L_{\max}^2} - 2k|\mathcal{Z}|^{2k} \cdot O\left(\left(\frac{m_n}{n}\right)^{1-\delta}\right) \right\}. \end{aligned} \quad (68)$$

Thus, if  $k$  grows sufficiently slowly that  $k|\mathcal{Z}|^{2k} = o\left(\left(\frac{n}{m_n}\right)^{1-\delta}\right)$ , then (68) becomes positive for sufficiently large  $n$ , and the right-hand side of (34) becomes summable. ■

## 8-F Proof of Theorem 4

First, denote the random variable  $A_{k,m}^n \triangleq L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x_{k+1}^{n-k}, Z^n) - D_{k,m}(x^n, Z^n)$ . Then, for part a), we have

$$L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x^n, Z^n) - D_{k,m}(x^n, Z^n) \leq \frac{2k\Lambda_{\max}}{n} + A_{k,m}^n \quad a.s.$$

Since the maximal rate for  $k$  is  $c_1 \log n$  as specified in Claim 1,  $\lim_{n \rightarrow \infty} \frac{2k\Lambda_{\max}}{n} = 0$ . Furthermore, from the summability condition on  $k$  and  $m$ , Theorem 3, and the Borel-Cantelli lemma, we get  $\lim_{n \rightarrow \infty} A_{k,m}^n = 0$  with probability 1, which proves part a). To prove part b), note that, for any  $\epsilon > 0$ ,

$$\begin{aligned} & E[L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x^n, Z^n) - D_{k,m}(x^n, Z^n)] \\ &\leq \frac{2k\Lambda_{\max}}{n} + E(A_{k,m}^n) \\ &= \frac{2k\Lambda_{\max}}{n} + E(A_{k,m}^n | A_{k,m}^n \leq \epsilon) Pr(A_{k,m}^n \leq \epsilon) + E(A_{k,m}^n | A_{k,m}^n > \epsilon) Pr(A_{k,m}^n > \epsilon) \\ &\leq \frac{2k\Lambda_{\max}}{n} + \epsilon + \Lambda_{\max} \cdot Pr(A_{k,m}^n > \epsilon) \\ &\leq \frac{2k\Lambda_{\max}}{n} + \epsilon + \Lambda_{\max} \cdot (\text{right-hand side of (34)}). \end{aligned} \quad (69)$$

From the proof of Claim 1, the condition of Theorem 4 requires  $k = k_n$  and  $m = m_n$  to satisfy

$$\lim_{n \rightarrow \infty} k_n |\mathcal{Z}|^{2k_n} \left(\frac{m_n}{n}\right)^{1-\delta} = 0.$$

Therefore, if we set  $\epsilon^2 = \Theta(k_n |\mathcal{Z}|^{2k_n} \left(\frac{m_n}{n}\right)^{1-\delta})$  with sufficiently large constant then, from (68), we can see that the right-hand side of (34) will decay almost exponentially, which is much faster than  $\Theta(k_n |\mathcal{Z}|^{2k_n} \left(\frac{m_n}{n}\right)^{1-\delta})$ . Hence, from (69), we conclude that  $E(A_{k,m}^n) = O\left(\sqrt{k_n |\mathcal{Z}|^{2k_n} \left(\frac{m_n}{n}\right)^{1-\delta}}\right)$ , which results in part b). ■

## 8-G Proof of Theorem 5

The fact that  $m = \Theta(n)$  implies the existence of  $\alpha > 0$ , such that  $m \geq n\alpha$  for all sufficiently large  $n$ . Let  $\mathbf{X}$  be the process formed by concatenating i.i.d. blocks of length  $\lceil 1/\alpha \rceil$ , each block consisting of the same repeated symbol chosen uniformly from  $\mathcal{X}$ . The first observation to note is that, for all  $n$  large enough that  $m \geq n\alpha$ ,

$$D_{0,m}(X^n, Z^n) = 0 \quad a.s. \quad (70)$$

This is because, by construction,  $X^n$  is, with probability 1, piecewise constant with constancy sub-blocks of length, at least,  $\lceil 1/\alpha \rceil$ . Thus, a genie with access to  $X^n$  can choose a sequence of symbol-by-symbol

