

Semidefinite Programming Approach to Gaussian Sequential Rate-Distortion Trade-offs

Takashi Tanaka, Kwang-Ki K. Kim, Pablo A. Parrilo, and Sanjoy K. Mitter

Abstract—Sequential rate-distortion (SRD) theory provides a framework for studying the fundamental trade-off between data-rate and data-quality in real-time communication systems. In this paper, we consider the SRD problem for multi-dimensional time-varying Gauss-Markov processes under mean-square distortion criteria. We first revisit the sensor-estimator separation principle, which asserts that considered SRD problem is equivalent to a joint sensor and estimator design problem in which data-rate of the sensor output is minimized while the estimator’s performance satisfies the distortion criteria. We then show that the optimal joint design can be performed by semidefinite programming. A semidefinite representation of the corresponding SRD function is obtained. Implications of the obtained result in the context of zero-delay source coding theory and applications to networked control theory are also discussed.

Index Terms—Control over communications; LMIs; Optimization algorithms; Stochastic optimal control; Kalman filtering

I. INTRODUCTION

In this paper, we study a fundamental performance limitation of zero-delay communication systems using the sequential rate-distortion (SRD) theory. Suppose that \mathbf{x}_t is an \mathbb{R}^n -valued discrete time random process with known statistical properties. At every time step, the encoder observes a realization of the source \mathbf{x}_t and generates a binary sequence $\mathbf{b}_t \in \{0, 1\}^{l_t}$ of length l_t , which is transmitted to the decoder. The decoder produces an estimation \mathbf{z}_t of \mathbf{x}_t based on the messages \mathbf{b}_t received up to time t . Both encoder and decoder have infinite memories of the past. A zero-delay communication system is determined by a selected encoder-decoder pair, whose performance is analyzed in the trade-off between the rate (*viz.* the average number of *bits* that must be transmitted per time step) and the distortion (*viz.* the discrepancy between the source signal \mathbf{x}_t and the reproduced signal \mathbf{z}_t). The region in the rate-distortion plane achievable by a zero-delay communication system is referred to as the *zero-delay rate-distortion region*.¹

The *standard rate-distortion region* identified by Shannon only provides a conservative outer bound of the zero-delay rate-distortion region. This is because, in general, achieving the standard rate-distortion region requires the use of anticipative (non-causal) codes (e.g., [1, Theorem 10.2.1]). It is well known that the standard rate-distortion region can be

expressed by the *rate-distortion function*² for general sources. In contrast, description of the zero-delay rate-distortion region requires more case-dependent knowledge of the optimal source coding schemes. For scalar memoryless sources, it is shown that the optimal performance of zero-delay codes is achievable by a scalar quantizer [2]. Witsenhausen [3] showed that for the k -th order Markov sources, there exists an optimal zero-delay quantizer with memory structure of order k . Neuhoff and Gilbert considered entropy-coded quantizers within the class of *causal source codes* [4], and showed that for memoryless sources, the optimal performance is achievable by time-sharing memoryless codes. This result is extended to sources with memory in [5]. An optimal memory structure of zero-delay quantizers for partially observable Markov processes on abstract (Polish) spaces is identified in [6]. The rate of finite-delay source codes for general sources and general distortion measures is analyzed in [7]. Zero-delay or finite-delay joint source-channel coding problems have also been studied in the literature; [8]–[11] to name a few.

In [12], [13], Tatikonda et al. studied the zero-delay rate-distortion region using a quantity called *sequential rate-distortion function*,³ which is defined as the infimum of the Massey’s directed information [18] from the source process to the reproduction process subject to the distortion constraint. Although the SRD function does not coincide with the boundary of the zero-delay rate-distortion region in general, it is recently shown that the SRD function provides a tight outer bound of the zero-delay rate-distortion region achievable by uniquely decodable codes [16], [19]. This observation shows an intimate connection between the SRD function and the fundamental performance limitations of real-time communication systems. For this reason, we consider the SRD function as the main object of interest in this paper.

Closely related quantity to the SRD function was studied by Gorbunov and Pinsker [14] in the early 1970’s. Bucy [15] derived the SRD function for Gauss-Markov processes in a simple case. In his approach, the problem of deriving the SRD function for Gauss-Markov processes under mean-square distortion criteria (which henceforth will be simply referred to as the *Gaussian SRD problem*) is viewed as a sensor-estimator joint design problem to minimize the estimation error subject

T. Tanaka is with ACCESS Linnaeus Center, KTH Royal Institute of Technology, Stockholm, 10044 Sweden.

P. A. Parrilo, and S. K. Mitter are with the Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA, 02139 USA.

K.-K. K. Kim is with the Electronic Control Development Team of Hyundai Motor Company (HMC) Research & Development Division in South Korea.

¹Formal definition of the zero-delay rate-distortion region is given in Section VI-A.

²This quantity is defined by the infimum of the mutual information between the source and the reproduction subject to the distortion constraint [1, Theorem 10.2.1].

³Closely related or apparently equivalent notions to the sequential rate-distortion function have been given various names in the literature, including nonanticipatory ϵ -entropy [14], constrained distortion rate function [15], causal rate-distortion function [16], and nonanticipative rate-distortion function [17].

to the data-rate constraint. This approach is justified by the “sensor-estimator separation principle,” which asserts that an optimal solution (i.e., the optimal stochastic kernel, to be made precise in the sequel) to the Gaussian SRD problem is realizable by a two-stage mechanism with a linear-Gaussian memoryless sensor and the Kalman filter. Although this fact is implicitly shown in [12], [13], for completeness, we reproduce a proof in this paper based on a technique used in [12], [13].

The sensor-estimator separation principle gives us a structural understanding of the Gaussian SRD problem. In particular, based on this principle, we show that the Gaussian SRD problem can be formulated as a semidefinite programming problem (Theorem 1), which is the main contribution of this paper. We derive a computationally accessible form (namely a semidefinite representation⁴ [20]) of the SRD function, and provide an efficient algorithm to solve Gaussian SRD problems numerically.

The semidefinite representation of the SRD function may be compared with an alternative analytical approach via Duncan’s theorem, which states that “twice the mutual information is merely the integration of the trace of the optimal mean square filtering error” [21]. Duncan’s result was significantly generalized as the “I-MMSE” relationships in non-causal [22] and causal [23] estimation problems. Our SDP-based approaches are applicable to the cases with multi-dimensional and time-varying Gauss-Markov sources to which the existing I-MMSE formulas cannot be applied straightforwardly. Although we focus on the Gaussian SRD problems in this paper, we note that the standard RD and SRD problems for general sources and distortion measures in abstract (Polish) spaces are discussed in [24] and [17], respectively.

This paper is organized as follows. In Section II, we formally introduce the Gaussian SRD problem, which is the main problem considered in this paper. In Section III, we show that the Gaussian SRD problem is equivalent to what we call the linear-Gaussian sensor design problem, which formally establishes the sensor-estimator separation principle. Then, in Section IV, we show that the linear-Gaussian sensor design problem can be reduced to an SDP problem, which thus provides us an SDP-based solution synthesis procedure for Gaussian SRD problems. Extensions to stationary and infinite horizon problems are given in Section V. In Section VI, we consider applications of SRD theory to real-time communication systems and networked control systems. Simple simulation results will be presented in Section VII. We conclude in Section VIII.

Notation: Let \mathcal{X} be an Euclidean space, and $\mathcal{B}_{\mathcal{X}}$ be the Borel σ -algebra on \mathcal{X} . Let $(\Omega, \mathcal{F}, \mathcal{P})$ be a probability space, and $\mathbf{x} : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ be a random variable. Throughout the paper, we use lower case boldface symbols such as \mathbf{x} to denote random variables, while $x \in \mathcal{X}$ is a realization of \mathbf{x} . We denote by $q_{\mathbf{x}}$ the probability measure of \mathbf{x} defined by $q_{\mathbf{x}}(A) = \mathcal{P}(\{\omega : \mathbf{x}(\omega) \in A\})$ for every $A \in \mathcal{B}_{\mathcal{X}}$. When no confusion occurs, this measure will be also denoted by $q_{\mathbf{x}}(x)$ or $q(x)$. For a Borel measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$,

we write $\mathbb{E}f(\mathbf{x}) \triangleq \int f(x)q_{\mathbf{x}}(dx)$. For a random vector, we write $\mathbf{x}^t \triangleq (\mathbf{x}_0, \dots, \mathbf{x}_t)$ or $\mathbf{x}^t \triangleq (\mathbf{x}_1, \dots, \mathbf{x}_t)$ depending on the initial index, and $\mathbf{x}_s^t \triangleq (\mathbf{x}_s, \dots, \mathbf{x}_t)$. Let Θ be a real symmetric matrix of size $n \times n$. Notations $\Theta \succ 0$ or $\Theta \in \mathbb{S}_{++}^n$ (resp. $\Theta \succeq 0$ or $\Theta \in \mathbb{S}_+^n$) mean that Θ is a positive definite (resp. positive semidefinite) matrix. For a positive semidefinite matrix Θ , we write $\|x\|_{\Theta} \triangleq \sqrt{x^{\top}\Theta x}$.

II. PROBLEM FORMULATION

We begin our discussion with an estimation-theoretic interpretation of a simple rate-distortion trade-off problem. Recall that a rate-distortion problem for a scalar Gaussian random variable $\mathbf{x} \sim \mathcal{N}(0, 1)$ with the mean square distortion constraint is an optimization problem of the following form:

$$\begin{aligned} \min \quad & I(\mathbf{x}; \mathbf{z}) \\ \text{s.t.} \quad & \mathbb{E}(\mathbf{x} - \mathbf{z})^2 \leq D. \end{aligned} \quad (1)$$

Here, \mathbf{z} is a reproduction of the source \mathbf{x} , and $I(\mathbf{x}; \mathbf{z})$ denotes the mutual information between \mathbf{x} and \mathbf{z} . The minimization is over the space of reproduction policies, i.e., stochastic kernels $q(dz|x)$. The optimal value of (1) is known as the rate-distortion function, $R(D)$, and can be explicitly obtained [1] as

$$R(D) = \max \left\{ 0, \frac{1}{2} \log \left(\frac{1}{D} \right) \right\}.$$

It is also possible to write the optimal reproduction policy $q(dz|x)$ explicitly. To this end, consider a linear sensor

$$\mathbf{y} = c\mathbf{x} + \mathbf{v} \quad (2)$$

where $\mathbf{v} \sim \mathcal{N}(0, \sigma^2)$ is a Gaussian noise independent of \mathbf{x} . Also, let

$$\mathbf{z} = \mathbb{E}(\mathbf{x}|\mathbf{y}) \quad (3)$$

be the least mean square error estimator of \mathbf{x} given \mathbf{y} . Notice that the right hand side of (3) is given by $\frac{c}{c^2 + \sigma^2} \mathbf{y}$. Then, it can be shown that an optimal solution $q(dz|x)$ to (1) is a composition of (2) and (3), provided that the signal-to-noise ratio of the sensor (2) is chosen to be

$$\text{SNR} \triangleq \frac{c^2}{\sigma^2} = \max \left\{ 0, \frac{1}{D} - 1 \right\}. \quad (4)$$

This gives us the following notable observations:

- **Fact 1:** A “sensor-estimator separation principle” holds for the Gaussian rate-distortion problem (1), in the sense that an optimal reproduction policy $q(dz|x)$ can be written as a two-stage mechanism with a linear sensor mechanism (2) and a least mean square error estimator (3).
- **Fact 2:** The original infinite dimensional optimization problem (1) with respect to $q(dz|x)$ is reduced to a simple optimization problem in terms of a scalar parameter SNR. Moreover, for a given $D > 0$, the optimal choice of SNR is given by a closed-form expression (4).

These facts can be significantly generalized, and serve as a guideline to develop a solution synthesis for Gaussian SRD problems in this paper.

⁴To be precise, we show that the *exponentiated* SRD function for multidimensional Gauss-Markov source is semidefinite representable by (27).

