# Universal Filtering Via Hidden Markov Modeling

Taesup Moon, *Student Member, IEEE*, and Tsachy Weissman, *Senior Member, IEEE*

*Abstract*—The problem of discrete universal filtering, in which the components of a discrete signal emitted by an unknown source and corrupted by a known discrete memoryless channel (DMC) are to be causally estimated, is considered. A family of filters are derived, and are shown to be universally asymptotically optimal in the sense of achieving the optimum filtering performance when the clean signal is stationary, ergodic, and satisfies an additional mild positivity condition. Our schemes are comprised of approximating the noisy signal using a hidden Markov process (HMP) via maximum-likelihood (ML) estimation, followed by the use of the forward recursions for HMP state estimation. It is shown that as the data length increases, and as the number of states in the HMP approximation increases, our family of filters attains the performance of the optimal distribution-dependent filter. An extension to the case of channels with memory is also established.

*Index Terms*—Finite alphabet, forward–backward recursion state estimation, hidden Markov process (HMP), maximum-likelihood (ML) parameter estimation, randomized scheme, stochastic setting, universal filtering.

## I. INTRODUCTION

**T**HE problem of estimating a discrete-time, finite-alphabet source signal $\{X_t\}_{t \in T}$ from the entire observation of a noisy signal $\{Z_t\}_{t \in T}$, which has been corrupted by a known discrete memoryless channel (DMC), has been thoroughly studied recently in [21]. It has been shown that even though the source distribution is unknown, an algorithm called DUDE can universally achieve the asymptotically optimal performance. This result has been extended in various directions such as the case of channel uncertainty [9], the case where the channel has memory [22], the case of nondiscrete noisy signal components [6], and the case where the reconstruction is required to depend causally on the noisy signal [18], [19]. In this paper, we revisit the last case, taking a different approach from [18], [19].

The case where we estimate $X_t$ causally based on observation of the noisy signal $Z^t = (Z_1, \ldots, Z_t)$, is referred to as *filtering*. The filter can be either deterministic or randomized (a concept that will be explained in detail later). In this paper, we will only focus on the *stochastic setting*, where we assume $\{X_t\}$ is a stationary and ergodic stochastic process. With the stochastic setting assumption, and under the same performance

criterion of [21], i.e., minimizing the expected normalized cumulative loss, knowledge of the conditional distribution of $X_t$ given $Z^t$ at each time $t$ is required to achieve the optimal performance. Also, by the same argument as in [21, Sec. III], this conditional distribution can be obtained by the conditional distribution of $Z_t$ given $Z^{t-1}$ when the invertible DMC is known. (We call a channel "invertible" if its transition probability matrix is of full row rank.)

However, for the *universal* filtering setting, where the probability distribution of the source is unknown, the conditional distribution of $Z_t$ given $Z^{t-1}$ is also not known and needs to be learned from the observed noisy signal. Therefore, if we can learn this conditional distribution accurately as the observation length increases, we can hope to build the universal filtering scheme that achieves the asymptotically optimal performance from the estimated conditional distribution. To pursue this goal, [18], [19] adopt the universal prediction [15] approach. That is, they first get an estimate of the conditional distribution of $Z_t$ given $Z^{t-1}$ by employing a universal predictor for the observed noisy signal, and then, by inverting the known DMC, obtain an estimate of the conditional distribution of $X_t$ given $Z^t$.

Unlike the approach of [18], [19], in this work, we turn our attention to the rich theory of hidden Markov process (HMP) models to directly obtain a different kind of estimate of the conditional distribution of $X_t$ given $Z^t$, without going through the channel inversion stage.[1]

Generally, HMPs are defined as a family of stochastic processes that are outputs of a memoryless channel whose inputs are finite-state Markov chains. As can be seen in [7], these HMP models arise in many areas, such as information theory, communications, statistics, learning, and speech recognition. Among these applications of HMPs, there are many situations where the state of the underlying Markov chain need be estimated based on the observed HMP. If the exact parameters of the HMP, namely, the state transition probability of the Markov chain and the channel transition density, as well as the order, the number of states, of the Markov chain are known, then this problem can be easily solved via well-known forward–backward recursions which were discovered by [4] and [2]. Especially, when we are estimating the state based on the causal observation of the HMP, we only need the forward recursion formula. In addition, much work has been done for the state estimation, where the order is known, but the parameters of the HMP are unknown. In this case, the parameters are first estimated via maximum-likelihood (ML) estimation or the expectation–minimization (EM) algorithm, then the state is estimated by using the estimated parameters in the recursion formula. A detailed explanation of this approach and the property of the ML parameter estimation can be found in [2], [3], [12], [8]. Furthermore, this was extended to the

---

[1]A part of this work was presented in [17].

case where the order of the Markov chain is also not known, but the upper bound on the order is known. In this case, the order estimation is first performed before the parameter and state estimation, and the above process is repeated. The references for the order estimation are given in [11], [13], [20]. There also has been work for the case where even the knowledge of the upper bound on the order of the Markov chain is not required [8], [23].

From these rich theories for the state, parameter, and order estimation of HMPs, we can see that it is possible to build a universal filtering scheme if the underlying source process is known to be a Markov process. That is, since the channel is memoryless and fixed in our setting, if our source $\{X_t\}$ is a Markov process, then obviously, $\{Z_t\}$ is an HMP, and we can first estimate the order[2] of the Markov process, then estimate the parameter, and finally perform forward recursion to learn the conditional distribution of $X_t$ given $Z^t$. From the consistency results of order estimation and parameter estimation, this conditional distribution will be an accurate estimate of the true one, and we can use it to build the universal filtering scheme.

Now, in our work, we extend this approach to the case where our source $\{X_t\}$ is a general stationary and ergodic process (with some benign conditions), which need not be a Markov process at all, and show that we can still build a universal filtering scheme that achieves asymptotically optimal performance. The skeleton of our scheme is the following: We first "model" our source as a Markov process with a certain order, or equivalently, model the noisy observed signal $\{Z_t\}$ as an HMP in a certain class. Then, we estimate the parameters of the HMP that "approximates" the noisy signal best in that class. We will show that from the consistency result about the ML parameter estimation for the mismatched model [8], these estimated parameters will give an accurate estimation of the conditional distribution of $X_t$ given $Z^t$, as the observation length increases and the HMP class gets richer. Then, this result will guarantee that our universal filter using this conditional distribution will attain the asymptotically optimal performance. In practice, this approach of HMP modeling has been heuristically employed in many applications, such as speech recognition [25], target tracking [26], and DNA sequence analysis [27], without theoretical justification. Additional samples of these practical applications can be found in [28]. Therefore, we focus on pursuing the rationale of the existing practical methodologies, and the main contribution of this work is in providing theoretical justification of the HMP modeling based approach to universal filtering.

The remainder of the paper is organized as follows. Section II introduces some notation and preliminaries that are needed for setting up the problem. In Section III, the universal filtering problem is defined explicitly. In Section IV, our universal filtering scheme is devised, the main theorem is stated, and proved. Section V extends our approach to the case where the channel has memory. Section VI gives discussions on our filter, and Section VII concludes the paper with some related future direc-

tions. Detailed technical proofs that are needed in the course of proving our main results are given in the Appendices.

## II. NOTATION AND PRELIMINARIES

### A. General Notation

We assume that the clean, noisy, and reconstruction signal components take their values in the same finite $M$-ary alphabet $\mathcal{A} = \{0, \ldots, M-1\}$. The simplex of $M$-dimensional column probability vectors will be denoted as $\mathcal{M}$.

The DMC is known to the filter and is denoted by its transition probability matrix $\mathbf{\Pi} = \{\Pi(i,j)\}_{i,j \in \mathcal{A}}$. Here, $\Pi(i,j)$ denotes the probability of channel output symbol $j$ when the input is $i$. We assume $\Pi(i,j) > 0 \, \forall i,j$, and let $\Pi_{\min} = \min_{i,j} \Pi(i,j)$. We assume this channel matrix is invertible and denote the inverse as $\mathbf{\Pi}^{-1}$. Let $\Pi_i^{-1}$ denote the $i$th column of $\mathbf{\Pi}^{-1}$. We also assume a given loss function (fidelity criterion) $\Lambda : \mathcal{A}^2 \to [0, \infty)$, represented by the loss matrix $\mathbf{\Lambda} = \{\Lambda(i,j)\}_{i,j \in \mathcal{A}}$, where $\Lambda(i,j)$ denotes the loss incurred when estimating the symbol $i$ with the symbol $j$. The maximum single-letter loss will be denoted by $\Lambda_{\max} = \max_{i,j \in \mathcal{A}} \Lambda(i,j)$, and $\boldsymbol{\lambda}_j$ will denote the $j$th column of $\mathbf{\Lambda}$.

As in [21], we define the extended Bayes response associated with the loss matrix $\mathbf{\Lambda}$ to any column vector $\boldsymbol{V} \in \mathbb{R}^M$ as

$$B(\boldsymbol{V}) = \arg\min_{\hat{x} \in \mathcal{A}} \lambda_{\hat{x}}^T \boldsymbol{V}$$

where $\arg\min_{\hat{x} \in \mathcal{A}}$ denotes the minimizing argument, resolving ties by taking the letter in the alphabet with the lowest index.

We let $P$ denote the true joint probability law of the clean and noisy signal, and $E(\cdot)$ denote expectation with respect to $P$. Throughout the paper, every almost sure convergence is with respect to $P$, and any equalities or inequalities between random variables should be understood in almost sure sense. If we need to refer to the probability law of clean or noisy signal induced by $P$, we denote $P_X$ and $P_Z$, respectively. If $P$ is written in a bold face, $\boldsymbol{P}$, with a subscript, it stands for a simplex vector in $\mathcal{M}$ for the corresponding distribution of the subscript. For example, $\boldsymbol{P}_{X_t|z^t}$ is a column $M$-vector whose $i$th component is $P(X_t = i | Z^t = z^t)$.

When we have some other probability law denoted as $Q$, and want to measure its difference from $P$, a natural choice of such a measure is the relative entropy rate. First, denote the $n$th-order relative entropy between $P$ and $Q$ as

$$D_n(P\|Q) = \sum_{z^n} P(z^n) \log \frac{P(z^n)}{Q(z^n)} = E\left(\log \frac{P(Z^n)}{Q(Z^n)}\right).$$

Then, the relative entropy rate (also known as Kullback–Leibler divergence rate) is defined as

$$\boldsymbol{D}(P\|Q) \triangleq \lim_{n \to \infty} \frac{1}{n} D_n(P\|Q)$$

if the limit exists. When $Q$ is a probability law in a certain class of HMPs, this limit always exists and the relative entropy rate is well defined. A more detailed discussion about this limit will be given in Lemma 2. This relative entropy rate will play a central role in analyzing our universal filtering scheme.

### B. Hidden Markov Processes (HMPs)

*1) Definition:* As stated in the Introduction, the HMPs are defined as a family of stochastic processes that are outputs of

---

[2]We slightly abuse the term "order" here. Generally, the order of a "finite-state Markov chain" stands for the number of states, but we also refer by the "order of a Markov process" to the length of the memory of the process. Hence, once we know the order of a Markov process, we also know the order of the associated finite-state Markov chain induced from the Markov process.

a memoryless channel whose inputs are finite-state Markov chains. Let us denote a general HMP as $\{Y_t\}$ and the underlying finite-state Markov chain as $\{S_t\}$. The corresponding alphabet sizes of each component are denoted as $|\mathcal{Y}|$ and $|\mathcal{S}|$, respectively. Then, there are three parameters that determine the probability laws of $\{Y_t\}$: $\pi \in \mathbb{R}^{1 \times |\mathcal{S}|}$, the initial distribution of finite-state Markov chain; $A \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, the probability transition matrix of finite-state Markov chain, and $C \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{Y}|}$, the probability transition matrix of the memoryless channel. The triplet $\{\pi, A, C\}$ is referred to as the parameter of HMP. Let $\Theta$ be a set of all $\theta$'s where $\theta := \{\pi_\theta, A_\theta, C_\theta\}$. For each $\theta$ and each realization $y^n$, we can calculate the likelihood function

$$Q_\theta(y^n) = \pi_\theta \prod_{t=1}^{n} (\hat{\boldsymbol{C}}_{\theta,t} A_\theta) \mathbf{1}$$

where $\hat{\boldsymbol{C}}_{\theta,t}$ is $|\mathcal{S}| \times |\mathcal{S}|$ diagonal matrix whose $(j,j)$th entry is the $(j, y_t)$th entry of $C_\theta$, and $\mathbf{1}$ is the $|\mathcal{S}| \times 1$ vector with all entries equal to 1.

Now, suppose $\{Z_t\}$ is an output of a DMC $\boldsymbol{\Pi}$, when the input is a stationary Markov process $\{X_t\}$. For simplicity, we assume that the alphabet of each component of $\{Z_t\}$ and $\{X_t\}$ are finite and equal, i.e., $\mathcal{Z} = \mathcal{X} = \mathcal{A}$. When the order of the underlying Markov process is $k$, we can associate the state $S_t$ of underlying finite-state Markov chain with $X_{t-k}^{t-1}$, which has alphabet size $M^k$. Then, clearly, $\{Z_t\}$ is also a stationary HMP, and the parameter set of such HMP is denoted as $\Theta_k \subset \Theta$. Note that when $\theta \in \Theta_k$, $\pi_\theta \in \mathbb{R}^{1 \times M^k}$, $A_\theta \in \mathbb{R}^{M^k \times M^k}$, and $C_\theta \in \mathbb{R}^{M^k \times M}$. Furthermore, for some $\delta > 0$, we define a set $\Theta_k^\delta \subset \Theta_k$ as the set of $\theta \in \Theta_k$ that has following properties: for the $i$th $k$-tuple state $x_1^k(i)$ and the $j$th $k$-tuple state $x_1^k(j)$

- $a_{ij,\theta} \geq \delta$, if $x_2^k(i)$ is equal to $x_1^{k-1}(j)$;
- $a_{ij,\theta} = 0$, otherwise;
- $c_{ir,\theta} = \boldsymbol{\Pi}(x_k(i), r)$, for all $1 \leq i \leq M^k$ and $1 \leq r \leq M$;

where $a_{ij,\theta}$ is the $(i,j)$th entry of $A_\theta$, and $c_{ir,\theta}$ is the $(i,r)$th entry of $C_\theta$. In particular, if $\theta \in \Theta_k^\delta$ then: 1) the stochastic matrix $A_\theta$ is irreducible and aperiodic; thus, since the Markov chain is stationary, $\pi_\theta$ is the stationary distribution of the Markov chain, and is uniquely determined from $A_\theta$, 2) $C_\theta$ is the same for all $\theta$, and, therefore, $\theta$ is completely specified by $A_\theta$. For notational brevity, we omit the subscript $\theta$ and denote the probability law $Q \in \Theta_k^\delta$, if $Q = Q_\theta$, and $\theta \in \Theta_k^\delta$.