schemes (in fact, ignoring the noisy sequence), with less than  $n\alpha$  (and, therefore, less than  $m$ ) switches, that perfectly recover  $X^n$  (and, therefore, by our assumption on the loss function, suffers zero loss). On the other hand, the assumptions on the loss function and the channel imply that, for the process  $\mathbf{X}$  just constructed,

$$\limsup_{n \rightarrow \infty} \min_{\hat{\mathbf{X}}^n} EL_{\hat{\mathbf{X}}^n}(X^n, Z^n) > 0, \quad (71)$$

since even the Bayes-optimal scheme for this process incurs a positive loss, with a positive probability, on each  $\lceil 1/\alpha \rceil$  super-symbol. Thus, we get

$$E \left\{ \limsup_{n \rightarrow \infty} E [L_{\hat{\mathbf{X}}^n}(X^n, Z^n) - D_{0,m}(X^n, Z^n) | X^n] \right\} \quad (72)$$

$$\geq \limsup_{n \rightarrow \infty} E [L_{\hat{\mathbf{X}}^n}(X^n, Z^n) - D_{0,m}(X^n, Z^n)] \quad (73)$$

$$= \limsup_{n \rightarrow \infty} EL_{\hat{\mathbf{X}}^n}(X^n, Z^n) \quad (74)$$

$$\geq \limsup_{n \rightarrow \infty} \min_{\hat{\mathbf{X}}^n} EL_{\hat{\mathbf{X}}^n}(X^n, Z^n) > 0, \quad (75)$$

where (73) follows from Fatou's lemma; (74) follows from (70); and (75) follows from (71). In particular, there must be one particular individual sequence  $\mathbf{x} \in \mathcal{X}^\infty$  for which the expression inside the curled brackets of (72) is positive, i.e.,

$$\limsup_{n \rightarrow \infty} E [L_{\hat{\mathbf{X}}^n}(X^n, Z^n) - D_{0,m}(X^n, Z^n) | X^n = x^n] > 0, \quad (76)$$

which is equivalent to (37).  $\blacksquare$

## 8-H Proof of Theorem 6

First, by adding and subtracting the same terms, we obtain

$$\begin{aligned} & EL_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \\ = & \underbrace{EL_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(X^n, Z^n)}_{(i)} + \underbrace{\min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi})}_{(ii)}. \end{aligned} \quad (77)$$

We will consider term (i) and term (ii) separately. For term (i),

$$\begin{aligned} (i) &= EL_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(X^n, Z^n) \\ &\leq \frac{2k\Lambda_{\max}}{n} + \frac{n-2k}{n} \cdot \left[ EL_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X_{k+1}^{n-k}, Z^n) - \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(X_{k+1}^{n-k}, Z^n) \right] \end{aligned} \quad (78)$$

$$\leq \frac{2k\Lambda_{\max}}{n} + \frac{n-2k}{n} \cdot E \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X_{k+1}^{n-k}, Z^n) - \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} L_{\hat{\mathbf{X}}^{n,\mathbf{S}}}(X_{k+1}^{n-k}, Z^n) \right] \quad (79)$$

$$\leq \frac{2k\Lambda_{\max}}{n} + E \left[ L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X_{k+1}^{n-k}, Z^n) - D_{k,m}(X^n, Z^n) \right], \quad (80)$$

where (78) follows from upper bounding and omitting the losses for time instances  $t \leq k$  and  $t > n - k$  in the first and second terms of (i), respectively; (79) follows from exchanging the minimum with the expectation, and (80) follows from the definition (31) and  $\frac{n-2k}{n} \leq 1$ .

For term (ii), we bound the first term in (ii) as

$$\begin{aligned} & \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^n, \mathbf{S}}(X^n, Z^n) \\ & \leq \frac{2k(m+1)\Lambda_{\max}}{n} + \frac{1}{n} \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} E \left[ E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} \Lambda(X_j, s_{k,j}(Z_{j-k}^{j+k})) \middle| A^n \right] \right], \end{aligned} \quad (81)$$

by upper bounding the losses with  $\Lambda_{\max}$  on the boundary of the shifting points. Now, let  $\mathbf{P}_{X_j|Z_i^l, A^n} \in \mathbb{R}^{|\mathcal{X}|}$  denote the  $|\mathcal{X}|$ -dimensional probability vector whose  $x$ -th component is  $Pr(X_j = x|Z_i^l, A^n)$ . Then, we can bound the second term in (81) by the following chain of inequalities:

$$\frac{1}{n} \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} E \left[ E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} \Lambda(X_j, s_{k,j}(Z_{j-k}^{j+k})) \middle| A^n \right] \right] \quad (82)$$

$$= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} \min_{s_k \in \mathcal{S}_k} E \left[ \Lambda(X_j, s_k(Z_{j-k}^{j+k})) \middle| A^n \right] \right] \quad (83)$$

$$= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} \sum_{z_{-k}^k \in \mathcal{Z}^{2k+1}} P(Z_{j-k}^{j+k} = z_{-k}^k | A^n) \min_{\hat{x} \in \mathcal{X}} E \left[ \Lambda(X_j, \hat{x}) | Z_{j-k}^{j+k} = z_{-k}^k, A^n \right] \right] \quad (84)$$

$$= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} \sum_{z_{-k}^k \in \mathcal{Z}^{2k+1}} P(Z_{j-k}^{j+k} = z_{-k}^k | A^n) U_{\Lambda}(\mathbf{P}_{X_j|Z_{j-k}^{j+k}=z_{-k}^k, A^n}) \right] \quad (85)$$

$$= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} E[U_{\Lambda}(\mathbf{P}_{X_j|Z_{j-k}^{j+k}, A^n}) | A^n] \right]$$

$$= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+k+1}^{\tau_i-k} E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-k}^k}^{(A_{\tau_i})}) | A^n] \right] \quad (86)$$

$$\leq \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-k}^k}^{(A_{\tau_i})}) | A^n] \right], \quad (87)$$

where (83) follows from the stationarity of the distribution in each block as well as the fact that the combination of the best  $k$ -th order sliding window denoiser for each block is in  $\mathcal{S}_{k,m}^n$  and achieves the minimum in (82); (84) follows from conditioning; (85) follows from the definition (2); (86) follows from the stationarity of the distribution in each  $i$ -th block; and (87) follows from adding more nonnegative terms.

For the second term in (ii), we first define

$$n_i(A^n) \triangleq \tau_i(A^n) - \tau_{i-1}(A^n)$$

as the length of the  $i$ -th block, for  $1 \leq i \leq r(A^n) + 1$ . Obviously,  $n_i(A^n)$  also depends on  $A^n$ , and, thus, is a random variable, but we again suppress  $A^n$  for brevity and denote it as  $n_i$ . Then, similar to the first

term above, we obtain

$$\begin{aligned}
\mathbb{D}(P_{X^n}, \mathbf{\Pi}) &= \min_{\hat{\mathbf{X}}^n \in \mathcal{D}_n} EL_{\hat{\mathbf{X}}^n}(X^n, Z^n) \\
&= \frac{1}{n} \min_{\hat{\mathbf{X}}^n \in \mathcal{D}_n} E \left[ E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} \Lambda(X_j, \hat{X}_j(Z^n)) \middle| A^n \right] \right] \\
&= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} \min_{\hat{X}: Z^n \rightarrow \hat{X}} E \left[ \Lambda(X_j, \hat{X}(Z^n)) \middle| A^n \right] \right] \\
&= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} \min_{\hat{X}: Z^{\tau_i} \rightarrow \hat{X}} E \left[ \Lambda(X_j, \hat{X}(Z^{\tau_i})) \middle| A^n \right] \right] \tag{88} \\
&= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} E[U_\Lambda(\mathbf{P}_{X_j|Z^{\tau_i}, A^n}) \middle| A^n] \right] \\
&= \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} E[U_\Lambda(\mathbf{P}_{X_0|Z_{1-j}^{n_i}}^{(A_{\tau_i})}) \middle| A^n] \right] \tag{89} \\
&\geq \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} E[U_\Lambda(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(A_{\tau_i})}) \middle| A^n] \right], \tag{90}
\end{aligned}$$

where (88) follows from the conditional independence between different blocks, given  $A^n$ ; (89) follows from the stationarity of the distribution in each block, and (90) follows from [1, Lemma 4(1)]. Therefore, from (81), (87), and (90), we obtain