A. Gaussian SRD problem

The Gaussian SRD problem can be viewed as a generalization of (1). Let $\{\mathbf{x}_t\}$ be an \mathbb{R}^{n_t} -valued Gauss-Markov process

$$\mathbf{x}_{t+1} = A_t \mathbf{x}_t + \mathbf{w}_t, \quad t = 0, 1, \dots, T-1 \quad (5)$$

where $\mathbf{x}_0 \sim \mathcal{N}(0, P_0)$, $P_0 \succ 0$ and $\mathbf{w}_t \sim \mathcal{N}(0, W_t)$, $W_t \succ 0$ for $t = 0, 1, \dots, T-1$ are mutually independent Gaussian random variables. The Gaussian SRD problem is formulated as

$$\text{(P-SRD): } \min_{\gamma \in \Gamma} I(\mathbf{x}^T \rightarrow \mathbf{z}^T) \quad (6a)$$

$$\text{s.t. } \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \leq D_t \quad (6b)$$

where (6b) is imposed for every $t = 1, \dots, T$. Here, $\{\mathbf{z}_t\}$ is an \mathbb{R}^{n_t} -valued reproduction of $\{\mathbf{x}_t\}$. The minimization (6a) is over the space Γ of zero-delay reproduction policies of \mathbf{z}_t given \mathbf{x}^t and \mathbf{z}^{t-1} , i.e., the sequences of causal stochastic kernels⁵ $\gamma = \otimes_{t=1}^T q(dz_t | x^t, z^{t-1})$. The term $I(\mathbf{x}^T \rightarrow \mathbf{z}^T)$ is known as *directed information*, introduced by Massey [18] following Marko's earlier work [25], and is defined by

$$I(\mathbf{x}^T \rightarrow \mathbf{z}^T) \triangleq \sum_{t=1}^T I(\mathbf{x}^t; \mathbf{z}_t | \mathbf{z}^{t-1}). \quad (7)$$

The Gaussian SRD problem is visualized in Fig. 1.

Remark 1: Directed information measures the amount of information flow from $\{\mathbf{x}_t\}$ to $\{\mathbf{z}_t\}$ and is not symmetric, i.e., $I(\mathbf{x}^T \rightarrow \mathbf{z}^T) \neq I(\mathbf{z}^T \rightarrow \mathbf{x}^T)$ in general. However, when the process $\{\mathbf{z}_t\}$ is causally dependent on $\{\mathbf{x}_t\}$ and $\{\mathbf{x}_t\}$ is not affected by $\{\mathbf{z}_t\}$, it can be shown [26] that $I(\mathbf{x}^T \rightarrow \mathbf{z}^T) = I(\mathbf{x}^T; \mathbf{z}^T)$. By definition of our source process (5), there is no information feedback from $\{\mathbf{z}_t\}$ to $\{\mathbf{x}_t\}$, and thus $I(\mathbf{x}^T \rightarrow \mathbf{z}^T) = I(\mathbf{x}^T; \mathbf{z}^T)$ holds in our setup. Hence, $I(\mathbf{x}^T; \mathbf{z}^T)$ can be equivalently used as an objective in (P-SRD). However, we choose to use $I(\mathbf{x}^T \rightarrow \mathbf{z}^T)$ for the future considerations (e.g., [27]) in which $\{\mathbf{x}_t\}$ is a controlled stochastic process and is dependent on $\{\mathbf{z}_t\}$. In such cases, $I(\mathbf{x}^T; \mathbf{z}^T)$ and $I(\mathbf{x}^T \rightarrow \mathbf{z}^T)$ are not equal, and the latter is a more meaningful quantity in many applications.

Since (P-SRD) is an infinite dimensional optimization problem, it is difficult to apply numerical methods directly. Hence, we first need to develop a structural understanding of its solution. It turns out that the sensor-estimator separation principle still holds for (P-SRD), and this observation plays an important role in the subsequent sections. We are going to establish the following facts:

- **Fact 1'**: A sensor-estimator separation principle holds for the Gaussian SRD problem. That is, an optimal policy $\otimes_{t=1}^T q(dz_t | x^t, z^{t-1})$ for (P-SRD) can be realized as a composition of a sensor mechanism

$$\mathbf{y}_t = C_t \mathbf{x}_t + \mathbf{v}_t, \quad t = 1, 2, \dots, T \quad (8)$$

where $\mathbf{v}_t \sim \mathcal{N}(0, V_t)$, $V_t \succ 0$ are mutually independent Gaussian random variables, and the least mean square error estimator (Kalman filter)

$$\mathbf{z}_t = \mathbb{E}(\mathbf{x}_t | \mathbf{y}^t), \quad t = 1, 2, \dots, T. \quad (9)$$

⁵See Appendix A for a formal description of causal stochastic kernels.

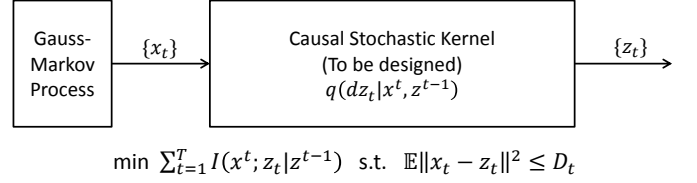


Fig. 1. The Gaussian sequential rate-distortion problem (P-SRD).

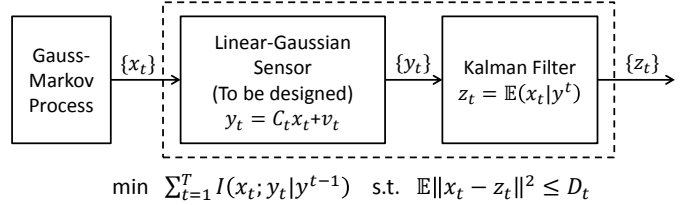


Fig. 2. The linear-Gaussian sensor design problem (P-LGS).

- **Fact 2'**: The original optimization problem (P-SRD) over an infinite-dimensional space Γ is reduced to an optimization problem over a finite-dimensional space of *matrix-valued* signal-to-noise ratios of the sensor (8), defined by

$$\text{SNR}_t \triangleq C_t^\top V_t^{-1} C_t \succeq 0, \quad t = 1, 2, \dots, T. \quad (10)$$

Moreover, the optimal $\{\text{SNR}_t\}_{t=1}^T$, which depends on $D_t > 0, t = 1, \dots, T$, can be obtained by SDP.

Unlike (4), an analytical expression of the optimal $\{\text{SNR}_t\}_{t=1}^T$ may not be available. Nevertheless, we will show that they can be easily obtained by SDP.

B. Linear-Gaussian sensor design problem

In Section III, we establish the sensor-estimator separation principle. To this end, we show that (P-SRD) is equivalent to what we call the *linear-Gaussian sensor design problem* (P-LGS) visualized in Fig. 2. Formally, (P-LGS) is formulated as

$$\text{(P-LGS): } \min_{\gamma \in \Gamma_{\text{LGS}}} \sum_{t=1}^T I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) \quad (11a)$$

$$\text{s.t. } \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \leq D_t \quad (11b)$$

where (11b) is imposed for every $t = 1, \dots, T$. We assume that $\{\mathbf{y}_t\}$ is produced by a linear-Gaussian sensor (8), and $\{\mathbf{z}_t\}$ is produced by the Kalman filter (9). In other words, the optimization domain $\Gamma_{\text{LGS}} \subset \Gamma$ is the space of causal stochastic kernels with a separation structure (8) and (9), which is parameterized by a sequence of matrices $\{C_t, V_t\}_{t=1}^T$. Intuitively, $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1})$ in (11a) can be understood as the amount of information acquired by the sensor (8) at time t . We call this problem a “sensor design problem” because our focus is on choosing an optimal sensing gain C_t in (8) and the noise covariance V_t . Notice that perfect observation with $C_t = I$ and $V_t = 0$ is trivially the best to minimize the estimation error in (11b) (in fact, $\mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 = 0$ is achieved), but it incurs significant information cost (i.e., $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) = +\infty$), and hence it is not an optimal solution to (P-LGS).

Remark 2: In (P-LGS), we search for the optimal $C_t \in \mathbb{R}^{r_t \times n}$ and $V_t \in \mathbb{S}_{++}^{r_t}$. However, the sensor dimension r_t is not given *a priori*, and choosing it optimally is part of the problem. In particular, if making no observation is the optimal sensing at some specific time instance t , we should be able to recover $r_t = 0$ as an optimal solution.

Although the objective functions (6a) and (11a) appear differently, it will be shown in Section III that they coincide in the domain Γ_{LGS} . Moreover, in the same section it will be shown that an optimal solution to (P-SRD) can always be found in the domain Γ_{LGS} . These observations imply that one can obtain an optimal solution to (P-SRD) by solving (P-LGS).

C. Stationary cases

We will also consider a time-invariant system

$$\mathbf{x}_{t+1} = A\mathbf{x}_t + \mathbf{w}_t, t = 0, 1, 2, \dots \quad (12)$$

where \mathbf{x}_t is an \mathbb{R}^n -valued random variable with $\mathbf{x}_0 \sim \mathcal{N}(0, P_0)$, and $\mathbf{w}_t \sim \mathcal{N}(0, W)$ is a stationary white Gaussian noise. We assume $P_0 \succ 0$ and $W \succ 0$. Stationary and infinite horizon version of the Gaussian SRD problem is formulated as

$$\min \limsup_{T \rightarrow \infty} \frac{1}{T} I(\mathbf{x}^T \rightarrow \mathbf{z}^T) \quad (13a)$$

$$\text{s.t.} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta}^2 \leq D. \quad (13b)$$

This is an optimization over the sequence of stochastic kernels $\otimes_{t \in \mathbb{N}} q(dz_t | x_t, z^{t-1})$. The optimal value of (13) as a function of the average distortion D is referred to as the *sequential rate-distortion function*, and is denoted by $R_{\text{SRD}}(D)$.

Similarly, a stationary and infinite horizon version of the linear-Gaussian sensor design problem is formulated as

$$\min \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) \quad (14a)$$

$$\text{s.t.} \quad \limsup_{t \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta}^2 \leq D. \quad (14b)$$

Here, we assume $\mathbf{y}_t = C_t \mathbf{x}_t + \mathbf{v}_t$ where $\mathbf{v}_t \sim \mathcal{N}(0, V_t)$, $V_t \succ 0$ is a mutually independent Gaussian stochastic process and $\mathbf{z}_t = \mathbb{E}(\mathbf{x}_t | \mathbf{y}^t)$. Design variables in (14) are $\{C_t, V_t\}_{t \in \mathbb{N}}$. Again, determining their dimensions is part of the problem.

D. Soft- vs. hard-constrained problems

Introducing Lagrange multipliers $\alpha_t > 0$, one can also consider a soft-constrained version of (P-SRD):

$$\min I(\mathbf{x}^T \rightarrow \mathbf{z}^T) + \frac{\alpha_t}{2} \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \quad (15)$$

Similarly to the Lagrange multiplier theorem (e.g., Proposition 3.1.1 in [28]), it is possible to show that there exists a set of multipliers such that an optimal solution to (15) is also an optimal solution to (P-SRD). We will prove this fact in Section IV after we establish that both (P-SRD) and (15) can be transformed as finite dimensional convex optimization problems. For this reason, we refer to both (P-SRD) and (15) as Gaussian SRD problems.

III. SENSOR-ESTIMATOR SEPARATION PRINCIPLE

Let f_{SRD}^* and f_{LGS}^* be the optimal values of (P-SRD) and (P-LGS) respectively. In this section, we show that $f_{\text{SRD}}^* = f_{\text{LGS}}^*$, and an optimal solution $\gamma \in \Gamma_{\text{LGS}}$ to (P-LGS) is also an optimal solution to (P-SRD). This result establishes the sensor-estimator separation principle (Fact 1'). We introduce another optimization problem (P-1), which serves as an intermediate step to establish this fact.

$$\begin{aligned} \text{(P-1):} \quad & \min_{\gamma \in \Gamma_1} \sum_{t=1}^T I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) \\ & \text{s.t.} \quad \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \leq D_t. \end{aligned}$$

The optimization is over the space Γ_1 of linear-Gaussian stochastic kernels $\gamma = \otimes_{t=1}^T q(dz_t | x_t, z^{t-1})$, where each stochastic kernel $q(dz_t | x_t, z^{t-1})$ is of the form

$$\mathbf{z}_t = E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}_{t-1} + \dots + F_{t,1} \mathbf{z}_1 + \mathbf{g}_t \quad (16)$$

where $E_t, F_{t,t-1}, \dots, F_{t,1}$ are some matrices with appropriate dimensions, and \mathbf{g}_t is a zero-mean, possibly degenerate Gaussian random variable that is independent of $\mathbf{x}_0, \mathbf{w}^t, \mathbf{g}^{t-1}$. Notice that $\Gamma_{\text{LGS}} \subset \Gamma_1 \subset \Gamma$. The underlying Gauss-Markov process $\{\mathbf{x}_t\}$ is defined by (5). Let f_1^* be the optimal value of (P-1). The next lemma claims the equivalence between (P-SRD) and (P-1).

Lemma 1:

- (i) If there exists $\gamma \in \Gamma$ attaining a value $f_{\text{SRD}} < +\infty$ of the objective function in (P-SRD), then there exists $\gamma_1 \in \Gamma_1$ attaining a value $f_1 \leq f_{\text{SRD}}$ of the objective function in (P-1).
- (ii) Every $\gamma_1 \in \Gamma_1 (\subset \Gamma)$ attaining $f_1 < +\infty$ in (P-1) also attains $f_{\text{SRD}} = f_1$ in (P-SRD).

