

Sliding Finite Window Codes: Near-Optimality and Q-Learning for Zero-Delay Coding

Liam Cregg¹, Fady Alajaji², *Senior Member, IEEE*, and Serdar Yüksel³, *Senior Member, IEEE*

Abstract—We study the problem of zero-delay coding for the transmission of a Markov source over a noisy channel with feedback and present a reinforcement learning solution which is guaranteed to approach optimality. To this end, we formulate the problem as a Markov decision process (MDP) where the state is a probability-measure valued predictor/belief and the actions are quantizer maps. This MDP formulation has been used to show the optimality of certain classes of encoder policies in prior work, but their computation is prohibitively complex due to the uncountable nature of the constructed state space. Based on recent results for partially observed MDPs, we present an approximation of the belief MDP using a sliding finite window of channel outputs and quantizers. Under an appropriate notion of predictor stability, we show that the lowest distortion achievable by such a sliding finite window policy approaches the true lowest distortion as the window length increases. We give sufficient conditions for predictor stability to hold. Finally, we propose a Q-learning algorithm which provably converges to the optimal policy and provide a detailed comparison of the sliding finite window scheme with another approximation scheme which quantizes the belief MDP in a nearest neighbor fashion, as well as other coding schemes from the literature.

Index Terms—Zero-delay coding, reinforcement learning, Q-learning, stochastic control, Markov decision processes, quantization, source-channel coding, noisy channels with feedback.

I. INTRODUCTION

THE zero-delay coding problem involves compressing and transmitting an information source at a fixed rate over a noisy channel with feedback and without delay, while minimizing the expected distortion at the receiver. This zero-delay restriction is of practical relevance in many applications, including live-streaming [1], [2], [3], [4], real-time tracking and estimating processes over erasure links [5] and real-time sensor networks [6], [7]. However, this restriction means that classical Shannon-theoretic methods [8], which require collecting large sequences of source symbols and compressing them at once, are not viable as they induce a large delay.

Received 26 February 2026; accepted 20 March 2026. Date of publication 6 April 2026; date of current version 21 May 2026. This work was supported by the Natural Sciences and Engineering Research Council of Canada. An earlier version of this paper was presented in part at the 2023 Conference on Decision and Control and in part at the 2024 American Control Conference [DOI: 10.1109/CDC49753.2023.10383642 and DOI: 10.1109/CDC49753.2023.10383642]. (*Corresponding author: Liam Cregg.*)

Liam Cregg is with the Department of Information Technology and Electrical Engineering, 8092 ETH Zürich, Switzerland (e-mail: lcregg@ethz.ch).

Fady Alajaji and Serdar Yüksel are with the Department of Mathematics and Statistics, Queen's University, Kingston, ON K7L 3N6, Canada (e-mail: fa@queensu.ca; yuksel@queensu.ca).

Communicated by V. Kostina, Associate Editor for Communications.

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TIT.2026.3681183>.

Digital Object Identifier 10.1109/TIT.2026.3681183

From an information-theoretic perspective, several strategies have been used to approach this problem, including mutual information constraints, entropy coding, and Shannon lower bounding techniques. Studies to this end include [9], [10], [11]. Within the context of linear systems, [12], [13], [14], [16] use sequential rate-distortion theory. Some of these works give applicable codes for zero-delay coding for Gaussian sources over additive-noise Gaussian channels, and some give upper and/or lower performance bounds; see [17] and [18] for further studies, and see [19], [20], [21] for studies in the context of channels with delay and other communication constraints.

Furthermore, learning theoretic methods have attracted significant interest in source-channel coding theory both in the classical literature and the recent literature, see for example [22], [23], [24] for the noiseless channel (quantization) case among several classical results, although usually restricted to independent and identically distributed (i.i.d.) sources. We note that our results are directly applicable for i.i.d. sources as well, since an optimal zero-delay code for an i.i.d. source is a memoryless code [25], [26], [27] (see also, for related discussions in a different causal coding context [28], [29], [30], [31]). More recently, deep learning is employed to construct powerful joint source-channel codes (see [32], [33], [34]), and reinforcement learning is used as a tool to estimate feedback capacity in [35] and [36]. Although effective in practice, these machine learning methods are generally experimental and do not provide a formal proof of convergence or optimality. Conversely, our reinforcement learning approach will be rigorously shown to converge to *near-optimality*, in the sense that, for any $\epsilon > 0$, there exists some window size N such that the resulting policy achieves distortion no further than ϵ from the optimal distortion (a formal definition is given in Definition 1).

There have been several studies about the zero-delay coding problem using stochastic control techniques. In particular, [25], [26], [37] consider Markov sources with finite alphabets and finite time horizons and show optimality of structured classes of policies. Similar optimality and existence results are presented for infinite time horizons in [27] (with feedback) and [38] (without feedback). The continuous-alphabet infinite-horizon case is examined in [39], although only over a noiseless channel. These results often rely on formulating the problem as a Markov decision process (MDP) in order to utilize existing results from stochastic control theory, such as dynamic programming and value iteration methods (see [40], [41] for detailed information on such methods). However, in the formulation of the MDP, these results utilize a state space

that is probability measure-valued (this state is often called the “predictor” in the literature) and an action space involving quantizers. These spaces are computationally difficult to work with, both in terms of complexity and implementation. Thus, while numerous existence and structural results have been established for this problem, the explicit development of effective coding schemes for a given zero-delay coding problem is still an open problem.

In this paper, we present an approximation method to simplify the resulting MDP, and we use this approximation to obtain near-optimal coding schemes for the zero-delay coding problem via a reinforcement learning approach. We emphasize that we provide guaranteed approximation and convergence results. In particular, we build on methods developed recently in [42], which were originally used to study partially observed Markov decision processes (POMDPs). The analysis requires significant adaptation for our setup as the measurement channels are selected by an encoding policy and are not constant (depending on the realization of the predictor measures). Based on these methods, we introduce a practical sliding finite window method and present several mathematical and algorithmic results on its near-optimal performance. We also compare this method with the approach used in [43], which uses a nearest neighbor quantization of the probability measure-valued state space. Although [43] only studied the noiseless channel (quantization) version of the problem, extensions to the noisy channel case follow with little additional effort as we discuss in the paper.

We emphasize that our new sliding finite window method has several advantages over the belief quantization scheme. In particular, it is easier to implement, computationally less complex, and valid for any initialization. This comes at a cost of some additional Dobrushin coefficient conditions on the source and channel, but we note that the conditions provided here are sufficient, not necessary. These differences can be seen by comparing Theorems 4 and 6, and are further detailed in Section V.

The sliding finite window code that we propose can be viewed as a finite state code [44], [45] where the states are the realizations of a finite window of encoder maps and outputs. We are not aware of such a construction in the literature with rigorous performance guarantees in the context of optimal real-time coding and in particular on an associated learning theoretic study. We also note the similarities to trellis or convolutional codes [46], [47], [48]; our method can be in some sense be seen as a zero-delay analog of these codes, but our analysis is based on stochastic control and thus entirely different than that of trellis codes.

A. Main Contributions

- To our knowledge, there is no prior study on near-optimality of sliding window coding schemes for the zero-delay coding problem with an explicit rate of convergence (involving all window lengths). To this end, in Theorem 2, we provide an approximation result for the zero-delay coding problem which bounds the sub-optimality gap of sliding finite window policies by a predictor stability term.

- Under Dobrushin coefficient conditions, we then show that these sliding finite window policies become near-optimal for the zero-delay coding problem in Corollary 2, with an explicit performance bound. We note that these conditions hold when the channel is “noisy enough” or when the source has some mixing conditions. To our knowledge, this gives the first rigorous proof of optimality of sliding finite window policies for the zero-delay coding problem.
- We provide a (reinforcement) Q-learning algorithm which allows for the computation of these sliding finite window policies and rigorously show its convergence in Theorem 4.
- Finally, we provide both theoretical and experimental comparisons of our sliding finite window method with a nearest neighbor approximation scheme presented in [43] (and adapted and generalized to the noisy channel setting). We also compare our algorithm to several other zero-delay schemes in the literature.

Notation. In general, we will denote random variables by capital letters and their realizations by lowercase letters. There are a few exceptions to this; in particular we will always use lowercase π and uppercase Q in order to avoid a conflict of notation with existing results in the literature. It will be clear from the context for these variables whether we are referring to a random variable or its realization. To denote the set of probability measures over a measurable space $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, we use $\mathcal{P}(\mathcal{X})$, and to denote a contiguous tuple of random variables (X_0, X_1, \dots, X_n) we will use the notation $X_{[0,n]}$ (and its realization by $x_{[0,n]}$). Probabilities and expectations will be denoted by P and \mathbf{E} , respectively. When the relevant distributions depend on some parameters, we include these in the superscript and/or subscript. Also note that, even for a finite set \mathcal{Y} , we sometimes write for consistency of notation $\sum_{\mathcal{Y}} f(y)P(y|x) = \int_{\mathcal{Y}} f(y)P(dy|x)$, where we use the counting measure over \mathcal{Y} .

II. PRELIMINARIES: OPTIMAL CODING PROBLEM AND ITS MDP FORMULATION

A. Optimal Zero-Delay Coding and Existence of an Optimal Policy

Let our information source be a Markov process $(X_t)_{t \geq 0}$ taking values in \mathcal{X} , which we assume is finite. Let $T(x'|x) := P(X_{t+1} = x'|X_t = x)$, $x, x' \in \mathcal{X}$, be its transition kernel. We assume that T is *irreducible*; that is, for every $x, x' \in \mathcal{X}$, there exists $t \in \mathbb{Z}_+$ such that $T^t(x'|x) > 0$, where

$$T^t(x'|x) := P(X_t = x'|X_0 = x) = \sum_{y \in \mathcal{X}} T(y|x)T^{t-1}(x'|y).$$

We say that a probability measure ζ is *invariant* for T if

$$\mu(x') = \sum_{x \in \mathcal{X}} T(x'|x)\zeta(x) \quad \forall x' \in \mathcal{X}.$$

One classical result in Markov chain theory is that an irreducible Markov chain on a finite space has a unique invariant measure [49], and thus our source $(X_t)_{t \geq 0}$ has a unique invariant measure ζ . Let $X_0 \sim \pi_0$ (we also call π_0

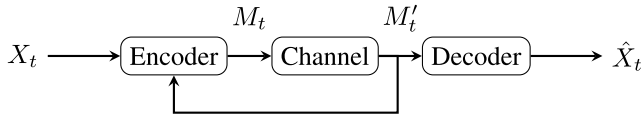


Fig. 1. Source-channel coding with feedback.

the prior). Let \mathcal{M} and \mathcal{M}' be the input and output alphabets of the memoryless channel, which we assume are finite, and let $(M_t)_{t \geq 0}$ and $(M'_t)_{t \geq 0}$ be the respective processes. We denote the channel kernel by $O(m'|m) := P(M'_t = m' | M_t = m)$ for $m' \in \mathcal{M}'$ and $m \in \mathcal{M}$. Finally, let $\hat{\mathcal{X}}$ be some finite set of reconstruction values, and let $(\hat{X}_t)_{t \geq 0} \subset \hat{\mathcal{X}}$ be the corresponding sequence of reproductions.

Consider sequences of functions $(\gamma_t^e)_{t \geq 0}$, which we call the encoder policy, and $(\gamma_t^d)_{t \geq 0}$, which we call the decoder policy. In addition to the current source symbol, the encoder has access to all past source symbols and channel inputs, and all past channel outputs in the form of feedback. In addition to the current channel output, the decoder has access to all previous channel outputs. That is, $(\gamma_t^e)_{t \geq 0}$ and $(\gamma_t^d)_{t \geq 0}$ are such that

$$\gamma_t^e : \mathcal{X}^{t+1} \times \mathcal{M}^t \times (\mathcal{M}')^t \rightarrow \mathcal{M} \quad \gamma_t^d : (\mathcal{M}')^{t+1} \rightarrow \hat{\mathcal{X}} \quad (\text{II.1})$$

$$(X_{[0,t]}, M_{[0,t-1]}, M'_{[0,t-1]}) \mapsto M_t \quad M'_{[0,t]} \mapsto \hat{X}_t. \quad (\text{II.2})$$

Our setup can be seen in Figure 1. Note that the initial distribution π_0 and the policies $(\gamma_t^e, \gamma_t^d)_{t \geq 0}$ induce a joint distribution on $(X_t, M_t, M'_t, \hat{X}_t)_{t \geq 0}$; this follows from the Ionescu-Tulcea Theorem, see e.g., [40, Proposition C.10].

We consider two performance criteria for the zero-delay coding problem. We wish to find encoder and decoder policies such that one of the following distortion quantities is minimized: the discounted distortion,

$$J_\beta(\pi_0, \gamma^e, \gamma^d) := \mathbf{E}_{\pi_0}^{\gamma^e, \gamma^d} \left[\sum_{t=0}^{\infty} \beta^t d(X_t, \hat{X}_t) \right], \quad (\text{II.3})$$

or the average distortion,

$$J(\pi_0, \gamma^e, \gamma^d) := \limsup_{T \rightarrow \infty} \mathbf{E}_{\pi_0}^{\gamma^e, \gamma^d} \left[\frac{1}{T} \sum_{t=0}^{T-1} d(X_t, \hat{X}_t) \right], \quad (\text{II.4})$$

where $d : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathbb{R}_+$ is a given distortion function and $\beta \in (0, 1)$ is a given discount factor.

We refer to the minimization of (II.3) as the discounted distortion problem and of (II.4) as the average distortion problem. For a fixed encoder policy γ^e , it is straightforward to show that the optimal decoder policy for both problems is given by:

$$\hat{X}_t = \gamma_t^{d*}(M'_{[0,t]}) := \operatorname{argmin} \hat{x} \in \hat{\mathcal{X}} \mathbf{E}_{\pi_0}^{\gamma^e} [d(X_t, \hat{x}) | M'_{[0,t]}]. \quad (\text{II.5})$$

Accordingly, we assume that we use an optimal decoder policy for a given encoder policy. We then denote by γ the resulting encoder-decoder policy (γ^e, γ^{d*}) , and the set of all such policies as Γ . Then (II.3) and (II.4) become $J_\beta(\pi_0, \gamma)$ and $J(\pi_0, \gamma)$, respectively. We denote the minimal distortions by

$$J_\beta^*(\pi_0) := \inf_{\gamma} J_\beta(\pi_0, \gamma) \quad (\text{II.6})$$

$$J^*(\pi_0) := \inf_{\gamma} J(\pi_0, \gamma). \quad (\text{II.7})$$

We will also consider policies which obtain the above infima within some arbitrary threshold $\epsilon > 0$, which we call *near-optimal*, as follows:

Definition 1: We say that a set of policies $\{\gamma\}$ depending on some parameter set is *near-optimal* for the discounted distortion problem (respectively, average distortion problem) if for any $\epsilon > 0$, there is some choice of parameters such that the resulting policy γ satisfies $J_\beta(\pi_0, \gamma) \leq J_\beta^*(\pi_0) + \epsilon$ (respectively, $J(\pi_0, \gamma) \leq J^*(\pi_0) + \epsilon$).

