# A Training-Based Mutual Information Lower Bound for Large-Scale Systems

Xiangbo Meng, *Student Member, IEEE,* Kang Gao, *Student Member, IEEE,* and Bertrand M. Hochwald, *Fellow, IEEE*

## Abstract

We provide a mutual information lower bound that can be used to analyze the effect of training in models with unknown parameters. For large-scale systems, we show that this bound can be calculated using the difference between two derivatives of a conditional entropy function. The bound does not require explicit estimation of the unknown parameters. We provide a step-by-step process for computing the bound, and provide an example application. A comparison with known classical mutual information bounds is provided.

## Index Terms

information rates, training, entropy, large-scale systems

## I. INTRODUCTION

Many systems have unknown parameters that are estimated during a training-phase with the help of known prescribed training signals. This phase is followed by a data phase, where knowledge of the estimated parameters is used to process the data. It is generally assumed that the parameters are constant during these two phases, the total duration of which is called the coherence time. It is often of great interest to optimize the training time for a given coherence

time, since time in the training phase, while useful for parameter estimation, generally takes away from time in the data phase.

In a communication system, the parameters of interest often include the channel, which is typically unknown and learned at the receiver with the help of pilot signals sent by the transmitter. For example, [1] analyzes a multi-antenna model and a capacity lower-bound is obtained by using the minimum mean-square estimate (MMSE) of the channel, and the residual channel error is treated as Gaussian noise. This lower bound is maximized over various parameters, including the fraction of the coherence time that should be dedicated to training. A similar optimization is considered in [2], where the power allocation and training duration are chosen to achieve the maximum sum-rate in a multiuser system. Such "one-shot learning", where the parameters are estimated only during the training phase, can be augmented by further refinement during the data phase [3], [4]. However, this refinement can suffer from error propagation [5], and we do not consider this herein.

Many of these previous efforts to analyze training assume that the unknown parameters appear linearly in the system model [1]–[4], [6], or appear in a linearized version of the model [7], [8], often by employing the Bussgang decomposition [9]. We develop a framework to analyze one-shot training that does not require the parameters to appear linearly in the model, nor does it require additive Gaussian noise; rather, it requires the system to be time-invariant and memoryless, and a certain entropy to be computed in the large-scale system limit. Herein, large-scale refers to long block lengths (time durations) or large dimensional inputs and outputs, or both. The fact that the large-scale system entropy can sometimes be computed even when the small-scale system entropy cannot is exploited for our training analysis.

*A. Problem setup and statement*

Consider a system model that has input and output processes as $\mathcal{X} = (\mathbf{x}_1, \mathbf{x}_2, \ldots)$ and $\mathcal{Y} = (\mathbf{y}_1, \mathbf{y}_2, \ldots)$, which comprise vectors $\mathbf{x}_t$ and $\mathbf{y}_t$ whose dimensions are $M$ and $N$ respectively. The input and output are connected through a conditional distribution parameterized by $\mathbf{g}_T$, whose value is unknown. We assume that $\mathbf{g}_T$ is constant during a coherence time block $T$, and then changes independently in the next block (same length $T$), and so on. The system is supplied with

known inputs during a "training phase" to learn the parameters, after which the system is used during its "data phase". The unknown parameters are assumed to have a known distribution, and the number of unknown parameters is allowed to be a function of $T$. *Problem statement:* We wish to determine the optimum amount of training.

To analyze the effects of training, a lower bound on the mutual information between the input and output is often used

$$\frac{1}{T}I(X_T; Y_T) \geq \frac{T - \tau T}{T}I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1}|X_{\tau T}, Y_{\tau T}), \tag{1}$$

where $\mathbf{x}_t$ and $\mathbf{y}_t$ are the $t$th vector input and output of the system, $X_t = [\mathbf{x}_1, \cdots, \mathbf{x}_t]$, $Y_t = [\mathbf{y}_1, \cdots, \mathbf{y}_t]$, and $\tau T$ is the number of training symbols in one coherence block. We assume $0 < \tau < 1$ is the fraction of the blocklength devoted to training, and $\tau T$ is integer for convenience. (We choose this in favor of using $\lceil \tau T \rceil$ throughout.) The optimal training fraction, in the sense of maximizing this lower bound, is then

$$\tau_{\text{opt}} = \underset{\tau}{\text{argmax}} \, (1 - \tau)I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1}|X_{\tau T}, Y_{\tau T}). \tag{2}$$

Such an analysis appears, for example, in [1], [2], [7], [8], but the right-hand side of (2) can be difficult to compute and is itself often approximated or lower bounded. For example, in [1], a wireless communication system with Rayleigh block-fading channel and additive Gaussian noise is considered, and the mutual information in (2) is lower bounded by treating the estimation error of the MMSE estimate of the channel as independent additive Gaussian noise. However, this form of analysis is often intractable when the parameters appear nonlinearly, or the additive noise is non-Gaussian, since explicit estimates of the unknown parameters are unavailable.

By considering a large-scale limit of the conditional mutual information in (2), we provide a method to revisit this computation. Let $T \to \infty$, and define the ratios

$$\alpha = \frac{N}{M}, \quad \beta = \frac{T}{M}. \tag{3}$$

It is possible, although not required, that $M$ and $N$ also grow to infinity with $T$, so that $\beta$ is finite. The large-scale limit of the conditional mutual information $I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1}|X_{\tau T}, Y_{\tau T})$ in

(2) is

$$\mathcal{I}'(\mathcal{X};\mathcal{Y}) = \lim_{T\to\infty} \frac{1}{N} I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | X_{\tau T}, Y_{\tau T}). \tag{4}$$

The normalization by $1/N$ is needed to keep this quantity finite if $N \to \infty$, and this limit (assuming that it exists) typically depends on $\alpha$, $\beta$, and $\tau$. The optimal training time in (2) then becomes

$$\tau_{\text{opt}} = \underset{\tau}{\text{argmax}}\ (1 - \tau)\mathcal{I}'(\mathcal{X};\mathcal{Y}), \tag{5}$$

and the corresponding optimal rate becomes

$$\mathcal{R}_{\text{opt}} = (1 - \tau_{\text{opt}})\mathcal{I}'(\mathcal{X};\mathcal{Y})\big|_{\tau=\tau_{\text{opt}}}.$$

The value of this analysis depends on our ability to compute $\mathcal{I}'(\mathcal{X};\mathcal{Y})$, and we show that this quantity can be computed as the derivative of a certain entropy.

## II. MAIN RESULTS

### A. Assumptions and definitions of useful quantities

Before we introduce the main results, we first make some assumptions and definitions. The bound in (1) is fully determined by the distribution of the triple $(X_T, Y_T, \mathbf{g}_T)$, and we make the following assumption:

$$\text{A1:}\ \ p(Y_T|X_T; \mathbf{g}_T) = \prod_{t=1}^{T} p(\mathbf{y}_t|\mathbf{x}_t; \mathbf{g}_T), \tag{6}$$

$$p(X_T) = p(X_{\tau T}) \prod_{t=\tau T+1}^{T} p(\mathbf{x}_t), \tag{7}$$

where $p(\mathbf{y}_t|\mathbf{x}_t; \mathbf{g}_T)$ is a fixed conditional distribution for all $t = 1, 2, \ldots, T$ and $p(\mathbf{x}_t)$ is a fixed distribution for all $t = \tau T + 1, \tau T + 2, \ldots, T$.