*2) Maximum-Likelihood (ML) Estimation:* Generally, suppose a probability law $Q$ is in a certain class $\Omega$. Then, the $n$th-order ML estimator in $\Omega$ for the observed sequence $z^n$, is defined as

$$\hat{Q}[z^n] = \arg \max_{Q \in \Omega} Q(z^n)$$

resolving ties arbitrarily. Now, if $Q \in \Theta_k^\delta$, then there is an algorithm called *expectation–maximization* (EM) [4] that iteratively updates the parameter estimates to maximize the likelihood. Thus, when $Q$ is in the class of probability laws of a HMP, the ML estimate can be efficiently attained.[3] We denote the ML estimator in $\Theta_k^\delta$ based on $z^n$ by

$$\hat{Q}_{k,\delta}[z^n] = \arg \max_{Q \in \Theta_k^\delta} Q(z^n).$$

---

[3]We neglect issues of convergence of the EM algorithm and assume that the ML estimation is performed perfectly.

Obviously, when the $n$-tuple $Z^n$ is random, $\hat{Q}_{k,\delta}[Z^n]$ is also a random probability law that is a function of $Z^n$.

*3) Consistency of ML Estimator:* When $P_Z \in \Theta_k^\delta$, an ML estimator $\hat{Q}_{k,\delta}[Z^n]$ is said to be *strongly consistent* if

$$\lim_{n \to \infty} \hat{Q}_{k,\delta}[Z^n] = P_Z \quad \text{a.s.}$$

The strong consistency of the ML estimator $\hat{Q}_{k,\delta}[Z^n]$ of the parameter of a finite-alphabet stationary ergodic HMP was proved in [1]. For the case of a general stationary ergodic HMP, the strong consistency was proved in [12].

We also have a sense of strong consistency for the case where $P_Z$ is a general stationary and ergodic process. By the similar argument as in [8, Theorem 2.2.1], we have the consistency in the sense that if the observed noisy signal is not necessarily an HMP, and we still perform the ML estimation in $\Theta_k^\delta$, then we get

$$\lim_{n \to \infty} \hat{Q}_{k,\delta}[Z^n] \in \mathcal{N} \quad \text{a.s.} \qquad (1)$$

where[4]

$$\mathcal{N} \triangleq \{Q \in \Theta_k^\delta : \boldsymbol{D}(P \| Q) = \min_{Q' \in \Theta_k^\delta} \boldsymbol{D}(P \| Q')\}.$$

This second consistency result is the key result that we will use in devising and analyzing our universal filtering scheme.

## III. THE UNIVERSAL FILTERING PROBLEM

As mentioned in the Introduction, we will assume a stochastic setting, that is, the underlying clean signal is an output of some stationary and ergodic process whose probability law is $P_X$. From $P_X$ and $\boldsymbol{\Pi}$, we can get the true joint probability law $P$ and the corresponding probability law of noisy observed signal, $P_Z$. That is

$$P(X^n = x^n, Z^n = z^n) = P_X(X^n = x^n) \prod_{t=1}^{n} \Pi(x_i, z_i)$$

and

$$P_Z(Z^n = z^n) = \sum_{x^n} P(X^n = x^n, Z^n = z^n).$$

A *filter* is a sequence of probability distributions $\hat{\boldsymbol{X}} = \{\hat{X}_t\}$, where $\hat{X}_t : \mathcal{A}^t \to \mathcal{M}$. The interpretation is that, upon observing $z^t$, the reconstruction for the underlying, unobserved $x_t$ is represented by the symbol $\hat{x}$ with probability $\hat{X}_t(z^t)[\hat{x}]$. The notation $\hat{X}_t(z^t)[\hat{x}]$ stands for the $\hat{x}$th element of a vector $\hat{X}_t(z^t)$. A filter is called *deterministic* if $\hat{X}_t(z^t)$ is a standard basis vector in $\mathbb{R}^M$ for all $t$ and $z^t$, and *randomized* if $\hat{X}_t(z^t)$ can be any vectors in $\mathcal{M}$ other than standard basis vectors for some $t$ and $z^t$. The normalized cumulative loss of the scheme $\hat{\boldsymbol{X}}$ on the individual pair $(x^n, z^n)$ is defined by

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

---

[4]Just as in [8, Theorem 2.2.1], the notion of a.s. set convergence is used. For any subset $\mathcal{E} \subset \Theta$, define $\mathcal{E}_\epsilon \triangleq \{Q \in \Theta : d(Q, \mathcal{E}) < \epsilon\}$, where $d$ is the Euclidean distance. Then, $\lim_{n \to \infty} \hat{Q}[Z^n] \in \mathcal{E}$ a.s. if for all $\epsilon > 0$, $\exists N(\epsilon, \omega)$ such that for all $n \geq N(\epsilon, \omega)$, $\hat{Q}[Z^n] \in \mathcal{E}_\epsilon$.

where

$$\ell(x_t, \hat{X}_t(z^t)) = \sum_{\hat{x} \in \hat{\mathcal{X}}} \Lambda(x_t, \hat{x}) \hat{X}_t(z^t)[\hat{x}].$$

Then, the goal of a filter is to minimize the expected normalized cumulative loss $E\big(L_{\hat{\boldsymbol{X}}}(X^n, Z^n)\big)$.

The optimal performance of the $n$th-order filter is defined as

$$\phi_n(P_X, \boldsymbol{\Pi}) = \min_{\hat{\boldsymbol{X}} \in \mathcal{F}} E\Big(L_{\hat{\boldsymbol{X}}}(X^n, Z^n)\Big)$$

where $\mathcal{F}$ denotes the class of all filters. Subadditivity arguments similar to those in [21] imply

$$\lim_{n \to \infty} \phi_n(P_X, \boldsymbol{\Pi}) = \inf_{n \geq 1} \phi_n(P_X, \boldsymbol{\Pi}) \triangleq \boldsymbol{\Phi}(P_X, \boldsymbol{\Pi}).$$

By definition, $\boldsymbol{\Phi}(P_X, \boldsymbol{\Pi})$ is the (distribution-dependent) optimal asymptotic filtering performance attainable when the clean signal is generated by the law $P_X$ and corrupted by $\boldsymbol{\Pi}$. This $\boldsymbol{\Phi}(P_X, \boldsymbol{\Pi})$ can be achieved by the optimal filter $\hat{\boldsymbol{X}}_P = \{\hat{X}_{P,t}\}$ where

$$\hat{X}_{P,t}(z^t)[\hat{x}] = \Pr(B(\boldsymbol{P}_{X_t|z^t}) = \hat{x}).$$

For brevity of notation, we denote $\hat{X}_P(z^t) = \hat{X}_{P,t}(z^t)$. Note that this is a *deterministic* filter, i.e., for a given $z^t$, the filter is a standard basis vector in $\mathbb{R}^M$ for all $t$. We can easily see that this filter is optimal since it minimizes $E(\ell(X_t, \hat{X}(Z^t))$ for all $t$, and thus, it minimizes $E\big(L_{\hat{\boldsymbol{X}}}(X^n, Z^n)\big)$ for all $n$.

As can be seen, $\hat{X}_P(z^t)$ needs the exact knowledge of $\boldsymbol{P}_{X_t|z^t}$, and thus, is dependent on the distribution of the underlying clean signal. The universal filtering problem is to construct (possibly a sequence of) filter(s), $\hat{\boldsymbol{X}}_{\text{univ}}$, that is independent of the distribution of underlying clean signal $P_X$, and yet asymptotically achieving $\boldsymbol{\Phi}(P_X, \boldsymbol{\Pi})$. We describe our sequence of universal filters in the next section.

## IV. Universal Filtering Based on Hidden Markov Modeling

### A. Description of the Filter

Before describing our sequence of universal filters, we make the following assumption on the source.

*Assumption 1:* There exists a sequence of positive reals $\{\delta_k\}$, such that $\delta_k \downarrow 0$ as $k \to \infty$, and $P_X$ satisfies

$$P_X(X_0|X_{-k}^{-1}) \geq \delta_k \quad \text{a.s.} \quad \forall k \in \mathbb{N}. \tag{2}$$

For any probability law $Q$, we construct a *randomized* filter as follows: For $\epsilon > 0$, denote $L_2$ $\epsilon$-ball in $\mathbb{R}^M$ as $B_\epsilon = \{\boldsymbol{V} \in \mathbb{R}^M : \|\boldsymbol{V}\|_2 \leq \epsilon\}$. Then, we define a filter for fixed $\epsilon$ as

$$\hat{X}_{Q,t}^\epsilon(z^t)[\hat{x}] = \Pr(B(\boldsymbol{Q}_{X_t|z^t} + \boldsymbol{U}) = \hat{x}) \tag{3}$$

where $\boldsymbol{U} \in \mathbb{R}^M$ is a random vector, uniformly distributed in $B_\epsilon$. For brevity of notation, we denote $\hat{X}_Q^\epsilon(z^t) = \hat{X}_{Q,t}^\epsilon(z^t)$. This filter is randomized since depending on $Q$ and $z^t$, $\hat{X}_Q^\epsilon(z^t)$ can be a probability simplex vector in $\mathcal{M}$ that is not a standard

basis vector. The reason we needed this randomization will be explained in proving Lemma 3.

To devise our filter, let us first consider an increasing sequence of positive integers $\{m_i\}_{i \geq 1}$ that satisfies following conditions:

$$\lim_{i \to \infty} \frac{m_{i-1}}{m_i} = 0, \quad \lim_{i \to \infty} m_i = \infty. \tag{4}$$

Now, define

$$i(t) \triangleq \max\{i : m_i \leq t\}.$$

Then, given that our source distribution satisfies (2), and for fixed $k$, define a random probability law

$$Q_k^t \triangleq \hat{Q}_{k,\delta_k}[Z^{m_{i(t)}}] = \arg\max_{Q \in \Theta_k^{\delta_k}} Q(Z^{m_{i(t)}}). \tag{5}$$

That is, $Q_k^t$ is the ML estimator in $\Theta_k^{\delta_k}$ based on $Z^{m_{i(t)}}$. As discussed in Section II-B.1, we only need to estimate the state transition probabilities of the underlying Markov chain to obtain this ML estimator, and this can be done by the EM algorithm. Performing the EM algorithm requires iterative forward–backward recursions, and is the most expensive part of our scheme in terms of the complexity. However, since the recursions are efficiently implemented by linear complexity dynamic programming (a.k.a. the Bahl–Cocke–Jelinek–Raviv (BCJR) algorithm), which is described in detail in [4], the overall complexity of our scheme is still linear in data length. Once we get $Q_k^t$, we can then calculate $\boldsymbol{Q}_{k X_t|z^t}^t$, which stands for the simplex vector in $\mathcal{M}$ whose $i$th component is $Q_k^t(X_t = i|Z^t = z^t)$. This vector can be obtained again from using the forward-recursion formula. Note that we get this conditional distribution directly, not by first estimating the output distribution, and then inverting the channel, as was done in [18], [19], [21].

Finally, we take as our sequence of universal filtering schemes, indexed by $k$ and $\epsilon$

$$\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon = \{\hat{X}_{Q_k^t,t}^\epsilon\}.$$

The following theorem states the main result of this paper.

*Theorem 1:* Let $\boldsymbol{X}^\infty \in \mathcal{A}^\infty$ be a stationary, ergodic process emitted by the source $P_X$ which satisfies Assumption 1. Let $\boldsymbol{Z}^\infty \in \mathcal{A}^\infty$ be the output of the DMC, $\boldsymbol{\Pi}$, whose input is $\boldsymbol{X}^\infty$. Then

a) $\displaystyle \lim_{\epsilon \to 0} \lim_{k \to \infty} \limsup_{n \to \infty} L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) \leq \boldsymbol{\Phi}(P_X, \boldsymbol{\Pi})$ a.s.

b) $\displaystyle \lim_{\epsilon \to 0} \lim_{k \to \infty} \limsup_{n \to \infty} E\big(L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n)\big) = \boldsymbol{\Phi}(P_X, \boldsymbol{\Pi}).$

*Remark:* In defining our universal filter $\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon$, one might intuitively think that it would be better to use the ML estimator that updates every point of time, i.e., to define $Q_k^t = \hat{Q}_{k,\delta_k}[Z^t]$. However, our definition of $Q_k^t$, namely, using the same ML estimator throughout each block, is crucial for our proof of the above theorem, especially for proving Corollary 1 that follows below. Besides this technical reason, updating the ML estimator every time would require higher complexity than our scheme requires.

## B. Intuition Behind the Scheme and Proof Sketch

The intuition behind our scheme parallels that of the universal compression and universal prediction problems in the stochastic setting. In the $n$th-order problem of both cases [5], [14], the excess expected codeword length per symbol, and the excess expected normalized cumulative loss incurred by using the wrong probability law $Q$ in place of the true probability law $P$ could be upper-bounded by the normalized $n$th-order relative entropy $\frac{1}{n}D_n(P\|Q)$. Then, to achieve the asymptotically optimum performance, the compressor and the predictor try to find and use some data-dependent $Q$ that makes $\frac{1}{n}D_n(P\|Q) \to 0$ as $n \to \infty$, that is, makes $\boldsymbol{D}(P\|Q)$ zero.

We follow the same intuition in our universal filtering problem. For fixed $k$ and $\epsilon$, our scheme, as can be seen from (5), divides the noisy observed signal into subblocks of length $(m_i - m_{i-1})$. Since $\frac{m_{i-1}}{m_i}$ tends to zero as $i \to \infty$, the length of each subblock grows faster than exponential. Now, to filter each subblock, it plugs the ML estimator in $\Theta_k^{\delta_k}$ obtained from the entire observation of noisy signal up to the previous subblock. From (1), we know that as the observation length $n$ increases, this ML estimator will converge to the parameter that minimizes the relative entropy rate between the true output probability law $P_Z$. Then, to show that this scheme achieves the asymptotically optimum performance, we bound the excess expected normalized cumulative loss with this relative entropy rate, and show that the bound goes to zero as the HMP parameter set becomes richer, that is, $k$ increases.

To be more specific, we briefly sketch the proof of our main theorem. Part b) of Theorem 1 states that our scheme is asymptotically optimal. As described in the proof of Theorem 1, it is not hard to show that Part b) follows directly from Part a) and Fatou's lemma. Therefore, proving Part a) is the key in proving the theorem. Part a) states that in the limit, the normalized cumulative loss of our scheme, for almost every realization, is less than or equal to the asymptotically optimum performance.

To prove Part a), we first fix $k$ and $\epsilon$, and get the following inequality:

$$\limsup_{n\to\infty} \left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \right)$$
$$\leq F\left( \limsup_{t\to\infty} \boldsymbol{D}(P_Z\|Q_k^t), \epsilon \right) \quad \text{a.s.} \quad (6)$$

where $F(x, y)$ is some function such that $F(x, y) \to 0$ as $x \downarrow 0$, and then $y \downarrow 0$.[5] There are two keys in getting this inequality. The first one is to show the concentration of $L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n)$ to its expectation which will be shown in Lemma 3 and Corollary 1. The second is to get the explicit upper bound function $F(x, y)$ which will be based on Lemma 4. Once establishing this inequality, we show that

$$\lim_{k\to\infty} \limsup_{t\to\infty} \boldsymbol{D}(P_Z\|Q_k^t) = 0 \quad \text{a.s.} \quad (7)$$

from Lemma 5 and then send $\epsilon \downarrow 0$ to get Part a). Keeping this proof sketch in mind, let us move on to the detailed proof in the next section.

[5]Note that $Q_k^t$ in $\boldsymbol{D}(P_Z\|Q_k^t)$ is a function of $Z^{m_i(t)}$, and thus, is random. A more formal definition of relative entropy rate between true and the random probability law like this case will be given after Lemma 4.