Note that in the zero-delay coding problem, we are usually concerned with the average distortion problem. However, it can be shown that as $\beta \rightarrow 1$, policies which are near-optimal for the discounted distortion problem is also near-optimal for the average distortion problem (see [41, Theorem 7.3.6]). Furthermore, the discounted distortion problem is generally easier to study from a reinforcement learning standpoint. Thus, we will target the discounted distortion problem throughout the majority of the paper and then make connections with the average distortion problem by taking $\beta \rightarrow 1$.

For fixed $(x_{[0,t-1]}, m_{[0,t-1]}, m'_{[0,t-1]})$, consider the function $\gamma(\cdot, x_{[0,t-1]}, m_{[0,t-1]}, m'_{[0,t-1]}) : \mathcal{X} \rightarrow \mathcal{M}$. Such a function (that is, a mapping from \mathcal{X} to \mathcal{M}) is called a *quantizer*. We denote the set of all quantizers by \mathcal{Q} . Thus we can view a policy γ as selecting a quantizer $Q_t \in \mathcal{Q}$ based on the information $(X_{[0,t-1]}, M_{[0,t-1]}, M'_{[0,t-1]})$, then generating the channel input M_t as $Q_t(X_t)$, as in [50].

Recall that we used $O(m'|m)$ to denote our channel transition kernel. Let $O_Q(m'|x)$ denote the kernel induced by a quantizer $Q \in \mathcal{Q}$; that is, $O_Q(m'|x) = O(m'|Q(x))$. Also, let $\pi_t, \bar{\pi}_t \in \mathcal{P}(\mathcal{X})$ be defined as

$$\pi_t(\cdot) = P_{\pi_0}^\gamma(X_t \in \cdot | M'_{[0,t-1]}) \quad (\text{II.8})$$

$$*\bar{\pi}_t(\cdot) = P_{\pi_0}^\gamma(X_t \in \cdot | M'_{[0,t]}), \quad (\text{II.9})$$

recalling that $X_0 \sim \pi_0$. We have dropped the γ for notational simplicity, but it should be noted that π_t and $\bar{\pi}_t$ are policy-dependent. In the literature, π_t is called the *predictor*, as it is the predictive probability of the next symbol X_t . $\bar{\pi}_t$ is called the *filter* due to its role in estimation and non-linear filtering of partially observable systems [51]. With a slight abuse of notation, we also let the source transition kernel T act as an operator on probability measures as follows:

$$T : \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{X}) \quad (\text{II.10})$$

$$\pi(x) \mapsto \sum_{x' \in \mathcal{X}} T(x'|x)\pi(x'). \quad (\text{II.11})$$

Then given π_0 , the above measures can be computed in a recursive manner as follows (see [52, Proposition 3.2.5]).

$$\begin{aligned} \bar{\pi}_t(x) &= \frac{O_{Q_t}(M'_t|x)\pi_t(x)}{\sum_{x'} O_{Q_t}(M'_t|x')\pi_t(x')}, \\ \pi_{t+1} &= T(\bar{\pi}_t). \end{aligned} \quad (\text{II.12})$$

Using the above update equations, one can compute π_t given $(M'_{[0,t-1]}, Q_{[0,t-1]})$, so that policies of the form $Q_t = \gamma_t(\pi_t)$ are valid. These policies form a special class.

Definition 2 [27]: We say a policy $\gamma = \{\gamma_t\}_{t \geq 0}$ is of the *Walrand-Varaiya type* if, at time t , γ selects a quantizer $Q_t = \gamma_t(\pi_t)$ and M_t is generated as $M_t = Q_t(X_t)$. Such a policy is

called *stationary* if it does not depend on t (that is, $\gamma_t = \bar{\gamma}$ for some $\bar{\gamma}$ and all $t \geq 0$). The set of all stationary Walrand-Varaiya policies is denoted by Γ_{WS} .

The following are key results, originally from Walrand and Varaiya [25] for a finite time horizon and extended to the infinite-horizon case in [27].

Proposition 1 [27, Proposition 2]: For any $\beta \in (0, 1)$, there exists $\gamma^* \in \Gamma_{\text{WS}}$ that solves the discounted distortion problem (that is, it minimizes (II.3)) for all priors $\pi_0 \in \mathcal{P}(\mathcal{X})$.

Proposition 2 [27, Theorem 3]: There exists $\gamma^* \in \Gamma_{\text{WS}}$ that solves the average distortion problem (that is, it minimizes (II.4)) for all priors $\pi_0 \in \mathcal{P}(\mathcal{X})$.

B. Regularity Properties of the Markov Decision Process

Utilized in the above results is the fact that, under any $\gamma \in \Gamma_{\text{WS}}$, the zero-delay coding problem can be viewed as a Markov decision process (MDP), which we now formally define.

Definition 3: We define a *Markov decision process* (MDP) as a 4-tuple $(\mathcal{Z}, \mathcal{U}, P, c)$, where:

- 1) \mathcal{Z} is the *state space*, which we assume is Polish (a Borel subset of a complete, separable metric space).
- 2) \mathcal{U} is the *action space*, also Polish.
- 3) $P : \mathcal{Z} \times \mathcal{U} \rightarrow \mathcal{P}(\mathcal{Z})$ is the *transition kernel*, such that $(z, u) \mapsto P(dz'|z, u)$.
- 4) $c : \mathcal{Z} \times \mathcal{U} \rightarrow [0, \infty)$ is the *cost function*.

The objective is to minimize $J_\beta(z_0, \gamma) := \mathbf{E}_{z_0}^\gamma \left[\sum_{t=0}^{\infty} \beta^t c(Z_t, U_t) \right]$ (the discounted cost problem) or $J(z_0, \gamma) := \limsup_{T \rightarrow \infty} \mathbf{E}_{z_0}^\gamma \left[\frac{1}{T} \sum_{t=0}^{T-1} c(Z_t, U_t) \right]$ (the average cost problem), over all γ , where $\gamma = (\gamma_t)_{t \geq 0}$ and $U_t = \gamma_t(Z_{[0,t]}, U_{[0,t-1]})$.

Proposition 3: Restricted to $\gamma \in \Gamma_{\text{WS}}$, the zero-delay coding problem is an MDP, where:

- 1) $\mathcal{Z} = \mathcal{P}(\mathcal{X})$.
- 2) $\mathcal{U} = \mathcal{Q}$.
- 3) $P = P(d\pi'|\pi, Q)$ induced by the update equations in (II.12).
- 4) $c(\pi, Q) = \sum_{\mathcal{M}'} \min_{\hat{x} \in \hat{\mathcal{X}}} \sum_x d(x, \hat{x}) O_Q(m'|x) \pi(x)$.

This follows directly from the update equations in (II.12) and the fact that, under any $\gamma \in \Gamma_{\text{WS}}$, π_t completely determines Q_t . The restriction to Γ_{WS} is without loss of optimality due to Propositions 1 and 2. The choice of c is due to the following result.

Lemma 1: If an optimal decoder is used, the expected distortion at the encoder (that is, before sending M_t) is given by

$$c(\pi_t, Q_t) = \sum_{m'} \min_{\hat{x} \in \hat{\mathcal{X}}} \sum_x d(x, \hat{x}) O_{Q_t}(m'|x) \pi_t(x). \quad (\text{II.13})$$

Proof: Under any $\gamma \in \Gamma_{\text{WS}}$ and given π_0 , the optimal decoder can compute the filter $\bar{\pi}_t$, and thus the optimal decoder chooses \hat{X}_t according to

$$\operatorname{argmin}_{\hat{x}} \hat{\mathbf{E}} \left[d(X_t, \hat{x}) | M'_{[0,t]} \right] = \operatorname{argmin}_{\hat{x}} \sum_x d(x, \hat{x}) \bar{\pi}_t(x). \quad (\text{II.14})$$

By the update equations in (II.12), we have

$$\bar{\pi}_t(x) = \frac{O_{Q_t}(M'_t|x) \pi_t(x)}{\sum_{x'} O_{Q_t}(M'_t|x') \pi_t(x')}, \quad (\text{II.15})$$

Thus at the decoder, the expected distortion is given by

$$\min_{\hat{x}} \sum_x d(x, \hat{x}) \frac{O_{Q_t}(M'_t|x) \pi_t(x)}{\sum_{x'} O_{Q_t}(M'_t|x') \pi_t(x')}. \quad (\text{II.16})$$

However, at the encoder we must take the further expectation over M'_t (conditioned on $M'_{[0,t-1]}$), since we do not yet have access to M'_t . Thus, at the encoder the expected distortion is

$$\sum_{m'} \min_{\hat{x}} \sum_x d(x, \hat{x}) \frac{O_{Q_t}(m'|x) \pi_t(x)}{\sum_{x'} O_{Q_t}(m'|x') \pi_t(x')} \quad (\text{II.17})$$

$$\cdot P_{\pi_0}^\gamma(M'_t = m' | M'_{[0,t-1]}) \\ = \sum_{m'} \min_{\hat{x}} \sum_x d(x, \hat{x}) \frac{O_{Q_t}(m'|x) \pi_t(x)}{\sum_{x'} O_{Q_t}(m'|x') \pi_t(x')} \quad (\text{II.18})$$

$$\cdot \sum_{x'} P_{\pi_0}^\gamma(X_t = x' | M'_{[0,t-1]}) P_{\pi_0}^\gamma(M'_t = m' | X_t = x', M'_{[0,t-1]}) \\ = \sum_{m'} \min_{\hat{x}} \sum_x d(x, \hat{x}) \frac{O_{Q_t}(m'|x) \pi_t(x)}{\sum_{x'} O_{Q_t}(m'|x') \pi_t(x')} \quad (\text{II.19})$$

$$\cdot \sum_{x'} O_{Q_t}(m'|x') \pi_t(x') \\ = \sum_{m'} \min_{\hat{x}} \sum_x d(x, \hat{x}) O_{Q_t}(m'|x) \pi_t(x), \quad (\text{II.20})$$

where second equality holds by the definition of π_t and since $P_{\pi_0}^\gamma(M'_t = m' | X_t = x', M'_{[0,t-1]}) = O_{Q_t}(m'|x')$ by the memoryless property of the channel and $\gamma(\pi_t) = Q_t$. ■

By this lemma, we have that the expected distortion at the encoder (assuming an optimal decoder), satisfies

$$\mathbf{E}_{\pi_0}^\gamma \left[\sum_{t=0}^{T-1} c(\pi_t, Q_t) \right] = \mathbf{E}_{\pi_0}^\gamma \left[\sum_{t=0}^{T-1} d(X_t, \hat{X}_t) \right]. \quad (\text{II.21})$$

Thus, this choice of c ensures that solving the MDP defined in Proposition 3 over all $\gamma \in \Gamma_{\text{WS}}$ (that is, minimizing $J_\beta(\pi_0, \gamma)$ or $J(\pi_0, \gamma)$) is equivalent to solving the zero-delay coding problem. Accordingly, we hereafter consider the discounted and average cost problems for this MDP (rather than the original discounted and average distortion problems). This allows us to use strategies from the literature of stochastic control; however, several complexities have been introduced:

- While the source alphabet \mathcal{X} is finite, the state space of the MDP, $\mathcal{P}(\mathcal{X})$, is uncountable. Furthermore, while our source process $(X_t)_{t \geq 0}$ is finite and irreducible (and hence has a unique invariant measure), there is no a priori reason for the MDP state process $(\pi_t)_{t \geq 0}$ to inherit these properties; in particular, irreducibility is too demanding.
- While we assume knowledge of the source transition kernel T , the calculation of the transition kernel $P(d\pi'|\pi, Q)$ is computationally demanding.

Thus even if one can approximate the MDP state space $\mathcal{P}(\mathcal{X})$ by some finite one, implementation of traditional MDP methods such as dynamic programming is difficult for this problem. This motivates the use of learning methods which learn the optimal policy empirically. We will cover this method in detail in Section IV. Finally, although explicit computation

of $P(d\pi'|\pi, Q)$ is difficult, the following key structural result was obtained in [50].

Lemma 2 [50, Lemma 11]: The transition kernel $P(d\pi'|\pi, Q)$ is weakly continuous. That is,

$$\int_{\pi'} f(\pi')P(d\pi'|\pi, Q) \quad (\text{II.22})$$

is continuous in (π, Q) for any continuous and bounded $f : \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$.

Here, we endow $\mathcal{P}(\mathcal{X})$ with the weak convergence topology and \mathcal{Q} with the Young topology (see [50]). Alternatively, since \mathcal{Q} is finite here the discrete topology would also suffice. MDPs with weakly continuous transition kernels as above are often called *weak Feller*.

C. Filter and Predictor Stability

A key property that we use is *filter/predictor stability* (recall from Definition II.8 that the predictor is given by π_t and the filter by $\bar{\pi}_t$).

Definition 4: The total variation distance between two probability measures μ, ν defined over \mathcal{X} is given by

$$\|\mu - \nu\|_{\text{TV}} := \sup_{\|f\|_{\infty} \leq 1} \left| \int_{\mathcal{X}} f(x)\mu(dx) - \int_{\mathcal{X}} f(x)\nu(dx) \right|, \quad (\text{II.23})$$

where the supremum is over all measurable real functions such that $\|f\|_{\infty} = \sup_{x \in \mathcal{X}} |f(x)| \leq 1$.

The total variation distance is equivalent to the L_1 metric when \mathcal{X} is finite. Recall that π_t , through the update equations (II.12), is sensitive to the value of π_0 . Accordingly, we use π_t^μ to denote the predictor when $\pi_0 = \mu$.

Definition 5: The predictor process $(\pi_t)_{t \geq 0}$ is *stable in total variation almost surely* under a policy $\gamma \in \Gamma_{\text{WS}}$ if, for all $\mu, \nu \in \mathcal{P}(\mathcal{X})$ such that $\mu \ll \nu$, we have

$$\lim_{t \rightarrow \infty} \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}} = 0 \text{ } P_\mu^\gamma\text{-a.s.} \quad (\text{II.24})$$

That is, the predictor process is insensitive to its initialization. In some cases, we are interested in a different form of stability, which we define next.

Definition 6: The predictor process $(\pi_t)_{t \geq 0}$ is *exponentially stable in total variation in expectation* under a policy $\gamma \in \Gamma_{\text{WS}}$ if there exists a coefficient $\alpha \in (0, 1)$ such that for all $\mu, \nu \in \mathcal{P}(\mathcal{X})$ (with $\mu \ll \nu$) and for all $t \geq 0$, we have

$$\mathbf{E}_\mu^\gamma [\|\pi_{t+1}^\mu - \pi_{t+1}^\nu\|_{\text{TV}}] \leq \alpha \mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}]. \quad (\text{II.25})$$

We can make equivalent definitions for the filter process; in fact, the problem of filter stability (in various senses) is a classical problem in probability and statistics, where it is typically established in two ways: (i) The transition kernel of the underlying state is in some sense *sufficiently ergodic*, so that regardless of the observations, the filter process inherits this ergodicity and forgets its prior over time. (ii) The observations are in some sense *sufficiently informative*, so that, regardless of the prior, the filter process tracks the true state process. For a detailed review of these filter stability methods, see [53]. However, we will need slightly more general results in our case, since it is usually assumed in the filter stability problem that the observation kernel is time-invariant; here, O_{Q_t} depends on Q_t and hence changes with time, and accordingly additional analysis is needed.