$$\begin{aligned}
(b) &= \min_{\mathbf{S} \in \mathcal{S}_{k,m}^n} EL_{\hat{\mathbf{X}}^n, \mathbf{S}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \\
&\leq \frac{2k(m+1)\Lambda_{\max}}{n} + \frac{1}{n} E \left[ \sum_{i=1}^{r+1} \sum_{j=\tau_{i-1}+1}^{\tau_i} E[U_\Lambda(\mathbf{P}_{X_0|Z_{-k}^k}^{(A_{\tau_i})}) \middle| A^n] - E[U_\Lambda(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(A_{\tau_i})}) \middle| A^n] \right] \\
&= \frac{2k(m+1)\Lambda_{\max}}{n} + E \left[ \sum_{i=1}^{r+1} \frac{n_i}{n} \cdot \left\{ E[U_\Lambda(\mathbf{P}_{X_0|Z_{-k}^k}^{(A_{\tau_i})}) \middle| A^n] - E[U_\Lambda(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(A_{\tau_i})}) \middle| A^n] \right\} \right]. \tag{91}
\end{aligned}$$

Now, observe that, regardless of  $A^n$ , the sequence of numbers  $\{\frac{n_i}{n}\}_{i=1}^{r+1}$  form a probability distribution, since  $\sum_{i=1}^{r+1} \frac{n_i}{n} = 1$  and  $\frac{n_i}{n} \geq 0$  for all  $i$ , with probability 1. Then, based on the fact that the average is less than the maximum, we obtain the further upper bound

$$(91) \leq \frac{2k(m+1)\Lambda_{\max}}{n} + E \left[ \max_{i \in \{1, \dots, M\}} \left\{ E[U_\Lambda(\mathbf{P}_{X_0|Z_{-k}^k}^{(i)})] - E[U_\Lambda(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(i)})] \right\} \right]. \tag{92}$$

The remaining argument to prove the theorem is to show that the upper bounds (80) and (92) converge to 0 as  $n$  tends to infinity. First, from the given condition on  $k = k_n$  and  $m = m_n$ , the maximal allowable growth rate for  $k$  is  $k = c_1 \log n$ , which leads to  $\lim_{n \rightarrow \infty} \frac{2k\Lambda_{\max}}{n} = 0$ . In addition, the condition requires  $m = o(n)$ , and  $k$  to be sufficiently slow, such that  $k|\mathcal{Z}|^{2k} = o\left(\left(\frac{n}{m}\right)^{1-\delta}\right)$ , which implies  $k = o\left(\frac{n}{m}\right)$ . Therefore,  $\lim_{n \rightarrow \infty} \frac{2k(m+1)\Lambda_{\max}}{n} = 0$ . Furthermore, from conditioning on  $X^n$ , bounded convergence theorem, and part

b) of Theorem 4, we obtain  $\lim_{n \rightarrow \infty} E[L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X_{k+1}^{n-k}, Z^n) - D_{k,m}(X^n, Z^n)] = 0$ . Thus, we have

$$\begin{aligned}
& \limsup_{n \rightarrow \infty} \left[ E L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi}) \right] \\
& \leq \limsup_{n \rightarrow \infty} E \left[ \max_{i \in \{1, \dots, M\}} \left\{ E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-k}^k}^{(i)})] - E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(i)})] \right\} \right] \\
& \leq E \left[ \limsup_{n \rightarrow \infty} \max_{i \in \{1, \dots, M\}} \left\{ E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-k}^k}^{(i)})] - E[U_{\Lambda}(\mathbf{P}_{X_0|Z_{-\infty}^{\infty}}^{(i)})] \right\} \right] \tag{93} \\
& = 0, \tag{94}
\end{aligned}$$

where (93) follows from the reverse Fatou's lemma, and (94) follows from [1, Lemma 4(2)] and  $M$  being finite. Since it is clear that  $\liminf_{n \rightarrow \infty} [E L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - \mathbb{D}(P_{X^n}, \mathbf{\Pi})] \geq 0$  by definition of  $\mathbb{D}(P_{X^n}, \mathbf{\Pi})$ , the theorem is proved. ■