Equation (6) says that the system is memoryless and time-invariant (given the input and parameters) and (7) says that the input $\mathbf{x}_t$ is *iid* and independent of $X_{\tau T}$ for all $t > \tau T$. We use the common convention of writing $p(X_{\tau T})$ and $p(\mathbf{x}_t)$ when we mean $p_{X_{\tau T}}(\cdot)$ and $p_{\mathbf{x}_t}(\cdot)$, even

4

though these functions can differ. Under A1, the distributions of $(X_T, Y_T, \mathbf{g}_T)$ are described by the set of known distributions

$$\mathcal{P}(T, \tau) = \{p(\mathbf{y}|\mathbf{x}; \mathbf{g}_T)), p(\mathbf{g}_T)), p(X_{\tau T}), p(\mathbf{x}_{\tau T+1})\}. \tag{8}$$

These distributions are used to calculate all of the entropies and mutual informations throughout. The entropies and mutual informations are "ergodic" in the sense that they are averaged over independent realizations of $\mathbf{g}_T$.

Define:

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \lim_{T \to \infty} \frac{1}{N} H(\mathbf{y}_{\varepsilon\tau T+1}|X_{\delta\tau T}, Y_{\varepsilon\tau T}), \tag{9}$$

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\delta+}) = \lim_{T \to \infty} \frac{1}{N} H(\mathbf{y}_{\varepsilon\tau T+1}|X_{\delta\tau T+1}, Y_{\varepsilon\tau T}), \tag{10}$$

with $\delta, \varepsilon \in [0, \frac{1}{\tau})$, again assuming these limits exist. Notice that here we treat $\delta\tau T, \varepsilon\tau T$ again as integers to avoid excessive use of the ceiling or floor notation. We drop the subscripts $\varepsilon$ and $\delta$ in $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ and $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\delta+})$ when $\varepsilon = 1$ or $\delta = 1$. For example, $\mathcal{H}'(\mathcal{Y}|\mathcal{X}) = \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)|_{\delta=1, \varepsilon=1}$ and $\mathcal{H}'(\mathcal{Y}|\mathcal{X}_+) = \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\delta+})|_{\delta=1, \varepsilon=1}$. With (9) and (10), we further make the following assumptions:

$$\text{A2:} \qquad \mathcal{H}'(\mathcal{Y}|\mathcal{X}_+) = \lim_{\varepsilon \searrow 1} \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+}), \tag{11}$$

$$\mathcal{H}'(\mathcal{Y}|\mathcal{X}) = \lim_{\varepsilon \searrow 1} \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}). \tag{12}$$

Assumptions A1–A2 in (6), (7), (11), and (12), are important for the main result (Theorem 1). A1 is often met in practice for a memoryless and time-invariant system with *iid* input in the data phase, independent of the input and output during training. However, we do not have a complete characterization of the processes $\mathcal{X}$ and $\mathcal{Y}$ that meet Assumptions A2. Nevertheless, A2 may be verified on a case-by-case basis by examining expressions of $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+})$ and $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X})$ with $\varepsilon \geq 1$ using Corollary 1(c) in Appendix A; see the Example.

Define

$$\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \lim_{T \to \infty} \frac{1}{N\tau T} H(Y_{\varepsilon\tau T}|X_{\delta\tau T}), \tag{13}$$

$$\mathcal{I}'(\mathcal{X}_\varepsilon; \mathcal{Y}_\varepsilon) = \lim_{T \to \infty} \frac{1}{N} I(\mathbf{x}_{\varepsilon\tau T+1}; \mathbf{y}_{\varepsilon\tau T+1} | X_{\varepsilon\tau T}, Y_{\varepsilon\tau T}), \tag{14}$$

$$\mathcal{I}(\mathcal{X}_\varepsilon; \mathcal{Y}_\varepsilon) = \lim_{T \to \infty} \frac{1}{NT} I(X^{\varepsilon\tau T+1}; Y^{\varepsilon\tau T+1} | X_{\varepsilon\tau T}, Y_{\varepsilon\tau T}), \tag{15}$$

where $X^t = [\mathbf{x}_t, \mathbf{x}_{t+1}, \cdots, \mathbf{x}_T]$, and $Y^t = [\mathbf{y}_t, \mathbf{y}_{t+1}, \cdots, \mathbf{y}_T]$. Similarly, we drop the subscripts $\varepsilon$ and $\delta$ in $\mathcal{I}'(\mathcal{X}_\varepsilon; \mathcal{Y}_\varepsilon)$, $\mathcal{I}(\mathcal{X}_\varepsilon; \mathcal{Y}_\varepsilon)$, and $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ when $\varepsilon = 1$ or $\delta = 1$.

*B. Main result*

**Theorem 1.** *Under Assumption A1,*

$$\mathcal{I}(\mathcal{X}_0; \mathcal{Y}_0) \geq (1 - \tau)\mathcal{I}'(\mathcal{X}; \mathcal{Y}) \tag{16}$$

$$\geq (1 - \tau) \lim_{T \to \infty} I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | \hat{\mathbf{g}}_T), \tag{17}$$

*where $\hat{\mathbf{g}}_T$ is any estimate of $\mathbf{g}_T$ that is a function of $(X_{\tau T}, Y_{\tau T})$. When A2 is also met,*

$$\mathcal{I}'(\mathcal{X}; \mathcal{Y}) = \lim_{\varepsilon \searrow 1} \frac{\partial \mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X})}{\partial \varepsilon} - \lim_{\varepsilon \searrow 1} \frac{\partial \mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\varepsilon)}{\partial \varepsilon}. \tag{18}$$

*Proof.* Please see Appendix A. □

Hence, the mutual information limit with one-shot learning $(1 - \tau)\mathcal{I}'(\mathcal{X}; \mathcal{Y})$ is a lower bound of the mutual information without any training, and is an upper bound of the mutual information with any estimate of the unknown parameters. The expression in (16) can be calculated as a derivative using (18) as long as $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ is available. The next section shows how $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ may be computed, and the section that follows provides operational significance to (16) in the form of a channel coding theorem. An example application of the theorem appears in Section III.

*C. Computation of $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$*

An expression for $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ may be derived from $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X})$ when this latter quantity is available. In some cases $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X})$ can be obtained through methods employed in statistical mechanics by treating the conditional entropy as free energy in a large-scale system. Free energy is a fundamental quantity [10], [11] that has been analyzed through the powerful "replica

method", and this, in turn, has been applied to entropy calculations in machine learning [12]–[15] and wireless communications [16]–[18], in both linear and nonlinear systems.

The entropy $\mathcal{H}(\mathcal{Y}|\mathcal{X})$ (equivalent to $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ where $\varepsilon = 1$) is considered in [12]–[15], where the input is multiplied by an unknown vector as an inner product and then passes through a nonlinearity to generate a scalar output. In [12], [13], [15], the inputs are *iid*, while orthogonal inputs are considered in [14]. The entropy $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ for MIMO systems is considered in [16]–[18], where the inputs are *iid* in the training phase and are *iid* in the data phase, but the distributions in the two phases can differ. In [16], a linear system is considered where the output is the result of the input multiplied by an unknown matrix, plus additive noise, while in [17], [18] uniform quantization is added at the output.