## C. Proof of Theorem 1

Before proving the theorem, we introduce several lemmas as building blocks. Lemmas 1 and 2 below give some general results for the HMPs that we are considering. Our lemmas are similar to [8, Lemma 2.3.4] and [8, Theorem 2.3.3]. The latter assumed that all the parameters are lower-bounded by $\delta > 0$, whereas in $\Theta_k^\delta$, some parameters can be zero. We take this into account in proving Lemmas 1 and 2. Lemma 3 shows the uniform concentration property of the normalized cumulative loss on $\Theta_k^\delta$, which is an important property that we need to prove the main theorem. Lemma 4 provides a key step to get the upper bound described in (6), and Lemma 5, which needs three additional definitions, enables to show (7). After building up the lemmas, we give the proof of the main theorem, which is merely an application of the lemmas.

*Lemma 1:* Suppose $Q \in \Theta_k^\delta$ and fix $\delta > 0$. Then, for all $\omega$, $Q(Z_0|Z_{-t}^{-1})$ converges to a limit $Q(Z_0|Z_{-\infty}^{-1})$ uniformly on $\Theta_k^\delta$.

*Proof:* To prove this lemma, we need three more lemmas in Appendix A, which are variations on those found in [1]. Let us denote $f_t := Q(Z_0|Z_{-t}^{-1})$ and $f_0 = 0$. Then, the sequence $\{f_t\}$ uniformly converges on $\Theta_k^\delta$, if the following $k$ subsequences:

$$\{f_{jk+l}, j = 0, 1, 2, \ldots, \}, \qquad 0 \leq l \leq k-1$$

uniformly converge on $\Theta_k^\delta$, and have the same limit.

First, the uniform convergence of each subsequence $\{f_{jk+l}\}$ can be shown by showing the series $\sum_{j=0}^t (f_{(j+1)k+l} - f_{jk+l})$ converges uniformly. From Lemma 8 in Appendix A, and setting $m = k$

$$\sum_{j=0}^t |f_{(j+1)k+l} - f_{jk+l}|$$

$$= \sum_{x_0} Q(Z_0|x_0) \sum_{j=1}^t |Q(x_0|Z_{-(j+1)k-l}^{-1}) - Q(x_0|Z_{-jk-l}^{-1})|$$

$$\leq M \sum_{j=1}^t (\rho_{\delta,k,k})^{j+1}$$

where $\rho_{\delta,k,k} < 1, M < \infty$ and $\rho_{\delta,k,k}$ does not depend on $Q, \omega$, and $l$. Therefore, the series $\sum_{j=0}^t (f_{(j+1)k+l} - f_{jk+l})$ converges absolutely regardless of $Q$, and, hence, we conclude that each subsequences $\{f_{jk+l}\}$ converges uniformly on $\Theta_k^\delta$.

Now, to show that the $k$ subsequences have the same limit, construct another subsequence, $\{f_{j(k+1)+1}, j = 0, 1, 2, \ldots, \}$. Since this subsequence contains infinitely many terms from all $k$ subsequences, if this subsequence converges uniformly on $\Theta_k^\delta$, we can conclude that the $k$ subsequences have the same limit. The derivation of the uniform convergence of this subsequence is the same as that described above, but setting $m = k + 1$ in Lemma 8. Therefore, the original sequence $\{f_t\}$ converges to its limit uniformly on $\Theta_k^\delta$. $\square$

The remarkable fact of this lemma is that the convergence is not only uniform on $\Theta_k^\delta$, but also in $\omega$. That is, the convergence holds uniformly on every realization of $z_{-\infty}^0$.

*Lemma 2:* For the distribution of the observed noisy process $\{Z_t\}$, $P_Z$, and every $Q \in \Theta_k^\delta$

$$\boldsymbol{D}(P_Z\|Q) \triangleq \lim_{n\to\infty} \frac{1}{n} D_n(P_Z\|Q) = E\left( \log \frac{P_Z(Z_0|Z_{-\infty}^{-1})}{Q(Z_0|Z_{-\infty}^{-1})} \right).$$

Moreover, uniformly on $\Theta_k^\delta$

$$\lim_{n \to \infty} \frac{1}{n} \log \frac{P_Z(Z^n)}{Q(Z^n)} = \boldsymbol{D}(P_Z \| Q) \qquad \text{a.s.}$$

*Proof:* This lemma consists of three parts. The first part is to show the existence of the first limit in the lemma so that the definition of $\boldsymbol{D}(P_Z \| Q)$ is valid. The second part is to show that the value of the limit is indeed $E\left( \log \frac{P_Z(Z_0 | Z_{-\infty}^{-1})}{Q(Z_0 | Z_{-\infty}^{-1})} \right)$. Finally, the last part is to show the uniform convergence of normalized log-likelihood ratio to the relative entropy rate. The first two parts and the pointwise convergence of the third part is a generalization of the Shannon–McMillan–Breiman theorem. The proof of these parts is identical to those in [8, Theorem 2.3.3] even for the case where some parameters in $\Theta_k^\delta$ can be zero.

The uniform convergence in the third part of the lemma is crucial in that it enables to obtain the second consistency result (1) as in [8, Theorem 2.2.1]. We take into account our parameter set, and repeat the argument of [8, Lemma 2.4.1]. To show the uniform convergence, we need to show

$$\lim_{n \to \infty} \frac{1}{n} \log Q(Z^n) = E\left( \log Q(Z_0 | Z_{-\infty}^0) \right) \quad \text{a.s.}$$

uniformly on $\Theta_k^\delta$. Since the pointwise convergence can be shown and the parameter set $\Theta_k^\delta$ is compact, it is enough to show that $\frac{1}{n} \log Q(Z^n)$ is an equicontinuous sequence by Ascoli's theorem. That is, we need to show that for all $\epsilon > 0$, there exists $\delta(\epsilon) > 0$ such that if $\|Q - Q'\|_1 < \delta(\epsilon)$, then

$$\left| \frac{1}{n} \log Q(Z^n) - \frac{1}{n} \log Q'(Z^n) \right| \le \epsilon, \quad \text{for all } n \qquad (8)$$

where

$$\|Q - Q'\|_1 \triangleq \sum_{i,j} |a_{ij} - a'_{ij}|$$

is defined to be the $L_1$ distance between the two parameters defining $Q$ and $Q'$. This equicontinuity can be proved by observing that a process $\{S_t = (X_{t-(k-1)}^t, Z_t)\}$ is a Markov process under any $Q \in \Theta_k^\delta$, where $\{S_t\}$ has a state space $\mathcal{S} = \mathcal{A}^k \times \mathcal{A}$. This is true since

$$
\begin{aligned}
Q(S_{t+1} | S^t) &= Q(X_{t+1}, Z_{t+1} | X^t, Z^t) \\
&= Q(X_{t+1} | X^t, Z^t) Q(Z_{t+1} | X^{t+1}, Z^t) \\
&= Q(X_{t+1} | X_{t-(k-1)}^t) \Pi(X_{t+1}, Z_{t+1}) \\
&= Q(S_{t+1} | S_t).
\end{aligned}
$$

Let $\{x_1^k(i) : i = 1, \ldots, M^k\}$ denote the set of all possible $k$-tuples of $\{X_t\}$, and let $s = (x_1^k(i), z)$, $\bar{s} = (x_1^k(j), \bar{z})$. Then, the transition matrix $T$ of $\{S_t\}$ has elements $t_{s\bar{s}} \triangleq Q(S_{t+1} = \bar{s} | S_t = s) = a_{ij} \Pi(x_k(j), \bar{z})$. Since all $A$ that are in $\Theta_k^\delta$ are irreducible and aperiodic, and $\Pi(x_k(j), \bar{z}) > 0$ for all $x_k(j)$ and $\bar{z}$, $T$ is also irreducible and aperiodic. Hence, $T$ has the unique stationary distribution $\tau$. In addition, from Assumption 1, we can observe that for all $Q \in \Theta_k^\delta$, $Q(S^n) > 0$ a.s. (with respect to $P$).

Since $\{S_t\}$ is also stationary, we have

$$Q(S^n) = \tau_{S_1} \prod_{t=k}^{n-1} t_{S_t S_{t+1}} = \tau_{S_1} \prod_{(s,\bar{s})} t_{s\bar{s}}^{N_{s\bar{s}}},$$

where

$$N_{s\bar{s}} \triangleq \sum_{t=k}^{n-1} \mathbf{1}(S_t = s, S_{t+1} = \bar{s}).$$

For another probability law $Q' \in \Theta_k^\delta$, we have

$$
\begin{aligned}
&\left| \frac{1}{n} \log Q(S^n) - \frac{1}{n} \log Q'(S^n) \right| \\
&\le \left| \frac{1}{n} \log \tau_{S_1} - \frac{1}{n} \log \tau'_{S_1} \right| + \left| \frac{1}{n} \sum_{(s,\bar{s})} N_{s\bar{s}} \log t_{s\bar{s}} - \frac{1}{n} \sum_{(s,\bar{s})} N_{s\bar{s}} \log t'_{s\bar{s}} \right| \\
&\le \left| \log \tau_{S_1} - \log \tau'_{S_1} \right| + \sum_{(s,\bar{s})} \left| \log t_{s\bar{s}} - \log t'_{s\bar{s}} \right| \qquad (9) \\
&= \left| \log \tau_{S_1} - \log \tau'_{S_1} \right| + \sum_{(i,j)} \left| \log a_{ij} - \log a'_{ij} \right| \qquad (10)
\end{aligned}
$$

where (9) follows from the fact that $\frac{1}{n} \le 1$, $\frac{N_{s\bar{s}}}{n} \le 1$, with probability 1, and (10) follows from the fact that DMC $\boldsymbol{\Pi}$ is equal for $Q$ and $Q'$. The summations are over the pairs that have nonzero transition probabilities.

Since the function $f(x) = \log x$ is a uniformly continuous function for $\delta \le x < 1$ and $a_{ij} \ge \delta$ that occur in the summation, we have for $\epsilon > 0$

$$\sum_{(i,j)} |\log a_{ij} - \log a'_{ij}| < \frac{\epsilon}{2} \quad \text{if} \quad \|Q - Q'\|_1 < \delta_1(\epsilon).$$

In addition, we know that all the elements of the stationary distribution of $T$ are bounded away from zero, since the largest element of the stationary distribution of $T$ is lower-bounded by $\frac{1}{M^{k+1}}$, and any state can be reached by a finite number of steps whose transition probabilities are bounded away from zero. Therefore, for some $C_1 < \infty$

$$|\log \tau_{S_1} - \log \tau'_{S_1}| < C_1 |\tau_{S_1} - \tau'_{S_1}|.$$

Then, from the result of the sensitivity of the stationary distribution of a Markov chain [10], for some $C_2 < \infty$

$$|\tau_{S1} - \tau'_{S1}| \le C_2 \sum_{(s,\bar{s})} |t_{s\bar{s}} - t'_{s\bar{s}}| = C_2 \sum_{(i,j)} |a_{ij} - a'_{ij}|.$$

Hence, for $\epsilon > 0$, we obtain,

$$|\log \tau_{S_1} - \log \tau'_{S_1}| < \frac{\epsilon}{2} \quad \text{if} \quad \|Q - Q'\|_1 < \delta_2(\epsilon).$$

Therefore, by letting $\delta(\epsilon) = \min(\delta_1(\epsilon), \delta_2(\epsilon))$, we have

$$\left| \frac{1}{n} \log Q(S^n) - \frac{1}{n} \log Q'(S^n) \right| < \epsilon \quad \text{if} \quad \|Q - Q'\|_1 < \delta(\epsilon).$$

Let us now go back to the original process $Z$. From

$$\left| \frac{1}{n} \log Q(S^n) - \frac{1}{n} \log Q'(S^n) \right| < \epsilon$$

we have

$$Q'(X^n, Z^n) < \exp(n\epsilon) Q(X^n, Z^n)$$

thus

$$
\begin{aligned}
Q'(Z^n) = \sum_{x^n} Q'(x^n, Z^n) &< \exp(n\epsilon) \sum_{x^n} Q(x^n, Z^n) \\
&= \exp(n\epsilon) Q(Z^n)
\end{aligned}
$$

where the summations are again over the sequences that have nonzero probabilities. By changing the role of $Q$ and $Q'$, we get the result (8), namely, $\frac{1}{n} \log Q(Z^n)$ is an equicontinuous sequence. Therefore, we have the uniform convergence of the lemma.                                                                        $\square$

*Lemma 3 (Uniform Concentration):* Suppose $Q \in \Theta_k^\delta$ for some fixed $\delta > 0$. Let $\hat{\boldsymbol{X}}_Q^\epsilon$ be the randomized filter defined in (3). Then

$$\lim_{n\to\infty} \left( L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - E\left( L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) \right) \right) = 0 \quad \text{a.s.}$$

uniformly on $\Theta_k^\delta$.

*Proof:* This lemma shows the uniform concentration property of $L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n)$. The randomization of the filter is needed to deal with ties that occur in deciding the Bayes response. A detailed proof of this lemma is given in Appendix B.        $\square$

*Lemma 4 (Continuity):* Consider a single letter filtering setting. Suppoes $Q$ is some other joint probability law of $X$ and $Z$. Define single letter filters $\hat{X}_P(z)$ and $\hat{X}_Q^\epsilon(z)$ as

$$
\begin{aligned}
\hat{X}_P(z)[\hat{x}] &= \Pr(B(\boldsymbol{P}_{X|z}) = \hat{x}) \\
\hat{X}_Q^\epsilon(z)[\hat{x}] &= \Pr(B(\boldsymbol{Q}_{X|z} + \boldsymbol{U}) = \hat{x})
\end{aligned}
$$

where $\boldsymbol{U} \in \mathbb{R}^M$ is a uniform random vector in $B_\epsilon$ as before. Then

$$
\begin{aligned}
E\left(\ell(X, \hat{X}_Q^\epsilon(Z))\right) &- E\left(\ell(X, \hat{X}_P(Z))\right) \\
&\leq \Lambda_{\max} K_{\boldsymbol{\Pi}} \cdot \|\boldsymbol{P}_Z - \boldsymbol{Q}_Z\|_1 + C_{\boldsymbol{\Lambda}} \cdot \epsilon
\end{aligned}
$$

where the expectations on the left-hand side of the inequality are under $P$ and $K_{\boldsymbol{\Pi}} = \sum_{i=1}^M \|\Pi_i^{-1}\|_2$, and

$$C_{\boldsymbol{\Lambda}} = \max_{a,b \in \mathcal{A}} \|\boldsymbol{\lambda}_a - \boldsymbol{\lambda}_b\|_2.$$

*Remark:* This lemma states that the excess expected loss of a randomized filter optimized for a mismatched probability law can be upper-bounded by the $L_1$ difference between the true and the mismatched probability laws of output symbol, plus a small constant term which diminishes with the randomization probability. This is somewhat analogous to a result for the prediction problem which was derived in [14, eq. (20)].