III. SLIDING FINITE WINDOW APPROXIMATION OF THE BELIEF MDP

Here we describe how to approximate π_t using a sliding finite window. The analysis in this section is inspired by [42], which used a similar construction to study sliding finite window policies for partially observed Markov decision processes (POMDPs). We note that although the proof methods are similar, our setup differs from [42] in two ways: 1) our problem is not strictly a POMDP as studied in [42], because the induced observation kernel changes with the choice of quantizer Q_t , 2) accordingly, our belief state is the predictor rather than the filter, and thus the Dobrushin coefficient terms we arrive at are different.

First, we must define a slightly different MDP than that defined in Proposition 3. Fix some window size $N \in \mathbb{Z}_+$. Recall the channel outputs $(M'_t)_{t \geq 0}$ and the quantizers $(Q_t)_{t \geq 0}$. We define the following:

$$I_t := (M'_{[t-N, t-1]}, Q_{[t-N, t-1]}) \quad (\text{III.1})$$

$$W_t := (\pi_{t-N}, I_t). \quad (\text{III.2})$$

We can compute π_t given W_t by applying the update equations in (II.12) N times. Denote this mapping by

$$\varphi : \mathcal{W} \rightarrow \mathcal{P}(\mathcal{X}) \quad (\text{III.3})$$

$$W_t \mapsto \pi_t \quad (\text{III.4})$$

where $\mathcal{W} = \mathcal{P}(\mathcal{X}) \times (\mathcal{M}')^N \times \mathcal{Q}^N$, endowed with the product topology, and we use the weak convergence topology on $\mathcal{P}(\mathcal{X})$ and standard coordinate topologies on \mathcal{M}' and \mathcal{Q} . We omit the dependence of φ on N for simplicity. We call W_t the sliding finite window belief term, and call policies of the form $Q_t = \gamma_t(W_t)$ *sliding finite window belief policies* (with window size N). If it does not depend on t , we call it stationary. Denote the set of all stationary sliding finite window belief policies by Γ_{FS} .

Remark: This approach assumes that we start at time $t \geq N$. Although for discounted MDPs the first N steps may be significant, for the zero-delay coding problem we are interested in taking $\beta \rightarrow 1$, so these first N steps will not be crucial. Accordingly, we assume that N steps have already been completed with some arbitrary $\gamma \in \Gamma_{\text{WS}}$. For notational simplicity, we assume that these steps have occurred from $t = -N, \dots, -1$ and thus the process starts from W_0 (and the prior would now be π_{-N}).

A. The Sliding Finite Window Belief MDP

This sliding finite window belief construction inherits the MDP properties of the original setup.

Proposition 4: Under any $\gamma \in \Gamma_{FS}$, the zero-delay coding problem is an MDP, where:

- 1) $\mathcal{Z} = \mathcal{W}$.
- 2) $\mathcal{U} = \mathcal{Q}$.
- 3) $P = P(dw'|w, Q)$.
- 4) $c(w, Q) = \sum_{\mathcal{M}'} \min_{\hat{x}} \sum_x d(x, \hat{x}) O_Q(m'|x) \varphi(w)(x)$.

Note that the cost function is exactly the cost function we had in Proposition 3, by simply replacing $\pi = \varphi(w)$. Then

an analog to Lemma 1 holds, and we have that solving the MDP from Proposition 4 is equivalent to solving the zero-delay coding problem. That is, we can equivalently consider $J_\beta^*(w_0) = \inf_{\gamma \in \Gamma_{FS}} J_\beta(w_0, \gamma)$.

B. Sliding Finite Window Approximation

The above representation is not particularly useful, as it still requires one to compute π_{t-N} . Instead, fix the first coordinate to ζ (the invariant distribution of the source) and let

$$\hat{W}_t = (\zeta, I_t) \quad (\text{III.5})$$

$$\hat{\pi}_t = \varphi(\hat{W}_t). \quad (\text{III.6})$$

That is, we obtain $\hat{\pi}_t$ by applying the update equations N times, but starting from an incorrect (fixed) prior ζ . Equivalently,

$$\pi_t(x) = P_{\pi_{t-N}}^\gamma(X_t = x | M'_{[t-N, t-1]}, Q_{[t-N, t-1]}) \quad (\text{III.7})$$

$$\hat{\pi}_t(x) = P_\zeta^\gamma(X_t = x | M'_{[t-N, t-1]}, Q_{[t-N, t-1]}). \quad (\text{III.8})$$

The key idea, which will be discussed in detail later, is that under predictor stability the correct predictor $\pi_t = \varphi(W_t)$ and the incorrect predictor $\hat{\pi}_t = \varphi(\hat{W}_t)$ will be close for large enough N , since the predictor will be insensitive to the prior.

The benefits of such an approximation are evident: rather than deal with all of \mathcal{W} , which is uncountable due to $\mathcal{P}(\mathcal{X})$, we only have to deal with the finite set $\mathcal{W}_N := \{\zeta\} \times (\mathcal{M}')^N \times \mathcal{Q}^N$. Furthermore, we no longer need to compute π_{t-N} , which can save significant computation resources especially when the relevant alphabets are large.

Consider the following transition kernel,

$$P_N(\hat{w}_1 | \hat{w}, Q) := P(\mathcal{P}(\mathcal{X}), i_1 | \hat{w}, Q), \quad (\text{III.9})$$

where P is the transition kernel of the sliding finite window belief MDP and $\hat{w}_1 = (\zeta, i_1)$, and cost function

$$c_N(\hat{w}, Q) := \sum_{m'} \min_{\hat{x} \in \hat{\mathcal{X}}} \sum_x d(x, \hat{x}) O_Q(m' | x) \hat{\pi}(x). \quad (\text{III.10})$$

Then, our approximate MDP becomes $\text{MDP}_N = (\mathcal{W}_N, \mathcal{Q}, P_N, c_N)$. Denote the discounted cost for this MDP under a given policy $\hat{\gamma}_N$ (from \mathcal{W}_N to \mathcal{Q}) by $\hat{J}_\beta(\hat{w}_0, \hat{\gamma}_N)$, and the optimal discounted cost by $\hat{J}_\beta^*(\hat{w}_0)$, with minimizing policy $\hat{\gamma}_N^*$. We are only looking at stationary policies, so the subscript N on the policy should not be confused with a time index. We can extend these function and policies from \mathcal{W}_N to \mathcal{W} by simply making them constant over $\mathcal{P}(\mathcal{X})$. We denote these extensions by \hat{J}_β^* and $\hat{\gamma}_N^*$.

The following is a key loss term:

$$L_t^N := \sup_{\gamma \in \Gamma_{WS}} \mathbf{E}_{\pi_{t-N}}^\gamma [\|\pi_t - \hat{\pi}_t\|_{\text{TV}}]. \quad (\text{III.11})$$

We now present our main results for this approximation scheme, which give a bound on the performance loss when using the given window length N . Note that here we take an expectation with respect to some policy acting on the previous N steps (and hence it generates W_0), and with respect to some prior π_{t-N} . Also, define $\|d\|_\infty := \max_{x, \hat{x}} d(x, \hat{x})$, where we recall that d is the distortion measure for the zero-delay coding problem.

Theorem 1: For any $\gamma \in \Gamma_{WS}$ acting on N time steps to generate W_0 and any prior $\pi_{t-N} \in \mathcal{P}(\mathcal{X})$, we have

$$\mathbf{E}_{\pi_{t-N}}^\gamma [|\tilde{J}_\beta^*(W_0) - J_\beta^*(W_0)|] \leq \frac{\|d\|_\infty}{1-\beta} \sum_{t=0}^{\infty} \beta^t L_t^N. \quad (\text{III.12})$$

Theorem 2: For any $\gamma \in \Gamma_{WS}$ acting on N time steps to generate w_0 and any prior $\pi_{t-N} \in \mathcal{P}(\mathcal{X})$, we have

$$\mathbf{E}_{\pi_{t-N}}^\gamma [|J_\beta(W_0, \tilde{\gamma}_N^*) - J_\beta^*(W_0)|] \leq \frac{2\|d\|_\infty}{1-\beta} \sum_{t=0}^{\infty} \beta^t L_t^N. \quad (\text{III.13})$$

That is, we can bound the sub-optimality of the optimal policy for the approximate MDP $\tilde{\gamma}_N^*$ if we can bound the differences in the *predictors*. This naturally connects to the question of predictor stability (recall Section II-C).

Theorems 1 and 2 build on the mathematical program considered [42, Theorems 3.2, 3.3]; however, in our setting the observation kernel is policy-dependent and we use the predictor as the effective state variable, instead of the filter in our formulation as well as the loss term. These entail further mathematical analysis in both the construction of the equivalent Markov model and the error analysis. The proofs are given in the Appendix.

C. Bounds on the Loss Term

The loss term L_t^N in the previous theorems is related to the question of predictor stability (recall Section II-C). Indeed, the term within the supremum is exactly

$$\mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\gamma\|_{\text{TV}}] \quad (\text{III.14})$$

when $\mu = \pi_{t-N}$ and $\nu = \zeta$, and under some $\gamma \in \Gamma_{WS}$. Thus bounding this term over all γ will give us a bound on L_t^N . We note that any notion of predictor stability could be used to give a bound on L_t^N (and thus on the performance of $\tilde{\gamma}_N^*$). In the following section, we give one such condition which is sufficient (but by no means necessary) to apply the above theorems.

Dobrushin Coefficient Conditions: The following results are inspired by the analysis in [54], which uses joint contraction properties of the state and observation kernels to bound (III.14). First we introduce some notation. For standard Borel spaces $\mathcal{A}_1, \mathcal{A}_2$ and some kernel $P : \mathcal{A}_1 \rightarrow \mathcal{P}(\mathcal{A}_2)$, we define the Dobrushin coefficient as

$$\delta(P) := \inf \sum_{i=1}^n \min(K(B_i|x), K(B_i|y)), \quad (\text{III.15})$$

where the infimum is over $x, y \in \mathcal{A}_1$ and all partitions $\{B_i\}_{i=1}^n$ of \mathcal{A}_2 . In particular, for finite spaces, the Dobrushin coefficient is equivalent to summing the minimum elements between every pair of rows, then taking the minimum of these sums. For example, take

$$P = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \\ \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{pmatrix}. \quad (\text{III.16})$$

Between the first and second rows, the sum of the minimum elements gives $\frac{2}{3}$, between the third and fourth gives $\frac{3}{4}$, etc.

One can verify that the minimum of such sums is $\frac{1}{2}$, so $\delta(P) = \frac{1}{2}$ (note that $\delta(P) \leq 1$ by definition).

The following is then a counterpart of [54, Theorem 3.6]. Note that in our case, the channel is not time-invariant, unlike in [54], but the analysis follows similarly. The proof is provided in the Appendix.

Theorem 3: For any $\gamma \in \Gamma_{\text{WS}}$ and for any $\mu \ll \nu$,

$$\begin{aligned} & \mathbf{E}_\mu^\gamma [\|\pi_{t+1}^\mu - \pi_{t+1}^\nu\|_{\text{TV}}] \\ & \leq (1 - \delta(T))(2 - \tilde{\delta}(O)) \mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}], \end{aligned}$$

where $\tilde{\delta}(O) = \min_{Q \in \mathcal{Q}} (\delta(O_Q))$.

We can get a simpler bound by using $\delta(O)$ directly rather than $\tilde{\delta}(O)$. For a given quantizer Q , the kernel $O_Q(m'|x) = O(m'|Q(x))$ only contains rows from the kernel O , so $\delta(O) \leq \delta(O_Q)$ for all Q . Thus we arrive at the following corollary.

Corollary 1: Assume $\alpha := (1 - \delta(T))(2 - \delta(O)) < 1$. Then for any $\gamma \in \Gamma_{\text{WS}}$ and for any $\mu \ll \nu$, we have

$$\mathbf{E}_\mu^\gamma [\|\pi_{t+1}^\mu - \pi_{t+1}^\nu\|_{\text{TV}}] \leq \alpha \mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}]. \quad (\text{III.17})$$

That is, the predictor process is exponentially stable in total variation in expectation. Furthermore, if $\delta(T) > \frac{1}{2}$, then the above is true with $\alpha = 1 - \delta(T)$ regardless of the channel O .

Applying this to the L_t^N term, we have

$$\begin{aligned} L_t^N &= \sup_{\gamma \in \Gamma_{\text{WS}}} \mathbf{E}_{\pi_{t-N}}^\gamma [\|\pi_t - \hat{\pi}_t\|_{\text{TV}}] \\ &\leq \alpha^N \|\pi_{t-N} - \pi'\|_{\text{TV}} \leq 2\alpha^N. \end{aligned} \quad (\text{III.18})$$

The requirement $(1 - \delta(T))(2 - \delta(O)) < 1$ places conditions on the source dynamics and the channel. A universal condition to make the result applicable over all channels is obtained when one considers the special case where the channel is noiseless: in this case, we have $\delta(O) = 0$ and $\delta(T) > \frac{1}{2}$ is to hold. Combining Theorem 2 and the bound in (III.18), we obtain the following result.

Corollary 2: Assume $\alpha := (1 - \delta(T))(2 - \delta(O)) < 1$. Then for any $\gamma \in \Gamma_{\text{WS}}$ which acts on N time steps to generate W_0 and any prior π_{-N} , we have

$$\mathbf{E}_{\pi_{-N}}^\gamma [|J_\beta(W_0, \tilde{\gamma}_N^*) - J_\beta^*(W_0) |] \leq \frac{4\|d\|_\infty}{(1 - \beta)^2} \alpha^N. \quad (\text{III.19})$$

IV. COMPUTATION OF OPTIMAL POLICY FOR THE APPROXIMATE MDP

We now turn our attention to actually calculating an optimal policy for the approximate MDP in Section III, which will be near-optimal for the zero-delay coding problem for sufficiently large parameter N . That is, the remaining problem is to find for some large N the $\hat{\gamma}_N^*$ which minimizes

$$\mathbf{E}_{\pi_{-N}}^{\hat{\gamma}_N^*} \left[\sum_{t \geq 0} \beta^t c_N(\hat{w}_t, Q_t) \right]. \quad (\text{IV.1})$$

This can be done a number of ways using tools from stochastic control. We present two methods here, 1) a model-free Q-learning algorithm, and 2) a model learning stage followed by a value iteration algorithm (see [40, Lemma 4.2.8]).