As we now show, the expression for $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ for $\varepsilon \geq 1$ can be leveraged to compute $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ for all $\varepsilon, \delta > 0$. We consider the case when the input $\mathbf{x}_t$ are *iid* for all $t$, and the distribution set $\mathcal{P}(T, \tau)$ defined in (8) can therefore be simplified as

$$\mathcal{P}(T, \tau) = \{p(\mathbf{y}|\mathbf{x}; \mathbf{g}_T), p(\mathbf{g}_T), p(\mathbf{x})\}. \tag{19}$$

The following theorem assumes that we have $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ available as a function of $(\tau, \varepsilon)$ for all $\varepsilon \geq 1$.

**Theorem 2.** *Assume that Assumption A1 is met, $\mathbf{x}_t$ are iid for all $t$, $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ exists and is continuous in $\tau$ and $\varepsilon$ for $\tau \in (0, 1)$ and $\varepsilon \in (0, \frac{1}{\tau}]$. Define*

$$F(\tau, \varepsilon) = \mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}), \tag{20}$$

*where $\varepsilon \geq 1$ and $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ is defined in (13). Then*

$$\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = u \cdot F\left(u\tau, \frac{\varepsilon - u}{\delta} + 1\right), \tag{21}$$

*for all $\varepsilon, \delta \in (0, \frac{1}{\tau}]$, where $u = \min(\varepsilon, \delta)$.*

*Proof.* According to (13) and (20), we have

$$F(\tau, \varepsilon) = \lim_{T \to \infty} \frac{1}{N\tau T} H(Y_{\varepsilon\tau T}|X_{\tau T}),$$

7

which is computed using $\mathcal{P}(T, \tau)$ defined in (19). When $\delta \geq \varepsilon > 0$, we have

$$\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta) = \lim_{T \to \infty} \frac{1}{N\tau T} H(Y_{\varepsilon\tau T} | X_{\delta\tau T}) \tag{22}$$

$$= \lim_{T \to \infty} \frac{1}{N\tau T} H(Y_{\varepsilon\tau T} | X_{\varepsilon\tau T}) = \lim_{T \to \infty} \frac{\varepsilon}{N\tilde{\tau} T} H(Y_{\tilde{\tau} T} | X_{\tilde{\tau} T}). \tag{23}$$

where $\tilde{\tau} = \varepsilon\tau$. Therefore, (20) and (23) yield

$$\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta) = \varepsilon \cdot F(\tilde{\tau}, 1) = \varepsilon \cdot F(\varepsilon\tau, 1). \tag{24}$$

When $\varepsilon > \delta > 0$, let $\tilde{\tau} = \delta\tau$, and then (22) yields

$$\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta) = \lim_{T \to \infty} \frac{\delta}{N\tilde{\tau} T} H(Y_{\varepsilon\tilde{\tau} T/\delta} | X_{\tilde{\tau} T})$$

$$= \delta \cdot F\left(\tilde{\tau}, \frac{\varepsilon}{\delta}\right) = \delta \cdot F\left(\delta\tau, \frac{\varepsilon}{\delta}\right). \tag{25}$$

By combining (24) and (25), we obtain (21). □

### D. Channel coding theorem

We now provide an operational description of the mutual information inequality (16). We consider a communication system where the channel is constant for blocklength $T$, and then changes independently and stays constant for another blocklength, and so on. The first $\tau T$ symbols of each block are used for training with known input and output. Under Assumption A1, the communication system is memoryless, is time-invariant within each block, and the input is *iid* independent of $X_{\tau T}$ after training. The system is retrained with every block, and the message to be transmitted is encoded over the data phase of multiple blocks.

A $(2^{nRNT}, n, T)$-code for a block-constant channel with blocklength $T$ is defined as an encoder that maps a message $S \in \{1, 2, \ldots, 2^{nRNT}\}$ to the input in the data phase $X^{\tau T+1}$ among $n$ blocks, and a decoder that maps $X_{\tau T}$, and the entire output $Y_T$ for $n$ blocks to $\hat{S} \in \{1, 2, \ldots, 2^{nRNT}\}$, where $N = \frac{\alpha T}{\beta}$. The code rate $R$ has units "bits per transmission per receiver", and the maximum probability of error of the code is defined as

$$P_{\mathrm{e}}(n, T) = \max_S P(\hat{S} \neq S). \tag{26}$$

8

The channel coding theorem is shown below.

**Theorem 3.** *Assume A1 is met, with a channel that is constant with blocklength $T$, whose conditional distribution is parameterized by $\mathbf{g}_T$ and is independent of the input. If $\mathcal{I}'(\mathcal{X}; \mathcal{Y})$ exists, then for every $R$ that satisfies*

$$R < (1 - \tau)\mathcal{I}'(\mathcal{X}; \mathcal{Y}),$$

*there exists $T_0 > 0$, so that for all $T > T_0$, we can find a code $(2^{nRNT}, n, T)$ with maximum probability of error $P_e(n, T) \to 0$ as $n \to \infty$.*

*Proof.* Define

$$\mathcal{R}_T = \frac{1}{TN} I(X^{\tau T+1}; Y_T | X_{\tau T}).$$

For any finite $T$, according to the classical channel coding theorem [19]–[21], for every $R < \mathcal{R}_T$, there exists a code $(2^{nRNT}, n, T)$ with maximum probability of error $P_e(n, T) \to 0$ as $n \to \infty$.

It is clear that $X^{\tau T+1}$ is independent of $(X_{\tau T}, Y_{\tau T})$. Therefore, we have

$$\mathcal{R}_T = \frac{1}{TN} I(X^{\tau T+1}; Y^{\tau T+1} | X_{\tau T}, Y_{\tau T}).$$

Since $\mathbf{x}_{\tau T+1}, \mathbf{x}_{\tau T+2}, \ldots, \mathbf{x}_T$ are *iid*, and $p(\mathbf{y}_t | \mathbf{x}_t; \mathbf{g}_T)$ is a fixed conditional distribution for all $t = 1, 2, \ldots, T$, we have

$$\mathcal{R}_T \geq \frac{(1 - \tau)}{N} I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | X_{\tau T}, Y_{\tau T}). \tag{27}$$

According to the definition in (14),

$$\mathcal{I}'(\mathcal{X}; \mathcal{Y}) = \lim_{T \to \infty} \frac{1}{N} I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | X_{\tau T}, Y_{\tau T}).$$

Therefore, for any $\kappa > 0$, there exists a number $T_0 > 0$ so that when $T > T_0$, we have

$$\frac{1}{N} I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | X_{\tau T}, Y_{\tau T}) > \mathcal{I}'(\mathcal{X}; \mathcal{Y}) - \kappa,$$

9

and (27) yields

$$\mathcal{R}_T > (1-\tau)(\mathcal{I}'(\mathcal{X};\mathcal{Y}) - \kappa),$$

which means any rate $R \le (1-\tau)(\mathcal{I}'(\mathcal{X};\mathcal{Y}) - \kappa)$ is achievable.

By taking the limit $\kappa \searrow 0$, we finish the proof. $\qquad\square$

This theorem shows that rates below $(1-\tau)\mathcal{I}'(\mathcal{X};\mathcal{Y})$ are achievable when $T$ is chosen large enough. Only an achievability statement is given here since $(1-\tau)\mathcal{I}'(\mathcal{X};\mathcal{Y})$ is a lower bound on $\mathcal{R}_T$ for large $T$.

