

Control with Minimum Communication Cost per Symbol

Justin Pearson, João P. Hespanha, Daniel Liberzon

Abstract—We address the problem of stabilizing a continuous-time linear time-invariant process under communication constraints. We assume that the sensor that measures the state is connected to the actuator through a finite capacity communication channel over which an encoder at the sensor sends symbols from a finite alphabet to a decoder at the actuator. We consider a situation where one symbol from the alphabet consumes no communication resources, whereas each of the others consumes one unit of communication resources to transmit. This paper explores how the imposition of limits on an encoder’s bit-rate and average resource consumption affect the encoder/decoder/controller’s ability to keep the process bounded. The main result is a necessary and sufficient condition for a bounding encoder/decoder/controller which depends on the encoder’s average bit rate, its average resource consumption, and the unstable eigenvalues of the process.

I. INTRODUCTION

This paper addresses the problem of stabilizing a continuous-time linear time-invariant process under communication constraints. As in [1–7], we assume that the sensor that measures the state is connected to the actuator through a finite capacity communication channel. At each sampling time, an encoder sends a symbol through the channel. The problem of determining whether or not it is possible to bound the state of the process under this type of encoding scheme is not new; it was established in [2–4] that a necessary and sufficient condition for stability can be expressed as a simple relationship between the unstable eigenvalues of A and the average communication bit-rate.

We expand upon this result by considering the notion that encoders can effectively save communication resources by not transmitting information, while noting that the absence of an explicit transmission nevertheless conveys information. To capture this, we suppose that one symbol from the alphabet consumes no communication resources to transmit, whereas each of the others consumes one unit of communication resources. We then proceed to define the *average cost per symbol* of an encoder, which is essentially the average fraction of non-free symbols emitted. This paper’s main technical contribution is a necessary and sufficient condition for the existence of an encoder/decoder/controller that bounds the state of the process. This condition depends on the channel’s average bit rate, the encoder’s average cost per symbol, and the unstable eigenvalues of A .

Our result extends [2] in the sense that as the constraint on the average cost per symbol is allowed to increase (becomes looser), our necessary and sufficient condition becomes the condition from [2]. As with [2–4], our result is constructive in the sense that we describe a family of encoder/decoder pairs that bound the process when our condition holds. A

counterintuitive corollary to our main result shows that if the process may be bounded with average bit-rate r , then there exists a bounding encoder/decoder/controller with average bit-rate r which uses no more than 50% non-free symbols in its symbol-stream.

The remainder of this paper is organized as follows. Section III contains the main negative result of the paper, namely that boundedness is not possible when our condition does not hold. To prove this result we actually show that it is not possible to bound the process with a large class of encoders — which we call M -of- N encoders — that includes all the encoders with average cost per symbol not exceeding a given threshold. Section IV contains the positive result of the paper, showing that when our condition *does* hold, there is an encoder/decoder pair that can bound the process; we provide the encoding scheme.

II. PROBLEM STATEMENT

Consider a stabilizable linear time-invariant process

$$\dot{x} = Ax + Bu, \quad x \in \mathbb{R}^n, u \in \mathbb{R}^m, \quad (1)$$

for which it is known that $x(0)$ belongs to some bounded set $\mathcal{X}_0 \subset \mathbb{R}^n$. A sensor that measures the state $x(t)$ is connected to the actuator through a finite-data-rate, error-free, and delay-free communication channel. An *encoder* collocated with the sensor samples the state once every T time units, and from this sequence of measurements $\{x(kT) : k \in \mathbb{N}_{>0}\}$ causally constructs a sequence of symbols $\{s_k \in \mathcal{A} : k \in \mathbb{N}_{>0}\}$ from a finite alphabet \mathcal{A} . The encoder sends this symbol sequence through the channel at a rate of 1 symbol every T time units to a *decoder/controller* collocated with the actuator, which causally constructs the control signal $u(t)$, $t \geq 0$ from the sequence of symbols $\{s_k \in \mathcal{A} : k \in \mathbb{N}_{>0}\}$ that arrive at the decoder. We assume the channel faithfully transmits each symbol without error.