A. Q-Learning

The Q-learning algorithm was introduced in [55] for finite sets and generalized (along with rigorous convergence and near-optimality guarantees) in [56]. The following is the Q-learning algorithm for our finite-window MDP in Section III, given by $\text{MDP}_N = (\mathcal{W}_N, \mathcal{Q}, P_N, c_N)$. Let $(Q_t)_{t \geq 0}$ be generated according to a uniform policy, that is $P(Q_t = Q) = \frac{1}{|\mathcal{Q}|}$ for all $t \geq 0$. For an arbitrary π_{-N} under this policy, consider the resulting sequence of finite windows $(\hat{W}_t)_{t \geq 0}$. Then we define the following update equations, where V_t and α_t are both functions from $\mathcal{W}_N \times \mathcal{Q}$ to \mathbb{R}_+ :

$$\begin{aligned} V_{t+1}(\hat{w}, Q) &= V_t(\hat{w}, Q) \quad \text{for all } (\hat{w}, Q) \neq (\hat{W}_t, Q_t), \\ V_{t+1}(\hat{W}_t, Q_t) &= (1 - \alpha_t(\hat{W}_t, Q_t)) V_t(\hat{W}_t, Q_t) \\ &\quad + \alpha_t(\hat{W}_t, Q_t) \left[c_N(\hat{W}_t, Q_t) + \beta \min_{Q \in \mathcal{Q}} V_t(\hat{W}_{t+1}, Q) \right], \\ \alpha_t(\hat{w}, Q) &= \frac{1}{1 + \sum_{k=0}^t \mathbf{1}(\hat{W}_k = \hat{w}, Q_k = Q)}, \end{aligned} \quad (\text{IV.2})$$

for some arbitrary V_0 , and where $\mathbf{1}(\hat{W}_k = \hat{w}, Q_k = Q)$ is 1 when $(\hat{W}_k = \hat{w}, Q_k = Q)$ and 0 otherwise. We show in the Appendix that the necessary assumptions of [56, Theorem 2.1] hold for MDP_N , which gives us the following result.

Theorem 4: [56, Theorem 2.1] The sequence $(V_t)_{t \geq 0}$ converges almost surely to a limit V^* , where

$$V^*(\hat{w}, Q) = c_N(\hat{w}, Q) + \beta \sum_{\hat{w}_1} \min_{Q_1} V^*(\hat{w}_1, Q_1) P_N(\hat{w}_1 | \hat{w}, Q) \quad (\text{IV.3})$$

Equation (IV.3) is exactly the classic Bellman optimality equation for MDP_N , so $\hat{\gamma}_N^*(\hat{w}) = \text{argmin}_{Q \in \mathcal{Q}} V^*(\hat{w}, Q)$ is the optimal policy for MDP_N .

B. Model Learning and Value Iteration

While Q-learning learns the function V^* empirically without explicitly learning P_N , this approach first learns P_N empirically and then performs value iteration on MDP_N . Again under a uniform policy and any prior π_{-N} , the same assumptions of [56, Theorem 2.1] hold and thus the kernel $P_N(\hat{w}_1 | \hat{w}, Q)$, is the almost sure limit of the empirical average (see also [57, Section 3.2]):

$$P_N(\hat{w}_1 | \hat{w}, Q) = \lim_{T \rightarrow \infty} \frac{\sum_{t=0}^{T-1} \mathbf{1}(\hat{W}_k = \hat{w}, Q_k = Q, \hat{W}_{k+1} = \hat{w}_1)}{\sum_{t=0}^{T-1} \mathbf{1}(\hat{W}_k = \hat{w}, Q_k = Q)}. \quad (\text{IV.4})$$

We can then perform value iteration [40, Lemma 4.2.8] as follows with $V_0(\hat{w}) \equiv 0$:

$$V_t(\hat{w}) = \min_Q \left\{ c_N(\hat{w}, Q) + \beta \sum_{\hat{w}_1} V_{t-1}(\hat{w}_1) P_N(\hat{w}_1 | \hat{w}, Q) \right\}.$$

Then an analogous result to 4 holds:

Theorem 5: The sequence $(V_t)_{t \geq 0}$ converges almost surely to a limit V^* , where

$$V^*(\hat{w}) = \min_Q \left\{ c_N(\hat{w}, Q) + \beta \sum_{\hat{w}_1} V^*(\hat{w}_1) P_N(\hat{w}_1 | \hat{w}, Q) \right\}. \quad (\text{IV.5})$$

Again choosing $\hat{\gamma}_N^*(\hat{w}) = Q_t$ as the minimizer of the above yields an optimal policy for MDP_N . The benefit of the value

iteration approach is that it can have better sample complexity for certain settings [58], but the policies obtained from either method will be the same.

V. AN ALTERNATIVE APPROXIMATION SCHEME AND COMPARISON

A. Nearest Neighbor Approximation

Another approximation scheme for the zero-delay coding problem was proposed in [43], along with a Q-learning algorithm to compute the corresponding near-optimal policy. Although this scheme was only formally analyzed in the case of a noiseless channel, the results go through in the case of a noisy channel with some modifications. Accordingly, we present the noisy channel analogs of the results in [43] without proof, but we do provide a detailed comparison of the two schemes.

First, we approximate π_t using a nearest neighbor quantization scheme with a finite number of bins. This approximation can be done efficiently using [59, Algorithm 1], which approximates $\mathcal{P}(\mathcal{X})$ by the following finite set,

$$\mathcal{P}_K(\mathcal{X}) = \left\{ \hat{\pi} \in \mathcal{P}(\mathcal{X}) : \hat{\pi} = \left[\frac{k_1}{K}, \dots, \frac{k_{|\mathcal{X}|}}{K} \right], \right. \\ \left. k_i = 0, \dots, K, i = 1, \dots, |\mathcal{X}| \right\}, \quad (\text{V.1})$$

where $\sum_{i=1}^K k_i = K$. Note that in the literature, such distributions are also called types or empirical distributions [60]. We let $\hat{\pi}$ be the nearest neighbor (in Euclidean distance) of π in $\mathcal{P}_K(\mathcal{X})$, and we clearly have that $\max_{\pi} d(\pi, \hat{\pi}) \rightarrow 0$ as $K \rightarrow \infty$.

Similarly to Section III, one can then define a transition kernel and cost function over $\mathcal{P}_K(\mathcal{X})$. We omit the details, but they are essentially the averages of the true transition kernel $P(d\pi'|\pi, Q)$ and cost function $c(\pi, Q)$ over the nearest neighbor bins defined by (V.1), with respect to an appropriate reference measure. We denote the resulting approximate MDP by $\text{MDP}_K = (\mathcal{P}_K, \mathcal{Q}, P_K, c_K)$.

Corollary 3: [43, Theorem 1] Let $\hat{\gamma}_K^* \in \Gamma_{\text{WS}}$ be optimal for MDP_K and let $\tilde{\gamma}_K^*(\pi) = \hat{\gamma}_K^*(\hat{\pi})$. Then for all $\pi_0 \in \mathcal{P}(\mathcal{X})$ and $\beta \in (0, 1)$,

$$\lim_{K \rightarrow \infty} |J_{\beta}(\pi_0, \tilde{\gamma}_K^*) - J_{\beta}^*(\pi_0)| = 0. \quad (\text{V.2})$$

That is, $\tilde{\gamma}_K^*$ is near-optimal for the zero-delay coding problem for sufficiently large K .

An analogous Q-learning result follows, using a uniform policy for quantizer selection to obtain $(\pi_t)_{t \geq 0}$ and applying the nearest neighbor approximation scheme onto $\mathcal{P}_K(\mathcal{X})$ to obtain $(\hat{\pi}_t)_{t \geq 0} \subset \mathcal{P}_K(\mathcal{X})$. Then the appropriate iterations are:

$$V_{t+1}(\hat{\pi}, Q) = V_t(\hat{\pi}, Q) \quad \text{for all } (\hat{\pi}, Q) \neq (\hat{\pi}_t, Q_t), \\ V_{t+1}(\hat{\pi}_t, Q_t) = (1 - \alpha_t(\hat{\pi}_t, Q_t)) V_t(\hat{\pi}_t, Q_t) \\ + \alpha_t(\hat{\pi}_t, Q_t) \left[c(\pi_t, Q_t) + \beta \min_{Q \in \mathcal{Q}} V_t(\hat{\pi}_{t+1}, Q) \right], \\ \alpha_t(\hat{\pi}, Q) = \frac{1}{1 + \sum_{k=0}^t \mathbf{1}(\hat{\pi}_k = \hat{\pi}, Q_k = Q)}, \quad (\text{V.3})$$

for some arbitrary V_0 . In the following result, we require that the prior π_0 is the *source invariant distribution* ζ , while in Theorem 4 it could be arbitrary. This is due to the recurrence

properties that are required on the predictor process, see details in [43].

Theorem 6 [43, Theorem 1]: Assume that π_0 in (V.3) is the unique invariant distribution of the source, ζ . Then V_t converges almost surely to a limit V^* . Furthermore, consider the policy

$$\hat{\gamma}_K^*(\hat{\pi}) = \operatorname{argmin}_{Q \in \mathcal{Q}} V^*(\hat{\pi}, Q), \quad (\text{V.4})$$

and let $\tilde{\gamma}_K^*$ be the extension of $\hat{\gamma}_K^*$ to $\mathcal{P}(\mathcal{X})$ by making it constant over the belief quantization bins, as in Corollary 3. Then as $K \rightarrow \infty$,

$$J_{\beta}(\zeta, \tilde{\gamma}_K^*) \rightarrow J_{\beta}^*(\zeta). \quad (\text{V.5})$$

That is, $\tilde{\gamma}_K^*$ is near-optimal for the discounted distortion zero-delay coding problem, provided that the source starts from its invariant distribution.

B. Comparison of the Two Schemes

In this section, we provide a detailed comparison noting explicit benefits and drawbacks of each of the two learning schemes to attain near-optimality. In particular, we emphasize the easier implementation and lower computational complexity of our proposed sliding window method, although there are some cases where the belief quantization method of [43] may be preferable.

- *Assumptions for convergence of Q-learning:* In Theorem 6, one can see that the convergence of V_t depends on the initial distribution used during learning (in particular, π_0 must be the source's unique invariant distribution ζ). Conversely, in Theorem 4, the convergence of V_t happens regardless of the initial distribution used during learning.
- *Conditions on initialization during implementation:* For the belief quantization scheme, the policy from Q-learning is near-optimal when $\pi_0 = \zeta$. We note that in the special case of the noiseless channel, valid initializations also include $\pi_0 = \delta_x$, i.e., deterministically starting from any $x \in \mathcal{X}$. This can be shown using [43, Corollary 2] and by noting that in this special case the Dirac masses δ_x are recurrent. However, in the noisy channel case, this argument fails and we require $\pi_0 = \zeta$; conversely, for the sliding finite window scheme, the policy obtained in Theorem 4 is near-optimal for almost every initial window.
- *Filter stability assumptions:* Because the predictor stability required for the sliding finite window scheme is rather strong, the near-optimality of approximations was only established for those source-channel combinations satisfying the Dobrushin coefficient conditions. However, we note that the Dobrushin coefficient conditions are not necessary; any filter stability result which gives a bound on the loss term in Theorem 2 could be used. To this end, Hilbert metric approaches may also be relevant, for example those in [61] and [62].
- *Computational complexity and need for Bayesian updates:* In the belief quantization scheme, one must compute the true belief process $(\pi_t)_{t \geq 0}$ using the update equations, then quantize this to the set $\mathcal{P}_K(\mathcal{X})$. This

increases the computational complexity of this scheme (and requires explicit knowledge of the system model), both during the Q-learning algorithm and during implementation of the learned policy. Conversely, the sliding finite window scheme uses the approximate predictor from a fixed prior and a given history. Since there are only finitely many such histories, one can compute these offline before running the Q-learning and before implementation of the learned policy. They can then be accessed in a lookup table fashion.

- *Model knowledge and a data-driven approach:* In the belief quantization scheme, both the encoder and the decoder must track the true belief π_t . The encoder needs π_t to compute the proper value of $\hat{\pi}_t$ and apply the learned policy, while the decoder needs it to compute the optimal reproduction value \hat{X}_t . Thus, the belief quantization approach requires knowledge, or at least a good estimate, of the underlying model (in particular, the initial distribution and the kernels T and O). Conversely, in the sliding finite window scheme the encoder policy is a constant function in $\mathcal{P}(\mathcal{X})$, so it can directly use the sliding finite window to apply the learned policy. For the decoder, if one has knowledge of the model, it can simply take the form of (II.5), as we did in our implementation. However, theoretically one could also apply learning at the decoder: treat the decoder as a map $\gamma^d : W_t, Q_t, M_t \rightarrow \hat{\mathcal{X}}$ and add this as another dimension in the Q-table, so that one now learns a joint encoder-decoder policy. Under the same Dobrushin coefficient conditions as for the encoder, this will become near-optimal for large N . Of course, this comes at a significant memory and runtime cost for the Q-learning algorithm, so if one has model knowledge it is still beneficial to incorporate this at the decoder.
- *Rate of convergence to near-optimality:* For the sliding finite window result in Theorem 2, the convergence is exponential (note that although we only bound the expectation, since there are only finitely many initial memories, the convergence is also exponential for each initial finite window). Conversely, the belief quantization result in Corollary 3 is only asymptotic.
- *On quantization:* For the belief quantization result in Corollary 3, the quantization of $\mathcal{P}(\mathcal{X})$ does not have to be uniform. Theoretically, this allows a more efficient quantization, although in our implementation we always use [59, Algorithm 1] (which gives a uniform quantization). The sliding finite window can be seen as a non-uniform quantization of \mathcal{W} (since it is constant over the belief coordinate). However, it is uniform over the product space $(\mathcal{M}')^N \times \mathcal{Q}^N$, since the sliding finite window scheme uses every possible finite window.

We summarize some of the key differences between the schemes in Table I.

C. Implications for the Average Cost Problem

Under mild assumptions, it can be shown that as $\beta \rightarrow 1$, a near-optimal policy for the discounted cost problem becomes near-optimal for the average cost problem

TABLE I
COMPARISON OF THE TWO APPROXIMATION SCHEMES

	Belief Quantization	Sliding Finite Window
Convergence of Q-learning and near-optimality	✓	✓
Stability conditions not needed	✓	✗
Insensitive to initialization	✗	✓
Exponential convergence of performance	✗	✓
Bayesian update not needed	✗	✓
Lookup table implementation	✓	✓

(see [41, Theorem 7.3.6]). It is straightforward to show that these assumptions hold for the zero-delay coding MDP; indeed, this same set of assumptions was shown to hold in [27] and used to show the existence of optimal Walrand-Varaiya codes for the average cost problem. Then through [41, Theorem 7.3.6] and our previous results we obtain the following corollaries. Note that the dependence on β is only asymptotic; in practice we set $\beta = 0.9999$ and this seems to give policies very close to the optimum. For a fixed β , the dependence on ϵ is as given in Theorems 4 and 6 (exponential for sliding finite window, asymptotic for belief quantization).