## III. STEPS FOR COMPUTING OPTIMAL TRAINING TIME AND AN EXAMPLE

We summarize the process to compute the optimal training time $\tau_{\mathrm{opt}}$ for a memoryless, time-invariant system with unknown parameters. We assume that the input dimension $M$, the output dimension $N$, and the coherence time (block of symbols) $T$ have the ratios defined in (3). The unknown parameters of the system are constant within the block, and change independently in the next block. The first $\tau T$ symbols of each block are used for training and the remaining $T - \tau T$ are for data. We assume $T \to \infty$, and solve (5) as an approximation of (2). The input $\mathbf{x}_t$ are *iid* for all $t = 1, \ldots, T$.

The process includes the following seven steps:

1) Verify Assumption A1 (6)–(7) based on the set of distributions in (8).
2) Compute $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ defined in (13) for $\varepsilon \ge 1$ based on (8), and express it as a function of $\tau$ and $\varepsilon$ as $F(\tau, \varepsilon)$ (20).
3) Compute $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ defined in (13), for all $\varepsilon, \delta \in (0, \frac{1}{\tau}]$ by using Theorem 2 and $F(\tau, \varepsilon)$.
4) Compute $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ and $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+})$ defined in (9)-(10) by taking the derivative of $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ and $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\varepsilon)$ (Corollary 1(a)–(b) in Appendix).
5) Verify Assumption A2 (11) by examining the expressions of $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+})$ and verify (12) with Corollary 1(c).
6) Compute $\mathcal{I}'(\mathcal{X};\mathcal{Y})$ by using (18).
7) Solve $\tau_{\mathrm{opt}}$ using (5).

The following simple example applies these steps.

*Example: Bit flipping through random channels*

Let

$$y_t = x_t \oplus g_{k_t}, \quad t = 1, \ldots, T, \tag{28}$$

where, since $M = N = 1$, the binary input $x_t$ and output $y_t$ are scalars, and $x_t$ is XOR'ed with a random bit $g_{k_t}$ intended to model the unknown "state" of the channel $k_t$. Thus, each channel either lets the input bit directly through, or inverts it. The $x_t$ are *iid* equally likely to be zero or one, Bernoulli($\frac{1}{2}$) random variables. Let $a > 0$ be a parameter, where $a \cdot T$ is the (integer) number of unique channels whose states are stored in the vector $\mathbf{g}_T = [g_1, g_2, \cdots, g_{a \cdot T}]^\mathsf{T}$ comprising *iid* Bernoulli($\frac{1}{2}$) random variables that are independent of the input. The channel selections $\mathbf{k}_T = [k_1, \ldots, k_T]$ are chosen as an *iid* uniform sample from $\{1, 2, \cdots, a \cdot T\}$ (with possible repetitions), and the choices are known to the receiver. We wish to send training signals through these channels to learn $\mathbf{g}_T$; the more entries of this vector that we learn, the more channels become useful for sending data, but the less time we have to send data before the blocklength $T$ runs out and $\mathbf{g}_T$ changes. We want to determine the optimum $\tau$ as $T \to \infty$ using (5). We therefore follow the steps above.

1) From (28), we have

$$p(\mathbf{y}_T | \mathbf{x}_T; \mathbf{g}_T) = \prod_{t=1}^{T} p(y_t | x_t; \mathbf{g}_T),$$

where $p(y_t | x_t; \mathbf{g}_T) = \mathbb{1}_{(y_t = x_t \oplus g_{k_t})}$ for all $t$. Here the notation is slightly abused, since now $M = N = 1$, we use $\mathbf{y}_t$ and $\mathbf{x}_t$ to denote $[y_1, \ldots, y_t]^\mathsf{T}$ and $[x_1, \ldots, x_t]^\mathsf{T}$. It is clear that Assumption A1 is met and $x_t$ are *iid* for all $t$ independent of $\mathbf{g}_T$.

2) By definition, $\mathcal{H}(\mathcal{Y}|\mathcal{X}) = \lim\limits_{T \to \infty} \frac{1}{\tau T} H(\mathbf{y}_{\tau T} | \mathbf{x}_{\tau T})$. The model (28) yields

$$H(\mathbf{y}_{\tau T} | \mathbf{x}_{\tau T}) \overset{(a)}{=} H(\{g_{k_1}, \ldots, g_{k_{\tau T}}\} | \mathbf{x}_{\tau T}) \overset{(b)}{=} H(\{g_{k_1}, \ldots, g_{k_{\tau T}}\}) \overset{(c)}{=} \mathbb{E}_{\mathbf{k}_{\tau T}} |A_{\tau T}|$$

$$\overset{(d)}{=} \sum_{i=1}^{aT} \mathbb{E}(\mathbb{1}_{(i \in A_{\tau T})}) = aT(1 - (1 - \frac{1}{aT})^{\tau T}).$$

where $A_{\tau T} = \{k_1, \ldots, k_{\tau T}\}$, $^{(a)}$ uses $g_{k_t} = y_t \oplus x_t$, $^{(b)}$ uses the independence between $\mathbf{x}_{\tau T}$ and $g_t$, $^{(c)}$ uses the independence between $g_t$, $g_k$ when $t \neq k$ and $^{(d)}$ uses $|A_{\tau T}| =$

$\sum_{i=1}^{aT} \mathbb{1}_{(i \in A_{\tau T})}$, where $\mathbb{1}_{(\cdot)}$ is the indicator function. Therefore,

$$\mathcal{H}(\mathcal{Y}|\mathcal{X}) = \lim_{T \to \infty} \frac{a}{\tau}(1 - (1 - \frac{1}{aT})^{\tau T}) = \frac{a}{\tau}(1 - e^{-\frac{\tau}{a}}). \tag{29}$$

By the chain rule for entropy, for $\varepsilon > 1$, we have

$$\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}) = \mathcal{H}(\mathcal{Y}|\mathcal{X}) + \lim_{T \to \infty} \frac{1}{\tau T} H(y_{\tau T+1}, y_{\tau T+2}, \ldots, y_{\varepsilon \tau T}|\mathbf{x}_{\tau T}, \mathbf{y}_{\tau T}).$$

Since

$$\varepsilon \tau T - \tau T \geq H(y_{\tau T+1}, y_{\tau T+2}, \ldots, y_{\varepsilon \tau T}|\mathbf{x}_{\tau T}, \mathbf{y}_{\tau T})$$

$$\geq H(x_{\tau T+1}, x_{\tau T+2}, \ldots, x_{\varepsilon \tau T}|\mathbf{x}_{\tau T}, \mathbf{y}_{\tau T}, \mathbf{g}_T)$$

$$= H(x_{\tau T+1}, x_{\tau T+2}, \ldots, x_{\varepsilon \tau T}) = \varepsilon \tau T - \tau T,$$

we conclude that

$$F(\tau, \varepsilon) = \mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}) = \frac{a}{\tau}(1 - e^{-\frac{\tau}{a}}) + \varepsilon - 1. \tag{30}$$

3) Theorem 2 yields

$$\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \begin{cases} \frac{a}{\tau}(1 - e^{-\frac{\tau}{a}\varepsilon}), & \varepsilon \leq \delta; \\ \frac{a}{\tau}(1 - e^{-\frac{\tau}{a}\delta}) + (\varepsilon - \delta), & \delta < \varepsilon, \end{cases}$$

for $\varepsilon, \delta \in (0, \frac{1}{\tau})$.