Lemma 1 is the most significant result in this section, which essentially guarantees the linearity of an optimal solution to the Gaussian SRD problems. The proof of Lemma 1 can be found in Appendix B. The basic idea of proof relies on the well-known fact that Gaussian distribution maximizes entropy when covariance is fixed. This proposition appears as Lemma 4.3 in [13], but we modified the proof using the Radon-Nikodym derivatives so that the proof does not require the existence of probability density functions. The next lemma establishes the equivalence between (P-1) and (P-LGS).

Lemma 2:

- (i) If there exists $\gamma_1 \in \Gamma_1$ attaining a value $f_1 < +\infty$ of the objective function in (P-1), then there exists $\gamma_{\text{LGS}} \in \Gamma_{\text{LGS}}$ attaining a value $f_{\text{LGS}} \leq f_1$ of the objective function in (P-LGS).
- (ii) Every $\gamma_{\text{LGS}} \in \Gamma_{\text{LGS}} (\subset \Gamma_1)$ attaining $f_{\text{LGS}} < +\infty$ in (P-LGS) also attains $f_1 \leq f_{\text{LGS}}$ in (P-1).

Proof of Lemma 2 is in Appendix C. Combining the above two lemmas, we obtain the following consequence, which is the main proposition in this section. It guarantees that we can alternatively solve (P-LGS) in order to solve (P-SRD).

Proposition 1: Suppose $f_{\text{SRD}}^* < +\infty$. Then there exists an optimal solution $\gamma_{\text{LGS}} \in \Gamma_{\text{LGS}} (\subset \Gamma)$ to (P-LGS). Moreover, an optimal solution to (P-LGS) is also an optimal solution to (P-SRD), and $f_{\text{SRD}}^* = f_{\text{LGS}}^*$.

IV. SDP-BASED SYNTHESIS

In this section, we develop an efficient numerical algorithm to solve (P-LGS). Due to the preceding discussion, this is equivalent to developing an algorithm to solve (P-SRD). Let (5) be given. Assume temporarily that (8) is also fixed. The Kalman filtering formula for computing $\mathbf{z}_t = \mathbb{E}(\mathbf{x}_t|\mathbf{y}^t)$ is

$$\begin{aligned}\mathbf{z}_t &= \mathbf{z}_{t|t-1} + P_{t|t-1}C_t^\top(C_tP_{t|t-1}C_t^\top + V_t)^{-1}(\mathbf{y}_t - C_t\mathbf{z}_{t|t-1}) \\ \mathbf{z}_{t|t-1} &= A_{t-1}\mathbf{z}_{t-1}\end{aligned}$$

where $P_{t|t-1}$ is the covariance matrix of $\mathbf{x}_t - \mathbb{E}(\mathbf{x}_t|\mathbf{y}^{t-1})$, which can be recursively computed as

$$P_{t|t-1} = A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1} \quad (17a)$$

$$P_{t|t} = (P_{t|t-1}^{-1} + \text{SNR}_t)^{-1} \quad (17b)$$

for $t = 1, \dots, T$ with $P_{0|0} = P_0$. The variable SNR_t is defined by (10). Using these quantities, mutual information terms in (11a) can be explicitly written as

$$\begin{aligned}I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) &= h(\mathbf{x}_t | \mathbf{y}^{t-1}) - h(\mathbf{x}_t | \mathbf{y}^t) \\ &= \frac{1}{2} \log \det(A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1}) - \frac{1}{2} \log \det P_{t|t}.\end{aligned}$$

Note that $W_t \succ 0$ and $V_t \succ 0$ guarantee that both differential entropy terms are finite. Hence, (P-LGS) is equivalent to the following optimization problem in terms of the variables $\{\text{SNR}_t, P_{t|t}\}_{t=1}^T$:

$$\min \sum_{t=1}^T \frac{1}{2} \log \det(A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1}) - \frac{1}{2} \log \det P_{t|t} \quad (18a)$$

$$\text{s.t. } P_{t|t}^{-1} = (A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1})^{-1} + \text{SNR}_t \quad (18b)$$

$$\text{SNR}_t \geq 0, \text{Tr}(\Theta_t P_{t|t}) \leq D_t. \quad (18c)$$

Equality (18b) is obtained by eliminating $P_{t|t-1}$ from (17). At this point, one may note that (18) can be viewed as an optimal control problem with state $P_{t|t}$ and control input SNR_t . Naturally, dynamic programming approach has been proposed in the literature in similar contexts [10]–[12], [15]. Alternatively, we next propose a method to transform (18) into an SDP problem. This allows us to solve (P-SRD) using standard SDP solvers, which is now a mature technology.

A. SRD optimization as max-det problem

Now we show that (18) can be converted to a determinant maximization problem [29] subject to linear matrix inequality constraints. The first step is to transform (18) into an optimization problem in terms of $\{P_{t|t}\}_{t=1}^T$ only. This is possible by simply replacing the nonlinear equality constraint (18b) with a linear inequality constraint

$$0 \prec P_{t|t} \preceq A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1}.$$

This replacement eliminates SNR_t from (18) giving us:

$$\min \sum_{t=1}^T \frac{1}{2} \log \det(A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1}) - \frac{1}{2} \log \det P_{t|t} \quad (19a)$$

$$\text{s.t. } 0 \prec P_{t|t} \preceq A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1} \quad (19b)$$

$$\text{Tr}(\Theta_t P_{t|t}) \leq D_t. \quad (19c)$$

Note that (18) and (19) are mathematically equivalent, since eliminated SNR variables can be easily constructed from $\{P_{t|t}\}_{t=1}^T$ through

$$\text{SNR}_t = P_{t|t}^{-1} - (A_{t-1}P_{t-1|t-1}A_{t-1}^\top + W_{t-1})^{-1}. \quad (20)$$

The second step is to rewrite the objective function (19a). Regrouping terms, (19a) can be written as a summation of the initial cost $\frac{1}{2} \log \det(A_0P_{0|0}A_0^\top + W_0)$, the final cost $-\frac{1}{2} \log \det P_{T|T}$, and stage-wise costs

$$\frac{1}{2} \log \det(A_tP_{t|t}A_t^\top + W_t) - \frac{1}{2} \log \det P_{t|t} \quad (21)$$

for $t = 1, \dots, T-1$. Applying the matrix determinant lemma (e.g., Theorem 18.1.1 in [30]), (21) can be rewritten as

$$\frac{1}{2} \log \det W_t - \frac{1}{2} \log \det(P_{t|t}^{-1} + A_t^\top W_t^{-1} A_t)^{-1}. \quad (22)$$

Due to the monotonicity of the determinant function, (22) is equal to the optimal value of

$$\min \frac{1}{2} \log \det W_t - \frac{1}{2} \log \det \Pi_t \quad (23a)$$

$$\text{s.t. } 0 \prec \Pi_t \preceq (P_{t|t}^{-1} + A_t^\top W_t^{-1} A_t)^{-1}. \quad (23b)$$

Applying the matrix inversion lemma, (23b) is equivalent to $0 \prec \Pi_t \preceq P_{t|t} - P_{t|t}A_t^\top(W_t + A_tP_{t|t}A_t^\top)^{-1}A_tP_{t|t}$, which is further equivalent to

$$\begin{bmatrix} P_{t|t} - \Pi_t & P_{t|t}A_t^\top \\ A_tP_{t|t} & W_t + A_tP_{t|t}A_t^\top \end{bmatrix} \succeq 0, \quad \Pi_t \succ 0. \quad (24)$$

Note that (24) is a linear matrix inequality (LMI) condition. The above discussion leads to the following conclusion.

Theorem 1: An optimal solution to (P-LGS) can be constructed by solving the following determinant maximization problem with decision variables $\{P_{t|t}, \Pi_t\}_{t=1}^T$:

$$\min - \sum_{t=1}^T \frac{1}{2} \log \det \Pi_t + c \quad (25a)$$

$$\text{s.t. } \Pi_t \succ 0, \quad t = 1, \dots, T \quad (25b)$$

$$P_{t+1|t+1} \preceq A_tP_{t|t}A_t^\top + W_t, \quad t = 0, \dots, T-1 \quad (25c)$$

$$\begin{bmatrix} P_{t|t} - \Pi_t & P_{t|t}A_t^\top \\ A_tP_{t|t} & W_t + A_tP_{t|t}A_t^\top \end{bmatrix} \succeq 0, \quad t = 1, \dots, T-1 \quad (25d)$$

$$\text{Tr}(\Theta_t P_{t|t}) \leq D_t, \quad t = 1, \dots, T \quad (25e)$$

$$P_{T|T} = \Pi_T, \quad (25f)$$

where $c = \frac{1}{2} \log \det(A_0P_{0|0}A_0^\top + W_0) + \sum_{t=1}^{T-1} \frac{1}{2} \log \det W_t$ is a constant. The optimal sequence $\{\text{SNR}_t\}_{t=1}^T$ can be reconstructed from (20), from which $\{C_t, V_t\}_{t=1}^T$ satisfying (10) can be reconstructed via the singular value decomposition. An optimal solution to (P-LGS) is obtained as a composition of (8) and (9).

Remark 3: Under the assumption that $W_t \succ 0, D_t > 0$ for every $t = 1, \dots, T$, the max-det problem (25) is always strictly feasible and there exists an optimal solution.⁶ Invoking Proposition 1, we have thus shown by construction that there always exists an optimal solution to (P-SRD) under this assumption.

Remark 4: As we mentioned in Remark 2, choosing an appropriate dimension r_t of the sensor output (8) is part of (P-LGS). It can be easily seen from Theorem 1 that the minimum sensor dimension to achieve the optimality in (P-LGS) is given by $r_t = \text{rank}(\text{SNR}_t)$.

Using the same technique, the soft-constrained version of the problem (15) can be formulated as:

$$\min \sum_{t=1}^T \left(\frac{\alpha_t}{2} \text{Tr}(\Theta_t P_{t|t}) - \frac{1}{2} \log \det \Pi_t \right) + c \quad (26a)$$

$$\text{s.t. } \Pi_t \succ 0, \quad t = 1, \dots, T \quad (26b)$$

$$P_{t+1|t+1} \preceq A_t P_{t|t} A_t^\top + W_t, \quad t = 0, \dots, T-1 \quad (26c)$$

$$\begin{bmatrix} P_{t|t} - \Pi_t & P_{t|t} A_t^\top \\ A_t P_{t|t} & W_t + A_t P_{t|t} A_t^\top \end{bmatrix} \succeq 0, \quad t = 1, \dots, T-1 \quad (26d)$$

$$P_{T|T} = \Pi_T \quad (26e)$$

The next proposition claims that (25) and (26) admit the same optimal solution provided Lagrange multipliers $\alpha_t, t = 1, \dots, T$, are chosen correctly. This further implies that, with the same choice of α_t , two versions of the Gaussian SRD problems (P-SRD) and (15) are equivalent.

Proposition 2: Suppose $W_t \succ 0, D_t > 0$ for $t = 1, \dots, T$. Then, there exist $\alpha_t, t = 1, \dots, T$ such that an optimal solution to (25) is also an optimal solution to (26).

Proof: Both (26) and (25) are strictly feasible. The result follows using the fact that the Slater's constraint qualification is satisfied for this problem, which guarantees that strong duality holds and the dual optimum is attained [31]. ■

B. Max-det problem as SDP

Strictly speaking, optimization problems (25) and (26) are in the class of determinant maximization problems [29], but not in the standard form of the SDP.⁷ However, they can be considered as SDPs in a broader sense for the following reasons. First, the hard constrained version (25) can be indeed transformed into a standard SDP problem. This conversion is possible by following the discussion in Chapter 4 of [32]. Second, sophisticated and efficient algorithms based on the interior-point method for SDP can almost directly be applied to max-det problems as well. In fact, off-the-shelf SDP solvers such as SDPT3 [33] have built-in functions to handle log-determinant terms directly.

Recall that (P-LGS) and (P-SRD) have a common optimal solution. Hence, Proposition 1 shows that both (P-LGS) and

⁶To see the strict feasibility, consider $P_{t|t} = \delta I$ for $t = 1, \dots, T-1$ and $\Pi_t = \delta^2 I$ for $t = 1, \dots, T$ with sufficiently small $\delta > 0$. The constraint set defined by (25b)-(25f) can be made compact by replacing (25b) with $\Pi_t \succeq \epsilon I$ without altering the result. Thus the existence of an optimal solution is guaranteed by the Weierstrass theorem.

⁷In the standard form, SDP is an optimization problem of the form $\min \langle C, X \rangle$ s.t. $\mathcal{A}(X) = B, X \succeq 0$.