Corollary 4: Let $\epsilon > 0$, let $\tilde{\gamma}_N^*$ be the policy given in Theorem 4, and assume the conditions of Theorem 4 are met. Then there exists $\beta \in (0, 1)$ and $N(\epsilon, \beta)$ sufficiently large such that for any γ which acts on the first N time steps to generate W_0 and any π_{-N} ,

$$\mathbf{E}_{\pi_{-N}}^\gamma [J(W_0, \tilde{\gamma}_N^*) - J^*(W_0)] < \epsilon. \quad (\text{V.6})$$

Corollary 5: Let $\epsilon > 0$, let $\tilde{\gamma}_K^*$ be the policy given in Theorem 6. Then there exists $\beta \in (0, 1)$ and $K(\epsilon, \beta)$ sufficiently large such that

$$J(\zeta, \tilde{\gamma}_K^*) - J^*(\zeta) < \epsilon. \quad (\text{V.7})$$

VI. SIMULATIONS

We now give some examples of zero-delay coding problems and simulate the performance of the policies resulting from Theorems 6 and 4. In all of the following, we use a discount factor $\beta = 0.9999$ and the distortion function $d(x, \hat{x}) = (x - \hat{x})^2$. The average distortion is calculated over $t = 0$ to $t = 10^5$. The initial distribution is ζ (that is, the invariant distribution for the source), so that the belief quantization scheme can be used. We observe that in each case with a known optimum, our Q-learning algorithm indeed nearly reached this optimum, validating our theoretical results. In the case with an unknown optimum, we observe convergence rates which agree with the bounds in the theorems, and we compare with existing schemes where applicable.

Remark: For the implementation of the Q-learning algorithms, the theoretical upper bound for the possible number of

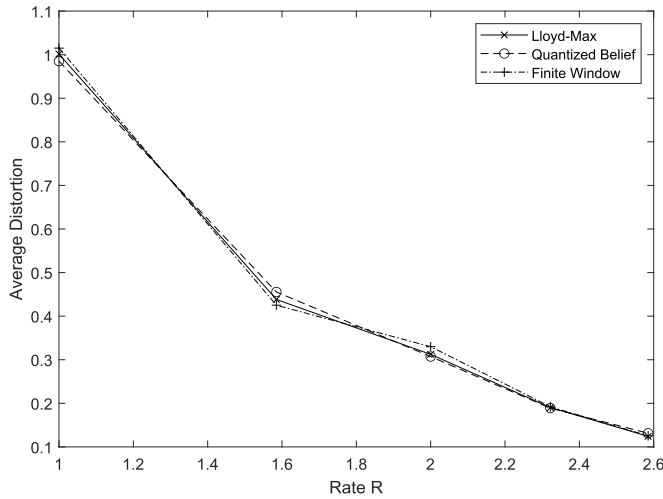


Fig. 2. Comparison with Lloyd-Max with i.i.d. source.

states in each scheme grows very quickly in their respective parameters. In particular, for the belief quantization approach we have $|\mathcal{P}_K(\mathcal{X})| = \binom{K+|\mathcal{X}|+1}{|\mathcal{X}|-1}$ [59], and for the sliding finite window approach we have $|\mathcal{W}_N| = |\mathcal{M}' \times \mathcal{Q}|^N$. However, the sets actually visited during Q-learning may be much smaller. Thus, one may wish to add entries to V_t as they are visited by the Q-learning algorithm in a dynamic fashion.

Furthermore, note that for certain problems it may be possible to significantly shrink the set of quantizers \mathcal{Q} with no loss of optimality. For example, for a noiseless channel one can omit those quantizers with empty bins, or for an i.i.d. source those with non-convex bins.

For a stopping criteria in the algorithms, any measure of the convergence of V_t would be suitable; in our results, we stop when the relative pointwise difference is sufficiently small. For example, $\max_{\hat{\pi}, Q} ((V_{t+1}(\hat{\pi}, Q) - V_t(\hat{\pi}, Q)) / V_t(\hat{\pi}, Q)) < \epsilon$ for some small $\epsilon > 0$.

A. Comparison With Lloyd-Max Quantizer for i.i.d. Source and Noiseless Channel

Let $\mathcal{X} = \{0, \dots, 7\}$ and consider an i.i.d. source $(X_t)_{t \geq 0}$, such that for all x ,

$$T(\cdot|x) = \left(\frac{1}{4} \frac{1}{8} \frac{1}{8} \frac{1}{16} \frac{1}{16} \frac{1}{16} \frac{1}{4} \frac{1}{16} \right). \quad (\text{VI.1})$$

Note that in the i.i.d. case, we trivially have $\delta(T) = 1$, so Corollary 2 is applicable. Indeed, here we have that $\alpha = 0$, so that for all $N \geq 1$,

$$\mathbf{E}_{\pi_{-N}}^y \left[|J_{\beta}(W_0, \tilde{\gamma}^*) - J_{\beta}^*(W_0)| \right] = 0. \quad (\text{VI.2})$$

That is, the optimal policy for the finite window representation is optimal (not just near-optimal) for the original problem for any N . This is not surprising given the i.i.d. nature of the source; the approximation of π_{t-N} to ζ is without any loss since π_{t-N} can be immediately recovered.

Similarly for the quantization approach, $K = 1$ is sufficient since $\pi_t = \zeta$ for all $t \geq 0$, so increasing K does not change the resulting policy. Accordingly, we let $N = K = 1$ and

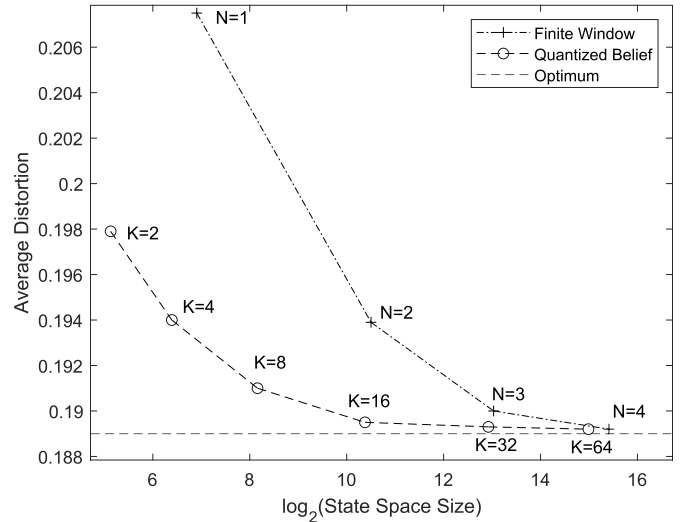


Fig. 3. Comparison with memoryless encoding, low Dobrushin coefficient.

compare the performance of both approaches against a Lloyd-Max quantizer when the channel is noiseless. We plot the performance for several sizes of \mathcal{M} . The rate is calculated as $R := \log_2 |\mathcal{M}|$. As expected by the above discussion, our algorithm matches with this quantizer in each case, which can be seen in Figure 2.

B. Comparison With Memoryless Encoding and Effect of Dobrushin Coefficient

Consider now a Markov source with transition kernel given by

$$T = \begin{pmatrix} \frac{1}{2} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \\ \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{pmatrix} \quad (\text{VI.3})$$

and where the channel is a 4-ary symmetric channel with error probability 0.06:

$$O = \begin{pmatrix} 0.94 & 0.02 & 0.02 & 0.02 \\ 0.02 & 0.94 & 0.02 & 0.02 \\ 0.02 & 0.02 & 0.94 & 0.02 \\ 0.02 & 0.02 & 0.02 & 0.94 \end{pmatrix}. \quad (\text{VI.4})$$

We have that $\delta(T) = \frac{2}{3} > \frac{1}{2}$ and $\delta(O) = 0.08$, so we can apply Corollary 2. In such a setup (where $\mathcal{X} = \mathcal{M}$ and the channel is symmetric), it was shown in [25] that “memoryless” encoding (that is, where $M_t = X_t$) is optimal. We compare our algorithms against such an encoding policy, shown in Figure 3. On the x-axis, we plot the number of states used by each policy on a \log_2 scale; that is, the number of states seen during learning. As remarked at the beginning of the section, in practice this is much smaller than the theoretical maximum number of states. We label each data point with the corresponding approximation parameter, which are $N = \{1, 2, 3, 4\}$ for the window length and $K = \{2, 4, 8, 16, 32, 64\}$ for the belief quantization.

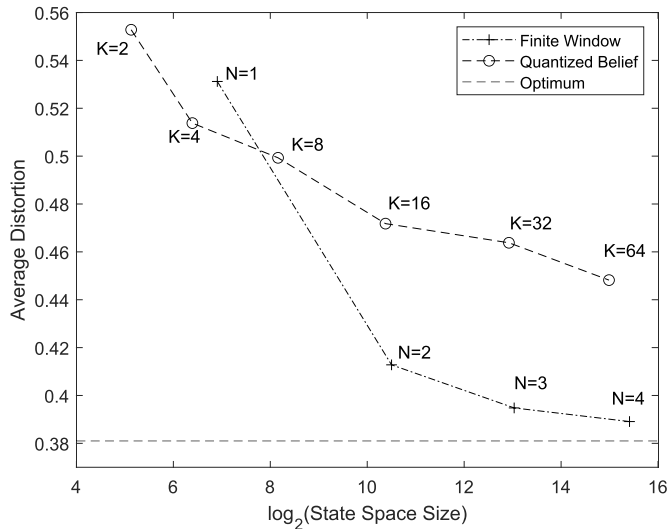


Fig. 4. Comparison with memoryless encoding, high Dobrushin coefficient.

We also examine the role of the Dobrushin coefficient in the performance of these policies by adding more channel noise while keeping T the same; we consider

$$O = \begin{pmatrix} 0.7 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.7 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.7 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.7 \end{pmatrix}. \quad (\text{VI.5})$$

We now have that $\delta(O) = 0.4$ is larger. While the distortion increases compared to Figure 3, as expected with the higher channel noise, the *rate* at which the finite window method approaches the optimum appears to increase, in line with our exponential bounds in Section III. Meanwhile the belief quantization does not appear to benefit as much from the increased noise. This is seen in Figure 4.

C. Quantization of Discretized Gauss-Markov Source

Although our method is only directly applicable for finite sources, we include an experiment with a Gauss-Markov source so that we can compare to existing schemes in the literature. We consider the following source:

$$Z_{t+1} = 0.9Z_t + \eta_t, \quad (\text{VI.6})$$

where $\eta_t \sim \mathcal{N}(0, 1)$ are i.i.d. and $Z_0 \sim \mathcal{N}(0, 1)$.

We approximate the Gauss-Markov source using a finite state source. The values of this finite state source are given by the reconstruction values of a Lloyd-Max quantizer designed for the stationary distribution of the Gauss-Markov source. For each reconstruction value, its row in the finite-state transition matrix is given by the (cumulative) distribution function of the Gauss-Markov source over the corresponding bin. In our experiments, we use 8 bins for the discretization, leading to the following finite-state source (rounded to two decimal places for visual clarity, we use higher precision in simulation):

$$\mathcal{X} = \{-2.19, -1.37, -0.77, -0.25, 0.25, 0.77, 1.37, 2.19\},$$

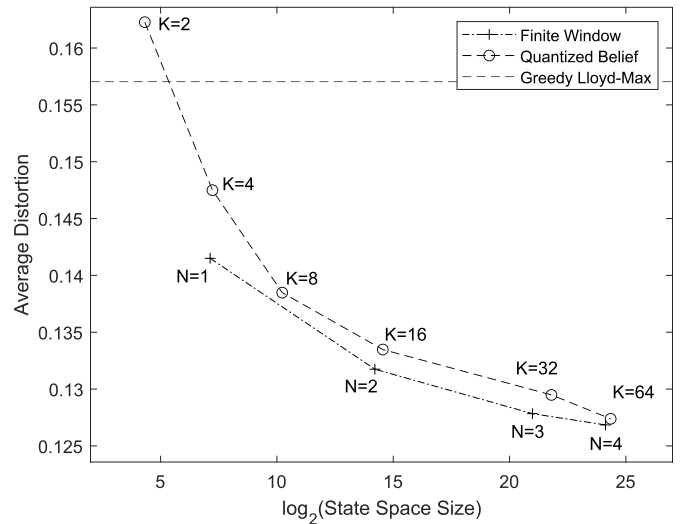


Fig. 5. Discrete Gauss-Markov quantization.

$$T = \begin{pmatrix} 0.11 & 0.19 & 0.21 & 0.19 & 0.15 & 0.09 & 0.04 & 0.01 \\ 0.07 & 0.15 & 0.19 & 0.20 & 0.17 & 0.13 & 0.07 & 0.02 \\ 0.05 & 0.13 & 0.18 & 0.20 & 0.19 & 0.14 & 0.08 & 0.03 \\ 0.04 & 0.11 & 0.17 & 0.20 & 0.19 & 0.16 & 0.10 & 0.03 \\ 0.03 & 0.10 & 0.16 & 0.19 & 0.20 & 0.17 & 0.11 & 0.04 \\ 0.03 & 0.08 & 0.14 & 0.19 & 0.20 & 0.18 & 0.13 & 0.05 \\ 0.02 & 0.07 & 0.13 & 0.17 & 0.20 & 0.19 & 0.15 & 0.07 \\ 0.01 & 0.04 & 0.09 & 0.15 & 0.19 & 0.21 & 0.19 & 0.11 \end{pmatrix}.$$

We wish to send this source over a noiseless channel with $\mathcal{M} = \mathcal{M}' = \{0, 1, 2, 3\}$ (that is, a quantization problem with four bins). Here we have $\delta(T) \approx 0.58$, so the Dobrushin coefficient conditions of Theorem 4 are met for this setup.

As a comparison, we use a greedy Lloyd-Max quantization scheme in which a Lloyd-Max quantizer is designed for the predictor π_t at each time step. This scheme is studied in [63] and corresponds to a greedy policy for zero-delay coding. It is shown that in the high-rate limit, this scheme also becomes globally optimal. We compare these schemes on two different setups:

- 1) Take as our true source the finite-state source, with the T matrix above and with X_0 distributed according to the invariant distribution of T . Here our reproduction alphabet is equal to our source alphabet, $\hat{\mathcal{X}} = \mathcal{X}$. For the greedy Lloyd-Max, since it can output any real number, we take the closest value in $\hat{\mathcal{X}}$ as the reproduction. We use the squared-error distortion $(X_t - \hat{X}_t)^2$; this is Figure 5.
- 2) Take as our true source the continuous Gauss-Markov source Z_t . At each time step, quantize it to the nearest $X_t \in \mathcal{X}$, and apply our finite-state schemes above. We use $\hat{\mathcal{X}} = \mathbb{R}$. Note that although we had assumed $\hat{\mathcal{X}}$ was finite in our scheme for simplicity, we can easily modify the decoder (II.5) and the cost function (II.13) to take any $\hat{\mathcal{X}}$. The distortion is given by $(Z_t - \hat{X}_t)^2$; this is Figure 6.

For the latter setup, the finite-state source $(X_t)_{t \geq 0}$ is not Markov as it is a quantization of the actual Gauss-Markov source, but our scheme treats it as if it were Markov (with the

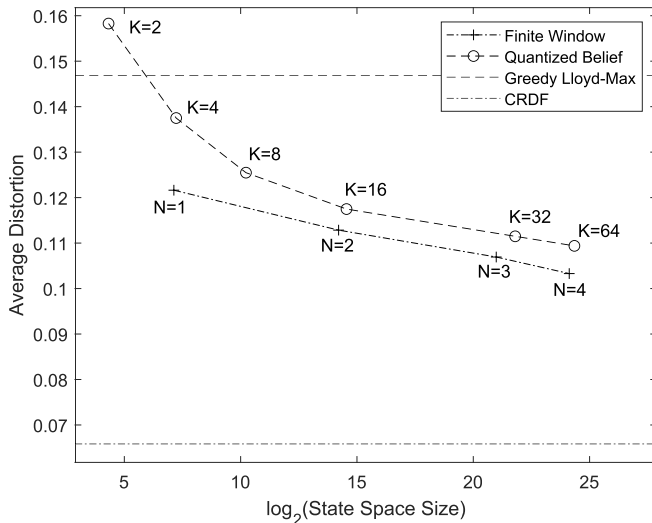


Fig. 6. Continuous Gauss-Markov quantization.