4) Then, Corollary 1(a)-(b) yields

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+}) = \frac{\partial \mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\varepsilon)}{\partial \varepsilon} = e^{-\frac{\tau}{a}\varepsilon},$$

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \frac{\partial \mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)}{\partial \varepsilon} = \begin{cases} e^{-\frac{\tau}{a}\varepsilon}, & \varepsilon < \delta; \\ 1, & \varepsilon > \delta. \end{cases}$$

5) $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+})$ and Corollary 1(c) allow us to conclude that Assumption A2 also holds.

6) From Theorem 1, we obtain

$$\mathcal{I}'(\mathcal{X};\mathcal{Y}) = 1 - e^{-\frac{\tau}{a}},$$

7) Finally, (5) yields

$$\tau_{\text{opt}} = \underset{\tau}{\text{argmax}}(1 - \tau)(1 - e^{-\frac{\tau}{a}}), \tag{31}$$

or

$$\tau_{\text{opt}} = \begin{cases} -a \ln a, & a \to 0; \\ \frac{1}{2}, & a \to \infty; \\ \frac{1}{e}, & a = \frac{1}{e}. \end{cases}$$

When $a$ is small, $\tau_{\text{opt}}$ is larger than $a$; when $a$ is large, $\tau_{\text{opt}}$ saturates at $\frac{1}{2}$; and $a = \frac{1}{e}$ is the dividing line between $\tau_{\text{opt}} > a$ and $\tau_{\text{opt}} < a$. The corresponding rates are

$$\mathcal{R}_{\text{opt}} = \begin{cases} (1 + a \ln a)(1 - a), & a \to 0; \\ \frac{1}{2}(1 - e^{-\frac{1}{2a}}), & a \to \infty; \\ (1 - \frac{1}{e})^2, & a = \frac{1}{e}. \end{cases}$$

The optimum fraction of the blocklength $T$ that should be devoted to training varies as a function of the number of possible unique channels. When $a = 1$, the number of unique channels equals $T$, and the $\tau_{\text{opt}} \approx 0.44$. For a large number of unique channels relative to the blocklength ($a \to \infty$), the fraction of the training time saturates at $1/2$. When $a$ is small, the optimum fraction of the blocklength devoted to training decreases to zero, but more slowly than $a$.

In this example, a traditional finite-system information-theoretic analysis and simulation is possible (these calculations are omitted). Figure 1 shows the results as a function of $T$, where we can see that as $T$ grows, the resulting $\tau_{\text{opt}}$ quickly approaches the large-system results. The fact that we can use a large-system limit to approximate a finite-system limit is important when applying the Theorems in realistic scenarios.
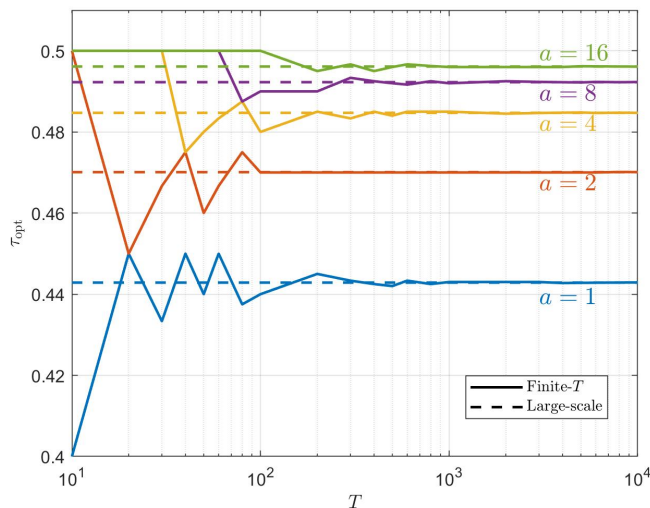
Fig. 1. Optimal training time $\tau_{\text{opt}}$ vs blocklength $T$, derived from finite-$T$ analysis (solid curves) and the proposed large-scale method (dashed curves). The large-scale analysis (31) shows excellent agreement with the finite-dimensional analysis (which is omitted) for even small values of $T$.

## IV. DISCUSSION AND CONCLUSION

### A. Number of unknowns and bilinear model

In general, a finite number of unknowns in the model leads to uninteresting results as $T \to \infty$. For example, consider a system modeled as

$$y_t = g x_t + v_t, \quad t = 1, 2, \ldots$$

where $g$ is the unknown gain of the system, $x_t$, $y_t$ are the input and corresponding output, $v_t$ is the additive noise, $\tau$ is the fraction of time used for training. This system is bilinear in the gain and the input. We assume that $v_t$ is modeled as *iid* Gaussian $\mathcal{N}(0, 1)$, independent of the input. The training signals are $x_t = 1$ for all $t = 1, 2, \ldots, \tau T$, and the data signals $x_t$ are modeled as *iid* Gaussian $\mathcal{N}(0, 1)$ for all $t = \tau T + 1, \tau T + 2, \ldots$ An analysis similar to the Example produces

$$\mathcal{I}(\mathcal{X}_0; \mathcal{Y}_0) \geq \frac{1 - \tau}{2} \mathbb{E}_g \log(1 + g^2),$$

and therefore $\tau_{\text{opt}} = 0$ maximizes this bound. This result reflects the fact that $g$ is learned perfectly for any $\tau > 0$ because there is only one unknown parameter for $\tau T$ training symbols

14

as $\tau T \to \infty$. Hence, trivially, it is advantageous to make $\tau$ as small as possible.

More interesting is the "large-scale" model

$$\mathbf{y}_t = \mathbf{f}(G\mathbf{x}_t + \mathbf{v}_t), \quad t = 1, 2, \ldots, \tag{32}$$

where $\mathbf{x}_t$ and $\mathbf{y}_t$ are the $t$th input and output vectors with dimension $M$ and $N$, $G$ is an $N \times M$ unknown random matrix that is not a function of $t$, $\mathbf{v}_1, \mathbf{v}_2, \ldots$ are *iid* unknown vectors with dimension $N$ and known distribution (not necessarily Gaussian), and $\mathbf{f}(\cdot)$ applies a possibly nonlinear function $f(\cdot)$ to each element of its input. The training interval $\tau T$ is used to learn $G$.

Let $M$ and $N$ increase proportionally to the blocklength $T$ with the ratios defined in (3); such a model can be used in large-scale wireless communication, signal processing, and machine learning applications. In wireless communication and signal processing [1], [3], [4], [7], [8], [16]–[18], [22], $\mathbf{x}_t$ and $\mathbf{y}_t$ are the transmitted signal and the received signal at time $t$ in a multiple-input-multiple-output (MIMO) system with $M$ transmitters and $N$ receivers, $G$ models the channel coefficients between the transmitters and receivers, $T$ is the coherence time during which the channel $G$ is constant, $\mathbf{v}_t$ is the additive noise at time $t$, $f(\cdot)$ models receiver effects such as quantization in analog-to-digital converters and nonlinearities in amplifiers. A linear receiver, $f(x) = x$, is considered in [3], [4], [16], [22]. Single-bit ADC's with $f(x) = \text{sign}(x)$ are considered in [7], [8], and low-resolution ADC's with $f(x)$ modeled as a uniform quantizer are considered in [17], [18]. The training and data signals can be chosen from different distributions, as in [1], [7], [8]. Conversely, the training and data signals can both be *iid*, as in [4], [16]–[18].