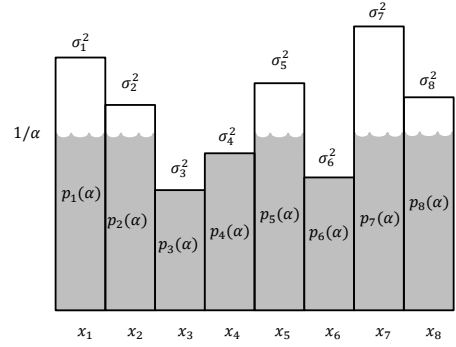


Fig. 3. Reverse water-filling solution to the Gaussian rate-distortion problem.

(P-SRD) are essentially solvable via SDP, which is much stronger than merely saying that they are convex problems. Note that convexity alone does not guarantee the existence of an efficient optimization algorithm.

C. Complexity analysis

In this section, we briefly consider the arithmetic complexity (i.e., the worst case number of arithmetic operations needed to obtain an ϵ -optimal solution) of problem (25), and how it grows as the horizon length T grows when the dimensions of the Gauss-Markov process (5) are fixed to $n_t = n \forall t = 1, \dots, T$. For a preliminary analysis, it would be natural for us to resort to the existing interior-point method literature (e.g., [32], [34]). Interior-point methods for the determinant maximization problem are already considered in [29], [35], [36]. The most computationally expensive step in the interior-point method is the Cholesky factorization involved in the Newton steps, which requires $\mathcal{O}(T^3)$ operations in general. However, it is possible to exploit the sparsity of coefficient matrices in the SDPs to reduce operation counts [37]–[39]. By exploiting the structure of our SDP formulation (25), it is theoretically expected that there exists a specialized interior-point method algorithm for (25) whose arithmetic complexity is $\mathcal{O}(T \log(1/\epsilon))$. However, more careful study and computational experiments are needed to verify this conjecture.

D. Single stage problem

When $T = 1$, the result of Proposition 1 recovers the well-known “reverse water-filling” solution for the standard Gaussian rate-distortion problem [1]. To see this, notice that $T = 1$ reduces problem (26) to

$$\begin{aligned} \min \quad & \text{Tr} P - \frac{1}{\alpha} \log \det P \\ \text{s.t.} \quad & 0 \preceq P \preceq \text{diag}(\sigma_1^2, \dots, \sigma_n^2). \end{aligned}$$

Here, we have already assumed $\Theta = I$ and $AP_0A^\top + W = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \succeq 0$. This does not result in loss of generality, since otherwise a change of variables $P \leftarrow U\Theta^{\frac{1}{2}}P\Theta^{\frac{1}{2}}U^\top$, where U is an orthonormal matrix that makes $U\Theta^{\frac{1}{2}}(AP_0A^\top + W)\Theta^{\frac{1}{2}}U^\top$ diagonal, converts the problem into the above form. For any positive definite matrix P , Hadamard's inequality (e.g., [1]) states that $\det P \leq \prod_i P_{ii}$ and the equality holds

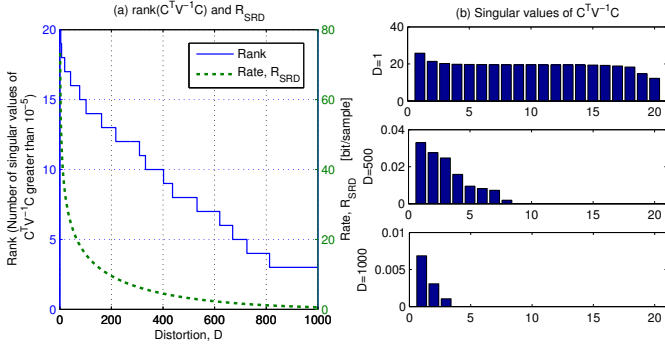


Fig. 4. Numerical experiments on rank monotonicity. 20 dimensional Gaussian process is randomly generated and $\text{SNR} = C^T V^{-1} C$ is constructed for various D . Observe $\text{rank}(\text{SNR})$ tends to decrease as D increase.

if and only if the matrix is diagonal. Hence, if diagonal elements of P are fixed, $\det P$ is maximized by setting all off-diagonal entries zero. Thus, the optimal solution to the above problem is diagonal. Writing $P = \text{diag}(p_1, \dots, p_n)$, the problem is decomposed as n independent optimization problems, each of which minimizes $p_i - \frac{1}{\alpha} \log p_i$ subject to $0 \leq p_i \leq \sigma_i^2$. It is easy to see that the optimal solution is $p_i = \min(1/\alpha, \sigma_i^2)$. This is the closed-form solution to (P-LGS) with $T = 1$, and its pictorial interpretation is shown in Fig. 3. This solution also indicates the optimal sensing formula is given by $\mathbf{y} = C\mathbf{x} + \mathbf{v}$, $\mathbf{v} \sim \mathcal{N}(0, V)$, where C and V satisfy

$$\begin{aligned} C^T V^{-1} C &= P^{-1} - (AP_0 A^T + W)^{-1} \\ &= \text{diag}_{1 \leq i \leq n} \left(\max \left\{ 0, \alpha - \frac{1}{\sigma_i^2} \right\} \right). \end{aligned}$$

In particular, we have $\dim(\mathbf{y}) = \text{rank}(C^T V^{-1} C) = \text{card}\{i : \sigma_i^2 > \frac{1}{\alpha}\}$, indicating that the optimal dimension of \mathbf{y} monotonically decreases as the ‘‘price of information’’ $1/\alpha$ increases.

V. STATIONARY PROBLEMS

A. Sequential rate-distortion function

We are often interested in infinite-horizon Gaussian SRD problems (13). Assuming that (A, Θ) is a detectable pair, it can be shown that (13) is equivalent to the infinite-horizon linear-Gaussian sensor design problem (14) [40]. Moreover, [40] shows that (13) and (14) admit an optimal solution that can be realized as a composition of a *time-invariant* sensor mechanism $\mathbf{y}_t = C\mathbf{x}_t + \mathbf{v}_t$ with i.i.d. process $\mathbf{v}_t \sim \mathcal{N}(0, V)$ and a *time-invariant* Kalman filter. Hence, it is enough to minimize the average cost per stage, which leads to the following simpler problem.

$$R_{\text{SRD}}(D) = \min -\frac{1}{2} \log \det \Pi + \frac{1}{2} \log \det W \quad (27a)$$

$$\text{s.t. } \Pi \succ 0 \quad (27b)$$

$$P \preceq APA^T + W \quad (27c)$$

$$\text{Tr}(\Theta P) \leq D \quad (27d)$$

$$\begin{bmatrix} P - \Pi & PA^T \\ AP & APA^T + W \end{bmatrix} \succeq 0. \quad (27e)$$

To confirm (27) is compatible with the existing result, consider a scalar system with $A = a, W = w, P = p$ and $\Theta = 1$. In

this case, a closed-form expression of the SRD function is known in the literature [12] [16], which is given by

$$R_{\text{SRD}}(D) = \min \left\{ 0, \frac{1}{2} \log \left(a^2 + \frac{w}{D} \right) \right\}. \quad (28)$$

For a scalar system, (27) further simplifies to

$$\min \log \left(a^2 + \frac{w}{p} \right) \quad (29a)$$

$$\text{s.t. } 0 < p \leq a^2 p + w, \quad p \leq D. \quad (29b)$$

It is elementary to verify that the optimal value of (29) is $\log(a^2 + \frac{w}{D})$ if $1 - \frac{w}{D} \leq a^2$, while it is 0 if $0 \leq a^2 < 1 - \frac{w}{D}$. Hence, it can be compactly written as $\min\{0, \frac{1}{2} \log(a^2 + \frac{w}{D})\}$, and the result recovers (28). Alternative representations of the SRD function (27) for stationary multi-dimensional Gauss-Markov processes when $\Theta = I$ are reported in [13, Section IV-B] and [17, Section VI].

B. Rank monotonicity

Using an optimal solution to (27) the optimal sensing matrices C and V are recovered from $C^T V^{-1} C = P^{-1} - (APA^T + W)^{-1}$. In particular, $\dim(\mathbf{y}) = \text{rank}(C^T V^{-1} C)$ determines the optimal dimension of the measurement vector. Similarly to the case of single stage problems, this rank has a tendency to decrease as D increases. A typical numerical behavior is shown in Figure 4. We do not attempt to prove the rank monotonicity here.

VI. APPLICATIONS AND RELATED WORKS

A. Zero-delay source coding

SRD theory plays an important role in the rate analysis of zero-delay source coding schemes. For each $t = 1, 2, \dots$, let

$$\mathcal{B}_t \subset \{0, 1, 00, 01, 10, 11, 000, \dots\}$$

be a set of variable-length uniquely decodable codewords. Assume that $\mathbf{b}_t \in \mathcal{B}_t$ for $t = 1, 2, \dots$, and let l_t be the length of \mathbf{b}_t . A zero-delay binary coder is a pair of a sequence of encoders $e_t(db_t|x^t, b^{t-1})$, i.e., stochastic kernels on \mathcal{B}_t given $\mathcal{X}^t \times \mathcal{B}^{t-1}$, and a sequence of decoders $d_t(dz_t|b^t, z^{t-1})$, i.e., stochastic kernels on \mathcal{Z}_t given $\mathcal{B}^t \times \mathcal{Z}^{t-1}$. The *zero-delay rate-distortion region* for the Gauss-Markov process (12) is the epigraph of the function

$$\begin{aligned} R_{\text{SRD}}^{\text{op}}(D) &= \inf_{\{\mathcal{B}_t, e_t, d_t\}_{t=1}^{\infty}} \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}(l_t) \\ &\text{s.t. } \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta}^2 \leq D. \end{aligned}$$

The SRD function is a lower bound of the achievable rate. Indeed, $R_{\text{SRD}}(D) \leq R_{\text{SRD}}^{\text{op}}(D) \quad \forall D > 0$ can be shown

straightforwardly as

$$I(\mathbf{x}^T \rightarrow \mathbf{z}^T) = I(\mathbf{x}^T; \mathbf{z}^T) \quad (30a)$$

$$\leq I(\mathbf{x}^T; \mathbf{b}^T) \quad (30b)$$

$$= H(\mathbf{b}^T) - H(\mathbf{b}^T | \mathbf{x}^T) \quad (30c)$$

$$\leq H(\mathbf{b}^T) \quad (30d)$$

$$\leq \sum_{t=1}^T H(\mathbf{b}_t) \quad (30e)$$

$$\leq \sum_{t=1}^T \mathbb{E}(l_t) \quad (30f)$$

where (30a) holds since there is no feedback from the process $\{\mathbf{z}_t\}$ to $\{\mathbf{x}_t\}$ (Remark 1), (30b) follows from the data processing inequality, (30d) holds since conditional entropy is non-negative, and (30e) is due to the chain rule for entropy. The final inequality (30f) holds since the expected length of a uniquely decodable code is lower bounded by its entropy [1, Theorem 5.3.1].

In general, $R_{\text{SRD}}(D)$ and $R_{\text{SRD}}^{\text{op}}(D)$ do not coincide. Nevertheless, by constructing an appropriate entropy-coded dithered quantizer (ECDQ), it is shown in [16] that $R_{\text{SRD}}^{\text{op}}(D)$ does not exceed $R_{\text{SRD}}(D)$ more than a constant due to the “space-filling loss” of the lattice quantizer and the loss of entropy coding.

B. Networked control theory

Zero-delay source/channel coding technologies are crucial in networked control systems [41]–[44]. Gaussian SRD theory plays an important role in the LQG control problems with information theoretic constraints [13]. It is shown in [19] that an LQG control problem in which observed data must be transmitted to the controller over a noiseless binary channel is closely related to the LQG control problem with directed information constraints. The latter problem is addressed in [27] using the SDP-based algorithm presented in this paper. In [27], the problem is viewed as a sensor-controller joint design problem in which directed information from the state process to the control input is minimized.⁸

C. Experimental design/Sensor scheduling

In this subsection, we compare the linear-Gaussian sensor design problem (P-LGS) with different types of sensor design/selection problems considered in the literature.

A problem of selecting the best subset of sensors to observe a random variable in order to minimize the estimation error and its convex relaxations are considered in [29]. A sensor selection problem for a linear dynamical system is considered in [47], where submodularity of the objective function is exploited. Dynamic sensor scheduling problems are also considered in the literature. In [48], an efficient algorithm to explore branches of the scheduling tree is proposed. In [49], a stochastic sensor selection strategy that minimizes the expected error covariance is considered.

The linear-Gaussian sensor design problem (P-LGS) is different from these sensor selection/scheduling problems in that it is essentially a continuous optimization problem (since

⁸The problem considered in [27] is different from the sensor-controller joint design problems considered in [45] and [46].