*Proof of Lemma 4:* Define $\hat{X}_Q(z)[\hat{x}] = \Pr(B(\boldsymbol{Q}_{X|z}) = \hat{x})$. Then

$$
\begin{aligned}
E\left(\ell(X, \hat{X}_Q^\epsilon(Z))\right) &- E\left(\ell(X, \hat{X}_P(Z))\right) \\
&= \sum_{x,z} P(x,z)\left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_P(z))\right)
\end{aligned}
$$

$$
\begin{aligned}
&= \sum_{x,z} \left(P(x,z) - Q(x,z)\right) \cdot \left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_P(z))\right) \\
&\quad + \sum_{x,z} Q(x,z) \cdot \left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_P(z))\right) \\
&\leq \Lambda_{\max} \sum_{x,z} |P(x,z) - Q(x,z)| \\
&\quad + \sum_{x,z} Q(x,z) \cdot \left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_P(z))\right) \qquad (11) \\
&\leq \Lambda_{\max} \sum_{x,z} |P(x,z) - Q(x,z)| \\
&\quad + \sum_{x,z} Q(x,z)\left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_Q(z))\right) \qquad (12)
\end{aligned}
$$

where (11) follows from Hölder's inequality, and (12) follows from the fact that

$$\sum_{x,z} Q(x,z)(\ell(x, \hat{X}_Q(z)) - \ell(x, \hat{X}_P(z))) \leq 0.$$

Now, let us bound the first term in (12)

$$
\begin{aligned}
\Lambda_{\max} &\sum_{x,z} |P(x,z) - Q(x,z)| \\
&= \Lambda_{\max} \sum_x |P(x) - Q(x)|\left(\sum_z \Pi(x,z)\right) \\
&= \Lambda_{\max} \sum_x |P(x) - Q(x)| \qquad (13) \\
&= \Lambda_{\max} \sum_i |(\boldsymbol{P}_Z - \boldsymbol{Q}_Z)^T \Pi_i^{-1}| \\
&\leq \Lambda_{\max} \sum_i \|\Pi_i^{-1}\|_2 \cdot \|\boldsymbol{P}_Z - \boldsymbol{Q}_Z\|_2 \qquad (14) \\
&\leq \Lambda_{\max} K_{\boldsymbol{\Pi}} \cdot \|\boldsymbol{P}_Z - \boldsymbol{Q}_Z\|_1 \qquad (15)
\end{aligned}
$$

where (13) follows from the fact that $\sum_z \Pi(x,z) = 1$; (14) follows from Cauchy–Schwartz inequality, and (15) follows from the fact that $L_2$-norm is less than or equal to $L_1$-norm.

The second term in (12) becomes

$$
\begin{aligned}
\sum_{x,z} &Q(x,z)\left(\ell(x, \hat{X}_Q^\epsilon(z)) - \ell(x, \hat{X}_Q(z))\right) \\
&= \sum_z Q(z) \sum_x Q(x|z) \sum_{\hat{x}} \Lambda(x,\hat{x}) \cdot \left(\hat{X}_Q^\epsilon(z)[\hat{x}] - \hat{X}_Q(z)[\hat{x}]\right) \\
&= \sum_z Q(z) \sum_{\hat{x}} \left(\hat{X}_Q^\epsilon(z)[\hat{x}] - \hat{X}_Q(z)[\hat{x}]\right) \sum_x \Lambda(x,\hat{x}) Q(x|z) \\
&= \sum_z Q(z) \sum_{\hat{x}} \left(\hat{X}_Q^\epsilon(z)[\hat{x}] - \hat{X}_Q(z)[\hat{x}]\right) \cdot \boldsymbol{\lambda}_{\hat{x}}^T \boldsymbol{Q}_{X|z}. \qquad (16)
\end{aligned}
$$

It is easy to see that the inner summation in (16) is always nonnegative since by definition, $\hat{X}_Q(z)$ assigns probability 1 to $B(\boldsymbol{Q}_{X|z})$. Now, for a given $Q$, define

$$\boldsymbol{U}_{\max} = \arg\max_{\boldsymbol{U} \in B_\epsilon} \left(\boldsymbol{\lambda}_{B(\boldsymbol{Q}_{X|z}+\boldsymbol{U})} - \boldsymbol{\lambda}_{B(\boldsymbol{Q}_{X|z})}\right)^T \boldsymbol{Q}_{X|z} \qquad (17)$$

resolving ties arbitrarily. Then, we have

$$
\begin{aligned}
\sum_{\hat{x}} &\left(\hat{X}_Q^\epsilon(z)[\hat{x}] - \hat{X}_Q(z)[\hat{x}]\right) \cdot \boldsymbol{\lambda}_{\hat{x}}^T \boldsymbol{Q}_{X|z} \\
&= \left(\sum_{\hat{x}} \left(\hat{X}_Q^\epsilon(z)[\hat{x}] \cdot \boldsymbol{\lambda}_{\hat{x}}\right) - \boldsymbol{\lambda}_{B(\boldsymbol{Q}(X|z))}\right)^T \boldsymbol{Q}_{X|z}
\end{aligned}
$$

$$\leq \left(\boldsymbol{\lambda}_{B(\boldsymbol{Q}(X|z)+\boldsymbol{U}_{\max})} - \boldsymbol{\lambda}_{B(\boldsymbol{Q}(X|z))}\right)^T \boldsymbol{Q}_{X|z} \qquad (18)$$

$$\leq \left(\boldsymbol{\lambda}_{B(\boldsymbol{Q}(X|z))} - \boldsymbol{\lambda}_{B(\boldsymbol{Q}_{X|z}+\boldsymbol{U}_{\max})}\right)^T \boldsymbol{U}_{\max} \qquad (19)$$

$$\leq \max_{a,b\in\mathcal{A}} \|\boldsymbol{\lambda}_a - \boldsymbol{\lambda}_b\|_2 \cdot \|\boldsymbol{U}_{\max}\|_2$$

$$\leq C_{\boldsymbol{\Lambda}} \cdot \epsilon \qquad (20)$$

where (18) follows from (17); (19) follows from the fact

$$\boldsymbol{\lambda}_{B(\boldsymbol{Q}_{X|z}+\boldsymbol{U}_{\max})}^T (\boldsymbol{Q}_{X|z} + \boldsymbol{U}_{\max}) \leq \boldsymbol{\lambda}_{B(\boldsymbol{Q}_{X|z})}^T (\boldsymbol{Q}_{X|z} + \boldsymbol{U}_{\max})$$

and (20) follows from the Cauchy–Schwartz inequality. Note that depending on $Q$ and $z$, (18) and (19) can be both zero and hold with equality. Together with (15), the lemma is proved. $\square$

Before moving on to Lemma 5, we need the following three definitions. In Lemma 2, we have seen that for $Q \in \Theta_k^\delta$, $\boldsymbol{D}(P_Z\|Q)$ is well defined. Now, let us consider the case where $Q \in \Theta_k^\delta$ is some function of the noisy observation $Z^n$ (denoted as $Q[Z^n]$). As mentioned in the footnote of Section IV-B, the notion of the relative entropy rate between $P_Z$ and that random $Q[Z^n]$ is defined in Definition 2 using Definition 1. Definition 3 is also needed for the inequality in Lemma 5.

*Definition 1:* Suppose $Q[Z^n] \in \Theta_k^\delta$ and $f$ is some function of $(X^\infty, Z^\infty, Q[Z^n])$ such that the expectation

$$E\Big(f(X^\infty, Z^\infty, Q[Z^n])\Big) = \int f(x^\infty, z^\infty, Q[z^n]) dP(x^\infty, z^\infty)$$

exists. Then, the notation $\hat{E}(\cdot)$ is defined as

$$\hat{E}\Big(f(X^\infty, Z^\infty, Q[Z^n])\Big) \triangleq \int f(x^\infty, z^\infty, Q[Z^n]) dP(x^\infty, z^\infty).$$

That is, in $\hat{E}\Big(f(X^\infty, Z^\infty, Q[Z^n])\Big)$, the Lebesgue integration with respect to the randomness of $Q[Z^n]$ is excluded. Moreover, suppose $\mathcal{T}$ is a set of some time indices, and $z_{\mathcal{T}}$ denotes a subsequence of $z^\infty$ with the time indices in $\mathcal{T}$. Then, the conditioning on $z_{\mathcal{T}}$ with respect to $\hat{E}$ is defined as

$$\hat{E}\Big(f(X^\infty, Z^\infty, Q[Z^n])\Big|Z_{\mathcal{T}} = z_{\mathcal{T}}\Big)$$
$$\triangleq \int f(x^\infty, z^\infty, Q[Z^n]) dP(x^\infty, z^\infty|z_{\mathcal{T}})$$

where $P(\cdot|z_{\mathcal{T}})$ denotes the conditional probability measure on $(x^\infty, z^\infty)$ given $z_{\mathcal{T}}$. Again, the randomness of $Q[Z^n]$ is excluded in the Lebesgue integration.

*Remark:* The above definitions of $\hat{E}(\cdot)$ and $\hat{E}(\cdot|\cdot)$ may seem subtle. However, the main point of two definitions is simple, namely, they exclude the randomness of $Q[Z^n]$ in calculating Lebesgue integrations.

*Definition 2:* Suppose $Q[Z^n] \in \Theta_k^\delta$. Then, the relative entropy rate between $P_Z$ and $Q[Z^n]$ is defined as

$$\boldsymbol{D}(P_Z\|Q[Z^n]) \triangleq \hat{E}\left(\log \frac{P_Z(Z_0|Z_{-\infty}^{-1})}{Q[Z^n](Z_0|Z_{-\infty}^{-1})}\right).$$

*Remark:* Note that $\boldsymbol{D}(P_Z\|Q[Z^n])$ is a function of $Z^n$, and still is a random variable.

*Definition 3:* Define the $k$th-order Markov approximation of $P_X$ for $n \geq k$ as

$$P_X^{(k)}(X^n) \triangleq P_X(X^k) \prod_{i=k+1}^n P_X(X_i|X_{i-k}^{i-1}).$$

Furthermore, denote $P_Z$ and $P_Z^{(k)}$ as the probability law of the output of DMC, $\boldsymbol{\Pi}$, when the probability law of input is $P_X$ and $P_X^{(k)}$, respectively.

*Remark:* Note that $P_Z^{(k)}$ is not the $k$th-order Markov approximation of $P_Z$, but is the distribution of the channel output whose input is $P_X^{(k)}$, the $k$th-order Markov approximation of the original input distribution $P_X$.

Now, we give the following lemma that upper-bounds the relative entropy rate between $P_Z$ and the ML estimator.

*Lemma 5:* For the given sequence $\{\delta_k\}$ defined in Section IV-A and for fixed $k$, we have

$$\lim_{n\to\infty} \boldsymbol{D}(P_Z\|\hat{Q}_{k,\delta_k}[Z^n]) \leq \boldsymbol{D}(P_X\|P_X^{(k)}) \quad \text{a.s.}$$

*Proof:* Recall that $\hat{Q}_{k,\delta_k}[Z^n]$ is an ML estimator in $\Theta_k^{\delta_k}$ based on the observation $Z^n$. From (1), we know that

$$\lim_{n\to\infty} \boldsymbol{D}(P_Z\|\hat{Q}_{k,\delta_k}[Z^n]) = \min_{Q\in\Theta_k^{\delta_k}} \boldsymbol{D}(P_Z\|Q) \quad \text{a.s.}$$

Also, (2) and Definition 3 assures that $P_Z^{(k)} \in \Theta_k^{\delta_k}$. Therefore, we have

$$\lim_{n\to\infty} \boldsymbol{D}(P_Z\|\hat{Q}_{k,\delta_k}[Z^n]) \leq \boldsymbol{D}(P_Z\|P_Z^{(k)}) \quad a.s..$$

This is the link where we needed Assumption 1. Now, let us denote $P^{(k)}$ as the joint probability law of $(X^n, Z^n)$ when the probability law of input process is $P_X^{(k)}$. Then, by the chain rule of relative entropy [5, eq. (2.67)], we have

$$E\left(\log \frac{P(X^n, Z^n)}{P^{(k)}(X^n, Z^n)}\right)$$
$$= D_n(P_X\|P_X^{(k)}) + E\left(\log \frac{P(Z^n|X^n)}{P^{(k)}(Z^n|X^n)}\right)$$
$$= D_n(P_Z\|P_Z^{(k)}) + E\left(\log \frac{P(X^n|Z^n)}{P^{(k)}(X^n|Z^n)}\right).$$

Since the DMC is fixed, we have $E\left(\log \frac{P(Z^n|X^n)}{P^{(k)}(Z^n|X^n)}\right) = 0$. Moreover, by the nonnegativity of relative entropy, $E\left(\log \frac{P(X^n|Z^n)}{P^{(k)}(X^n|Z^n)}\right) \geq 0$. Therefore, we get

$$D_n(P_Z\|P_Z^{(k)}) \leq D_n(P_X\|P_X^{(k)}).$$

Since $\boldsymbol{D}(P_X\|P_X^{(k)}) = \lim_{n\to\infty} \frac{1}{n} D_n(P_X\|P_X^{(k)})$ always exists by ergodicity, we have

$$\boldsymbol{D}(P_Z\|P_Z^{(k)}) \leq \boldsymbol{D}(P_X\|P_X^{(k)})$$

and the lemma is proved. $\square$

We are now finally in a position to prove our main theorem.

*Proof of Theorem 1:* As mentioned in Section IV-B, we first fix $k$ and $\epsilon$, and try to get the inequality in the form of (6) to prove Part a). To refresh, (6) is given again here

$$\limsup_{n \to \infty} \Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$\leq F\Big( \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t), \epsilon \Big) \quad \text{a.s.}$$

From the definition of $L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n)$

$$L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) = \frac{1}{n} \sum_{t=1}^n \ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z^t))$$

where from (5), we know that $Q_k^t$ is a function of $Z^{m_i(t)}$. Since $\ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z^t))$ is a function of $(X_t, Z^t, Q[Z^{m_i(t)}])$, we can define a quantity $\hat{E}(\ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z^t)))$ from Definition 1. From this, we also define

$$\hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) = \frac{1}{n} \sum_{t=1}^n \hat{E}\Big( \ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z^t)) \Big).$$

Now, we have following Corollary 1 from Lemma 3, whose proof is given in Appendix C. This corollary is a key step in proving the main theorem, since it provides a crucial link that enables to get the inequality in (6).