T matrix given above). For this setup, we also plot the causal rate distortion (CRDF) for the Gauss-Markov source, given in e.g., [15], [64] as $R(D) = \max \left\{ 0, \frac{1}{2} \log_2 \left(a^2 + \frac{\sigma^2}{D} \right) \right\}$. Note that in this low-rate regime the CRDF is not necessarily achievable [65]. Here each data point corresponds to $N = \{1, 2, 3, 4\}$ for the window length and $K = \{2, 4, 8, 16, 32, 64\}$ for the belief quantization, and again we plot against the number of states used in the policy.

We observe from Figures 5 and 6 that our method outperforms the greedy Lloyd-Max except for very small codebook sizes (in particular, for $K = 2$ in the belief quantization scheme). Furthermore, while our method takes much more offline computation, it is faster online as the policy is simply a lookup table, while the Lloyd-Max scheme must run the Lloyd-Max scheme to convergence at each iteration. However, our method is clearly sensitive to the discretization chosen for the Gauss-Markov source; if it is too coarse, it introduces additional error in the assumed dynamics, while if it is too fine the set of quantizers grows rapidly and thus the policy becomes more expensive to learn and implement.

Finally, although in this setting our state space grows very fast with N and K and thus becomes intractable for high values of these parameters, we note that we outperform the greedy Lloyd-Max method even for small parameter values, where the state space is still manageable. This is especially true in the finite-window scheme when the Dobrushin coefficients are high, as the exponential convergence gives us good performance for relatively small window sizes.

VII. CONCLUSION

We have provided two complementary approximation schemes for the zero-delay coding problem over a noisy channel with feedback, one in which the underlying belief space is quantized through a nearest-neighbor map, and the other in which the belief is approximated with a finite window of past observations. We then showed the convergence of a Q-learning algorithm to policies which are optimal for these approximations, yielding near-optimal policies for the zero-delay coding problem. Finally, we illustrated the convergence

of our algorithm to the optimum through simulations, and compared against other zero-delay schemes when the optimum is unknown. For future research, we wish to generalize our results to the case where our source-channel pair is continuous. We note that such a generalization should be mostly straightforward; under some recurrence assumptions on the source, one can still show the required regularity properties of the underlying MDP [39], and then most of the approximation and convergence results go through with minor technical changes. We also wish to find less stringent filter stability conditions for the finite window scheme than the Dobrushin coefficient ones given here; this may be possible via observability-type conditions, such as those in [66]. The controlled case, where one communicates over a channel to a controller which affects some dynamical system, is also a promising research direction.

APPENDIX A

PROOF OF THEOREMS 1 AND 2

Lemma 3: Consider the finite window approximation used in Section III, and recall $I_t = (M'_{[t-N,t-1]}, Q_{[t-N,t-1]})$, $W_t = (\pi_{t-N}, I_t)$, $\hat{W}_t = (\zeta, I_t)$, and $\hat{\pi}_t = P_\zeta^\gamma(X_t | m'_{[t-N,t-1]}, Q_{[t-N,t-1]})$. Then for any $w_t \in \mathcal{W}$, $\hat{w}_t \in \mathcal{W}_N$, and $Q_t \in \mathcal{Q}$ we have

$$\|P(M'_t \in \cdot | w_t, Q_t) - P(M'_t \in \cdot | \hat{w}_t, Q_t)\|_{\text{TV}} \leq \|\pi_t - \hat{\pi}_t\|_{\text{TV}}.$$

Proof: Let $f : \mathcal{M}' \rightarrow \mathbb{R}$ be measurable with $\|f\|_\infty \leq 1$. Then,

$$\begin{aligned} & \left| \sum_{\mathcal{M}'} f(m'_t) P(m'_t | w_t, Q_t) - \sum_{\mathcal{M}'} f(m'_t) P(m'_t | \hat{w}_t, Q_t) \right| \\ &= \left| \sum_{\mathcal{X}} \sum_{\mathcal{M}'} f(m'_t) P(m'_t | w_t, x_t, Q_t) P(x_t | w_t, Q_t) \right. \\ & \quad \left. - \sum_{\mathcal{X}} \sum_{\mathcal{M}'} f(m'_t) P(m'_t | \hat{w}_t, x_t, Q_t) P(x_t | \hat{w}_t, Q_t) \right| \\ &= \left| \sum_{\mathcal{X}} \sum_{\mathcal{M}'} f(m'_t) O_{Q_t}(m'_t | x_t) \pi_t(x_t) \right. \\ & \quad \left. - \sum_{\mathcal{X}} \sum_{\mathcal{M}'} f(m'_t) O_{Q_t}(m'_t | x_t) \hat{\pi}_t(x_t) \right| \\ &\leq \|\pi_t - \hat{\pi}_t\|_{\text{TV}}, \end{aligned}$$

where the third line follows from conditional independence of M'_t and W_t given (X_t, Q_t) , and since Q_t is a function of W_t under any $\gamma \in \Gamma_{\text{WS}}$. The last line follows from the fact that $g(x) := \sum_{\mathcal{M}'} f(m') O_{Q_t}(m' | x)$ is upper bounded by 1. Taking the supremum over all such f yields the result. ■

A. Proof of Theorem 1

We provide a proof for the case when $N = 1$, but an analogous proof follows for $N > 1$. It can be shown (see [40, Theorem 4.2.3]) that the functions $J_\beta(w_t, \gamma)$ and $J_\beta^*(w_t)$ satisfy the following fixed-point equations:

$$\begin{aligned} J_\beta(w_t, \gamma) &= c(w_t, \gamma(w_t)) \\ & \quad + \beta \int_{\mathcal{W}} J_\beta(w_{t+1}, \gamma) P(dw_{t+1} | w_t, \gamma(w_t)) \end{aligned}$$

$$J_\beta^*(w_t) = \min_{Q_t \in \mathcal{Q}} \left(c(w_t, Q_t) + \beta \int_{\mathcal{W}} J_\beta^*(w_{t+1}) P(dw_{t+1}|w_t, Q_t) \right),$$

for all $w_t \in \mathcal{W}$ and $\gamma \in \Gamma_{FS}$. Note that although the integral is over \mathcal{W} , which is uncountable, we can only reach finitely many elements from a given w_t since the observation space \mathcal{M}' is finite. In particular, when $N = 1$, we can write $w_t = (\pi_{t-1}, m'_{t-1}, Q_{t-1})$ and $w_{t+1} = (\pi_t, m'_t, Q_t)$, so the above becomes

$$J_\beta(w_t, \gamma) = c(w_t, \gamma(w_t)) + \beta \sum_{m'_t \in \mathcal{M}'} J_\beta((\pi_t, m'_t, \gamma(w_t)), \gamma) P(m'_t|w_t, \gamma(w_t)) \quad (\text{A.1})$$

$$J_\beta^*(w_t) = \min_{Q_t \in \mathcal{Q}} \left(c(w_t, Q_t) + \beta \sum_{m'_t \in \mathcal{M}'} J_\beta^*(\pi_t, m'_t, Q_t) P(m'_t|w_t, Q_t) \right). \quad (\text{A.2})$$

The functions $\hat{J}_\beta(\hat{w}_t, \hat{\gamma})$ and $\hat{J}_\beta^*(\hat{w}_t)$ satisfy equivalent fixed-point equations to (A.2), so that for $N = 1$,

$$\begin{aligned} \hat{J}_\beta(\hat{w}_t, \hat{\gamma}) &= c_1(\hat{w}_t, \hat{\gamma}(\hat{w}_t)) \\ &\quad + \beta \sum_{m'_t \in \mathcal{M}'} \hat{J}_\beta(\zeta, m'_t, \hat{\gamma}(\hat{w}_t)) P(m'_t|\hat{w}_t, \hat{\gamma}(\hat{w}_t)) \\ \hat{J}_\beta^*(\hat{w}_t) &= \min_{Q_t \in \mathcal{Q}} \left(c_1(\hat{w}_t, Q_t) \right. \\ &\quad \left. + \beta \sum_{m'_t \in \mathcal{M}'} \hat{J}_\beta^*(\zeta, m'_t, Q_t) P(m'_t|\hat{w}_t, Q_t) \right). \end{aligned} \quad (\text{A.3})$$

By definition of the extension \tilde{J}_β^* we have $\hat{J}_\beta^*(\hat{w}_1) = \tilde{J}_\beta^*(w_1)$. Thus,

$$\begin{aligned} &\beta \sum_{m'_0 \in \mathcal{M}'} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|w_0, Q_0) \\ &= \beta \sum_{m'_0 \in \mathcal{M}'} \tilde{J}_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|w_0, Q_0). \end{aligned}$$

We add and subtract the above term and use $\tilde{J}_\beta^*(w_0) = \hat{J}_\beta^*(\hat{w}_0)$ to obtain

$$\begin{aligned} &|\tilde{J}_\beta^*(w_0) - J_\beta^*(w_0)| \\ &= \left| \hat{J}_\beta^*(\hat{w}_0) - \beta \sum_{m'_0 \in \mathcal{M}'} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|w_0, Q_0) \right. \\ &\quad \left. + \beta \sum_{m'_0 \in \mathcal{M}'} \tilde{J}_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|w_0, Q_0) - J_\beta^*(w_0) \right|. \end{aligned}$$

Then applying the fixed-point equations (A.2) and (A.3) to the last and first terms respectively,

$$\begin{aligned} &|\tilde{J}_\beta^*(w_0) - J_\beta^*(w_0)| \\ &\leq \max_{Q_0 \in \mathcal{Q}} |c_1(\hat{w}_0, Q_0) - c(w_0, Q_0)| \\ &\quad + \beta \max_{Q_0 \in \mathcal{Q}} \left| \sum_{m'_0} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|\hat{w}_0, Q_0) \right. \\ &\quad \left. - \sum_{m'_0} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|w_0, Q_0) \right| \end{aligned}$$

$$\begin{aligned} &+ \beta \max_{Q_0 \in \mathcal{Q}} \left| \sum_{m'_0} \tilde{J}_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|w_0, Q_0) \right. \\ &\quad \left. - \sum_{m'_0} J_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|w_0, Q_0) \right|. \end{aligned}$$

We now bound these three terms in expectation. The expectation is on W_0 and \hat{W}_0 , with respect to the prior π_{-1} and some $\gamma \in \Gamma_{WS}$, but we omit these in the expectation for notational simplicity. For the first term, by definition of c and c_N we have

$$\begin{aligned} \mathbf{E} \left[|c_1(\hat{W}_0, Q_0) - c(W_0, Q_0)| \right] &\leq \|d\|_\infty \mathbf{E} [\|\hat{\pi}_0 - \pi_0\|_{TV}] \\ &\leq \|d\|_\infty L_0^1, \end{aligned} \quad (\text{A.4})$$

where $\|d\|_\infty = \max_{x, \hat{x}} d(x, \hat{x})$ and we recall the definition of L_t^N in (III.11); that is, $L_0^1 = \sup_{\gamma \in \Gamma_{WS}} \mathbf{E} [\|\pi_t - \hat{\pi}_t\|_{TV}]$. For the second term, we have

$$\begin{aligned} &\mathbf{E} \left[\max_{Q_0 \in \mathcal{Q}} \left| \sum_{m'_0} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|\hat{W}_0, Q_0) \right. \right. \\ &\quad \left. \left. - \sum_{m'_0} \hat{J}_\beta^*(\zeta, m'_0, Q_0) P(m'_0|W_0, Q_0) \right| \right] \\ &\leq \|\hat{J}_\beta^*\|_\infty \mathbf{E} \left[\max_{Q_0 \in \mathcal{Q}} \|P(m'_0|\hat{W}_0, Q_0) - P(m'_0|W_0, Q_0)\|_{TV} \right] \\ &\leq \|\hat{J}_\beta^*\|_\infty L_0^1, \end{aligned}$$

where the first inequality follows from the definition of the total variation (since $\hat{J}_\beta^*/\|\hat{J}_\beta^*\|_\infty$ is bounded by 1) and the second inequality is due to Lemma 3. Finally, since both sums in the last term are over the same measure $P(m'_0|W_0, Q_0)$, we have

$$\begin{aligned} &\mathbf{E} \left[\max_{Q_0 \in \mathcal{Q}} \left| \sum_{m'_0} \tilde{J}_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|W_0, Q_0) \right. \right. \\ &\quad \left. \left. - \sum_{m'_0} J_\beta^*(\pi_0, m'_0, Q_0) P(m'_0|W_0, Q_0) \right| \right] \\ &\leq \sup_{\gamma' \in \Gamma_{WS}} \mathbf{E} \left[|\tilde{J}_\beta^*(W_1) - J_\beta^*(W_1)| \right], \end{aligned}$$

where we have used $(\pi_0, M'_0, Q_0) = W_1$. Combining all three bounds (and multiplying by β where appropriate) gives us

$$\begin{aligned} &\mathbf{E} \left[|\tilde{J}_\beta^*(W_0) - J_\beta^*(W_0)| \right] \\ &\leq (\|d\|_\infty + \beta \|\hat{J}_\beta^*\|_\infty) L_0^1 + \beta \sup_{\gamma' \in \Gamma_{WS}} \mathbf{E} \left[|\tilde{J}_\beta^*(W_1) - J_\beta^*(W_1)| \right] \end{aligned}$$

We apply the same process to the final term, then recursively and by the fact that $\|\hat{J}_\beta^*\|_\infty \leq \frac{\|d\|_\infty}{1-\beta}$, we have

$$\mathbf{E} \left[|\tilde{J}_\beta^*(W_0) - J_\beta^*(W_0)| \right] \leq \frac{\|d\|_\infty}{1-\beta} \sum_{t=0}^{\infty} \beta^t L_t^1. \quad \blacksquare$$

B. Proof of Theorem 2

As before, we consider $N = 1$, but analogous arguments follow for $N > 1$. We apply a similar strategy, by using the fixed-point equations in the proof of Theorem 2. Also, let $Q_0^* := \tilde{\gamma}_1^*(w_0)$; that is, the action given by making

the optimal policy for MDP_1 constant over $\mathcal{P}(\mathcal{X})$. Then, by (A.1),

$$J_\beta(w_0, \tilde{\gamma}_1^*) = c(w_0, Q_0^*) + \beta \sum_{m'_0 \in \mathcal{M}'} J_\beta((\pi_0, m'_0, Q_0^*), \tilde{\gamma}_1^*) P(m'_0 | w_0, Q_0^*)$$

and using (A.3) and the fact that $\tilde{J}_\beta^*(w_0) = \hat{J}_\beta^*(\hat{w}_0)$,

$$\tilde{J}_\beta^*(w_0) = c_1(\hat{w}_0, Q_0^*) + \beta \sum_{m'_0 \in \mathcal{M}'} \tilde{J}_\beta^*(\pi_0, m'_0, Q_0^*) P(m'_0 | \hat{w}_0, Q_0^*).$$