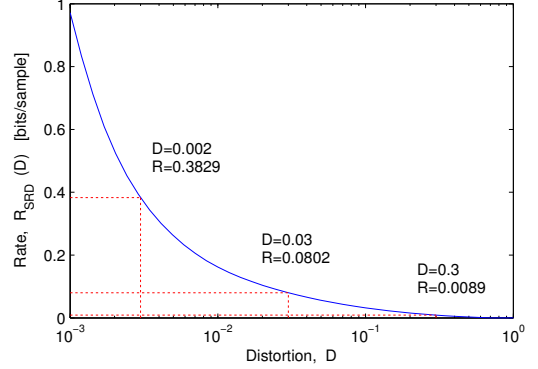


Fig. 5. Sequential rate-distortion function for the noisy double pendulum.

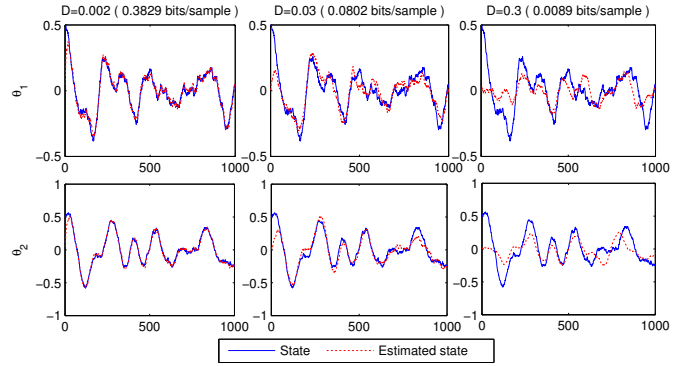


Fig. 6. Tracking performance of the Kalman filter under different distortion constraints. (Tested on the same sample path of the noisy double pendulum.)

matrices $\{C_t, V_t\}_{t=1}^T$ can be freely chosen), and the objective is to minimize an information-theoretic cost (11a).

VII. NUMERICAL SIMULATIONS

In this section, we consider two numerical examples to demonstrate how the SDP-based formulation of the Gaussian SRD problem can be used to calculate the minimal communication bandwidth required for the real-time estimation with desired accuracy.

A. Optimal sensor design for double pendulum

A linearized equation of motion of a double pendulum with friction and disturbance is given by

$$\begin{bmatrix} d\theta_1 \\ d\theta_2 \\ d\omega_1 \\ d\omega_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -\frac{(m_1+m_2)g}{m_1 l_1} & \frac{m_2 g}{m_1 l_1} & -c_1 & 0 \\ \frac{(m_1+m_2)g}{m_1 l_2} & -\frac{(m_1+m_2)g}{m_1 l_2} & 0 & -c_2 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \omega_1 \\ \omega_2 \end{bmatrix} dt + d\mathbf{b}$$

where \mathbf{b} is a Brownian motion. We consider a discrete time model of the above equation of motion obtained through the Tustin transformation. We are interested in designing a sensing model $\mathbf{y}_t = C\mathbf{x}_t + \mathbf{v}_t$, $\mathbf{v}_t \sim \mathcal{N}(0, V)$ that optimally trades-off information cost and distortion level.⁹ We solve the stationary

⁹In practice, it is often the case that \mathbf{x}_t is partially observable through a given sensor mechanism. In such cases, the framework discussed in this paper is not appropriate. Instead, one can formulate an SRD problem for *partially observable* Gauss-Markov processes. See [50] for details.

optimization problem (27) for this example with various values of D . The result is the sequential rate-distortion function shown in Figure 5. Finally, for every point on the trade-off curve, the optimal sensing matrices C and V are reconstructed, and the Kalman filter is designed based on them. Figure 6 shows the trade-off between the distortion level and the tracking performance of the Kalman filter. When the distortion constraint is strict ($D = 0.002$), the optimally designed sensor generates high rate information (0.3829 bits/sample) and the Kalman filter built on it tracks true state very well. When D is large ($D = 0.3$), the optimal sensing strategy chooses “not to observe much”, and the resulting Kalman filter shows poor tracking performance.

B. Minimum down-link bandwidth for satellite attitude determination

The equation of motion of the angular velocity vector of a spin-stabilized satellite linearized around the nominal angular velocity vector $(\omega_0, 0, 0)$ is

$$\begin{bmatrix} d\omega_1 \\ d\omega_2 \\ d\omega_3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & \frac{I_3 - I_1}{I_2} \omega_0 \\ 0 & \frac{I_1 - I_2}{I_3} \omega_0 & 1 \end{bmatrix} \begin{bmatrix} \omega_1 \\ \omega_2 \\ \omega_3 \end{bmatrix} dt + d\mathbf{b}$$

where \mathbf{b} is a disturbance. Again, the equation of motion is converted to a discrete time model in the simulation. Suppose that the satellite has on-board sensors that can accurately measure angular velocities, and the ground station needs to estimate them with some required accuracy (distortion) based on the transmitted data from the satellite. Our interest is to determine the minimum down-link bit-rate that makes it possible, and identify what information needs to be transmitted to achieve this. Assume that the distortion constraints D_t are time varying, but given *a priori*. (For instance, it must be kept small only when the satellite is in a mission.) The discussion so far indicates that the data to be transmitted is in the form of $\mathbf{y}_t = C_t \mathbf{x}_t + \mathbf{v}_t$ in order to minimize communication cost measured by $\sum_{t=1}^T I(\mathbf{x}_t, \mathbf{y}_t | \mathbf{y}^{t-1})$. In Figure 7, a result of the SDP (26) is plotted, when the scheduling horizon is $T = 120$ and a particular distortion constraint profile D_t is given (shown in red in (a)). The optimal down-link schedule shown in (b) requires no communication at all when the distortion constraint is met. As by-products of the SDP (26), the optimal scheduling of sensing matrices C_t and noise covariances V_t of \mathbf{v}_t can be also explicitly obtained.

VIII. CONCLUSION

In this paper, we revisited the “sensor-estimator separation principle” and showed that an optimal solution to the Gaussian SRD problem can be found by considering a related linear-Gaussian sensor design problem, which can be formulated as a determinant maximization problem with LMI constraints. The implication is that Gaussian SRD problems are efficiently solvable using standard SDP solvers. We have also considered several potential applications of the Gaussian SRD problem and its relationship to real-time communication theory, networked control theory, and sensor scheduling problems.

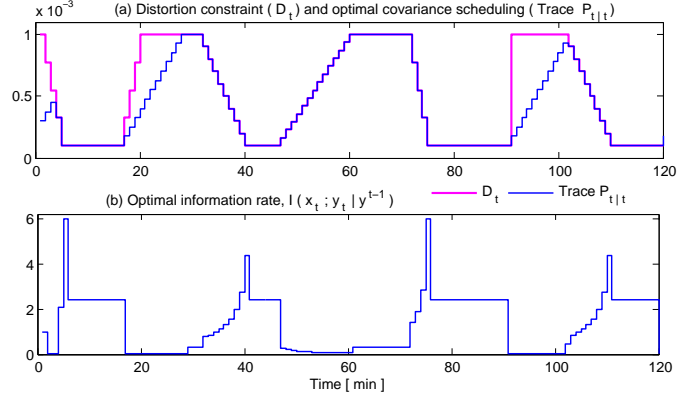


Fig. 7. Satellite attitude determination with time-varying distortion constraints.

APPENDIX A MATHEMATICAL PRELIMINARIES

A. Stochastic kernels

Let \mathcal{X}, \mathcal{Y} be Euclidean spaces. A (Borel-measurable) stochastic kernel on \mathcal{Y} given \mathcal{X} is a map $q_{\mathcal{Y}|\mathcal{X}} : \mathcal{B}_{\mathcal{Y}} \times \mathcal{X} \rightarrow [0, 1]$ such that $q_{\mathcal{Y}|\mathcal{X}}(\cdot|x)$ is a probability measure on $(\mathcal{Y}, \mathcal{B}_{\mathcal{Y}})$ for every $x \in \mathcal{X}$, and $q_{\mathcal{Y}|\mathcal{X}}(A|\cdot)$ is a Borel measurable function for every $A \in \mathcal{B}_{\mathcal{Y}}$. For simplicity, a stochastic kernel on \mathcal{Y} given \mathcal{X} will be denoted by $q(dy|x)$. The following results can be found in Propositions 7.27 and 7.28 in [51].

Lemma 3: Let \mathcal{X}, \mathcal{Y} be Euclidean spaces.

- (a) Let r be a probability measure on $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$, and $q(dy|x)$ be a Borel measurable stochastic kernel on \mathcal{Y} given \mathcal{X} . Then, there exists a unique probability measure p on $(\mathcal{X} \times \mathcal{Y}, \mathcal{B}_{\mathcal{X} \times \mathcal{Y}})$ such that

$$p(B_X \times B_Y) = \int_{B_X} q(B_Y|x)r(dx) \quad \forall B_X \in \mathcal{B}_{\mathcal{X}}, B_Y \in \mathcal{B}_{\mathcal{Y}}. \quad (31)$$

- (b) Let p be a probability measure on $(\mathcal{X} \times \mathcal{Y}, \mathcal{B}_{\mathcal{X} \times \mathcal{Y}})$. Then there exists a Borel-measurable stochastic kernel $q(dy|x)$ on \mathcal{Y} given \mathcal{X} such that (31) holds, where r is the marginal of p on \mathcal{X} .

Lemma 3 (a) guarantees the function p defined on the algebra of measurable rectangles by (31) has a unique extension to the σ -algebra $\mathcal{B}_{\mathcal{X} \times \mathcal{Y}}$. For simplicity, the joint probability measure defined this way is denoted by

$$p(dx, dy) = q(dy|x)r(dx). \quad (32)$$

Conversely, if the left hand side of (32) is given, Lemma 3 (b) guarantees the existence of the decomposition on the right hand side.

Definition 1: A stochastic kernel $q(dz^T|x^T)$ on \mathcal{Z}^T given \mathcal{X}^T is said to be zero-delay if it admits a factorization $q(dz^T|x^T) = \prod_{t=1}^T q(dz_t|z^{t-1}, x^t)$.

Once a zero-delay stochastic kernel is specified, successive applications of Lemma 3 (a) uniquely determine a joint probability measure by $q(dx^T, dz^T) = q(dx^T) \prod_{t=1}^T q(dz_t|z^{t-1}, x^t)$. The mutual information and the expectation in (P-SRD) is understood with respect to this joint probability measure.

Let p and q be probability measures on $\mathcal{X} = \mathbb{R}^n$. Whenever p is absolutely continuous with respect to q (denoted by $p \ll q$), $\frac{dp}{dq}$ denotes the Radon-Nikodym derivative.

Lemma 4: Let \mathcal{X} and \mathcal{Y} be Polish spaces.

- (a) If p, q, r are probability measures on \mathcal{X} such that $r \ll q$ and $q \ll p$, then $r \ll p$ and $\frac{dr}{dp} = \frac{dr}{dq} \frac{dq}{dp} p - a.e..$ If $q \ll p$ and $p \ll q$, then $\frac{dp}{dq} \frac{dq}{dp} = 1$ a.e..
- (b) Let $p_{\mathbf{x}, \mathbf{y}}$ be a joint probability measure on $\mathcal{X} \times \mathcal{Y}$, and $p_{\mathbf{x}}, p_{\mathbf{y}}$ be its marginals. Let $p_{\mathbf{x}|\mathbf{y}}$ be a Borel-measurable stochastic kernel such that

$$p_{\mathbf{x}, \mathbf{y}}(B_X \times B_Y) = \int_{B_Y} p_{\mathbf{x}|\mathbf{y}}(B_X | y) p_{\mathbf{y}}(dy) \quad (33)$$

for every $B_X \in \mathcal{B}_X, B_Y \in \mathcal{B}_Y$. If $p_{\mathbf{x}, \mathbf{y}} \ll p_{\mathbf{x}} \times p_{\mathbf{y}}$, then

$$\frac{dp_{\mathbf{x}, \mathbf{y}}}{d(p_{\mathbf{x}} \times p_{\mathbf{y}})} = \frac{dp_{\mathbf{x}|\mathbf{y}}}{dp_{\mathbf{x}}} p_{\mathbf{y}} - a.e.. \quad (34)$$

Proof: For (a), see Proposition 3.9 in [52]. To prove (b), let $f(x, y) = \frac{dp_{\mathbf{x}, \mathbf{y}}}{d(p_{\mathbf{x}} \times p_{\mathbf{y}})}$. By definition,

$$\begin{aligned} p_{\mathbf{x}, \mathbf{y}}(B_X \times B_Y) &= \int_{B_X \times B_Y} f(x, y) (p_{\mathbf{x}} \times p_{\mathbf{y}})(dx, dy) \\ &= \int_{B_Y} \left(\int_{B_X} f(x, y) p_{\mathbf{x}}(dx) \right) p_{\mathbf{y}}(dy) \end{aligned}$$