*Corollary 1:* For fixed $k$ and $\epsilon$, we have

$$\lim_{n \to \infty} \Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) - \hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) \Big) = 0 \quad \text{a.s.}$$

From Corollary 1, we have following equality:

$$\limsup_{n \to \infty} \Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$= \limsup_{n \to \infty} \Big( \hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) - \phi_n(P_X, \boldsymbol{\Pi}) \Big) \quad \text{a.s.}$$

Therefore , to get the inequality of the form of (6), we can equivalently show

$$\limsup_{n \to \infty} \Big( \hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$\leq F\Big( \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t), \epsilon \Big).$$

Now, let us consider following chain of inequalities:[6]

$$\hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) - \phi_n(P_X, \boldsymbol{\Pi})$$
$$= \frac{1}{n} \sum_{t=1}^n \Big( \hat{E}\Big( \ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z^t)) \Big) - \hat{E}\Big( \ell(X_t, \hat{X}_P(Z^t)) \Big) \Big)$$
$$= \frac{1}{n} \sum_{t=1}^n \hat{E}\Big( \hat{E}\big( \ell(X_t, \hat{X}^\epsilon_{Q_k^t}(Z_t, Z^{t-1})) | Z^{t-1} \big)$$
$$\qquad - \hat{E}\big( \ell(X_t, \hat{X}_P(Z_t, Z^{t-1})) | Z^{t-1} \big) \Big)$$
$$\leq \frac{K_{\boldsymbol{\Pi}} \Lambda_{\max}}{n} \sum_{t=1}^n \hat{E} \| \boldsymbol{P}_{Z_t|Z^{t-1}} - \boldsymbol{Q}^t_{k\, Z_t|Z^{t-1}} \|_1 + C_{\boldsymbol{\Lambda}} \cdot \epsilon \quad (21)$$

---

[6]All the equalities and inequalities between random variables in this proof should be understood in almost sure sense.

$$\leq \frac{\sqrt{2 \ln 2} K_{\boldsymbol{\Pi}} \Lambda_{\max}}{n} \sum_{t=1}^n \hat{E} \sqrt{\hat{E}\Big( \log \frac{P(Z_t|Z^{t-1})}{Q_k^t(Z_t|Z^{t-1})} \Big| Z^{t-1} \Big)}$$
$$+ C_{\boldsymbol{\Lambda}} \cdot \epsilon \qquad (22)$$
$$\leq \sqrt{2 \ln 2} K_{\boldsymbol{\Pi}} \Lambda_{\max} \sqrt{\frac{1}{n} \sum_{t=1}^n \hat{E}\Big( \log \frac{P(Z_t|Z^{t-1})}{Q_k^t(Z_t|Z^{t-1})} \Big)} + C_{\boldsymbol{\Lambda}} \cdot \epsilon.$$
$$(23)$$

The notation $\boldsymbol{Q}^t_{k\, Z_t|Z^{t-1}}$ in (21) stands for the simplex vector in $\mathcal{M}$ whose $i$th component stands for $Q_k^t(Z_t = i | Z^{t-1})$. The inequality in (21) is obtained from Lemma 4, since $\boldsymbol{\Pi}$ does not vary with $t$, and given $Z^{t-1}$, estimating $X_t$ based on $Z^t$ is equivalent to the single letter setting as in Lemma 4 with the corresponding conditional distribution. Furthermore, (22) follows from Pinsker's inequality [5, Lemma 12.6.1], and (23) follows from Jensen's inequality. By taking $\limsup$ on both sides, we have

$$\limsup_{n \to \infty} \Big( \hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$\leq \sqrt{2 \ln 2} K_{\boldsymbol{\Pi}} \Lambda_{\max} \sqrt{\limsup_{n \to \infty} \frac{1}{n} \sum_{t=1}^n \hat{E}\Big( \log \frac{P(Z_t|Z^{t-1})}{Q_k^t(Z_t|Z^{t-1})} \Big)}$$
$$+ C_{\boldsymbol{\Lambda}} \cdot \epsilon \quad a.s.,$$

since the square root function is a continuous function. For the expression inside the square root of the right-hand side of the inequality

$$\limsup_{n \to \infty} \frac{1}{n} \sum_{t=1}^n \hat{E}\Big( \log \frac{P(Z_t|Z^{t-1})}{Q_k^t(Z_t|Z^{t-1})} \Big)$$
$$= \limsup_{t \to \infty} \hat{E}\Big( \log \frac{P(Z_t|Z^{t-1})}{Q_k^t(Z_t|Z^{t-1})} \Big) \quad \text{a.s.} \quad (24)$$
$$= \limsup_{t \to \infty} \hat{E}\Big( \log \frac{P(Z_0|Z_{-\infty}^{-1})}{Q_k^t(Z_0|Z_{-\infty}^{-1})} \Big) \quad \text{a.s.} \quad (25)$$
$$= \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t) \quad \text{a.s.} \quad (26)$$

where (24) follows from Cesáro's mean convergence theorem; the numerator of (25) follows from the fact that $P$ is stationary and $P(Z_0|Z_{-t}^{-1}) \to P(Z_0|Z_{-\infty}^{-1})$ almost surely by martingale convergence theorem; and the denominator of (25) follows from the fact that $Q_k^t$ is also a stationary law, and with probability 1, for any $\epsilon > 0$, there exists $N$ such that for all $t \geq N$

$$|Q_k^t(Z_0|Z_{-t}^{-1}) - Q_k^t(Z_0|Z_{-\infty}^{-1})| < \epsilon$$

which is guaranteed by the uniform convergence result of Lemma 1. Finally, (26) follows from Definition 2. Therefore

$$\limsup_{n \to \infty} \Big( \hat{E}\Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) \Big) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$\limsup_{n \to \infty} \Big( L_{\hat{\boldsymbol{X}}^\epsilon_{\mathrm{univ},k}}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \Big)$$
$$\leq \sqrt{2 \ln 2} K_{\boldsymbol{\Pi}} \Lambda_{\max} \sqrt{\limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t)} + C_{\boldsymbol{\Lambda}} \cdot \epsilon \quad \text{a.s.}$$
$$(27)$$

which finally is in the form of (6). Now, we need to check if the right-hand side of (27) goes to zero if we let $k \to \infty$ and $\epsilon \downarrow 0$. To see this, consider the following further upper bounds:

$$\limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t)$$
$$= \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| \hat{Q}_{k,\delta_k}[Z^t]) \qquad (28)$$
$$\leq \boldsymbol{D}(P_{X} \| P_{X_k}), \qquad (29)$$

where (28) follows from the fact that $m_{i(t)} \to \infty$ as $t \to \infty$, and (29) follows from Lemma 5. The inequality (29) holds for every $k$, and by Shannon–McMillan–Breiman theorem [5, Ch. 15.7], we know that $\boldsymbol{D}(P_X \| P_{X_k}) \to 0$ as $k \to \infty$. Therefore

$$\lim_{k \to \infty} \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t) = 0$$

and thus

$$\lim_{k \to \infty} \limsup_{n \to \infty} \left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \right) \leq C_{\boldsymbol{\Lambda}} \cdot \epsilon \quad \text{a.s.}$$

Finally, sending $\epsilon$ to zero, Part a) of the theorem is proved. Part b) follows directly from a), and Fatou's lemma. That is

$$\lim_{k \to \infty} \limsup_{n \to \infty} \left( E\left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) \right) - \phi_n(P_X, \boldsymbol{\Pi}) \right)$$
$$= \lim_{k \to \infty} \limsup_{n \to \infty} E\left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \right)$$
$$\leq \lim_{k \to \infty} E\left( \limsup_{n \to \infty} \left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) - \phi_n(P_X, \boldsymbol{\Pi}) \right) \right)$$
$$\leq C_{\boldsymbol{\Lambda}} \cdot \epsilon.$$

Note that the expectation here is with respect to the randomness of probability law within the paranthesis, too. By sending $\epsilon$ to zero, Part b) is proved. $\square$

## V. EXTENSION: UNIVERSAL FILTERING FOR CHANNEL WITH MEMORY

Now, let us extend our result to the case where channel has memory. With the identical assumption on $\{X_t\}$, now suppose $\{Z_t\}$ is expressed as

$$Z_t = X_t \oplus N_t \qquad (30)$$

where $\oplus$ denotes modulo-$M$ addition, and $\{N_t\}$ is an $\mathcal{A}$-valued noise process which is not necessarily memoryless. We assume we have complete knowledge of the probability law of $\{N_t\}$. Specifically, let us consider the case where $\{N_t\}$ is an HMP, that is, it is an output of an invertible memoryless channel $\boldsymbol{\Gamma} = \{\Gamma(i,j)\}_{i,j \in \mathcal{A}}$ whose input is irreducible, aperiodic $\ell$th-order Markov process $\{S_t\}$, which is independent of $\{X_t\}$. Let $\Gamma_{\min} = \min_{i,j \in \mathcal{A}} \{\Gamma(i,j)\}$, and suppose $\Gamma_{\min} > 0$. For simplicity, assume that the alphabet size of $\{S_t\}$ is also $\mathcal{A}$.

In this model, the channel between $X_t$ and $Z_t$ at time $t$ is an $M$-ary symmetric channel, which is specified by the $S_t$th row of $\boldsymbol{\Gamma}$. Define an $M \times M$ matrix $\boldsymbol{\Pi}_t$ whose $(x_t, z_t)$th element is

$$\Pi_t(x_t, z_t) = P_{N_t}(z_t \ominus x_t)$$
$$= \Pr(Z_t = z_t | X_t = x_t)$$
$$= \sum_{s_t} \Pr(Z_t = z_t | X_t = x_t, S_t = s_t) \Pr(S_t = s_t)$$

where $\ominus$ denotes modulo-$M$ subtraction. Now, let us make following assumptions on the noise process:
- $\{N_t\}$ is stationary, i.e., $\boldsymbol{\Pi}_t$ is identical for all $t$;
- $\boldsymbol{\Pi}_t$ is invertible;
- for all $S_{t-\ell}^t(\omega)$, there exists an $\alpha$ such that $\Pr(S_t | S_{t-\ell}^{t-1}) \geq \alpha > 0$.

As stated in [22, Sec. 2-A], the first and the second assumptions are rather benign. Especially, for the second assumption, it can be shown that under benign conditions on the parametrization, almost all parameter values except for those in a set of Lebesgue measure zero, give rise to a process satisfying this assumption. In addition, since this only corresponds to the case when $k = 0$ in [22, Assumption 1], it is a much weaker assumption. The third assumption is a similar positivity assumption as Assumption 1, which enables our universal filtering scheme.

Under these assumptions on the noise process, we can extend our scheme to do the universal filtering for this channel. First, we can convert this channel to the equivalent memoryless channel $\boldsymbol{\Xi} = \{\xi((i,j), h)\}_{i,j,h \in \mathcal{A}}$, where the input process is $\{(X_t, S_t)\}$ and the output is $\{Z_t\}$. Here, $\boldsymbol{\Xi}$ is $M^2 \times M$ matrix, and the channel transition probability is

$$\xi((i,j), h) = \Gamma(j, h \ominus i)$$

for all $i, j$, and $k$. To do the filtering, we apply our scheme to this equivalent memoryless channel. For fixed $k \geq \ell$, as in Section II-B.1, define a parameter set of HMPs, $\Theta_k$, whose Markov chain has $M^{k+\ell}$ states, and the memoryless channel has dimension $\mathbb{R}^{M^{(k+\ell)} \times M}$. The $k$th-order conditional probability of our new input process is

$$\Pr(X_t, S_t | X_{t-k}^{t-1}, S_{t-k}^{t-1})$$
$$= \Pr(X_t | X_{t-k}^{t-1}) \cdot \Pr(S_t | S_{t-\ell}^{t-1})$$
$$\geq \delta_k \cdot \alpha, \qquad (31)$$

where (31) follows from Assumption 1 and the third condition on the noise process. Let $\gamma_k = \delta_k \cdot \alpha$. Then, we can model $\{Z_t\}$ in $\Theta_k^{\gamma_k}$, or equivalently, model $(X_t, S_t)$ as $k$th-order Markov chain, and obtain $Q_k^t$, the ML estimator in $\Theta_k^{\gamma_k}$ based on $Z^{m_{i(t)}}$. By forward recursion, we can get $Q_k^t(X_t, S_t | Z^t)$, and by summing over $S_t$'s we can calculate $\boldsymbol{Q}_{k,X_t|Z^t}^t$, the simplex vector in $\mathcal{M}$ whose $i$th component is $Q_k^t(X_t = i | Z^t)$. Then, finally, we define our sequence of universal filtering schemes as

$$\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon = \{\hat{X}_{Q_k^t, t}^\epsilon\}$$

exactly the same as we proposed in Section IV-A.

The analysis of this scheme is identical to the one given in the proof of the main theorem. Equation (21), which is the only place where the invertibility of the $\boldsymbol{\Pi}$ is used, can also be obtained in this case due to the second assumption of the noise process. Thus, we again get

$$\limsup_{n \to \infty} \left( L_{\hat{\boldsymbol{X}}_{\text{univ},k}^\epsilon}(X^n, Z^n) - \phi_n(P_X, \Pi) \right)$$
$$\leq 2\sqrt{2 \ln 2} K_\Pi \Lambda_{\max} \sqrt{\limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t)} + C_{\boldsymbol{\Lambda}} \cdot \epsilon \quad \text{a.s.}$$

Since

$$\limsup_{t \to \infty} \boldsymbol{D}(P_Z \| Q_k^t)$$
$$= \limsup_{t \to \infty} \boldsymbol{D}(P_Z \| \hat{Q}_{k,\gamma_k}[Z^t]) \leq \boldsymbol{D}(P_X \| P_{X_k})$$

by the same argument as Lemma 5, we have the same result as Theorem 1. Thus, we can successfully extend our scheme to the case where the channel noise is an HMP with some mild assumptions.

## VI. DISCUSSION

Throughout the paper, we have only considered the case where the input, output, and reconstruction alphabets are equal. However, we can easily extend our result to the case where the alphabet sizes are different (but still finite). In that case, the condition on the channel, parallel to the invertibility condition, is that the channel transition matrix should have full row-rank. Since the argument of the extension would be rather straightforward, we omit the details in this paper.

The result that we attain in Theorem 1 can also be attained by the schemes devised in [18], [19]. Therefore, we have shown that a completely different approach can achieve the same goal in the universal filtering problem when the underlying signal is a stationary and ergodic process. In addition, our work gives the first theoretical justification of using HMP models for filtering, which is a prevalent approach in practice where the underlying signal need not be a Markov process. Furthermore, it is not clear how to extend the schemes in [18], [19] to the case of channels with memory, whereas the extension of our scheme to such cases is quite simple in some settings (e.g., when the noise is an HMP), as in Section V.

As described in Section IV-A, our filter is a randomized filter. The randomization is necessary in obtaining the continuity result of Lemma 9 Part a), which we use for proving our main theorem. Whether a deterministic version of our filter, i.e., a filter that is defined without $\boldsymbol{U}$ in (3), is universally optimal remains an open question. The filter devised in [18], [19] is also a randomized scheme, and a parallel discussion regarding the randomization is also given in [19, Sec. VI]. In contrast, a filter that appears in [24], which is equivalent to the scheme in [21] that only utilizes a one-sided context, is a deterministic scheme that can indeed achieve the asymptotically optimal performance in Theorem 1. Therefore, when the channel is memoryless, the randomization of a universal filter is not necessary in general to achieve the performance goal of Theorem 1. However, when the channel has memory as in Section V, we are not aware of any deterministic filter that can be universally optimal for any underlying stationary and ergodic process.

## VII. CONCLUDING REMARK AND FUTURE WORK

In this paper, we proved that, for the known, invertible DMC, a family of filters based on HMPs is universally asymptotically optimal for any general stationary and ergodic $\{X_t\}$ satisfying some mild positivity condition. That is, we showed that our sequence of schemes indexed by $k$ and $\epsilon$ achieves the best asymptotically optimal performance regardless of clean source distribution. We could also extend this scheme to the case where

channel has memory, especially where the channel noise process is an HMP.

The future direction of the work would be to ascertain the relationship between $k$ and $n$, such that we can devise a single scheme that grows $k$ with some rate related to $n$. Attempting to loosen the positivity assumption that we made in our main theorem and extending our discrete universal filtering schemes to discrete universal denoising schemes are additional future directions of our research.

## APPENDIX A
## THREE LEMMAS

Here, we revise three lemmas from [1] regarding probability law of HMP. These are needed to prove Lemma 1. For the following three lemmas, fix $k$ and $\delta$, and suppose $Q \in \Theta_k^\delta$. Also, fix some $m \in \mathbb{N}$, such that $m \geq k$. Proofs are similar to [1, the Appendix]. Note that $\{X_t\}$ is still our clean signal and $\{Z_t\}$ is the noisy observed signal (not necessarily an HMP).