Using $w_1 = (\pi_0, m'_0, Q_0^*)$, we add and subtract

$$\sum_{m'_0 \in \mathcal{M}'} \tilde{J}_\beta^*(w_1) P(m'_0 | w_0, Q_0^*),$$

to obtain

$$\begin{aligned} & |J_\beta(w_0, \tilde{\gamma}_1^*) - \tilde{J}_\beta^*(w_0)| \\ & \leq |c(w_0, Q_0^*) - c(\hat{w}_0, Q_0^*)| \\ & + \beta \sum_{m'_0 \in \mathcal{M}'} |\tilde{J}_\beta^*(w_1) P(m'_0 | \hat{w}_0, Q_0^*) - \tilde{J}_\beta^*(w_1) P(m'_0 | w_0, Q_0^*)| \\ & + \beta \sum_{m'_0 \in \mathcal{M}'} |\tilde{J}_\beta^*(w_1) - J_\beta(w_1, \tilde{\gamma}_1^*)| P(m'_0 | w_0, Q_0^*). \end{aligned}$$

Thus, using (A.4) and Lemma 3,

$$\begin{aligned} & \mathbf{E} [|J_\beta(W_0, \tilde{\gamma}_1^*) - \tilde{J}_\beta^*(W_0)|] \\ & \leq \|d\|_\infty L_0^1 \\ & + \beta \|\tilde{J}_\beta^*\|_\infty \sup_{\gamma \in \Gamma_{\text{ws}}} \mathbf{E} [|P(m'_0 | \hat{W}_0, Q_0^*) - P(m'_0 | W_0, Q_0^*)|_{\text{TV}}] \\ & + \beta \sup_{\gamma \in \Gamma_{\text{ws}}} \mathbf{E} [| \tilde{J}_\beta^*(W_1) - J_\beta(W_1, \tilde{\gamma}_1^*) |] \\ & \leq (\|d\|_\infty + \beta \|\tilde{J}_\beta^*\|_\infty) L_0^1 \\ & + \beta \sup_{\gamma} \mathbf{E} [|J_\beta(W_1, \tilde{\gamma}_1^*) - \tilde{J}_\beta^*(W_1)|]. \end{aligned}$$

Recursively, and using the fact that $\|\tilde{J}_\beta^*\|_\infty \leq \frac{\|d\|_\infty}{1-\beta}$,

$$\mathbf{E} [|J_\beta^*(W_0, \tilde{\gamma}_1^*) - J_\beta^*(W_0)|] \leq \frac{\|d\|_\infty}{1-\beta} \sum_{t=0}^{\infty} \beta^t L_t^1. \quad (\text{A.5})$$

Finally, we have

$$\begin{aligned} & \mathbf{E} [|J_\beta^*(W_0, \tilde{\gamma}_1^*) - J_\beta^*(W_0)|] \\ & \leq \mathbf{E} [|J_\beta^*(W_0, \tilde{\gamma}_1^*) - \tilde{J}_\beta^*(W_0)|] + \mathbf{E} [| \tilde{J}_\beta^*(W_0) - J_\beta^*(W_0) |] \\ & \leq \frac{2\|d\|_\infty}{1-\beta} \sum_{t=0}^{\infty} \beta^t L_t^1, \end{aligned}$$

where the final inequality follows from (A.5) and Theorem 1. ■

APPENDIX B PROOF OF THEOREM 3

Lemma 4: The following holds:

$$\begin{aligned} & \sum_x \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} O_Q(m'|x) \pi_t^\mu(x) \\ & \leq (2 - \tilde{\delta}(O)) \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}, \end{aligned}$$

where $\tilde{\delta}(O) = \min_{Q \in \mathcal{Q}} (\delta(O_Q))$.

Proof: The following argument is from [54, Lemma 3.5], adapted to our setup. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be measurable with $\|f\|_\infty \leq 1$. Recall the update equations for $\pi_t, \bar{\pi}_t$ given in (II.12), and let $N^\mu(M_t, Q_t)$ denote the normalizing term for the $(\pi_t^\mu)_{t \geq 0}$ process, $N^\mu(M_t', Q_t) := \sum_x O_{Q_t}(M_t'|x) \pi_t^\mu(x)$. Then we have for any $M_t' = m'$ and $Q_t = Q$,

$$\begin{aligned} & \left| \sum_x f(x) \bar{\pi}_t^\mu(x) - \sum_x f(x) \bar{\pi}_t^\nu(x) \right| \\ & = \left| \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\mu(x)}{N^\mu(m', Q)} - \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\nu(x)}{N^\nu(m', Q)} \right| \\ & \leq \left| \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\mu(x)}{N^\mu(m', Q)} - \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\nu(x)}{N^\mu(m', Q)} \right| \\ & \quad + \left| \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\nu(x)}{N^\mu(m', Q)} - \sum_x f(x) \frac{O_Q(m'|x) \pi_t^\nu(x)}{N^\nu(m', Q)} \right| \\ & = \frac{1}{N^\mu(m', Q)} \left| \sum_x f(x) O_Q(m'|x) \pi_t^\mu(x) - \sum_x f(x) O_Q(m'|x) \pi_t^\nu(x) \right| \\ & \quad + \left| \frac{1}{N^\mu(m', Q)} - \frac{1}{N^\nu(m', Q)} \right| \left| \sum_x f(x) O_Q(m'|x) \pi_t^\nu(x) \right| \\ & \leq \frac{1}{N^\mu(m', Q)} \sum_x O_Q(m'|x) |\pi_t^\mu - \pi_t^\nu|(x) \\ & \quad + \left| \frac{1}{N^\mu(m', Q)} - \frac{1}{N^\nu(m', Q)} \right| N^\nu(m', Q) \\ & \leq \frac{1}{N^\mu(m', Q)} \left(\sum_x O_Q(m'|x) |\pi_t^\mu - \pi_t^\nu|(x) + |N^\mu(m', Q) - N^\nu(m', Q)| \right), \end{aligned} \quad (\text{B.1})$$

where in the second last line we have used the notation $\sum_x |\pi_t^\mu - \pi_t^\nu|(x) = \sum_x (\mathbf{1}_{S^+} - \mathbf{1}_{S^-})(\pi_t^\mu - \pi_t^\nu)(x)$ with $S^+ = \{x | (\pi_t^\mu - \pi_t^\nu)(x) > 0\}$ and $S^- = \{x | (\pi_t^\mu - \pi_t^\nu)(x) \leq 0\}$. Note that $\sum_x |\pi_t^\mu - \pi_t^\nu|(x) = \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}$. Taking the supremum over all f gives

$$\begin{aligned} & \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} \\ & \leq \frac{1}{N^\mu(m', Q)} \left(\sum_x O_Q(m'|x) |\pi_t^\mu - \pi_t^\nu|(x) + |N^\mu(m', Q) - N^\nu(m', Q)| \right). \end{aligned} \quad (\text{B.1})$$

Thus, we have

$$\begin{aligned} & \sum_x \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} O_Q(m'|x) \pi_t^\mu(x) \\ & = \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} N^\mu(m', Q_t) \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{\mathcal{M}'} \left(\sum_{\mathcal{X}} O_{Q_t}(m'|x) |\pi_t^\mu - \pi_t^\nu|(x) \right. \\
&\quad \left. + |N^\mu(m', Q_t) - N^\nu(m', Q_t)| \right) \\
&\leq \sum_{\mathcal{X}} \left(\sum_{\mathcal{M}'} O_{Q_t}(m'|x) \right) |\pi_t^\mu - \pi_t^\nu|(x) \\
&\quad + \sum_{\mathcal{M}'} \left| \sum_{\mathcal{X}} O_{Q_t}(m'|x) (\pi_t^\mu - \pi_t^\nu)(x) \right| \\
&\leq \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}} + \sum_{\mathcal{M}'} |O_{Q_t}(\pi_t^\mu) - O_{Q_t}(\pi_t^\nu)|(m') \\
&= \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}} + \|O_{Q_t}(\pi_t^\mu) - O_{Q_t}(\pi_t^\nu)\|_{\text{TV}},
\end{aligned}$$

where in the second last line we have used the kernel operator notation $O_{Q_t}(\pi)(m') = \sum_{\mathcal{X}} O_{Q_t}(m'|x)\pi(x)$. It is shown in [67] that the Dobrushin coefficient acts as a contraction coefficient for kernel operators under total variation. In particular

$$\|O_{Q_t}(\pi_t^\mu) - O_{Q_t}(\pi_t^\nu)\|_{\text{TV}} \leq (1 - \delta(O_{Q_t})) \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}. \quad (\text{B.2})$$

Thus,

$$\begin{aligned}
&\sum_{\mathcal{X}} \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} O_{Q_t}(m'|x) \pi_t^\mu(x) \\
&\leq (2 - \delta(O_{Q_t})) \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}} \\
&\leq (2 - \tilde{\delta}(O)) \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}},
\end{aligned}$$

where $\tilde{\delta}(O) = \min_{Q \in \mathcal{Q}} (\delta(O_Q))$.

A. Proof of Theorem 3

Note that, in \mathbf{E}_μ^γ expectations of $\bar{\pi}_t^\mu$ and $\bar{\pi}_t^\nu$, it is enough to take the expectation over only $M'_{[0,t]}$, since under any $\gamma \in \Gamma_{\text{ws}}$, $Q_{[0,t]}$ are deterministic given μ and $M'_{[0,t-1]}$. Thus,

$$\begin{aligned}
&\mathbf{E}_\mu^\gamma [\|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}}] \\
&= \sum_{(\mathcal{M}')^{\gamma+1}} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} P_\mu^\gamma(m'_{[0,t]}) \\
&= \sum_{(\mathcal{M}')^\gamma} \sum_{\mathcal{X}} \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} P_\mu^\gamma(m'_t|x_t, m'_{[0,t-1]}) \\
&\quad \cdot P_\mu^\gamma(x_t|m'_{[0,t-1]}) P_\mu^\gamma(m'_{[0,t-1]}) \\
&= \sum_{(\mathcal{M}')^\gamma} \sum_{\mathcal{X}} \sum_{\mathcal{M}'} \|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}} O_{Q_t}(m'_t|x_t) \pi_t^\mu(x) P_\mu^\gamma(m'_{[0,t-1]}) \\
&\leq (2 - \tilde{\delta}(O)) \sum_{(\mathcal{M}')^\gamma} \|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}} P_\mu^\gamma(m'_{[0,t-1]}) \\
&= (2 - \tilde{\delta}(O)) \mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}],
\end{aligned}$$

where the third equality follows from the fact that Q_t is deterministic given μ and $M'_{[0,t-1]}$, and that given X_t and Q_t , M'_t depends only on the kernel O_{Q_t} . For the inequality we used Lemma 4. Finally, using the Dobrushin contraction property for kernels (as noted in the derivation of (B.2)), we have

$$\begin{aligned}
&\mathbf{E}_\mu^\gamma [\|\pi_{t+1}^\mu - \pi_{t+1}^\nu\|_{\text{TV}}] \\
&\leq (1 - \delta(T)) \mathbf{E}_\mu^\gamma [\|\bar{\pi}_t^\mu - \bar{\pi}_t^\nu\|_{\text{TV}}] \\
&\leq (1 - \delta(T)) (2 - \tilde{\delta}(O)) \mathbf{E}_\mu^\gamma [\|\pi_t^\mu - \pi_t^\nu\|_{\text{TV}}].
\end{aligned}$$

APPENDIX C

PROOF OF THEOREM 4

We will show certain ergodic behavior (in particular, [56, Assumption 2.1]) and then invoke [56, Theorem 2.1]. Indeed, we have trivially that

$$\frac{\sum_{k=0}^t c_N(\hat{w}_k, Q_k) \mathbf{1}(\hat{w}_k = \hat{w}, Q_k = Q)}{\sum_{k=0}^t \mathbf{1}(\hat{w}_k = \hat{w}, Q_k = Q)} = c_N(\hat{w}, Q),$$

so that [56, Assumption 2.2 (ii)] holds. Additionally, by positive Harris recurrence of $(X_t)_{t \geq 0}$, the marginals on $P(X_t \in \cdot)$ converge to ζ , so we have for any f that

$$\frac{\sum_{k=0}^t f(\hat{w}_{k+1}) \mathbf{1}(\hat{w}_k = \hat{w}, Q_k = Q)}{\sum_{k=0}^t \mathbf{1}(\hat{w}_k = \hat{w}, Q_k = Q)} \rightarrow \int f(\hat{w}_1) P_N(\hat{w}_1 | \hat{w}, Q)$$

almost surely as $t \rightarrow \infty$, where P_N is from (III.9). Thus [56, Assumption 2.1 (iii)] holds.

Finally, note that not every $\hat{w} \in \mathcal{W}_N$ has positive probability of being visited (certain sequences of channel outputs and quantizers are impossible depending on the source and channel), so we restrict ourselves to only those with positive probability. Then [56, Assumption 2.1 (i)] holds for (almost every) \hat{w} , so we apply [56, Theorem 2.1] to obtain that almost surely, as $N \rightarrow \infty$,

$$V_t(\hat{w}, Q) \rightarrow V^*(\hat{w}, Q),$$

where

$$V^*(\hat{w}, Q) = c_N(\hat{w}, Q) + \beta \sum_{\hat{w}_1 \in \mathcal{W}_N} \min_Q V^*(\hat{w}_1) P_N(\hat{w}_1 | \hat{w}, Q),$$

where P_N and c_N are from (III.9) and (III.10). Taking the minimum over Q gives us the classic Bellman optimality equation, and hence $\hat{J}_\beta(\hat{w}, \hat{\gamma}_N^*) = \hat{J}_\beta^*(\hat{w})$ for almost every $\hat{w} \in \mathcal{W}_N$. The result follows by applying Corollary 2. ■

REFERENCES

- [1] A. Badr, A. Khisti, W.-T. Tan, and J. Apostolopoulos, "Streaming codes for channels with burst and isolated erasures," in *Proc. IEEE INFOCOM*, Apr. 2013, pp. 2850–2858.
- [2] M. Rudow and K. V. Rashmi, "Streaming codes for variable-size messages," *IEEE Trans. Inf. Theory*, vol. 68, no. 9, pp. 5823–5849, Sep. 2022.
- [3] A. Khisti and S. C. Draper, "The streaming-DMT of fading channels," *IEEE Trans. Inf. Theory*, vol. 60, no. 11, pp. 7058–7072, Nov. 2014.
- [4] A. Badr, P. Patil, A. Khisti, W.-T. Tan, and J. Apostolopoulos, "Layered constructions for low-delay streaming codes," *IEEE Trans. Inf. Theory*, vol. 63, no. 1, pp. 111–141, Jan. 2017.
- [5] A. Khina, V. Kostina, A. Khisti, and B. Hassibi, "Tracking and control of Gauss–Markov processes over packet-drop channels with acknowledgments," *IEEE Trans. Control Netw. Syst.*, vol. 6, no. 2, pp. 549–560, Jun. 2019.
- [6] M. Gastpar, B. Rimoldi, and M. Vetterli, "To code, or not to code: Lossy source-channel communication revisited," *IEEE Trans. Inf. Theory*, vol. 49, no. 5, pp. 1147–1158, May 2003.
- [7] M. Gastpar, "Uncoded transmission is exactly optimal for a simple Gaussian 'sensor' network," *IEEE Trans. Inf. Theory*, vol. 54, pp. 177–182, 2007.
- [8] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, no. 3, pp. 379–423, Jul. 1948.
- [9] E. I. Silva, M. S. Derpich, and J. Ostergaard, "A framework for control system design subject to average data-rate constraints," *IEEE Trans. Autom. Control*, vol. 56, no. 8, pp. 1886–1899, Aug. 2011.
- [10] E. I. Silva, M. S. Derpich, J. Ostergaard, and M. A. Encina, "A characterization of the minimal average data rate that guarantees a given closed-loop performance level," *IEEE Trans. Autom. Control*, vol. 61, no. 8, pp. 2171–2186, Aug. 2016.