Let $N = 1$. In machine learning, (32) is a model of a single layer neural network (perceptron) [12]–[14] and $\mathbf{x}_t$ is the input to the perceptron with dimension $M$, $\mathbf{y}_t$ is the scalar decision variable at time $t$, $G$ holds the unknown weights of the perceptron, and $f(\cdot)$ is the nonlinear activation function. A perceptron is often used as a classifier, where the output of the perceptron is the class label of the corresponding input. In [12], [13], *iid* inputs are used to learn the weights, and orthogonal inputs are used in [14]. Binary class classifiers are considered in [12]–[14]. Training employs $\tau T$ labeled input-output pairs $(\mathbf{x}_t, \mathbf{y}_t)$, and the trained perceptron then classifies new inputs before it is retrained on a new dataset. Generally, both the training and data are modeled as having the same distribution.

To obtain optimal training results for (32), Theorems 1–2 show that a starting point for computing $\mathcal{I}'(\mathcal{X};\mathcal{Y})$ is $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ for $\varepsilon \geq 1$. Fortunately, $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ results can sometimes be found in the existing literature; for example, in [16]–[18], $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ is used to calculate the mean-square error of the estimated input signal, conditioned on the training. Our analysis shows how to leverage these same $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})$ results to derive the training-based mutual information.

## B. Models for which assumptions are superfluous

Assumption A2 is likely to be superfluous for certain common system models, such as when the distribution on $\mathbf{x}_t$ is *iid* through the training and data phases, and the transition probabilities can be written as a product as in Assumption A1. However, we have not yet characterized for which models A2 is automatically satisfied without additional assumptions on $\mathcal{H}'$, and think that this would be an interesting research topic for further work.

<div align="center">

APPENDIX A

PROOF OF THEOREM 1

</div>

### A. Proof of the inequalities (16)–(17)

Under Assumption A1, we have

$$I(X^{\tau T+1}; Y_T|X_{\tau T})$$
$$=I(X^{\tau T+1}; Y_{\tau T}|X_{\tau T}) + I(X^{\tau T+1}; Y^{\tau T+1}|X_{\tau T}, Y_{\tau T})$$
$$=I(X^{\tau T+1}; Y^{\tau T+1}|X_{\tau T}, Y_{\tau T}),$$

where the first equality uses the chain rule and the second uses that $X^{\tau T+1}$ is independent of $(X_{\tau T}, Y_{\tau T})$. Moreover,

$$I(X^{\tau T+1}; Y^{\tau T+1}|X_{\tau T}, Y_{\tau T})$$
$$= H(X^{\tau T+1}|X_{\tau T}, Y_{\tau T}) - H(X^{\tau T+1}|X_{\tau T}, Y_T)$$
$$\overset{(a)}{=} \sum_{t=\tau T+1}^{T} \left( H(\mathbf{x}_t|X_{t-1}, Y_{\tau T}) - H(\mathbf{x}_t|X_{t-1}, Y_T) \right)$$

<div align="center">16</div>

$$\overset{(b)}{\geq} \sum_{t=\tau T+1}^{T} \left( H(\mathbf{x}_t | X_{t-1}, Y_{t-1}) - H(\mathbf{x}_t | X_{t-1}, Y_t) \right)$$

$$= \sum_{t=\tau T+1}^{T} I(\mathbf{x}_t; \mathbf{y}_t | X_{t-1}, Y_{t-1}).$$

Here, $^{(a)}$ uses the chain rule, $^{(b)}$ uses conditioning to reduce entropy. Equality in $^{(b)}$ can be achieved when $\mathbf{g}_T$ is estimated perfectly from $(X_{\tau T}, Y_{\tau T})$. Since Assumption A1 implies that $\mathbf{x}_t$ is independent of $(\mathbf{x}_k, \mathbf{y}_k)$ when $k \neq t$ and $t \geq \tau T + 1$, for all $t \geq \tau T + 1$, we have

$$I(\mathbf{x}_{t+1}; \mathbf{y}_{t+1} | X_t, Y_t) - I(\mathbf{x}_t; \mathbf{y}_t | X_{t-1}, Y_{t-1})$$

$$\overset{(a)}{=} (H(\mathbf{x}_{t+1}) - H(\mathbf{x}_{t+1} | X_t, Y_{t+1}))$$

$$- (H(\mathbf{x}_t) - H(\mathbf{x}_t | X_{t-1}, Y_t))$$

$$= H(\mathbf{x}_t | X_{t-1}, Y_t) - H(\mathbf{x}_{t+1} | X_t, Y_{t+1})$$

$$\overset{(b)}{=} H(\mathbf{x}_{t+1} | X_{t-1}, Y_{t-1}, \mathbf{y}_{t+1}) - H(\mathbf{x}_{t+1} | X_t, Y_{t+1}) \overset{(c)}{\geq} 0.$$

Here, $^{(a)}$ uses the independence between $\mathbf{x}_t$ and $(X_{t-1}, Y_{t-1})$, $^{(b)}$ uses Assumption A1, $^{(c)}$ uses conditioning to reduce entropy. Thus, $I(\mathbf{x}_t; \mathbf{y}_t | X_{t-1}, Y_{t-1})$ is monotonically increasing with $t$ for all $t > \tau T$. Then, in the limit when $T \to \infty$, we get (16).

Also,

$$I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | X_{\tau T}, Y_{\tau T})$$

$$= H(\mathbf{x}_{\tau T+1}) - H(\mathbf{x}_{\tau T+1} | X_{\tau T}, Y_{\tau T}, \mathbf{y}_{\tau T+1})$$

$$\overset{(a)}{=} H(\mathbf{x}_{\tau T+1}) - H(\mathbf{x}_{\tau T+1} | X_{\tau T}, Y_{\tau T}, \hat{\mathbf{g}}_T, \mathbf{y}_{\tau T+1})$$

$$\overset{(b)}{\geq} H(\mathbf{x}_{\tau T+1}) - H(\mathbf{x}_{\tau T+1} | \hat{\mathbf{g}}_T, \mathbf{y}_{\tau T+1})$$

$$= I(\mathbf{x}_{\tau T+1}; \mathbf{y}_{\tau T+1} | \hat{\mathbf{g}}_T),$$

where $^{(a)}$ uses that $\hat{\mathbf{g}}_T$ is a function of $(X_{\tau T}, Y_{\tau T})$, and $^{(b)}$ uses conditioning to reduce entropy. By taking the limit $T \to \infty$, we have (17).

*B. Proof of* (18)

We first show the derivative relationship between $\mathcal{H}'(\mathcal{Y}_\varepsilon)$ and $\mathcal{H}(\mathcal{Y}_\varepsilon)$ defined below, and then generalize to the conditional entropies which directly lead to the conclusion (18). Define

$$\mathcal{H}'(\mathcal{Y}_\varepsilon) = \lim_{T \to \infty} H(\mathbf{y}_{\lceil \varepsilon \tau T \rceil + 1} | Y_{\lceil \varepsilon \tau T \rceil}), \tag{33}$$

$$\mathcal{H}(\mathcal{Y}_\varepsilon) = \lim_{T \to \infty} \frac{1}{\tau T} H(Y_{\lceil \varepsilon \tau T \rceil}), \tag{34}$$

which can be considered as $\mathcal{H}'(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ and $\mathcal{H}(\mathcal{Y}_\varepsilon | \mathcal{X}_\delta)$ with $\delta = 0$. For mathematical rigorousness, we keep the $\lceil \cdot \rceil$ notation here. We show that, under some conditions, $\mathcal{H}'(\mathcal{Y}_\varepsilon)$ is the derivative of $\mathcal{H}(\mathcal{Y}_\varepsilon)$.