*Lemma 6:* We have

$$Q(X_{t+m} = j | X_t = i, Z_{-\infty}^\infty) \geq \mu_{\delta,k,m}$$

where

$$\mu_{\delta,m,k} = \left(1 + \frac{M-1}{(\delta \cdot \Pi_{\min})^{m+k}}\right)^{-1}$$

is independent of $Q, Z_{-\infty}^\infty, i, j$.

*Proof:* From Markovity and conditioning

$$\frac{Q(X_{t+m} = j | X_t = i, Z_{-\infty}^\infty)}{Q(X_{t+m} = j' | X_t = i, Z_{-\infty}^\infty)}$$
$$= \frac{Q(X_{t+m} = j, Z_{-\infty}^\infty | X_t = i)}{Q(X_{t+m} = j', Z_{-\infty}^\infty | X_t = i)}$$
$$= \frac{Q(X_{t+m} = j, Z_{t+m+k+1}^\infty | X_t = i)}{Q(X_{t+m} = j', Z_{t+m+k+1}^\infty | X_t = i)}$$
$$\times \frac{Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j)}{Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j')}. \tag{32}$$

Now, let us bound the terms in (32). First

$$\frac{Q(X_{t+m} = j, Z_{t+m+k+1}^\infty | X_t = i)}{Q(X_{t+m} = j', Z_{t+m+k+1}^\infty | X_t = i)}$$
$$= \frac{\sum_{j_0} Q(X_{t+m+k} = j_0, X_{t+m} = j, Z_{t+m+k+1}^\infty | X_t = i)}{\sum_{j_0} Q(X_{t+m+k} = j_0, X_{t+m} = j', Z_{t+m+k+1}^\infty | X_t = i)}$$
$$= \frac{\sum_{j_0} a_{ij}^m a_{jj_0}^k Q(Z_{t+m+k+1}^\infty | X_{t+m+k} = j_0)}{\sum_{j_0} a_{ij'}^m a_{j'j_0}^k Q(Z_{t+m+k+1}^\infty | X_{t+m+k} = j_0)}$$

where $a_{rs}^p$ stands for the $(r, s)$th element of $A^p$, the $p$th power of the transition matrix $A$.

Note that $a_{ij}^m \geq \delta^m$ and $a_{jj_0}^k \geq \delta^k$, $\forall i, j, j_0$ from the assumption of $\Theta_k^\delta$. Let $Q(Z_{t+m+k+1}^\infty | X_{t+m+k} = j_0) = \alpha_{j_0}$. Then, the last expression is

$$\frac{a_{ij}^m}{a_{ij'}^m} \frac{\sum_{j_0} a_{jj_0}^k \alpha_{j_0}}{\sum_{j_0} a_{j'j_0}^k \alpha_{j_0}}. \tag{33}$$

Since

$$\frac{\sum_{j_0} a_{jj_0}^k \alpha_{j_0}}{\sum_{j_0} a_{j'j_0}^k \alpha_{j_0}} = \frac{\sum_{j_0} \alpha_{j_0} a_{j'j_0}^k \frac{a_{jj_0}^k}{a_{j'j_0}^k}}{\sum_{j_0} \alpha_{j_0} a_{j'j_0}^k} \leq \max_{j_0} \left( \frac{a_{jj_0}^k}{a_{j'j_0}^k} \right)$$

we have

$$(33) \leq \frac{a_{ij}^m}{a_{ij'}^m} \max_{j_0} \left( \frac{a_{jj_0}^k}{a_{j'j_0}^k} \right) \leq \max_{i,j,j',j_0} \left( \frac{a_{ij}^m a_{jj_0}^k}{a_{ij'}^m a_{j'j_0}^k} \right) \leq \frac{1}{\delta^{m+k}}. \quad (34)$$

Now, let us look at the second term in (32). Suppose $\mathcal{T} = \{t+1, \ldots, t+m+k\} \setminus \{t+m\}$, and let $x_{\mathcal{T}}$ stand for the sequence of $x_{t'}$'s where $t' \in \mathcal{T}$. Then

$$\frac{Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j)}{Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j')}$$
$$= \frac{\sum_{x_{\mathcal{T}}} Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j, X_{\mathcal{T}} = x_{\mathcal{T}}) \cdot Q(x_{\mathcal{T}} | i, j)}{\sum_{x_{\mathcal{T}}} Q(Z_{t+1}^{t+m+k} | X_t = i, X_{t+m} = j', X_{\mathcal{T}} = x_{\mathcal{T}}) \cdot Q(x_{\mathcal{T}} | i, j')}$$
$$\leq \frac{1}{(\Pi_{\min})^{m+k}} \quad (35)$$

where $Q(x_{\mathcal{T}} | i, j)$ stands for the conditional probability $Q(X_{\mathcal{T}} = x_{\mathcal{T}} | X_t = i, X_{t+m} = j)$. Thus, from (34) and (35)

$$(32) \leq \frac{1}{(\delta \cdot \Pi_{\min})^{m+k}}.$$

Let now $\rho_j \triangleq Q(X_{t+m} = j | X_t = i, Z_{-\infty}^\infty)$, then

$$1 = \rho_j + \sum_{j' \neq j} \rho_{j'} \leq \rho_j + (M-1) \frac{\rho_j}{(\delta \cdot \Pi_{\min})^{m+k}}$$

and thus, $\rho_j \geq (1 + \frac{M-1}{(\delta \cdot \Pi_{\min})^{m+k}})^{-1}$, which proves the lemma. $\square$

*Lemma 7:* Suppose when $\mathcal{T}$ is a set of time indices, $x_{\mathcal{T}}$ and $z_{\mathcal{T}}$ stand for the sequences of $x_{t'}$'s and $z_{t'}$'s where $t' \in \mathcal{T}$. Now, consider following two arbitrarily given sets:

$$C_t \in \mathcal{X}_t^\infty \triangleq \left\{ x_{\mathcal{T}} : \mathcal{T} \subseteq \mathbb{Z}_{\geq t} \cup \{\infty\} \right\}$$

and

$$D \in \mathcal{Z}_{-\infty}^\infty \triangleq \left\{ z_{\mathcal{T}} : \mathcal{T} \subseteq \mathbb{Z} \cup \{\infty, -\infty\} \right\}.$$

For $d \in \mathbb{N}$, define

$$M_d^+ \triangleq \max_i Q(C_t | X_{t-dm} = i, D)$$
$$M_d^- \triangleq \min_i Q(C_t | X_{t-dm} = i, D).$$

Then

$$M_d^+ - M_d^- \leq (\rho_{\delta,k,m})^{d-1}$$

where $\rho_{\delta,k,m} = 1 - 2\mu_{\delta,k,m}$.

*Proof:* From the argument of Lemma 6, it is easy to see that

$$Q(X_{t+m} = j | X_t = i, D) \geq \mu_{\delta,k,m}$$

is independent of $D$ as well. Now, define

$$\gamma_i(d) \triangleq Q(C_t | X_{t-dm} = i, D)$$

$$\beta_{ij}(d) \triangleq Q(X_{t-dm} = j | X_{t-(d+1)m} = i, D)$$
$$i^+(d) \triangleq \arg\max_i Q(C_t | X_{t-(d+1)m} = i, D)$$
$$i^-(d) \triangleq \arg\min_i Q(C_t | X_{t-dm} = i, D).$$

Since $\delta$, $k$, and $m$ are fixed, let us simply denote $\mu = \mu_{\delta,k,m}$. Also, let us omit $d$ and the parenthesis for the above four quantities to simplify notation. Then

$$M_{d+1}^+ = Q(C_t | X_{t-(d+1)m} = i^+, D) = \sum_j \gamma_j \beta_{i^+ j}$$

$$= \mu M_d^- + (\beta_{i^+ i^-} - \mu) M_d^- + \sum_{j \neq i^-} \gamma_j \beta_{i^+ j}$$

$$\leq \mu M_d^- + (\beta_{i^+ i^-} - \mu) M_d^+ + \sum_{j \neq i^-} \beta_{i^+ j} M_d^+ \quad (36)$$

$$= \mu M_d^- + (1 - \mu) M_d^+ \quad (37)$$

where (36) is possible from Lemma 6, since $\beta_{ij} \geq \mu$ for all $i, j$. By the similar argument, we get

$$M_{d+1}^- \geq \mu M_d^+ + (1 - \mu) M_d^-. \quad (38)$$

By subtracting (38) from (37), we get

$$M_{d+1}^+ - M_{d+1}^- \leq (1 - 2\mu)(M_d^+ - M_d^-) \leq \cdots \leq (1 - 2\mu)^d,$$

thus, proves the lemma. Note that since $\mu = \mu_{\delta,k,m} < \frac{1}{2}$, we know $0 < \rho_{\delta,k,m} < 1$. Also, the result does not depend on $Q$. $\square$

*Lemma 8:* For all $p, d \geq 1$, and $0 \leq l \leq m-1$

$$|Q(C_t | Z_{t-dm-l}^p) - Q(C_t | Z_{t-(d+1)m-l}^p)| \leq (\rho_{\delta,k,m})^{d+1}.$$

*Proof:* By conditioning

$$Q(C_t | Z_{t-(d+1)m-l}^p)$$
$$= \sum_j Q(C_t | Z_{t-(d+1)m-l}^p, X_{t-(d+2)m} = j)$$
$$\cdot Q(X_{t-(d+2)m} = j | Z_{t-(d+1)m-l}^p)$$

and, therefore

$$M_{d+2}^- \leq Q(C_t | Z_{t-(d+1)m-l}^p) \leq M_{d+2}^+.$$

On the other hand

$$Q(C_t | Z_{t-dm-l}^p)$$
$$= \sum_{z_{t-(d+1)m-l}^{t-dm-l-1}} Q(C_t | Z_{t-(d+1)m-l}^p)$$
$$\cdot Q(Z_{t-(d+1)m-l}^{t-dm-l-1} = z_{t-(d+1)m-l}^{t-dm-l-1} | Z_{t-dm-l}^p)$$

and, thus

$$M_{d+2}^- \leq Q(C_t | Z_{t-dm-l}^p) \leq M_{d+2}^+.$$

Therefore, from Lemma 7, we have

$$|Q(C_t | Z_{t-dm-l}^p) - Q(C_t | Z_{t-(d+1)m-l}^p)|$$
$$\leq M_{d+2}^+ - M_{d+2}^- \leq (\rho_{\delta,k,m})^{d+1}.$$

Note that the result does not depend on either $Q$ or $l$. $\square$

APPENDIX B
PROOF OF LEMMA 3

Before proving Lemma 3 we need the following lemma first. Parts b)-d) are crucial for Lemma 3, and Part a) enables Part b). The continuity result of $\hat{X}_Q^\epsilon(z_{-t}^0)$ in Part a) is the key reason why we need the randomization of the filter.

*Lemma 9:* Suppose $Q \in \Theta_k^\delta$ and fix $\delta > 0$.
a) We have

$$\|\hat{X}_Q^\epsilon(z_{-t_1}^0) - \hat{X}_Q^\epsilon(z_{-t_2}^0)\|_1 \\ \leq M^2 \cdot \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_1$$

where $t_1, t_2 > 0$ are arbitrary integers. That is, for any integer $t > 0$ and any individual sequence $z_{-t}^0$, $\hat{X}_Q^\epsilon(z_{-t}^0)$ is a Lipshitz continuous function in $\boldsymbol{Q}_{X_0|z_{-t}^0}$.
b) $\ell(X_0, \hat{X}_Q^\epsilon(Z_{-t}^0)) \rightarrow \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))$ a.s. uniformly on $\Theta_k^\delta$.
c) For all $Q \in \Theta_k^\delta$, and for all $\omega$, there exist $0 < \gamma < 1, \beta > 0$, such that $|Q(X_0|Z_{-t}^0) - Q(X_0|Z_{-\infty}^0)| < \beta\gamma^t$.
d) For fixed $t, \eta > 0$, there exists some finite set $\mathcal{F}_k(t, \eta) \subset \Theta_k^\delta$, such that

$$\max_{Q \in \Theta_k^\delta} \min_{Q' \in \mathcal{F}_k(t,\eta)} \max_{x_0, z_{-t}^0} |Q(x_0|z_{-t}^0) - Q'(x_0|z_{-t}^0)| \leq \eta.$$

*Proof:*
a) For given simplex vector $\boldsymbol{Q}$, fixed $\hat{x}$, and $B_\epsilon$ defined as in Section IV-A, we define followings:

- $S_{\hat{x}}(\boldsymbol{Q}) \triangleq \{\boldsymbol{W} \in B_\epsilon : B(\boldsymbol{Q} + \boldsymbol{W}) = \hat{x}\}$;
- $DP(\hat{x}) \triangleq \{\boldsymbol{c} : \boldsymbol{c} = \boldsymbol{\lambda}_{\hat{x}} - \boldsymbol{\lambda}_a, \text{ for all } a \in \mathcal{A}\setminus\{\hat{x}\}\}$;
- $\text{dist}(\boldsymbol{Q}, \boldsymbol{c}^T\boldsymbol{y} = 0) \triangleq \min_{\boldsymbol{y} \in \{\boldsymbol{y} \in \mathbb{R}^M : \boldsymbol{c}^T\boldsymbol{y}=0\}} \|\boldsymbol{Q} - \boldsymbol{y}\|_2$.