Since clearly $f \in L^1(p_{\mathbf{x}} \times p_{\mathbf{y}})$, the Fubini's theorem [52] is applicable in the second line. Substituting this expression into (33), we have $\int_{B_X} f(x, y) p_{\mathbf{x}}(dx) = p_{\mathbf{x}|\mathbf{y}}(B_X | y) p_{\mathbf{y}} - a.e..$ Thus $f(x, y) = \frac{dp_{\mathbf{x}|\mathbf{y}}}{dp_{\mathbf{x}}} p_{\mathbf{y}} - a.e..$ ■

B. Information theoretic quantities

The relative entropy, also known as the Kullback–Leibler divergence, from p to q is defined by

$$D_{\text{KL}}(p \| q) = \begin{cases} \int \log \frac{dp}{dq} dq & \text{if } p \ll q \\ +\infty & \text{otherwise.} \end{cases}$$

Relative entropy is always nonnegative. Given two stochastic kernels $p_{\mathbf{x}|\mathbf{y}}(dx|y)$ and $q_{\mathbf{x}|\mathbf{y}}(dx|y)$ on \mathcal{X} given \mathcal{Y} , and a probability measure $r_{\mathbf{y}}(dy)$, the conditional relative entropy is defined by

$$D_{\text{KL}}(p_{\mathbf{x}|\mathbf{y}} \| q_{\mathbf{x}|\mathbf{y}} | r_{\mathbf{y}}) = \int_{\mathcal{Y}} D_{\text{KL}}(p_{\mathbf{x}|\mathbf{y}}(dx|y) \| q_{\mathbf{x}|\mathbf{y}}(dx|y)) r_{\mathbf{y}}(dy).$$

Suppose $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ are Euclidean spaces, and $q_{\mathbf{x}, \mathbf{y}}$ is a joint probability measure on $\mathcal{X} \times \mathcal{Y}$. Let $q_{\mathbf{x}}, q_{\mathbf{y}}$ be its marginals, and $q_{\mathbf{x}} \times q_{\mathbf{y}}$ be the product measure. The mutual information between \mathbf{x} and \mathbf{y} is defined by $I(\mathbf{x}; \mathbf{y}) = D_{\text{KL}}(q_{\mathbf{x}, \mathbf{y}} \| q_{\mathbf{x}} \times q_{\mathbf{y}})$. Given a joint probability measure $q(dx, dy, dz)$, the conditional mutual information is defined by

$$I(\mathbf{x}; \mathbf{y} | \mathbf{z}) = D_{\text{KL}}(q_{\mathbf{x}, \mathbf{y} | \mathbf{z}} \| q_{\mathbf{x} | \mathbf{z}} \times q_{\mathbf{y} | \mathbf{z}} | q_{\mathbf{z}}).$$

Suppose $\mathcal{X} = \mathbb{R}^n$, and \mathbf{x} is a $(\mathcal{X}, \mathcal{B}_X)$ -valued random variable with probability measure $q_{\mathbf{x}}$. Let λ be the Lebesgue measure on \mathcal{X} restricted to \mathcal{B}_X . The differential entropy of \mathbf{x} is defined by

$$h(\mathbf{x}) = \begin{cases} - \int \log \frac{dq_{\mathbf{x}}}{d\lambda} dq_{\mathbf{x}} & \text{if } q_{\mathbf{x}} \ll \lambda \\ -\infty & \text{otherwise.} \end{cases}$$

APPENDIX B PROOF OF LEMMA 1

(i): Given a sequence of stochastic kernels $\gamma = \otimes_{t=1}^T q(dz_t | x^t, z^{t-1}) \in \Gamma$ attaining cost $f_{\text{SRD}} < +\infty$ in (P-SRD), we are going to construct a sequence of linear-Gaussian stochastic kernels of the form (16) that incurs no greater cost than f_{SRD} in (P-1). Let $q(dx^T, dz^T)$ be the joint probability measure generated by γ and the underlying Gauss-Markov process (5). Without loss of generality, we can assume $q(dx^T, dz^T)$ has zero-mean. Otherwise, it is possible to choose an alternative feasible policy $\tilde{\gamma} \in \Gamma$ by linearly shifting γ so that the resulting probability measure $\tilde{q}(dx^T, dz^T)$ has zero-mean. This operation does not increase the mutual information terms in the objective function.

Let $r(dx^T, dz^T)$ be a zero-mean, jointly Gaussian probability measure with the same covariance as $q(dx^T, dz^T)$. Let $E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}_{t-1} + \dots + F_{t,1} \mathbf{z}_1$ be the least mean square error estimate of \mathbf{z}_t given $\mathbf{x}_t, \mathbf{z}^{t-1}$ in $r(dx^T, dz^T)$, and let Γ_t be the covariance matrix of the corresponding estimation error. Let $\{\mathbf{g}_t\}$ be a sequence of Gaussian random vectors such that \mathbf{g}_t is independent of $\mathbf{x}_0, \mathbf{w}^t, \mathbf{g}^{t-1}$ and $\mathbf{g}_t \sim \mathcal{N}(0, \Gamma_t)$. For every $t = 1, \dots, T$, define a stochastic kernel $s(dz_t | x_t, z^{t-1})$ by

$$\mathbf{z}_t = E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}_{t-1} + \dots + F_{t,1} \mathbf{z}_1 + \mathbf{g}_t.$$

We set $\gamma_1 = \otimes_{t=1}^T s(dz_t | x_t, z^{t-1}) \in \Gamma_1$ as a candidate solution to (P-1). By construction of $s(dz_t | x_t, z^{t-1})$, the following relation holds for every $t = 1, \dots, T$:

$$r(dx_t, dz^t) = s(dz_t | x_t, z^{t-1}) r(dx_t, dz^{t-1}). \quad (35)$$

Let $s(dx^T, dz^T)$ be a jointly Gaussian measure defined by $\{s(dz_t | x_t, z^{t-1})\}_{t=1}^T$ and the process (5). That is, it is a joint measure recursively defined by

$$s(dx^t, dz^{t-1}) = q(dx_t | x_{t-1}) s(dx^{t-1}, dz^{t-1}) \quad (36a)$$

$$s(dx^t, dz^t) = s(dz_t | x_t, z^{t-1}) s(dx^t, dz^{t-1}). \quad (36b)$$

where $q(dx_t | x_{t-1})$ is a stochastic kernel defined by (5).

Notice the following fact about $r(dx^T, dz^T)$.

Proposition 3: For $t = 2, \dots, T$, let $r(dx_{t-1}, dz^{t-1})$ and $r(dx_{t-1}, dx_t, dz^{t-1})$ be marginals of $r(dx^T, dz^T)$. Then $r(dx_{t-1}, dx_t, dz^{t-1}) = q(dx_t | x_{t-1}) r(dx_{t-1}, dz^{t-1})$.

Proof: Since $\mathbf{z}^{t-1} - \mathbf{x}_{t-1} - \mathbf{x}_t$ forms a Markov chain in the measure $q(dx^T, dz^T)$, by Lemma 3.2 of [53], $\mathbf{z}^{t-1} - \mathbf{x}_{t-1} - \mathbf{x}_t$ forms a Markov chain under $r(dx^T, dz^T)$ as well. Hence under r , \mathbf{x}_t is independent of \mathbf{z}^{t-1} given \mathbf{x}_{t-1} , or $r(dx_t | x_{t-1}, z^{t-1}) = r(dx_t | x_{t-1})$. Moreover, since $q(dx_t, dx_{t-1})$ is a Gaussian distribution, and since r is defined to be a Gaussian distribution with the same covariance as q , $r(dx_t, dx_{t-1})$ and $q(dx_t, dx_{t-1})$ have the same joint distribution. Hence, $q(dx_t | x_{t-1}) = r(dx_t | x_{t-1})$. Thus, $r(dx_t | x_{t-1}, z^{t-1}) = q(dx_t | x_{t-1})$, proving the claim. ■

In general, $r(dx^T, dz^T)$ and $s(dx^T, dz^T)$ are different joint probability measures. However, we have the following result.

Proposition 4: For every $t = 1, \dots, T$, let $r(dx_t, dz^t)$ and $s(dx_t, dz^t)$ be marginals of $r(dx^T, dz^T)$ and $s(dx^T, dz^T)$ respectively. Then $r(dx_t, dz^t) = s(dx_t, dz^t)$.

Proof: By definitions,

$$\begin{aligned} r(dx_1, dz_1) &= s(dz_1|x_1)r(dx_1) \\ s(dx_1, dz_1) &= s(dz_1|x_1)q(dx_1). \end{aligned}$$

Since $r(dx_1) = q(dx_1)$, $r(dx_1, dz_1) = s(dx_1, dz_1)$ holds. So assume that the claim holds for $t = k - 1$. Then

$$\begin{aligned} s(dx_k, dz^k) &= s(dz_k|x_k, z^{k-1})s(dx_k, dz^{k-1}) \quad (37a) \end{aligned}$$

$$\begin{aligned} &= s(dz_k|x_k, z^{k-1}) \int_{\mathcal{X}_{k-1}} s(dx_{k-1}, dx_k, dz^{k-1}) \\ &= s(dz_k|x_k, z^{k-1}) \int_{\mathcal{X}_{k-1}} q(dx_k|x_{k-1})s(dx_{k-1}, dz^{k-1}) \quad (37b) \end{aligned}$$

$$= s(dz_k|x_k, z^{k-1}) \int_{\mathcal{X}_{k-1}} q(dx_k|x_{k-1})r(dx_{k-1}, dz^{k-1}) \quad (37c)$$

$$= s(dz_k|x_k, z^{k-1}) \int_{\mathcal{X}_{k-1}} r(dx_{k-1}, dx_k, dz^{k-1}) \quad (37d)$$

$$\begin{aligned} &= s(dz_k|x_k, z^{k-1})r(dx_k, dz^{k-1}) \\ &= r(dx_k, dz^k). \quad (37e) \end{aligned}$$

The first step (37a) follows from the definition (36b). Step (37b) also follows from the definition (36a). In (37c), the induction assumption $s(dx_{k-1}, dz^{k-1}) = r(dx_{k-1}, dz^{k-1})$ was used. The result of Proposition 3 was used in (37d). The final step (37e) is due to (35). ■

To prove that $\gamma_1 = \otimes_{t=1}^T s(dz_t|x_t, z^{t-1})$ incurs no greater cost than f_{SRD} in (P-1), notice that replacing $q(dz_t|x^t, z^{t-1})$ with $s(dz_t|x_t, z^{t-1})$ will not change the distortion:

$$\begin{aligned} \mathbb{E}_q \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 &= \int \|x_t - z_t\|_{\Theta_t}^2 q(dx_t, dz^t) \\ &= \int \|x_t - z_t\|_{\Theta_t}^2 r(dx_t, dz^t) \quad (38) \end{aligned}$$

$$\begin{aligned} &= \int \|x_t - z_t\|_{\Theta_t}^2 s(dx_t, dz^t) \quad (39) \\ &= \mathbb{E}_s \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2. \end{aligned}$$

Equality (38) holds since q and r have the same second order properties. The result of Proposition 4 was used in step (39).

Next, we show that the mutual information never increases by this replacement.

Proposition 5: If $I_q(\mathbf{x}^T; \mathbf{z}^T) < +\infty$, then $I_r(\mathbf{x}^T; \mathbf{z}^T) \leq I_q(\mathbf{x}^T; \mathbf{z}^T)$.