**Theorem 4.** *Suppose there exists a $\kappa > 0$ so that $H(\mathbf{y}_{t+1} | Y_t)$ is monotonic in $t$ when $t \in [\lfloor (\varepsilon - \kappa) \tau T \rfloor, \lceil (\varepsilon + \kappa) \tau T \rceil]$ as $T \to \infty$.*

*If $\mathcal{H}(\mathcal{Y}_\varepsilon)$ and its derivative with respect to $\varepsilon$ exist, we have*

$$\mathcal{H}'(\mathcal{Y}_\varepsilon) = \frac{\partial \mathcal{H}(\mathcal{Y}_\varepsilon)}{\partial \varepsilon}. \tag{35}$$

*If both $\mathcal{H}(\mathcal{Y}_\varepsilon)$ and $\mathcal{H}'(\mathcal{Y}_\varepsilon)$ exist, and there exists a $c > 0$ independent of $t$ and $\tau T$ so that $|H(\mathbf{y}_{t+1} | Y_t)| < c$, we have*

$$\mathcal{H}(\mathcal{Y}_\varepsilon) = \int_0^\varepsilon \mathcal{H}'(\mathcal{Y}_u) du. \tag{36}$$

*Proof.* Equation (36) is an integral equivalent of (35) and we only prove (35) for simplicity. Without loss of generality, we assume that $H(\mathbf{y}_{t+1} | Y_t)$ is monotonically decreasing. Using the definition of $\mathcal{H}(\mathcal{Y}_\varepsilon)$ in (34), we have

$$\frac{1}{\kappa} (\mathcal{H}(\mathcal{Y}_{\varepsilon + \kappa}) - \mathcal{H}(\mathcal{Y}_\varepsilon))$$

$$= \lim_{T \to \infty} \frac{H(Y_{\lceil (\varepsilon + \kappa) \tau T \rceil}) - H(Y_{\lceil \varepsilon \tau T \rceil})}{\kappa \tau T}$$

$$= \lim_{T \to \infty} \frac{\sum_{t = \lceil \varepsilon \tau T \rceil + 1}^{\lceil (\varepsilon + \kappa) \tau T \rceil} H(\mathbf{y}_t | Y_{t-1})}{\kappa \tau T}$$

$$\le \lim_{T \to \infty} \frac{(\kappa \tau T + 1) \cdot H(\mathbf{y}_{\lceil \varepsilon \tau T \rceil + 1} | Y_{\lceil \varepsilon \tau T \rceil})}{\kappa \tau T}$$

18

$$= \lim_{T\to\infty} H(\mathbf{y}_{\lceil\varepsilon\tau T\rceil+1}|Y_{\lceil\varepsilon\tau T\rceil}). \tag{37}$$

Similarly to (37), we also have

$$\lim_{T\to\infty} H(\mathbf{y}_{\lceil\varepsilon\tau T\rceil+1}|Y_{\lceil\varepsilon\tau T\rceil}) \leq \frac{1}{\kappa}(\mathcal{H}(\mathcal{Y}_\varepsilon) - \mathcal{H}(\mathcal{Y}_{\varepsilon-\kappa})). \tag{38}$$

Let $\kappa \searrow 0$ in both (37) and (38); because we assume that the derivative of $\mathcal{H}(\mathcal{Y}_\varepsilon)$ exists, these limits both equal this derivative. Then, the definition of $\mathcal{H}'(\mathcal{Y}_\varepsilon)$ in (33) yields (35). $\square$

Theorem 4 is a consequence of the entropy chain rule and letting an infinite sum converge to an integral (standard Riemann sum approximation). Such an analysis has also been used in the context of computing mutual information; for example [23]–[27]. Theorem 4 can be generalized to include conditioning on $\mathcal{X}$, thus leading to the following corollary, provided that $\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)$ and its derivative with respect to $\varepsilon$ exist.

**Corollary 1.** *Assume A1 holds. (a) For $\varepsilon > 1$,*

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}) = \frac{\partial\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X})}{\partial\varepsilon}, \tag{39}$$

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+}) = \frac{\partial\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\varepsilon)}{\partial\varepsilon}. \tag{40}$$

*(b) If $\mathbf{x}_t$ are iid for all $t$, then for all $\varepsilon, \delta > 0$ and $\varepsilon \neq \delta$,*

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \frac{\partial\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)}{\partial\varepsilon}, \tag{41}$$

$$\mathcal{H}'(\mathcal{Y}|\mathcal{X}_+) = \frac{\partial\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\varepsilon)}{\partial\varepsilon}\bigg|_{\varepsilon=1}. \tag{42}$$

*(c) If*

$$\lim_{\varepsilon\searrow 1}\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}) = \lim_{\delta\nearrow 1}\mathcal{H}'(\mathcal{Y}|\mathcal{X}_\delta), \tag{43}$$

*then Assumption A2 (12) is met.*

*Proof.* (a) Under A1, for all $\delta \geq 1$, we have

$$H(\mathbf{y}_{t+1}|X_{\lceil\delta\tau T\rceil},Y_t) \leq H(\mathbf{y}_{t+1}|X_{\lceil\delta\tau T\rceil},Y_{t-1})$$

$$= H(\mathbf{y}_t|X_{\lceil\delta\tau T\rceil},Y_{t-1}), \tag{44}$$

when $\lceil\tau T\rceil + 1 \leq t \leq \lceil\delta\tau T\rceil - 1$ or $t \geq \lceil\delta\tau T\rceil + 1$. Here, we use that the input is *iid* and the system is memoryless and time invariant; the inequality follows from the fact that conditioning reduces entropy. Therefore, $\forall \kappa \in (0, \varepsilon - 1)$, $H(\mathbf{y}_{t+1}|X_{\lceil\tau T\rceil},Y_t)$ is monotonically decreasing in $t$ for $t \in [\lfloor(\varepsilon - \kappa)\tau T\rfloor, \lceil(\varepsilon + \kappa)\tau T\rceil]$ when $\tau T > \frac{2}{\varepsilon-1-\kappa}$. Then, Theorem 4 yields (39).

Also, $\forall \kappa \in (0, \varepsilon - 1), \delta > 2\varepsilon - 1$, $H(Y_{t+1}|X_{\lceil\delta\tau T\rceil},Y_t)$ is monotonically decreasing in $t$ for $t \in [\lfloor(\varepsilon - \kappa)\tau T\rfloor, \lceil(\varepsilon + \kappa)\tau T\rceil]$ when $\tau T > \max(\frac{2}{\varepsilon-1-\kappa}, \frac{2}{\delta-\varepsilon-\kappa})$. Then, Theorem 4 yields

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \frac{\partial\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta)}{\partial\varepsilon}. \tag{45}$$

Assumption A1 yields

$$\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+}), \tag{46}$$

where $\mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+})$ is defined in (10), and

$$\mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\delta) = \mathcal{H}(\mathcal{Y}_\varepsilon|\mathcal{X}_\varepsilon). \tag{47}$$

Therefore, (45) becomes (40).

(b) If $\mathbf{x}_t$ are *iid* for all $t$, then (44) is valid for all $t \leq \lceil\delta\tau T\rceil - 1$ or $t \geq \lceil\delta\tau T\rceil + 1$. Therefore, Theorem 4 yields (41). By taking $\varepsilon = 1$ and $\delta > 1$, (41), (46), and (47) then yield (42).