In words, $S_{\hat{x}}(\boldsymbol{Q})$ is a set of vectors in $\epsilon$-ball, $B_\epsilon$, that makes the Bayes response $B(\boldsymbol{Q}+\boldsymbol{W})$ equal to $\hat{x}$; $DP(\hat{x})$ is a set of normal vectors that define the decision planes $\{\boldsymbol{y} \in \mathbb{R}^M : \boldsymbol{c}^T\boldsymbol{y} = 0, \boldsymbol{c} \in DP(\hat{x})\}$ which separate the reconstruction alphabet $\hat{x}$ and other alphabets, and $\text{dist}(\boldsymbol{Q}, \boldsymbol{c}^T\boldsymbol{y} = 0)$ is the shortest $L_2$ distance from a simplex vector $\boldsymbol{Q}$ to the plane $\{\boldsymbol{y} \in \mathbb{R}^M : \boldsymbol{c}^T\boldsymbol{y} = 0\}$. Then, for some fixed $t$, by definition

$$\hat{X}_Q^\epsilon(z_{-t}^0)[\hat{x}] = \frac{\text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t}^0}))}{\text{Vol}(B_\epsilon)}$$

where $\text{Vol}(\cdot)$ is a volume of a set. Since $\text{Vol}(B_\epsilon)$ is a constant, for any $t_1$ and $t_2$, we have

$$|\hat{X}_Q^\epsilon(z_{-t_1}^0)[\hat{x}] - \hat{X}_Q^\epsilon(z_{-t_2}^0)[\hat{x}]| \\ = \frac{|\text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_1}^0})) - \text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_2}^0}))|}{\text{Vol}(B_\epsilon)}. \quad (39)$$

For the numerator, as a crude bound, we obtain

$$|\text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_1}^0})) - \text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_2}^0}))| \\ \leq \text{Vol}(B_\epsilon^{M-1}) \sum_{\boldsymbol{c} \in DP(\hat{x})} \Big|\text{dist}(\boldsymbol{Q}_{X_0|z_{-t_1}^0}, \boldsymbol{c}^T\boldsymbol{y} = 0) \\ - \text{dist}(\boldsymbol{Q}_{X_0|z_{-t_2}^0}, \boldsymbol{c}^T\boldsymbol{y} = 0)\Big|, \quad (40)$$

where $B_\epsilon^{M-1} = \{\boldsymbol{U} \in \mathbb{R}^{M-1} : \|\boldsymbol{U}\|_2 \leq \epsilon\}$. Since

$$\text{dist}(\boldsymbol{Q}, \boldsymbol{c}^T\boldsymbol{y} = 0) = \frac{|\boldsymbol{c}^T\boldsymbol{Q}|}{\|\boldsymbol{c}\|_2}$$

we have

$$\text{dist}(\boldsymbol{Q}_{X_0|z_{-t_1}^0}, \boldsymbol{c}^T\boldsymbol{y} = 0) - \text{dist}(\boldsymbol{Q}_{X_0|z_{-t_2}^0}, \boldsymbol{c}^T\boldsymbol{y} = 0)$$
$$= \frac{|\boldsymbol{c}^T\boldsymbol{Q}_{X_0|z_{-t_1}^0}| - |\boldsymbol{c}^T\boldsymbol{Q}_{X_0|z_{-t_2}^0}|}{\|\boldsymbol{c}\|_2}$$
$$\leq \frac{\left|\boldsymbol{c}^T(\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0})\right|}{\|\boldsymbol{c}\|_2} \quad (41)$$
$$\leq \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_2 \quad (42)$$
$$\leq \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_1 \quad (43)$$

where (41) follows from the triangle inequality; (42) follows from Cauchy–Schwartz inequality, and (43) follows from the fact that the $L_2$-norm is less than or equal to the $L_1$-norm. Therefore, (40) becomes

$$|\text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_1}^0})) - \text{Vol}(S_{\hat{x}}(\boldsymbol{Q}_{X_0|z_{-t_2}^0}))| \\ \leq M \cdot \text{Vol}(B_\epsilon^{M-1}) \cdot \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_1$$

and, thus, (39) becomes

$$|\hat{X}_Q^\epsilon(z_{-t_1}^0)[\hat{x}] - \hat{X}_Q^\epsilon(z_{-t_2}^0)[\hat{x}]| \\ \leq M \cdot \frac{\text{Vol}(B_\epsilon^{M-1})}{\text{Vol}(B_\epsilon)} \cdot \|\boldsymbol{Q}(X_0|z_{-t_1}^0) - \boldsymbol{Q}(X_0|z_{-t_2}^0)\|_1 \\ \leq M \cdot \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_1.$$

Therefore, we have

$$\|\hat{X}_Q^\epsilon(z_{-t_1}^0) - \hat{X}_Q^\epsilon(z_{-t_2}^0)\|_1 \\ \leq M^2 \cdot \|\boldsymbol{Q}_{X_0|z_{-t_1}^0} - \boldsymbol{Q}_{X_0|z_{-t_2}^0}\|_1$$

and Part a) is proved.
b) By the exact same argument as in proving Lemma 1, we can easily know that $Q(X_0|Z_{-t}^0) \rightarrow Q(X_0|Z_{-\infty}^0)$ for all $\omega$, uniformly on $\Theta_k^{\delta_k}$. Since we have

$$\left|\ell(X_0, \hat{X}_Q^\epsilon(Z_{-t}^0)) - \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right| \\ = \left|\sum_{\hat{x}} \Lambda(X_0, \hat{x})\left(\hat{X}_Q^\epsilon(Z_{-t}^0)[\hat{x}] - \hat{X}_Q^\epsilon(Z_{-\infty}^0)[\hat{x}]\right)\right| \\ \leq \Lambda_{\max}\|\hat{X}_Q^\epsilon(Z_{-t}^0) - \hat{X}_Q^\epsilon(Z_{-\infty}^0)\|_1 \\ \leq \Lambda_{\max}M^2 \cdot \|\boldsymbol{Q}(X_0|Z_{-t}^0) - \boldsymbol{Q}(X_0|Z_{-\infty}^0)\|_1$$

we get the uniform convergence.
c) Again, let us follow the argument in the proof of Lemma 1. Suppose $t = jk + l$, where $j = \lfloor t/k \rfloor$, and $l = t \mod k$. Then

$$|Q(X_0|Z_{-t}^0) - Q(X_0|Z_{-\infty}^0)| \\ = |Q(X_0|Z_{-jk-l}^0) - Q(X_0|Z_{-\infty}^0)| \\ \leq \sum_{i=j}^\infty |Q(X_0|Z_{-ik-l}^0) - Q(X_0|Z_{-(i+1)k-l}^0)|$$

$$\leq \sum_{i=j}^{\infty} \rho^{i+1} \tag{44}$$

$$= \frac{\rho^{j+1}}{1-\rho} = \frac{\rho}{1-\rho} \rho^{\lfloor t/k \rfloor} = \frac{\rho^{1-\frac{t}{k}}}{1-\rho} (\rho^{1/k})^t \tag{45}$$

$$\leq \frac{1}{1-\rho} (\rho^{1/k})^t \tag{46}$$

where $\rho = \rho_{\delta,k,k}$ as defined in Lemma 7, and (44) follows from Lemma 8. By letting $\beta = \frac{1}{1-\rho}$, and $\gamma = \rho^{1/k}$, we have proved Part c).

d) We know that for the individual sequence pair $(x_0, z_{-t}^0)$

$$Q(x_0 | z_{-t}^0) = \frac{\sum_{x_{-t}^{-1}} Q(x_{-t}^0, z_{-t}^0)}{Q(z_{-t}^0)}$$

$$= \frac{\sum_{x_{-t}^{-1}} Q(x_{-t}^0, z_{-t}^0)}{\sum_{x_{-t}^0} Q(x_{-t}^0, z_{-t}^0)}$$

$$= \frac{\sum_{x_{-t}^{-1}} Q(x_{-t}^0) Q(z_{-t}^0 | x_{-t}^0)}{\sum_{x_{-t}^0} Q(x_{-t}^0) Q(z_{-t}^0 | x_{-t}^0)}$$

$$= \frac{\sum_{x_{-t}^{-1}} \left( Q(x_{-t}^0) \prod_{i=-t}^0 \Pi(x_i, z_i) \right)}{\sum_{x_{-t}^0} \left( Q(x_{-t}^0) \prod_{i=-t}^0 \Pi(x_i, z_i) \right)}.$$

For $Q \in \Theta_k^\delta$, $\mathbf{\Pi}$ is fixed and we can think of $\prod_{i=-t}^0 \Pi(x_i, z_i)$ as a constant for the individual sequence pair $(x_{-t}^0, z_{-t}^0)$. Since

$$Q(x_{-t}^0) = Q(x_{-t}^{k-1-t}) \prod_{j=k-t}^0 a_{x_{j-k}^{j-1} x_{j-k+1}^j}$$

$Q(x_0 | z_{-t}^0)$ is the ratio of two finite-order polynomials of $\{a_{ij}\}$, and as $\Theta_k^\delta$ is closed and bounded, $Q(x_0 | z_{-t}^0)$ is a uniformly continuous function of $\{a_{ij}\}$. Therefore, for given $\eta$, $\exists \epsilon(\eta)$ such that $\|Q - Q'\|_1 < \epsilon(\eta)$ implies

$$\max_{x_0, z_{-t}^0} |Q(x_0 | z_{-t}^0) - Q'(x_0 | z_{-t}^0)| \leq \eta$$

since there are only finite numbers of possible $(x_0, z_{-t}^0)$ pairs. Also, since $\Theta_k^\delta$ is compact, we can always find a finite set $\mathcal{F}_k(t, \eta)$ that for any $Q \in \Theta_k^\delta$, there exists at least one $Q' \in \mathcal{F}_k(t, \eta)$ that satisfies $\|Q - Q'\|_1 < \epsilon(\eta)$. Therefore, Part d) is proved. $\qquad \square$

*Proof of Lemma 3:* To prove Lemma 3, first consider following limit:

$$\lim_{n \to \infty} E\left( L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) \right)$$

$$= \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n E\left( \ell(X_t, \hat{X}_Q^\epsilon(Z^t)) \right)$$

$$= \lim_{t \to \infty} E\left( \ell(X_t, \hat{X}_Q^\epsilon(Z^t)) \right) \tag{47}$$

$$= \lim_{t \to \infty} E\left( \ell(X_0, \hat{X}_Q^\epsilon(Z_{-(t-1)}^0)) \right) \tag{48}$$

$$= E\left( \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0)) \right) \text{uniformly on } \Theta_k^\delta \tag{49}$$

where (47) is from Cesáro's mean convergence theorem, (48) is from stationarity, and (49) is from Lemma 9 Part b) and bounded convergence theorem. Thus, to complete the proof, we need to show that

$$\lim_{n \to \infty} L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) = E\left( \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0)) \right) \quad \text{a.s.} \tag{50}$$

uniformly on $\Theta_k^\delta$. Now, let us show the pointwise convergence in (50) without the uniformity by using ergodic theorem. For given $Q$, define

$$g_{t,Q}(X, Z) \triangleq \ell(X_0, \hat{X}_Q^\epsilon(Z_{-(t-1)}^0))$$

$$g_Q(X, Z) \triangleq \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))$$

and denote by $T$ the shift operator. Then, what we should prove becomes

$$\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n g_{t,Q}(T^t(X, Z)) = E\left( g_Q(X, Z) \right) \quad \text{a.s.}$$

while the ergodic theorem gives

$$\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^n g_Q(T^t(X, Z)) = E\left( g_Q(X, Z) \right) \quad \text{a.s.}$$

Observe that

$$\left| \frac{1}{n} \sum_{t=1}^n g_{t,Q}(T^t(X, Z)) - \frac{1}{n} \sum_{t=1}^n g_Q(T^t(X, Z)) \right|$$

$$\leq \frac{1}{n} \sum_{t=1}^n \left| g_{t,Q}(T^t(X, Z)) - g_Q(T^t(X, Z)) \right|$$

$$= \frac{1}{n} \sum_{t=1}^n \left| \ell(X_t, \hat{X}_Q^\epsilon(Z_1^t)) - \ell(X_t, \hat{X}_Q^\epsilon(Z_{-\infty}^t)) \right|.$$

Since Lemma 9 Part c) holds for all $\omega$, we can think that the lemma holds for all individual sequence pairs $(x_0, z_{-\infty}^0)$. Thus, it holds for all individual pairs $(x_t, z_{-\infty}^t)$ as well, and we can conclude that $Q(X_t | Z_1^t) \to Q(X_t | Z_{-\infty}^t)$ for all $\omega$ as $t \to \infty$. Hence, by exactly the same argument as Lemma 9 Part a) and Lemma 9 Part b), we conclude that $\ell(X_t, \hat{X}_Q^\epsilon(Z_1^t)) \to \ell(X_t, \hat{X}_Q^\epsilon(Z_{-\infty}^t))$ almost surely as $t \to \infty$. Now, by Cesáro's mean convergence theorem, we obtain

$$\frac{1}{n} \sum_{t=1}^n \left| \ell(X_t, \hat{X}_Q^\epsilon(Z_1^t)) - \ell(X_t, \hat{X}_Q^\epsilon(Z_{-\infty}^t)) \right| \to 0 \quad \text{a.s.}$$

Therefore, we get

$$L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) \to E\left( \ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0)) \right) \quad \text{a.s.}$$

Note that up to this point we cannot guarantee the uniformity of the convergence, since the ergodic theorem only gives the individual convergence for each $Q$. To show the uniformity of the convergence in (50), first define the following quantity for some fixed integer $t \in [1, n-1]$:

$$L_{\hat{\boldsymbol{X}}_{Q,t}^\epsilon}(X^n, Z^n)$$

$$= \frac{1}{n} \left( \sum_{i=1}^t \ell(X_i, \hat{X}_Q^\epsilon(Z^i)) + \sum_{i=t+1}^n \ell(X_i, \hat{X}_Q^\epsilon(Z_{i-t}^i)) \right).$$

From Lemma 9 Part d), for any $Q \in \Theta_k^\delta$ and fixed $t, \eta > 0$, we can pick some $Q' \in \mathcal{F}_k(t, \eta)$ such that $\|Q - Q'\|_1 < \epsilon(\eta)$, and thus

$$\max_{x_0, z_{-t}^0} |Q(x_0|z_{-t}^0) - Q'(x_0|z_{-t}^0)| \le \eta.$$

By adding and subtracting some common terms involving such $Q'$, and from the triangle inequality, we have

$$\left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right|$$
$$\le \left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q,t}^\epsilon}(X^n, Z^n) \right|$$
$$+ \left| L_{\hat{\boldsymbol{X}}_{Q,t}^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q',t}^\epsilon}(X^n, Z^n) \right|$$
$$+ \left| L_{\hat{\boldsymbol{X}}_{Q',t}^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q'}^\epsilon}(X^n, Z^n) \right|$$
$$+ \left| L_{\hat{\boldsymbol{X}}_{Q'}^\epsilon}(X^n, Z^n) - E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-\infty}^0))\right) \right|$$
$$+ \left| E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-\infty}^0))\right) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right|. \quad (51)$$

Now, the goal becomes to show that the terms in the right-hand side of the inequality converges to zero independent of $Q$ as $n$, $t$, and $\eta$ vary. First, we will bound each term, and send $n \to \infty$.

i)

$$\left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q,t}^\epsilon}(X^n, Z^n) \right|$$
$$\le \frac{1}{n} \sum_{i=t+1}^n \left| \ell(X_i, \hat{X}_Q^\epsilon(Z^i)) - \ell(X_i, \hat{X}_Q^\epsilon(Z_{i-t}^i)) \right|$$
$$\le \Lambda_{\max} \cdot \frac{1}{n} \sum_{i=t+1}^n \|\hat{X}_Q^\epsilon(Z^i) - \hat{X}_Q^\epsilon(Z_{i-t}^i)\|_1$$
$$\le \Lambda_{\max} M^2 \cdot \frac{1}{n} \sum_{i=t+1}^n \|\boldsymbol{Q}_{X_0|Z_{-i}^0} - \boldsymbol{Q}_{X_0|Z_{-t}^0}\|_1 \quad (52)$$
$$\le \Lambda_{\max} M^3 \cdot \frac{1}{n} \sum_{i=t+1}^n (\beta\gamma^t + \beta\gamma^i) \quad (53)$$
$$\to \Lambda_{\max} M^3 \beta\gamma^t \quad \text{a.s. uniformly on } \Theta_k^\delta \quad (54)$$

where (52) follows from stationarity and Lemma 9 Part a); (53) follows from Lemma 9 Part c), and (54) follows from the Cesáro's mean convergence theorem. Since (53) does not depend on $Q$, the limit is uniform on $\Theta_k^\delta$.