- [11] P. A. Stavrou, J. Østergaard, and C. D. Charalambous, “Zero-delay rate distortion via filtering for vector-valued Gaussian sources,” *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 5, pp. 841–856, Oct. 2018.
- [12] R. Bansal and T. Başar, “Simultaneous design of measurement and control strategies for stochastic systems with feedback,” *Automatica*, vol. 25, no. 5, pp. 679–694, Sep. 1989.
- [13] S. Tatikonda, A. Sahai, and S. Mitter, “Stochastic linear control over a communication channels,” *IEEE Trans. Autom. Control*, vol. 49, pp. 1549–1561, 2004.
- [14] T. Tanaka, K.-K.-K. Kim, P. A. Parrilo, and S. K. Mitter, “Semidefinite programming approach to Gaussian sequential rate-distortion trade-offs,” *IEEE Trans. Autom. Control*, vol. 62, no. 4, pp. 1896–1910, Apr. 2017.
- [15] M. S. Derpich and J. Ostergaard, “Improved upper bounds to the causal quadratic rate-distortion function for Gaussian stationary sources,” *IEEE Trans. Inf. Theory*, vol. 58, no. 5, pp. 3131–3152, May 2012.
- [16] P. A. Stavrou and M. Skoglund, “Asymptotic reverse waterfilling algorithm of NRDF for certain classes of vector Gauss–Markov processes,” *IEEE Trans. Autom. Control*, vol. 67, no. 6, pp. 3196–3203, Jun. 2022.
- [17] P. A. Stavrou, T. Tanaka, and S. Tatikonda, “The time-invariant multidimensional Gaussian sequential rate-distortion problem revisited,” *IEEE Trans. Autom. Control*, vol. 65, no. 5, pp. 2245–2249, May 2020.
- [18] V. Kostina and B. Hassibi, “Rate-cost tradeoffs in control,” *IEEE Trans. Autom. Control*, vol. 64, no. 11, pp. 4525–4540, Nov. 2019.
- [19] J. Chakravorty and A. Mahajan, “Fundamental limits of remote estimation of autoregressive Markov processes under communication constraints,” *IEEE Trans. Autom. Control*, vol. 62, no. 3, pp. 1109–1124, Mar. 2017.
- [20] N. Guo and V. Kostina, “Optimal causal rate-constrained sampling for a class of continuous Markov processes,” *IEEE Trans. Inf. Theory*, vol. 67, no. 12, pp. 7876–7890, Dec. 2021.
- [21] T. Soleymani, J. S. Baras, and K. H. Johansson, “State estimation over delayed and lossy channels: An encoder–decoder synthesis,” *IEEE Trans. Autom. Control*, vol. 69, no. 3, pp. 1568–1583, Mar. 2024.
- [22] D. Pollard, “Quantization and the method of k -means,” *IEEE Trans. Inf. Theory*, vol. IT-28, pp. 199–205, 1982.
- [23] T. Linder, G. Lugosi, and K. Zeger, “Rates of convergence in the source coding theorem, in empirical quantizer design, and in universal lossy source coding,” *IEEE Trans. Inf. Theory*, vol. 40, no. 6, pp. 1728–1740, Nov. 1994.
- [24] T. Linder, *Learning-theoretic Methods in Vector Quantization*. New York, NY, USA: Springer, 2002.
- [25] J. Walrand and P. Varaiya, “Optimal causal coding–decoding problems,” *IEEE Trans. Inf. Theory*, vol. IT-29, no. 6, pp. 814–820, Nov. 1983.
- [26] H. S. Witsenhausen, “On the structure of real-time source coders,” *Bell Syst. Tech. J.*, vol. 58, no. 6, pp. 1437–1451, Jul. 1979.
- [27] R. G. Wood, T. Linder, and S. Yüksel, “Optimal zero delay coding of Markov sources: Stationary and finite memory codes,” *IEEE Trans. Inf. Theory*, vol. 63, no. 9, pp. 5968–5980, Sep. 2017.
- [28] D. Neuhoff and R. Gilbert, “Causal source codes,” *IEEE Trans. Inf. Theory*, vol. IT-28, no. 5, pp. 701–713, Sep. 1982.
- [29] N. Gaarder and D. Slepian, “On optimal finite-state digital transmission systems,” *IEEE Trans. Inf. Theory*, vol. IT-28, no. 2, pp. 167–186, Mar. 1982.
- [30] T. Weissman and N. Merhav, “On causal source codes with side information,” *IEEE Trans. Inf. Theory*, vol. 51, no. 11, pp. 4003–4013, Nov. 2005.
- [31] H. Asnani and T. Weissman, “On real time coding with limited lookahead,” *IEEE Trans. Inf. Theory*, vol. 59, no. 6, pp. 3582–3606, Jun. 2013.
- [32] D. Gündüz, P. de Kerret, N. D. Sidiropoulos, D. Gesbert, C. R. Murthy, and M. van der Schaar, “Machine learning in the air,” *IEEE J. Sel. Areas Commun.*, vol. 37, no. 10, pp. 2184–2199, Oct. 2019.
- [33] N. Farsad, M. Rao, and A. Goldsmith, “Deep learning for joint source-channel coding of text,” in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Calgary, AB, Canada, Apr. 2018, pp. 2326–2330.
- [34] E. Boursoulatzé, D. B. Kurka, and D. Gündüz, “Deep joint source-channel coding for wireless image transmission,” *IEEE Trans. Cognit. Commun. Netw.*, vol. 5, no. 3, pp. 567–579, Sep. 2019.
- [35] Z. Aharoni, O. Sabag, and H. H. Permuter, “Computing the feedback capacity of finite state channels using reinforcement learning,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Jul. 2019, pp. 837–841.
- [36] Z. Aharoni, O. Sabag, and H. H. Permuter, “Feedback capacity of Ising channels with large alphabet via reinforcement learning,” *IEEE Trans. Inf. Theory*, vol. 68, no. 9, pp. 5637–5656, Sep. 2022.
- [37] D. Teneketzis, “On the structure of optimal real-time encoders and decoders in noisy communication,” *IEEE Trans. Inf. Theory*, vol. 52, no. 9, pp. 4017–4035, Sep. 2006.
- [38] A. Mahajan and D. Teneketzis, “Optimal design of sequential real-time communication systems,” *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 5317–5338, Nov. 2009.
- [39] M. Ghomi, T. Linder, and S. Yüksel, “Zero-delay lossy coding of linear vector Markov sources: Optimality of stationary codes and near optimality of finite memory codes,” *IEEE Trans. Inf. Theory*, vol. 68, no. 5, pp. 3474–3488, May 2022.
- [40] O. Hernandez-Lerma and J. B. Lasserre, *Discrete-Time Markov Control Processes: Basic Optimality Criteria*. New York, NY, USA: Springer, 1996.
- [41] S. Yüksel. (2023). *Optimization and Control of Stochastic Systems*. [Online]. Available: <https://mast.queensu.ca/yuksel/LectureNotesOnStochasticOptControl.pdf>
- [42] A. D. Kara and S. Yüksel, “Convergence of finite memory q learning for POMDPs and near optimality of learned policies under filter stability,” *Math. Oper. Res.*, vol. 48, no. 4, pp. 2066–2093, Nov. 2022.
- [43] L. Cregg, T. Linder, and S. Yüksel, “Reinforcement learning for near-optimal design of zero-delay codes for Markov sources,” *IEEE Trans. Inf. Theory*, vol. 70, no. 11, pp. 8399–8413, Nov. 2024.
- [44] F. Pollara, R. J. McEliece, and K. A. Ghaffar, “Finite-state codes,” *IEEE Trans. Inf. Theory*, vol. IT-34, no. 5, pp. 1083–1089, Sep. 1988.
- [45] J. Kieffer and J. Dunham, “On a type of stochastic stability for a class of encoding schemes,” *IEEE Trans. Inf. Theory*, vol. IT-29, no. 6, pp. 793–797, Nov. 1983.
- [46] G. Ungerboeck, “Channel coding with multilevel/phase signals,” *IEEE Trans. Inf. Theory*, vol. IT-28, no. 1, pp. 55–67, Jan. 1982.
- [47] A. Calderbank and J. Mazo, “A new description of trellis codes,” *IEEE Trans. Inf. Theory*, vol. IT-30, no. 6, pp. 784–791, Nov. 1984.
- [48] R. Johansson and K. S. Zigangirov, *Fundamentals of Convolutional Coding*. Hoboken, NJ, USA: Wiley, 2015.
- [49] O. Hernández-Lerma and J. B. Lasserre, *Markov Chains and Invariant Probabilities*. Basel, Switzerland: Birkhäuser, 2003.
- [50] T. Linder and S. Yüksel, “On optimal zero-delay coding of vector Markov sources,” *IEEE Trans. Inf. Theory*, vol. 60, no. 10, pp. 5975–5991, Oct. 2014.
- [51] R. van Handel. (2007). *Stochastic Calculus, Filtering, and Stochastic Control*. [Online]. Available: <http://www.princeton.edu/~rvan/acm217/ACM217.pdf>
- [52] O. Cappé, E. Moulines, and T. Rydén, *Inference in Hidden Markov Models*. Berlin, Germany: Springer, 2005.
- [53] P. Chigansky and R. Liptser, “Stability of nonlinear filters in non-mixing case,” *Ann. Appl. Probab.*, vol. 14, no. 4, pp. 2038–2056, Nov. 2004.
- [54] C. McDonald and S. Yüksel, “Exponential filter stability via Dobrushin’s coefficient,” *Electron. Commun. Probab.*, vol. 25, no. none, pp. 1–13, Jan. 2020.
- [55] C. J. Watkins and P. Dayan, “Q-learning,” *Mach. Learn.*, vol. 8, nos. 3–4, pp. 279–292, 1992.
- [56] A. D. Kara and S. Yüksel, “Q-learning for stochastic control under general information structures and non-Markovian environments,” *Trans. Mach. Learn. Res.*, May 2024.
- [57] O. Bicer, A. D. Kara, and S. Yüksel, “Quantizer design for finite model approximations, model learning, and quantized Q-learning for MDPs with unbounded spaces,” 2025, *arXiv:2510.04355*.
- [58] A. D. Kara and S. Yüksel, “Robustness to incorrect system models in stochastic control,” *SIAM J. Control Optim.*, vol. 58, no. 2, pp. 1144–1182, Jan. 2020.
- [59] Y. A. Reznik, “An algorithm for quantization of discrete probability distributions,” in *Proc. Data Compress. Conf.*, Mar. 2011, pp. 333–342.
- [60] I. Csizsar, “The method of types [information theory],” *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2505–2523, Oct. 1998.
- [61] F. Le Gland and N. Oudjane, “Stability and uniform approximation of nonlinear filters using the Hilbert metric and application to particle filters,” *Ann. Appl. Probab.*, vol. 14, no. 1, pp. 144–187, Feb. 2004.
- [62] Y. E. Demirci, A. D. Kara, and S. Yüksel, “Refined bounds on near optimality finite window policies in POMDPs and their reinforcement learning,” 2024, *arXiv:2409.04351*.

- [63] A. Khina, Y. Nakahira, Y. Su, and B. Hassibi, "Algorithms for optimal control with fixed-rate feedback," in *Proc. IEEE 56th Annu. Conf. Decis. Control (CDC)*, Dec. 2017, pp. 6015–6020.
- [64] A. Gorbunov and M. Pinsker, "Prognostic epsilon entropy of a Gaussian message and a Gaussian source (in russian)," *Problemy Peredachi Informatsii*, vol. 10, no. 2, pp. 5–25, 1974.
- [65] T. Linder and R. Zamir, "Causal coding of stationary sources and individual sequences with high resolution," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 662–680, Feb. 2006.
- [66] C. McDonald and S. Yüksel, "Stochastic observability and filter stability under several criteria," *IEEE Trans. Autom. Control*, vol. 69, no. 5, pp. 1–16, May 2024.
- [67] R. L. Dobrushin, "Central limit theorem for nonstationary Markov Chains. I," *Theory Probab. Its Appl.*, vol. 1, no. 1, pp. 65–80, Jan. 1956.

Liam Cregg received the B.A.Sc. and M.A.Sc. degrees in mathematics and engineering from Queen's University in 2022 and 2024, respectively. He is currently pursuing the Ph.D. degree in electrical engineering with the Institut für Automatik, ETH Zürich. His research interests include stochastic control theory, information theory, and probability.

Fady Alajaji (Senior Member, IEEE) received the B.E. degree (Hons.) in electrical engineering from American University of Beirut, Lebanon, in 1988, and the M.Sc. and Ph.D. degrees in electrical engineering from the University of Maryland, College Park, in 1990 and 1994, respectively. In 1994, he held a post-doctoral appointment with the Institute for Systems Research, University of Maryland. In 1995, he joined the Department of Mathematics and Statistics, Queen's University, Kingston, ON, Canada, where he is currently a Professor of mathematics and engineering. Since 1997, he has been cross appointed with the Department of Electrical and Computer Engineering, Queen's University. From 2013 to 2014, he was the Acting Head of the Department of Mathematics and Statistics. From 2003 to 2008 and from 2018 to 2019, he was the Chair of the Queen's Mathematics and Engineering Program. His research interests include information theory, digital communications, error control coding, data compression, joint source-channel coding, network epidemics, generative models, data privacy, and fairness in machine learning. He received the Premiers Research Excellence Award from the Province of Ontario. He served as an organizer and a technical program committee member for several international conferences and workshops. From 2019 to 2022, he was an Associate Editor of Shannon Theory and Information Measures for IEEE TRANSACTIONS ON INFORMATION THEORY. He served as an Area Editor (2008–2015) for Source-Channel Coding and Signal Processing and an Editor (2003–2012) for Source and Source-Channel Coding for IEEE TRANSACTIONS ON COMMUNICATIONS.

Serdar Yüksel (Senior Member, IEEE) received the B.Sc. degree in electrical and electronics engineering from Bilkent University in 2001 and the M.S. and Ph.D. degrees in electrical and computer engineering from the University of Illinois Urbana-Champaign in 2003 and 2006, respectively. He was a Post-Doctoral Researcher with Yale University before joining the Department of Mathematics and Statistics, Queen's University, as an Assistant Professor, where he is currently a Professor. His research interests include stochastic control theory, information theory, and probability.