(c) For $t \geq \lceil\tau T\rceil + 1$, we have

$$H(\mathbf{y}_{t+1}|X_{\lceil\tau T\rceil}, Y_t) \leq H(\mathbf{y}_t|X_{\lceil\tau T\rceil}, Y_{t-1}).$$

Therefore,

$$\lim_{\varepsilon\searrow 1} \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}) \leq \mathcal{H}'(\mathcal{Y}|\mathcal{X}).$$

Conditioning to reduce entropy again yields

$$H(\mathbf{y}_{\lceil \tau T \rceil + 1} | X_{\lceil \tau T \rceil}, Y_{\lceil \tau T \rceil}) \leq H(\mathbf{y}_{\lceil \tau T \rceil + 1} | X_{\lceil \delta \tau T \rceil}, Y_{\lceil \tau T \rceil}),$$

for any $\delta < 1$ and therefore,

$$\mathcal{H}'(\mathcal{Y}|\mathcal{X}) \leq \lim_{\delta \nearrow 1} \mathcal{H}'(\mathcal{Y}|\mathcal{X}_\delta).$$

Equation (43) then implies A2 (12). □

Corollary 1 can now be used to finish the proof of Theorem 1. By (14), and Assumptions A1–2

$$\mathcal{I}'(\mathcal{X}; \mathcal{Y}) = \mathcal{H}'(\mathcal{Y}|\mathcal{X}) - \mathcal{H}'(\mathcal{Y}|\mathcal{X}_+) = \lim_{\varepsilon \searrow 1} \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}) - \lim_{\varepsilon \searrow 1} \mathcal{H}'(\mathcal{Y}_\varepsilon|\mathcal{X}_{\varepsilon+}),$$

and together with Corollary 1(a), we have (18).

## REFERENCES

[1] B. Hassibi and B. M. Hochwald, "How much training is needed in multiple-antenna wireless links?" *IEEE Trans. Inf. Theory*, vol. 49, no. 4, pp. 951–963, 2003.

[2] R. Muharar, "Optimal power allocation and training duration for uplink multiuser massive MIMO systems with MMSE receivers," *IEEE Access*, vol. 8, pp. 23 378–23 390, 2020.

[3] K. Takeuchi, M. Vehkapera, T. Tanaka, and R. R. Muller, "Large-system analysis of joint channel and data estimation for MIMO DS-CDMA systems," *IEEE Trans. Inf. Theory*, vol. 58, no. 3, pp. 1385–1412, 2012.

[4] K. Takeuchi, R. R. Müller, M. Vehkaperä, and T. Tanaka, "On an achievable rate of large Rayleigh block-fading MIMO channels with no CSI," *IEEE Trans. Inf. Theory*, vol. 59, no. 10, pp. 6517–6541, 2013.

[5] N. I. Miridakis and T. A. Tsiftsis, "On the joint impact of hardware impairments and imperfect CSI on successive decoding," *IEEE Trans. Veh. Technol.*, vol. 66, no. 6, pp. 4810–4822, 2016.

[6] Z. Sheng, H. D. Tuan, H. H. Nguyen, and M. Debbah, "Optimal training sequences for large-scale MIMO-OFDM systems," *IEEE Trans. Signal Process.*, vol. 65, no. 13, pp. 3329–3343, 2017.

[7] Y. Li, C. Tao, L. Liu, A. Mezghani, and A. L. Swindlehurst, "How much training is needed in one-bit massive MIMO systems at low SNR?" in *IEEE GLOBECOM*, Washington, D.C., USA, 2016, pp. 1–6.

[8] Y. Li, C. Tao, G. Seco-Granados, A. Mezghani, A. L. Swindlehurst, and L. Liu, "Channel estimation and performance analysis of one-bit massive MIMO systems," *IEEE Trans. Signal Process.*, vol. 65, no. 15, pp. 4075–4089, 2017.

[9] J. Bussgang, "Crosscorrelation functions of amplitude-distorted Gaussian signals," *MIT, Cambridge, MA, USA, Tech. Rep.*, 1952.

[10] T. Castellani and A. Cavagna, "Spin-glass theory for pedestrians," *J. Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 05, p. P05012, 2005.

[11] M. Mezard and A. Montanari, *Information, physics, and computation.* New York, NY, USA: Oxford University Press, 2009.

[12] A. Engel and C. Van den Broeck, *Statistical Mechanics of Learning*. Cambridge, U.K.: Cambridge University Press, 2001.

[13] M. Opper and W. Kinzel, "Statistical mechanics of generalization," in *Models of Neural Networks III,* Klaus Schulten, E. Domany, and J. Leo van Hemmen, Eds., New York, NY, USA: Springer, 1996, ch. 5, pp. 151–209.

[14] T. Shinzato and Y. Kabashima, "Learning from correlated patterns by simple perceptrons," *J. Physics A: Mathematical and Theor.*, vol. 42, no. 1, p. 015005, 2008.

[15] S. Ha, K. Kang, J.-H. Oh, C. Kwon, and Y. Park, "Generalization in a perceptron with a sigmoid transfer function," in *IEEE IJCNN*, vol. 2, Nagoya, Japan, 1993, pp. 1723–1726.

[16] C.-K. Wen, Y. Wu, K.-K. Wong, R. Schober, and P. Ting, "Performance limits of massive MIMO systems based on Bayes-optimal inference," in *IEEE ICC*, London, U.K., 2015, pp. 1783–1788.

[17] C.-K. Wen, S. Jin, K.-K. Wong, C.-J. Wang, and G. Wu, "Joint channel-and-data estimation for large-MIMO systems with low-precision ADCs," in *IEEE ISIT*, Hong Kong, 2015, pp. 1237–1241.

[18] C.-K. Wen, C.-J. Wang, S. Jin, K.-K. Wong, and P. Ting, "Bayes-optimal joint channel-and-data estimation for massive MIMO with low-precision ADCs," *IEEE Trans. Signal Process.*, vol. 64, no. 10, pp. 2541–2556, 2016.

[19] T. M. Cover and J. A. Thomas, *Elements of information theory*. Hoboken, NJ, USA: John Wiley & Sons, 2012.

[20] R. W. Yeung, *Information Theory and Network Coding*. New York, NY, USA: Springer Science & Business Media, 2008.

[21] M. Effros, A. Goldsmith, and Y. Liang, "Generalizing capacity: New definitions and capacity theorems for composite channels," *IEEE Trans. Inf. Theory*, vol. 56, no. 7, pp. 3069–3087, 2010.

[22] K. Takeuchi, R. R. Müller, M. Vehkaperä, and T. Tanaka, "An achievable rate of large block-fading MIMO systems with no CSI via successive decoding," in *IEEE ISITA*, Taichung, Taiwan, 2010, pp. 519–524.

[23] S. Shamai and S. Verdú, "The impact of frequency-flat fading on the spectral efficiency of CDMA," *IEEE Trans. Inf. Theory*, vol. 47, no. 4, pp. 1302–1327, 2001.

[24] D. Guo and S. Verdú, "Randomly spread CDMA: Asymptotics via statistical physics," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 1983–2010, 2005.

[25] D. Guo, S. Shamai, and S. Verdú, "Mutual information and minimum mean-square error in Gaussian channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.

[26] D. Guo and C.-C. Wang, "Multiuser detection of sparsely spread CDMA," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 3, pp. 421–431, 2008.

[27] M. L. Honig *et al.*, *Advances in multiuser detection*. Hoboken, NJ, USA: John Wiley & Sons, 2009.