ii)

$$\left| L_{\hat{\boldsymbol{X}}_{Q,t}^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q',t}^\epsilon}(X^n, Z^n) \right|$$
$$\le \frac{1}{n} \sum_{i=t+1}^n |\ell(X_i, \hat{X}_Q^\epsilon(Z_{i-t}^i)) - \ell(X_i, \hat{X}_{Q'}^\epsilon(Z_{i-t}^i))| + \frac{t\Lambda_{\max}}{n}$$
$$\le \Lambda_{\max} \cdot \frac{1}{n} \sum_{i=t+1}^n \|\hat{X}_Q^\epsilon(Z_{i-t}^i) - \hat{X}_{Q'}^\epsilon(Z_{i-t}^i)\|_1 + \frac{t\Lambda_{\max}}{n}$$
$$\le \Lambda_{\max} M^2 \cdot \frac{1}{n} \sum_{i=t+1}^n \|\boldsymbol{Q}_{X_i|Z_{i-t}^i} - \boldsymbol{Q}'_{X_i|Z_{i-t}^i}\|_1 + \frac{t\Lambda_{\max}}{n} \quad (55)$$
$$\le \Lambda_{\max} M^3 \frac{n-t}{n} \cdot \eta + \frac{t\Lambda_{\max}}{n}$$
$$\to \Lambda_{\max} M^3 \eta \quad \text{a.s. uniformly on } \Theta_k^\delta \quad (56)$$

where (55) follows from Lemma 9 Part a), and (56) follows from Lemma 9 Part d). Since (56) does not depend on $Q$, the limit is also uniform on $\Theta_k^\delta$.

iii)

$$\left| L_{\hat{\boldsymbol{X}}_{Q',t}^\epsilon}(X^n, Z^n) - L_{\hat{\boldsymbol{X}}_{Q'}^\epsilon}(X^n, Z^n) \right|$$
$$\to \Lambda_{\max} M^3 \beta\gamma^t \quad \text{a.s.}$$

by following the same argument as i). Since $\mathcal{F}_k(t, \eta)$ is finite, this convergence is uniform on $\mathcal{F}_k(t, \eta)$.

iv)

$$\left| L_{\hat{\boldsymbol{X}}_{Q'}^\epsilon}(X^n, Z^n) - E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-\infty}^0))\right) \right| \to 0 \quad \text{a.s.}$$

from the proof of pointwise convergence above. As in iii), this convergence is also uniform on $\mathcal{F}_k(t, \eta)$.

v)

$$\left| E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-\infty}^0))\right) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right|$$
$$\le \left| E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-\infty}^0))\right) - E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-t}^0))\right) \right|$$
$$+ \left| E\left(\ell(X_0, \hat{X}_{Q'}^\epsilon(Z_{-t}^0))\right) - E[\ell(X_0, \hat{X}_Q^\epsilon(Z_{-t}^0))] \right|$$
$$+ \left| E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-t}^0))\right) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right|$$
$$\le \sum_{x_0, z_{-\infty}^0} P(x_0, z_{-\infty}^0) \left| \ell(x_0, \hat{X}_{Q'}^\epsilon(z_{-\infty}^0)) - \ell(x_0, \hat{X}_{Q'}^\epsilon(z_{-t}^0)) \right|$$
$$+ \sum_{x_0, z_{-t}^0} P(x_0, z_{-t}^0) \left| \ell(x_0, \hat{X}_{Q'}^\epsilon(z_{-t}^0)) - \ell(x_0, \hat{X}_Q^\epsilon(z_{-t}^0)) \right|$$
$$+ \sum_{x_0, z_{-\infty}^0} P(x_0, z_{-\infty}^0) \left| \ell(x_0, \hat{X}_Q^\epsilon(z_{-\infty}^0)) - \ell(x_0, \hat{X}_Q^\epsilon(z_{-t}^0)) \right|$$
$$\le \Lambda_{\max} M^3 \left(2\beta\gamma^t + \eta\right)$$

by similar argument as in i) and ii).

Therefore, by taking limit supremum on both side of (51), we get

$$\limsup_{n\to\infty} \left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right|$$
$$\le \Lambda_{\max} M^3 \left(4\beta\gamma^t + 2\eta\right) \quad \text{a.s.}$$

uniformly on $\Theta_k^\delta$. Since $t$ and $\eta$ are arbitrary, by sending $t \to \infty$ and $\eta \downarrow 0$, we have

$$\limsup_{n\to\infty} \left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - E\left(\ell(X_0, \hat{X}_Q^\epsilon(Z_{-\infty}^0))\right) \right| \le 0 \quad \text{a.s.}$$

uniformly on $\Theta_k^\delta$. Thus, the lemma is proved. $\square$

## APPENDIX C
## PROOF OF COROLLARY 1

*Proof of Corollary 1:* First note the subtle point that Corollary 1 does not directly follow from Lemma 3. Since the probability law $Q_k^t$ that we are using to filter each block is changing every block, whereas the uniform convergence in Lemma 3 is for the fixed $Q \in \Theta_k^{\delta_k}$ for all $t$, it is not enough to guarantee the corollary. However, since $Q_k^t$ remains the same within each block, we can still use the result of Lemma 3 if the block length gets long enough. Keeping this in mind, let us take a more
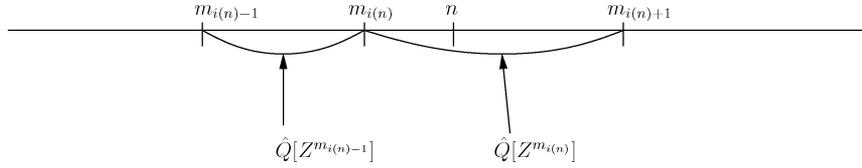
Fig. 1. The time line.

careful look at each block. In the proof, for brevity of notation, let us denote

$$\ell_t(Q) \triangleq \ell(X_t, \hat{X}_Q^\epsilon(Z^t))$$

since we are always dealing with the randomized filter, and there is no possibility of confusion. Now, fix any $\delta > 0$. Then, from (4), there exists some $I$, such that

$$\frac{m_{I-1}}{m_I} < \frac{\delta}{8\Lambda_{\max}}$$

and from Lemma 3, there exists some $N$ such that for all $n \geq N$

$$\max_{Q \in \Theta_K^{\delta_k}} \left| L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) - E L_{\hat{\boldsymbol{X}}_Q^\epsilon}(X^n, Z^n) \right| < \delta/4. \quad (57)$$

Recalling the definition $i(t) \triangleq \max\{i : m_i \leq t\}$, we let $I_0 = \max(I, i(N) + 1)$. Then, for any $n \geq m_{I_0}$ and $m_{i(n)} \leq n < m_{i(n)+1}$

$$\left| L_{\hat{\boldsymbol{X}}_{\mathrm{univ},k}^\epsilon}(X^n, Z^n) - \hat{E} L_{\hat{\boldsymbol{X}}_{\mathrm{univ},k}^\epsilon}(X^n, Z^n) \right|$$

$$\leq \frac{1}{n}\left| \sum_{t=1}^{m_{i(n)-1}} \left( \ell_t(Q_k^t) - \hat{E}(\ell_t(Q_k^t)) \right) \right|$$

$$+ \frac{1}{n}\left| \sum_{t=m_{i(n)-1}+1}^{m_{i(n)}} \left( \ell_t(\hat{Q}[Z^{m_{i(n)-1}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)-1}}])) \right) \right|$$

$$+ \frac{1}{n}\left| \sum_{t=m_{i(n)}+1}^{n} \left( \ell_t(\hat{Q}[Z^{m_{i(n)}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)}}])) \right) \right|. \quad (58)$$

Note that in the second and third terms, $Q_k^t$ is fixed to $\hat{Q}[Z^{m_{i(n)-1}}]$ and $\hat{Q}[Z^{m_{i(n)}}]$ from the definition of our filter. Fig. 1 summarizes the time line with the above notations.

Now, we can bound each term in (58). For the first term, since $n \geq m_{i(n)} \geq m_I$, we know that

$$\frac{m_{i(n)-1}}{n} \leq \frac{m_{i(n)-1}}{m_{i(n)}} < \frac{\delta}{8\Lambda_{\max}}.$$

Therefore

$$\frac{1}{n}\left| \sum_{t=1}^{m_{i(n)-1}} \left( \ell_t(Q_k^t) - \hat{E}(\ell_t(Q_k^t)) \right) \right| \leq \frac{\delta}{8\Lambda_{\max}} \cdot \Lambda_{\max} = \frac{\delta}{8}.$$

For the second term, since $n \geq m_{i(n)} \geq N$, and from (57)

$$\frac{1}{n}\left| \sum_{t=m_{i(n)-1}+1}^{m_{i(n)}} \left( \ell_t(\hat{Q}[Z^{m_{i(n)-1}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)-1}}])) \right) \right|$$

$$\leq \frac{m_{i(n)}}{n} \frac{1}{m_{i(n)}}\left| \sum_{t=1}^{m_{i(n)}} \left( \ell_t(\hat{Q}[Z^{m_{i(n)-1}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)-1}}])) \right) \right|$$

$$+ \frac{1}{n}\left| \sum_{t=1}^{m_{i(n)-1}} \left( \ell_t(\hat{Q}[Z^{m_{i(n)-1}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)-1}}])) \right) \right|$$

$$\leq \frac{\delta}{4} + \frac{\delta}{8\Lambda_{\max}} \cdot \Lambda_{\max} = \frac{3\delta}{8}.$$

Finally, for the last term

$$\frac{1}{n}\left| \sum_{t=m_{i(n)}+1}^{n} \left( \ell_t(\hat{Q}[Z^{m_{i(n)}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)}}])) \right) \right|$$

$$\leq \frac{1}{n}\left| \sum_{t=1}^{n} \left( \ell_t(\hat{Q}[Z^{m_{i(n)}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)}}])) \right) \right|$$

$$+ \frac{1}{n}\left| \sum_{t=1}^{m_{i(n)}} \left( \ell_t(\hat{Q}[Z^{m_{i(n)}}]) - \hat{E}(\ell_t(\hat{Q}[Z^{m_{i(n)}}])) \right) \right|$$

$$\leq \frac{\delta}{4} + \frac{\delta}{4} = \frac{\delta}{2}.$$

Therefore, for any $n \geq m_{I_0}$ and $m_{i(n)} \leq n \leq m_{i(n)+1}$, we have

$$\left| L_{\hat{\boldsymbol{X}}_{\mathrm{univ},k}^\epsilon}(X^n, Z^n) - \hat{E} L_{\hat{\boldsymbol{X}}_{\mathrm{univ},k}^\epsilon}(X^n, Z^n) \right| < \delta,$$

and since $\delta$ was arbitrary, we have the corollary. $\qquad\square$

## REFERENCES

[1] L. E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains," *Ann. Math. Statist.*, vol. 37, pp. 1554–1563, 1966.
[2] L. E. Baum, T. Petrie, G. Soules, and N. Weiss, "A maximization technique occuring in the statistical analysis of probabilistic functions of Markov chains," *Ann. Math. Statist.*, vol. 41, pp. 164–171, 1970.
[3] P. J. Bickel, Y. Ritov, and T. Rydén, "Asymptotic normality of the maxumum-likelihood estimator for general hidden Markov models," *Ann. Statist.*, vol. 26, no. 4, pp. 1614–1635, 1998.
[4] R. W. Chang and J. C. Hancock, "On receiver structures for channels having memory," *IEEE Trans. Inf. Theory*, vol. IT-12, no. 4, pp. 463–468, Oct. 1966.
[5] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
[6] A. Dembo and T. Weissman, "Universal denoising for the finite-input general-output channel," *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1507–1517, Apr. 2005.

[7] Y. Ephraim and N. Merhav, "Hidden Markov processes," *IEEE Trans. Inf. Theory*, vol. 48, no. 6, pp. 1518–1569, Jun. 2002.

[8] L. Finesso, "Consistent Estimation of the Order for Markov and Hidden Markov Chains," Ph.D. dissertation, Univ. Maryland, College Park, 1990.

[9] G. M. Gemelos, S. Sigurjonsson, and T. Weissman, "Universal minimax discrete denoising under channel uncertainty," *IEEE Trans. Inf. Theory*, vol. 52, no. 8, pp. 3476–3497, Aug. 2006.

[10] G. H. Golub and C. D. Meyer, "Using the QR factorization and group inversion to compute, differentiate, and estimate the sensitivity of stationary probabilities for Markov chains," *SIAM J. Alg. Discr. Meth.*, vol. 7, pp. 273–281, 1986.

[11] J. C. Kieffer, "Strongly consistent code-based identification and order estimation for constrained finite-state model classes," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 893–902, May 1993.

[12] B. G. Leroux, "Maximum-likelihood estimation for hidden Markov models," *Stochastic Processes Their Applic.*, vol. 40, pp. 127–143, 1992.

[13] C.-C. Liu and P. Narayan, "Order estimation and sequential universal data compression of a hidden Markov source by the model of mixtures," *IEEE Trans. Infm. Theory*, vol. 40, no. 4, pp. 1167–1180, Jul. 1994.

[14] N. Merhav and M. Feder, "Universal prediction," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2124–2147, Oct. 1998.

[15] N. Feder, N. Merhav, and M. Gutman, "Universal prediction for individual sequences," *IEEE Trans. Inf. Theory*, vol. 38, no. 4, pp. 1258–1270, Jul. 1992.

[16] L. Mevel and L. Finesso, "Asymptotical statistics of misspecified hidden Markov models," *IEEE Trans. Autom. Control*, vol. 49, no. 7, pp. 1123–1132, Jul. 2004.

[17] T. Moon and T. Weissman, "Discrete universal filtering via hidden Markov modeling," in *Proc IEEE Int. Symp. Information Theory*, Adelaide, Australia, Sep. 2005, pp. 1285–1289.

[18] E. Ordentlich, T. Weissman, M. Weinberger, A. Somekh-Baruch, and N. Merhav, "Discrete universal filtering through incremental parsing," in *Proc. Data Compression Conf. (DCC 2004)*, Snowbird, UT, Mar. 2004, pp. 352–361.

[19] T. Weissman, E. Ordentlich, M. Weinberger, A. Somekh-Baruch, and N. Merhav, "Universal filtering via prediction," *IEEE Trans. Inf. Theory*, vol. 53, no. 4, pp. 1253–1264, Apr. 2007.

[20] T. Rydén, "Estimating the order of hidden Markov models," *Statistics*, vol. 26, pp. 345–354, 1995.

[21] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdú, and M. Weinberger, "Universal discrete denoising: Known channel," *IEEE Trans. Inf. Theory*, vol. 51, no. 1, pp. 5–28, Jan. 2005.

[22] R. Zhang and T. Weissman, "Discrete denoising for channels with memory," *Commun. Inf. Syst.*, vol. 5, no. 2, pp. 257–288, 2005.

[23] J. Ziv and N. Merhav, "Estimating the number of states of a finite-state source," *IEEE Trans. Inf. Theory*, vol. 38, no. 1, pp. 61–65, Jan. 1992.

[24] S. B. Vardeman, "Admissible solutions of $k$-extended finite state set and the sequence compound decision problems," *J. Multiv. Anal.*, vol. 10, pp. 426–441, 1980.

[25] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proc. IEEE*, vol. 77, no. 2, pp. 257–286, Feb. 1989.

[26] S. Sitbon and J. M. Passerieux, "New efficient target tracking based upon hidden Markov model and probabilistic data association," in *Proc. 29th Asilomar Conf. Signals, Systems and Computers*, Pacific Grove, CA, Nov. 1995, p. 849.

[27] P. Boufounos, S. El-Difrawy, and D. Ehrlich, "Hidden Markov models for DNA sequencing," in *Proc. Workshop on Genomic Signal Processing and Statistics (GENSIPS 2002)*, Oct. 2002.

[28] Course Website for CS 294 Univ. Calif. Berkeley [Online]. Available: http://www.cs.berkeley.edu/~asimma/294-fall06/lectures/hidden-Markov.html