*Remark:* As in [1, Theorem 3], the convergence rate in (44) may depend on  $P_{\mathbf{X}}$ , and there is no vanishing upper bound on this rate that holds for all  $P_{\mathbf{X}} \in \mathcal{P}\{m_n\}$ . However, we can glean some insight into the convergence rate from (i) and (ii): whereas the term (i) is uniformly upper bounded for all  $P_{\mathbf{X}} \in \mathcal{P}\{m_n\}$ ,<sup>8</sup> the rate at which term (ii) vanishes depends on  $P_{\mathbf{X}}$ . In general, we observe that the slower the rate of increase of  $k = k_n$ , the faster the convergence in (i), but the convergence in (ii) is slower. With respect to the rate of increase of  $m_n$ , the slower it is, the faster the convergence in (i), but whether or not the convergence in (ii) is accelerated by a slower rate of increase of  $m_n$  may depend on the underlying process distribution  $P_{\mathbf{X}}$ .

## References

- [1] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdú, and M. Weinberger, "Universal discrete denoising: Known channel," *IEEE Trans. Inform. Theory*, 51(1):5-28, Jan 2005
- [2] T. Weissman, E. Ordentlich, M. Weinberger, A. Somekh-Baruch, and N. Merhav, "Universal filtering via prediction," *IEEE Trans. Inform. Theory*, 53(4):1253-1264, Apr 2007
- [3] E. Ordentlich, M. Weinberger, and T. Weissman, "Multi-directional context sets with applications to universal denoising and compression," *Proceedings of IEEE Int. Symp. Inform. Theory*, 1270-1274, Sep 2005
- [4] J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE Trans. Inform. Theory*, 24(5):530-536, Sep 1978
- [5] N. Merhav and M. Feder, "Universal prediction," *IEEE Trans. Inform. Theory*, 44(6):2124-2147, Oct 1998
- [6] D. Blackwell, "Controlled random walks", *Proceedings International Congress of Mathematicians*, 3:336-338, 1956. Amsterdam: North Holland
- [7] D. Blackwell, "An analog of the minimax theorem for vector payoffs", *Pacific J. Math*, 6:1-8, 1956

---

<sup>8</sup>Recall part b) of Theorem 4, where a uniform bound (uniform in the underlying individual sequence) on  $E[L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(x^n, Z^n) - D_{k,m}(x^n, Z^n)]$  was provided in the semi-stochastic setting. Clearly, in the stochastic setting the same bound holds on  $E[L_{\hat{\mathbf{X}}_{\text{univ}}^{n,k,m}}(X^n, Z^n) - D_{k,m}(X^n, Z^n)]$ , regardless of the distribution of  $X^n$ .