The positive time T between successive samplings is called the *sampling period* and has units of time units. The *average bit-rate* of an encoder/decoder pair is the amount of information that the encoder transmits in units of bits per time unit. For an encoder/decoder pair whose encoder transmits one symbol from an alphabet \mathcal{A} every T time units, the pair’s average bit-rate is given by

$$r := \frac{\log_2 |\mathcal{A}|}{T}. \quad (2)$$

We consider encoder/decoder pairs whose alphabets each contain one special symbol $0 \in \mathcal{A}$ that can be transmitted without consuming any communication resources, and S

additional symbols that each require one unit of communication resources per transmission. One can think of the “free” symbol 0 as the absence of an explicit transmission. The “communication resources” at stake may be energy, time, or any other resource that may be consumed in the course of the communication process. In order to capture the average rate that an encoder consumes communication resources, we define the *average cost per symbol* of an encoder as follows: We say an encoder has *average cost per symbol not exceeding* γ_{\max} if for every symbol sequence $\{s_k\}$ the encoder may generate, we have

$$\frac{1}{N-M+1} \sum_{k=M}^N I_{s_k \neq 0} \leq \gamma_{\max} + O\left(\frac{1}{N-M+1}\right) \quad (3)$$

for all integers N, M satisfying $N \geq M \geq 0$, where $I_{s_k \neq 0} := 1$ if the k th symbol is not the free symbol, and 0 if it is. The summation in (3) captures the total resources spent transmitting symbols s_M, s_{M+1}, \dots, s_N . Motivating this definition of average cost per symbol is the observation that the lefthand side has the intuitive interpretation as the average cost per transmitted symbol between symbols s_M and s_N . As $N - M \rightarrow \infty$, the rightmost “big-oh” term vanishes, leaving γ_{\max} as an upper bound on the average long-term cost per symbol of the symbol sequence. Note that the average cost per symbol of any encoder never exceeds 1 and does not depend on the sampling period T .

Whereas an encoder/decoder pair’s average bit-rate r only depends on its symbol alphabet \mathcal{A} and sampling period T , its average cost per symbol depends on every possible symbol sequence it may generate, and therefore depends on the encoder/decoder pair, the controller, the process (1), and the initial condition $x(0)$.

The specific question considered in this paper is: under what conditions on the bit-rate and average cost per symbol does there exist a controller and encoder/decoder pair that keep the state of process (1) bounded?

III. NECESSARY CONDITION FOR BOUNDEDNESS WITH LIMITED-COMMUNICATION ENCODERS

It is well known from [2–4] that it is possible to construct a controller and encoder/decoder pair that bounds process (1) with average bit-rate r only if

$$r \ln 2 \geq \sum_{i: \Re \lambda_i[A] \geq 0} \lambda_i[A], \quad (4)$$

where \ln denotes the base- e logarithm, and the summation is over all eigenvalues of A with nonnegative real part. The result that follows shows that a larger bit-rate may be needed when one poses constraints on the encoder’s average cost per symbol γ_{\max} . Specifically, when $\gamma_{\max} \geq S/(S+1)$ the (necessary) stability condition reduces to (4), but when $\gamma_{\max} < S/(S+1)$ a bit-rate larger than (4) is necessary to compensate for the “dilution” of information content in the transmitted symbols due to the constraint on the average cost per symbol.

Theorem 1: Consider an encoder/decoder pair with sampling period T , an alphabet with S nonfree symbols and one

free symbol, and an average cost per symbol not exceeding γ_{\max} . If this pair keeps the state of process (1) bounded for every initial condition $x_0 \in \mathcal{X}_0$, then we must have

$$r f(\gamma_{\max}, S) \ln 2 \geq \sum_{i: \Re \lambda_i[A] \geq 0} \lambda_i[A], \quad (5)$$

where the bit-rate r is related to S and T via Equation (2), and the function $f: [0, 1] \times [0, \infty) \rightarrow [0, \infty)$ is defined as

$$f(\gamma, S) := \begin{cases} \frac{H(\gamma) + \gamma \log_2 S}{\log_2(S+1)} & 0 \leq \gamma \leq \frac{S}{S+1} \\ 1 & \frac{S}{S+1} < \gamma \leq 1, \end{cases} \quad (6)$$

and $H(p) := -p \log_2(p) - (1-p) \log_2(1-p)$ is the base-2 entropy of a Bernoulli random variable with parameter p . \square

Remark 1: The function f is plotted in Figure 1 for several values of S . It is worth making two observations regarding f . The first observation is that the function $f(\gamma, S)$ is monotone nonincreasing in S for any fixed $\gamma \in [0, 1]$, which implies that smaller alphabets are preferable to large ones when trying to satisfy (5) with a given fixed bit-rate. The second observation is that the average cost per time unit, which is γ/T , can be made arbitrarily small in two different ways.