Proof: This can be directly verified as

$$\begin{aligned} I_q(\mathbf{x}^T; \mathbf{z}^T) - I_r(\mathbf{x}^T; \mathbf{z}^T) &= \int \log \frac{dq(x^T|z^T)}{dq(x^T)} q(dx^T, dz^T) \quad (40) \end{aligned}$$

$$- \int \log \frac{dr(x^T|z^T)}{dr(x^T)} r(dx^T, dz^T) \quad (41)$$

$$\begin{aligned} &= \int \log \frac{dq(x^T|z^T)}{dq(x^T)} q(dx^T, dz^T) \\ &\quad - \int \log \frac{dr(x^T|z^T)}{dr(x^T)} r(dx^T, dz^T) \quad (42) \end{aligned}$$

$$\begin{aligned} &= \int \log \left(\frac{dq(x^T|z^T)}{dq(x^T)} \cdot \frac{dr(x^T)}{dr(x^T|z^T)} \right) q(dx^T, dz^T) \\ &= \int \log \left(\frac{dq(x^T|z^T)}{dr(x^T|z^T)} \right) q(dx^T, dz^T) \quad (43) \end{aligned}$$

$$\begin{aligned} &= \int \left(\int \log \left(\frac{dq(x^T|z^T)}{dr(x^T|z^T)} \right) q(dx^T|z^T) \right) q(dz^T) \\ &= \int D_{\text{KL}}(q(x^T|z^T) \| r(x^T|z^T)) q(dz^T) \geq 0. \end{aligned}$$

(40) is by definition of mutual information and Lemma 4 (b). Since $q(dx^T)$ is a non-degenerate Gaussian probability measure, $I_q(\mathbf{x}^T; \mathbf{z}^T) < +\infty$ implies that $q(dx^T|z^T)$ admits a density $q(dz^T) - a.e.$. This further requires that a Gaussian measure $r(dx^T|z^T)$ admits a density everywhere in $\text{supp}(r(dz^T))$, i.e., the support of the probability measure $r(dz^T)$. Thus, the Radon-Nikodym derivative in (41) exists everywhere in $\text{supp}(r(dz^T))$. Since r is a Gaussian probability measure, $\log \frac{dr(x^T|z^T)}{dr(x^T)}$ is a quadratic function of x^T and z^T everywhere in $\text{supp}(r(dx^T, dz^T))$. Since it can be shown that $\text{supp}(q(dx^T, dz^T)) \subseteq \text{supp}(r(dx^T, dz^T))$, this allows us to replace $r(dx^T, dz^T)$ with $q(dx^T, dz^T)$ in (42) since they have the same second order moments. Lemma 33 (a) is applicable in (43) since $r(dx^T) = q(dx^T)$. ■

Finally,

$$\sum_{t=1}^T I_q(\mathbf{x}^t; \mathbf{z}_t | \mathbf{z}^{t-1}) = I_q(\mathbf{x}^T; \mathbf{z}^T) \quad (44)$$

$$\geq I_r(\mathbf{x}^T; \mathbf{z}^T) \quad (45)$$

$$= \sum_{t=1}^T I_r(\mathbf{x}^t; \mathbf{z}_t | \mathbf{z}^{t-1})$$

$$\geq \sum_{t=1}^T I_r(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1})$$

$$= \sum_{t=1}^T I_s(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) \quad (46)$$

See Remark 1 for the equality (44). The result of Proposition 5 was used in (45). Equality (46) follows from Proposition 4. Thus, using $\gamma = \otimes_{t=1}^T q(dz_t|x_t, z^{t-1}) \in \Gamma$ attaining cost f_{SRD} in (P-SRD), we have constructed $\gamma_1 = \otimes_{t=1}^T s(dz_t|x_t, z^{t-1}) \in \Gamma_1$ incurring smaller cost in (P-1) than f_{SRD} .

(ii): Let $\gamma_1 = \otimes_{t=1}^T q(dz_t|x_t, z^{t-1}) \in \Gamma_1$ be a sequence of linear-Gaussian stochastic kernels attaining $f_1 < +\infty$ in (P-1), and $q(dx^T, dz^T)$ be the resulting joint probability measure. Since $\mathbf{z}_t - (\mathbf{x}_t, \mathbf{z}^{t-1}) - \mathbf{x}^{t-1}$ forms a Markov chain

in $q(dx^T, dz^T)$, we have

$$\begin{aligned} I(\mathbf{x}^t; \mathbf{z}_t | \mathbf{z}^{t-1}) &= I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) + I(\mathbf{x}^{t-1}; \mathbf{z}_t | \mathbf{x}_t, \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}). \end{aligned} \quad (47)$$

Hence the mutual information terms in (P-1) can be replaced with the ones in (P-SRD) without increasing cost.

APPENDIX C PROOF OF LEMMA 2

(i): Suppose

$$\mathbf{z}_t = E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}_{t-1} + \dots + F_{t,1} \mathbf{z}_1 + \mathbf{g}_t, t = 1, \dots, T \quad (48)$$

is a linear-Gaussian stochastic kernel that attains $f_1 < +\infty$ in (P-1). It is sufficient for us to show that there exist nonnegative integers r_1, \dots, r_T and matrices $C_t \in \mathbb{R}^{r_t \times n_t}$, $V_t \in \mathbb{S}_{++}^{r_t}$, $t = 1, \dots, T$ such that $\{C_t, V_t\}_{t=1}^T$ attains a smaller cost than f_1 in (P-LGS). Let

$$\begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix} = \mathbb{E} \mathbf{g}_t \mathbf{g}_t^\top$$

with an orthonormal matrix $U = \begin{bmatrix} U_1 & U_2 \end{bmatrix}$ be a singular value decomposition of the covariance matrix of \mathbf{g}_t . If \mathbf{g}_t is nondegenerate, we understand that $U = U_1$, while if \mathbf{g}_t is a point mass at zero, then $U = U_2$. Clearly $\tilde{\mathbf{g}}_t = U_1^\top \mathbf{g}_t$ is a zero-mean, nondegenerate Gaussian random vector and $U_2^\top \mathbf{g}_t = 0$. Define

$$\begin{bmatrix} \tilde{\mathbf{z}}_t \\ \hat{\mathbf{z}}_t \end{bmatrix} = \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix} \mathbf{z}_t, \begin{bmatrix} \tilde{E}_t \\ \hat{E}_t \end{bmatrix} = \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix} E_t, \begin{bmatrix} \tilde{F}_{t,s} \\ \hat{F}_{t,s} \end{bmatrix} = \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix} F_{t,s}$$

for $s = 1, \dots, t-1$. Then multiplying (48) by U^\top from the left yields

$$\begin{bmatrix} \tilde{\mathbf{z}}_t \\ \hat{\mathbf{z}}_t \end{bmatrix} = \begin{bmatrix} \tilde{E}_t \\ \hat{E}_t \end{bmatrix} \mathbf{x}_t + \begin{bmatrix} \tilde{F}_{t,t-1} \\ \hat{F}_{t,t-1} \end{bmatrix} \mathbf{z}_{t-1} + \dots + \begin{bmatrix} \tilde{F}_{t,1} \\ \hat{F}_{t,1} \end{bmatrix} \mathbf{z}_1 + \begin{bmatrix} \tilde{\mathbf{g}}_t \\ 0 \end{bmatrix}. \quad (49)$$

Proposition 6: $\hat{E}_t = 0 \forall t = 1, \dots, T$ is necessary for $f_1 < +\infty$.

Proof: Focus on the mutual information terms in (P-1).

$$\begin{aligned} I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t, \hat{\mathbf{z}}_t | \mathbf{z}^{t-1}) \\ &\geq I(\mathbf{x}_t; \hat{\mathbf{z}}_t | \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \hat{E}_t \mathbf{x}_t + \hat{F}_{t,t-1} \mathbf{z}_{t-1} + \dots + \hat{F}_{t,1} \mathbf{z}_1 | \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \hat{E}_t \mathbf{x}_t | \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \hat{E}_t \mathbf{x}_t, \mathbf{z}^{t-1}) - I(\mathbf{x}_t; \mathbf{z}^{t-1}) \\ &\geq I(\mathbf{x}_t; \hat{E}_t \mathbf{x}_t) - I(\mathbf{x}_t; \mathbf{z}^{t-1}) \end{aligned}$$

Recall that \mathbf{x}_t is defined by (5) and is a nondegenerate Gaussian random vector. If $\hat{E}_t \mathbf{x}_t$ is a non-zero linear function of \mathbf{x}_t , then $I(\mathbf{x}_t; \hat{E}_t \mathbf{x}_t) = +\infty$, while $I(\mathbf{x}_t; \mathbf{z}^{t-1})$ is bounded. Therefore, $\hat{E}_t = 0$ is necessary for $I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1})$ to be bounded. ■

Proposition 6, together with (49), implies that $\hat{\mathbf{z}}_t$ is a linear function of \mathbf{z}^{t-1} . Hence, there exist some matrices

$H_{t,1}, \dots, H_{t,t-1}$ such that the first row of (49) can be rewritten as

$$\tilde{\mathbf{z}}_t = \tilde{E}_t \mathbf{x}_t + H_{t,t-1} \tilde{\mathbf{z}}_{t-1} + \dots + H_{t,1} \tilde{\mathbf{z}}_1 + \tilde{\mathbf{g}}_t. \quad (50)$$

It is also easy to see that \mathbf{z}^t can be fully reconstructed if $\tilde{\mathbf{z}}^t$ is given. In particular, this implies that the σ -algebras generated by \mathbf{z}^t and $\tilde{\mathbf{z}}^t$ are the same.

$$\sigma(\mathbf{z}^t) = \sigma(\tilde{\mathbf{z}}^t). \quad (51)$$

Proposition 7: $I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) = I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}^{t-1}) \forall t = 1, \dots, T$.

Proof: This can be directly verified as follows.

$$\begin{aligned} I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t, \hat{\mathbf{z}}_t | \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \mathbf{z}^{t-1}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}_{t-1}, \hat{\mathbf{z}}_{t-1}, \mathbf{z}^{t-2}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}_{t-1}, \mathbf{z}^{t-2}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}_{t-1}, \tilde{\mathbf{z}}_{t-2}, \mathbf{z}^{t-3}) \\ &\quad \vdots \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}^{t-1}) \end{aligned} \quad (52)$$

Equality (52) holds since $\hat{\mathbf{z}}_t$ is a linear function of \mathbf{z}^{t-1} . Similarly, (53) holds because $\hat{\mathbf{z}}_{t-1}$ is a linear function of \mathbf{z}^{t-2} . The remaining equalities can be shown by repeating the same argument. ■

Now, for every $t = 1, \dots, T$, set $C_t = \tilde{E}_t$ and $\mathbf{v}_t = \tilde{\mathbf{g}}_t$. Then, by construction, \mathbf{v}_t is a zero-mean, nondegenerate Gaussian random vector that is independent of $\mathbf{x}_0, \mathbf{w}^t, \mathbf{v}^{t-1}$. Hence $\mathbf{y}_t = C_t \mathbf{x}_t + \mathbf{v}_t$ is an admissible sensor equation for (P-LGS).

Proposition 8: $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) = I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}^{t-1}) \forall t = 1, \dots, T$.

Proof: By concatenating (50), it can be easily seen that an identity $\mathcal{H}_t \tilde{\mathbf{z}}^t = \mathbf{y}^t$ holds for every $t = 1, \dots, T$, where \mathcal{H}_t is an invertible matrix defined by

$$\mathcal{H}_t = \begin{bmatrix} I & 0 & \dots & 0 \\ -H_{2,1} & I & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ -H_{t,1} & \dots & -H_{t,t-1} & I \end{bmatrix}.$$

Hence,

$$\begin{aligned} I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) &= I(\mathbf{x}_t; \mathbf{y}_t | \tilde{\mathbf{z}}^{t-1}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t - H_{t,t-1} \tilde{\mathbf{z}}_{t-1} - \dots - H_{t,1} \tilde{\mathbf{z}}_1 | \tilde{\mathbf{z}}^{t-1}) \\ &= I(\mathbf{x}_t; \tilde{\mathbf{z}}_t | \tilde{\mathbf{z}}^{t-1}). \end{aligned}$$

Thus, starting from a sequence of linear-Gaussian stochastic kernels (48), we have constructed a sequence of sensor equations of the form $\mathbf{y}_t = C_t \mathbf{x}_t + \mathbf{v}_t$ such that $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) = I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1})$. The last equality is a consequence of Propositions 7 and 8. To complete the proof of the first statement of Lemma 2, it is left to show that

$$\mathbb{E} \|\mathbf{x}_t - \mathbf{z}'_t\|_{\Theta_t}^2 \leq \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \forall t = 1, \dots, T \quad (54)$$

where $\mathbf{z}'_t = \mathbb{E}(\mathbf{x}_t | \mathbf{y}^t)$. (Here, we refer to the variable “ \mathbf{z}_t ” in (P-LGS) as \mathbf{z}'_t in order to distinguish it from the variable \mathbf{z}_t in (P-1).) The inequality (54) can be verified by the following observation. Since $\mathcal{H}_t \tilde{\mathbf{z}}^t = \mathbf{y}^t$, we have $\sigma(\mathbf{y}^t) = \sigma(\tilde{\mathbf{z}}^t)$. Moreover, it follows from (51) that $\sigma(\mathbf{y}^t) = \sigma(\mathbf{z}^t)$. Thus, \mathbf{z}_t is $\sigma(\mathbf{y}^t)$ -measurable. However, since $\mathbf{z}'_t = \mathbb{E}(\mathbf{x}_t | \mathbf{y}^t)$, \mathbf{z}'_t minimizes the mean square estimation error among all $\sigma(\mathbf{y}^t)$ -measurable functions. Thus (54) must hold.