- [8] J. Hannan, "Approximation to Bayes risk in repeated play," *Contributions to the Theory of Games*, vol.III:97-139, Princeton, NJ: Princeton Univ. Press, 1957
- [9] E. Ordentlich and T. Cover, "The cost of achieving the best portfolio in hindsight," *Mathematics of Operations Research*, 23(4):960-982, Nov 1998
- [10] N. Littlestone and M. K. Warmuth, "The weighted majority algorithm," *Information and Computation*, 108(2):212-261 Feb 1994.
- [11] V. Vovk, "Aggregating strategies," *Proc. 3rd Annu. Workshop on Computational Learning Theory*, 371-382, Morgan Kaufman, San Mateo, CA 1990
- [12] A. Gyorgy, T. Linder, and G. Lugosi, "Efficient algorithms and minimax bounds for zero-delay lossy source coding," *IEEE Trans. Signal Processing*, 52(8):2337-2347, Aug 2004
- [13] S. Matloub and T. Weissman, "Universal zero-delay joint source-channel coding," *IEEE Trans. Inform. Theory*, 52(12):5240-5250, Dec 2006
- [14] N. Cesa-Bianchi and G. Lugosi, "Prediction, learning, and games," *Cambridge University Press*, 2006
- [15] M. Herbster and M. K. Warmuth, "Tracking the best expert," in the special issue on context sensitivity and concept drift of the *Journal of Machine Learning*, Vol. 32(2):151-178, Aug 1998
- [16] O. Bousquet and M. K. Warmuth, "Tracking a small set of experts by mixing past posteriors," in a special issue of *Journal of Machine Learning Research* for COLT 2001, vol. 3:363-396, Nov 2002
- [17] G. I. Shamir and N. Merhav, "Low-complexity sequential lossless coding for piecewise-stationary memoryless sources," *IEEE Trans. Inform. Theory*, 45(5):1498-1519, 1999.
- [18] F.M.J. Willems, "Coding for binary independent piecewise identically distributed source," *IEEE Trans. Inform. Theory*, 42(6):2210-2217, 1996
- [19] S.S. Kozat and A.C. Singer, "Universal piecewise constant and least squares prediction," to appear in *IEEE Trans. on Signal Processing*, 2007
- [20] Y. Singer, "Switching portfolios," *Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence*, pp.1498-1519, 1998
- [21] G. I. Shamir and D. J. Costello, Jr., "On the redundancy of universal lossless coding for general piecewise stationary sources", *Communications in Information and Systems*, Vol. 1(3):305-332, Sep 2001
- [22] A. György, T. Linder, and G. Lugosi, "Tracking the best quantizer," *Proc. Int. Symp. Inform. Theory*, 1163- 1167, Sep 2005
- [23] D. Siegmund, "Confidence sets in change-point problems," *Int. Statist. Review* 56, 31-48, 1989
- [24] D. Siegmund and E.S. Venkatraman, "Using the generalized likelihood ratio statistic for sequential detection of a change-point," *Ann. Statist.* 23, 255-271, 1995
- [25] S. M. Oh, J. M. Rehg, and F. Dellaert, "Parameterized duration modeling for switching linear dynamic systems," *IEEE Int. Conf. on Computer Vision and Pattern Recognition (CVPR)*, NYC, 2006
- [26] B. Mesot and D. Barber, "Switching linear dynamical systems for noise robust speech recognition," *IEEE Trans. on Audio, Speech and Language Processing*, 15(6):1850-1858, Aug 2007

- [27] P. Lancaster and M. Tismenetsky, “The theory of matrices,” *Orlando, FL: Academic*, 1985
- [28] Y. Ephraim and N. Merhav, “Hidden Markov processes,” *IEEE Trans. Inform. Theory*, vol. 48:1518-1569, June 2002.
- [29] A. Lempel and J. Ziv, “Compression of two-dimensional data,” *IEEE Trans. Inform. Theory*, vol. 32(1):2-8, Jan 1986
- [30] A. Cohen, N. Merhav, and T. Weissman, “Scanning and sequential decision making for multi-dimensional data - Part I: the noiseless Case,” to appear in *IEEE Trans. Inform. Theory*.
- [31] A. Cohen, N. Merhav, and T. Weissman, “Scanning and sequential decision making for multi-dimensional data - Part II: the noisy case,” submitted to *IEEE Trans. Inform. Theory*, May 2007
- [32] K. Sivaramakrishnan and T. Weissman, “Universal denoising of discrete-time continuous-amplitude signals,” submitted to *IEEE Trans. Inform. Theory*, available at [http://www.stanford.edu/~tsachy/ieee\\_it\\_draft.pdf](http://www.stanford.edu/~tsachy/ieee_it_draft.pdf)
- [33] G. M. Gemelos, S. Sigurjonsson, T. Weissman, “Universal minimax discrete denoising under channel uncertainty,” *IEEE Trans. Inform. Theory*, 52(8):3476-3497, 2006
- [34] R. Zhang and T. Weissman, “Discrete denoising for channels with memory,” *Communications in Information and Systems (CIS)*, 5(2):257.288, 2005
- [35] A. Dembo and T. Weissman, “Universal denoising for the finite-input-general-output channel,” *IEEE Trans. Inform. Theory*, 51(4):1507-1517, April 2005
- [36] S. Jalali, S. Verdú and T. Weissman, “A universal Wyner-Ziv scheme for discrete sources,” *Proc. IEEE Int. Symp. Inform. Theory*, Nice, France, July 2007
- [37] E. Ordentlich, G. Seroussi, S. Verdú, and K. Viswanathan, “Universal algorithms for channel decoding of uncompressed sources,” Submitted to *IEEE Trans. Inform. Theory*, 2006
- [38] J. Yu and S. Verdú, “Schemes for bidirectional modeling of discrete stationary sources,” *IEEE Trans. Inform. Theory*, 52(11):4789.4807, 2006