(ii): Let $\{C_t, V_t\}_{t=1}^T$ be a sequence of matrices that attains $f_{\text{LGS}} < +\infty$ in (P-LGS). Let \mathbf{y}_t be defined by (8), and $\mathbf{z}'_t = \mathbb{E}(\mathbf{x}_t | \mathbf{y}^t)$ be the least mean square error estimate of \mathbf{x}_t given \mathbf{y}^t obtained by the Kalman filter. From the Kalman filtering formula, we have

$$\begin{aligned} \mathbf{z}'_t &= A_{t-1} \mathbf{z}'_{t-1} + P_{t|t-1} C_t^\top (C_t P_{t|t-1} C_t^\top + V_t)^{-1} (\mathbf{y}_t - C_t A_{t-1} \mathbf{z}'_{t-1}) \\ &= E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}'_{t-1} + \cdots + F_{t,1} \mathbf{z}'_1 + \mathbf{g}_t \end{aligned}$$

where $E_t, F_{t,t-1}, \dots, F_{t,1}$ are some matrices (in fact, all $F_{t,t-2}, \dots, F_{t,1}$ are zero matrices) and \mathbf{g}_t is a zero-mean Gaussian random vector that is independent of $\mathbf{x}_0, \mathbf{w}^t$ and \mathbf{g}^{t-1} . Hence, by constructing a linear-Gaussian stochastic kernel for (P-1) by

$$\mathbf{z}_t = E_t \mathbf{x}_t + F_{t,t-1} \mathbf{z}_{t-1} + \cdots + F_{t,1} \mathbf{z}_1 + \mathbf{g}_t$$

using the same $E_t, F_{t,t-1}, \dots, F_{t,1}$ and \mathbf{g}_t , $(\mathbf{x}^T, \mathbf{z}^T)$ and $(\mathbf{x}^T, \mathbf{z}'^T)$ have the same joint distribution. Thus $\mathbb{E} \|\mathbf{x}_t - \mathbf{z}'_t\|_{\Theta_t}^2 = \mathbb{E} \|\mathbf{x}_t - \mathbf{z}_t\|_{\Theta_t}^2 \forall t = 1, \dots, T$. Hence, it remains to prove that

$$I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) \geq I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1}) \quad \forall t = 1, \dots, T.$$

Notice that $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{z}^{t-1}) \geq I(\mathbf{x}_t; \mathbf{z}_t | \mathbf{z}^{t-1})$ is immediate from the data-processing inequality. Moreover, an equality $I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{y}^{t-1}) = I(\mathbf{x}_t; \mathbf{y}_t | \mathbf{z}^{t-1})$ holds since the input sequence \mathbf{y}^{t-1} and the output sequence \mathbf{z}^{t-1} of the Kalman filter contain statistically equivalent information. Formally, this can be shown by proving that the Kalman filter is causally invertible [54], and thus one can construct \mathbf{y}^{t-1} from \mathbf{z}^{t-1} and *vice versa*.

ACKNOWLEDGMENT

The authors would like to thank Prof. Sekhar Tatikonda for valuable discussions.

REFERENCES

- [1] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York, NY, USA: Wiley-Interscience, 1991.
- [2] N. T. Gaarder and D. Slepian, “On optimal finite-state digital transmission systems,” *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 167–186, 1982.
- [3] H. Witsenhausen, “On the structure of real-time source coders,” *Bell System Technical Journal*, vol. 58, no. 6, pp. 1437–1451, 1979.
- [4] D. L. Neuhoff and R. K. Gilbert, “Causal source codes,” *IEEE Transactions on Information Theory*, vol. 28, no. 5, pp. 701–713, 1982.
- [5] T. Linder and R. Zamir, “Causal coding of stationary sources and individual sequences with high resolution,” *IEEE Transactions on Information Theory*, vol. 52, no. 2, pp. 662–680, 2006.
- [6] S. Yüksel, “On optimal causal coding of partially observed markov sources in single and multiterminal settings,” *IEEE Transactions on Information Theory*, vol. 59, no. 1, pp. 424–437, 2013.
- [7] V. Kostina and S. Verdú, “Fixed-length lossy compression in the finite blocklength regime,” *Information Theory, IEEE Transactions on*, vol. 58, no. 6, pp. 3309–3338, 2012.
- [8] J. C. Walrand and P. Varaiya, “Optimal causal coding-decoding problems,” *IEEE Transactions on Information Theory*, vol. 29, no. 6, pp. 814–820, 1983.
- [9] D. Teneketzis, “On the structure of optimal real-time encoders and decoders in noisy communication,” *IEEE Transactions on Information Theory*, vol. 52, no. 9, pp. 4017–4035, 2006.
- [10] A. Mahajan and D. Teneketzis, “Optimal design of sequential real-time communication systems,” *IEEE Transactions on Information Theory*, vol. 55, no. 11, pp. 5317–5338, 2009.
- [11] S. K. Gorantla and T. P. Coleman, “Information-theoretic viewpoints on optimal causal coding-decoding problems,” *CoRR*, vol. abs/1102.0250, 2011.
- [12] S. Tatikonda, “Control under communication constraints,” *PhD thesis, Massachusetts Institute of Technology*, 2000.
- [13] S. Tatikonda, A. Sahai, and S. Mitter, “Stochastic linear control over a communication channel,” *IEEE Transactions on Automatic Control*, vol. 49, no. 9, pp. 1549–1561, 2004.
- [14] A. Gorbunov and M. Pinsker, “Nonanticipatory and prognostic epsilon entropies and message generation rates,” *Problemy Peredachi Informat-sii*, vol. 9, no. 3, pp. 12–21, 1973.
- [15] R. Bucy, “Distortion-rate theory and filtering,” in *Advances in Communications*. Springer, 1980, pp. 11–15.
- [16] M. S. Derpich and J. Ostergaard, “Improved upper bounds to the causal quadratic rate-distortion function for Gaussian stationary sources,” *IEEE Transactions on Information Theory*, vol. 58, no. 5, pp. 3131–3152, 2012.
- [17] C. D. Charalambous, P. Stavrou, N. U. Ahmed *et al.*, “Nonanticipative rate distortion function and relations to filtering theory,” *IEEE Transactions on Automatic Control*, vol. 59, no. 4, pp. 937–952, 2014.
- [18] J. Massey, “Causality, feedback and directed information,” *International Symposium on Information Theory and Its Applications (ISITA)*, pp. 303–305, 1990.
- [19] E. I. Silva, M. S. Derpich, and J. Ostergaard, “A framework for control system design subject to average data-rate constraints,” *Automatic Control, IEEE Transactions on*, vol. 56, no. 8, pp. 1886–1899, 2011.
- [20] G. Blekherman, P. Parrilo, and R. Thomas, *Semidefinite Optimization and Convex Algebraic Geometry*, ser. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics, 2013.
- [21] T. E. Duncan, “On the calculation of mutual information,” *SIAM Journal on Applied Mathematics*, vol. 19, no. 1, pp. 215–220, 1970.
- [22] D. Guo, S. Shamai, and S. Verdú, “Mutual information and minimum mean-square error in Gaussian channels,” *IEEE Transactions on Information Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.
- [23] T. Weissman, Y.-H. Kim, and H. H. Permuter, “Directed information, causal estimation, and communication in continuous time,” *IEEE Transactions on Information Theory*, vol. 59, no. 3, pp. 1271–1287, 2013.
- [24] F. Rezaei, N. Ahmed, and C. D. Charalambous, “Rate distortion theory for general sources with potential application to image compression,” *International Journal of Applied Mathematical Sciences*, vol. 3, no. 2, pp. 141–165, 2006.
- [25] H. Marko, “The bidirectional communication theory—a generalization of information theory,” *IEEE Transactions on Communications*, vol. 21, no. 12, pp. 1345–1351, 1973.
- [26] J. L. Massey and P. C. Massey, “Conservation of mutual and directed information,” in *Proceedings. International Symposium on Information Theory, 2005*. IEEE, 2005, pp. 157–158.
- [27] T. Tanaka and H. Sandberg, “SDP-based joint sensor and controller design for information-regularized optimal LQG control,” *54th IEEE Conference on Decision and Control (CDC)*, 2015.
- [28] D. P. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1995.
- [29] L. Vandenberghe, S. Boyd, and S.-P. Wu, “Determinant maximization with linear matrix inequality constraints,” *SIAM journal on matrix analysis and applications*, vol. 19, no. 2, pp. 499–533, 1998.
- [30] D. A. Harville, *Matrix algebra from a statistician’s perspective*. Springer, 1997, vol. 1.
- [31] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.
- [32] A. Ben-Tal and A. Nemirovski, *Lectures on modern convex optimization: Analysis, algorithms, and engineering applications*. Philadelphia, PA, USA: SIAM, 2001, vol. 2.
- [33] R. H. Tütüncü, K. C. Toh, and M. J. Todd, “Solving semidefinite-quadratic-linear programs using SDPT3,” *Mathematical programming*, vol. 95, no. 2, pp. 189–217, 2003.
- [34] J. Renegar, *A mathematical view of interior-point methods in convex optimization*. Philadelphia, PA, USA: SIAM, 2001, vol. 3.

- [35] K.-C. Toh, "Primal-dual path-following algorithms for determinant maximization problems with linear matrix inequalities," *Computational Optimization and Applications*, vol. 14, no. 3, pp. 309–330, 1999.
- [36] T. Tsuchiya and Y. Xia, "An extension of the standard polynomial-time primal-dual path-following algorithm to the weighted determinant maximization problem with semidefinite constraints," *Pacific Journal of Optimization*, vol. 3, no. 1, pp. 165–182, 2007.
- [37] M. Fukuda, M. Kojima, K. Murota, and K. Nakata, "Exploiting sparsity in semidefinite programming via matrix completion I: General framework," *SIAM Journal on Optimization*, vol. 11, no. 3, pp. 647–674, 2001.
- [38] K. Nakata, K. Fujisawa, M. Fukuda, M. Kojima, and K. Murota, "Exploiting sparsity in semidefinite programming via matrix completion II: Implementation and numerical results," *Mathematical Programming*, vol. 95, no. 2, pp. 303–327, 2003.
- [39] L. Vandenberghe and M. Andersen, *Chordal Graphs and Semidefinite Optimization*. Now Publishers Incorporated, 2015.
- [40] T. Tanaka, "Semidefinite representation of sequential rate-distortion function for stationary Gauss-Markov processes," *IEEE Multi-Conference on Systems and Control (MSC)*, 2015.
- [41] G. N. Nair, F. Fagnani, S. Zampieri, and R. J. Evans, "Feedback control under data rate constraints: An overview," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 108–137, 2007.
- [42] J. Baillieul and P. J. Antsaklis, "Control and communication challenges in networked real-time systems," *Proceedings of the IEEE*, vol. 95, no. 1, pp. 9–28, 2007.
- [43] J. P. Hespanha, P. Naghshtabrizi, and Y. Xu, "A survey of recent results in networked control systems," *Proceedings of the IEEE*, vol. 95, no. 1, p. 138, 2007.
- [44] S. Yüksel and T. Başar, *Stochastic networked control systems*, ser. Systems & Control Foundations & Applications. New York, NY: Springer, 2013, vol. 10.
- [45] 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, 1989.
- [46] B. M. Miller and W. J. Runggaldier, "Optimization of observations: a stochastic control approach," *SIAM journal on control and optimization*, vol. 35, no. 3, pp. 1030–1052, 1997.
- [47] T. Summers, F. Cortesi, and J. Lygeros, "On submodularity and controllability in complex dynamical networks," *IEEE Transactions on Control of Network Systems*, vol. PP, no. 99, 2015.
- [48] M. P. Vitus, W. Zhang, A. Abate, J. Hu, and C. J. Tomlin, "On efficient sensor scheduling for linear dynamical systems," *Automatica*, vol. 48, no. 10, pp. 2482–2493, 2012.
- [49] V. Gupta, T. H. Chung, B. Hassibi, and R. M. Murray, "On a stochastic sensor selection algorithm with applications in sensor scheduling and sensor coverage," *Automatica*, vol. 42, no. 2, pp. 251–260, 2006.
- [50] T. Tanaka, "Zero-delay rate-distortion optimization for partially observable gauss-markov processes," *54th IEEE Conference on Decision and Control (CDC)*, 2015.
- [51] D. P. Bertsekas and S. E. Shreve, *Stochastic optimal control: The discrete time case*. Academic Press New York, 1978, vol. 139.
- [52] G. Folland, *Real analysis: Modern techniques and their applications*. Hoboken, NJ, USA: John Wiley & Sons, 1999.
- [53] Y.-H. Kim, "Feedback capacity of stationary Gaussian channels," *IEEE Transactions on Information Theory*, vol. 56, no. 1, pp. 57–85, 2010.
- [54] T. Kailath, "An innovations approach to least-squares estimation – Part I: Linear filtering in additive white noise," *IEEE Transactions on Automatic Control*, vol. 13, no. 6, pp. 646–655, 1968.