- 1) One could pick T very large, then leveraging the fact that $r f(\gamma, S)$ is monotone increasing in S , pick S large enough to satisfy (5). This approach has two downsides: First, with a large sampling period, the state, although remaining bounded, can grow quite large between transmissions. Second, large S means that the encoder/decoder pair must store and process a large symbol library, adding complexity to the pair’s implementation.
- 2) Alternatively, one can make γ/T arbitrary small by observing that the sequences

$$\gamma_k := e^{-k}, \quad T_k := e^{-k} \sqrt{k}, \quad k \in \mathbb{N}_{>0}$$

have the property that $\gamma_k \rightarrow 0$, $T_k \rightarrow 0$, and $\gamma_k/T_k \rightarrow 0$, but $H(\gamma_k)/T_k \rightarrow \infty$, and hence $r_k f(\gamma_k, S) \ln 2 \rightarrow \infty$ (where $r_k := \log_2(S+1)/T_k$). This means that one can find $k \in \mathbb{N}_{>0}$ to make the average cost per time unit γ_k/T_k arbitrarily small, and also satisfy the necessary condition (5). The drawback of this approach is that to achieve a small sampling period T in practice requires an encoder/decoder pair with a very precise clock.

Remark 2: The addition of the “free” symbol effectively increases the data rate without increasing the rate of resource consumption, as seen by the following two observations:

- Without the free symbols, the size of the alphabet would be S and the bit-rate would be $\log_2(S)/T$. It could happen that this bit-rate is too small to bound the plant, yet after the introduction of the free symbol, the condition (5) is satisfied.
- Since γ_{\max} is the fraction of non-free symbols, then the quantity $r \gamma_{\max}$ is the number of bits per time unit spent transmitting non-free symbols. But since $f(\gamma, S) \geq \gamma$, again we see that the free symbols help satisfy (5).

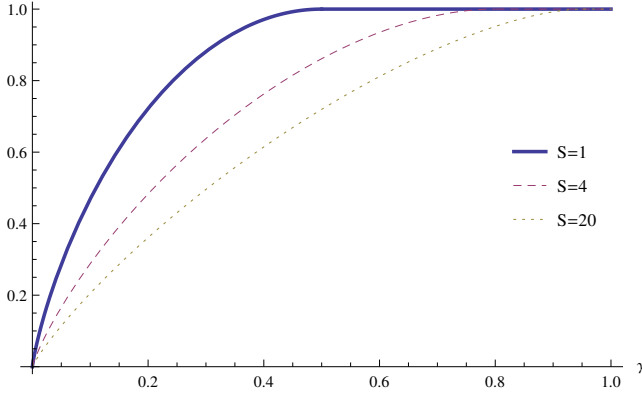


Fig. 1. A plot of $f(\gamma, S)$ for $S = 1, 4, 20$.

A. Theorem 1 Proof Setup

We lead up to the proof of Theorem 1 by first establishing three lemmas centered around a restricted large class of encoders called M -of- N encoders. We first define M -of- N encoders, which essentially partition their symbol sequences into N -length *codewords*, each with M or fewer non-free symbols. Lemma 1 demonstrates that every encoder with a limited average cost per symbol is an M -of- N encoder for appropriate N and M . Next, in Lemma 2 we establish a relationship between the number of codewords available to an M -of- N encoder and the function f as defined in (6). Then, in Lemma 3 we leverage previous work to establish a necessary condition for an M -of- N encoder to bound the state of the process. Finally, the proof of Theorem 1 leverages these three results.

To introduce the class of M -of- N encoders, we define an N -symbol *codeword* to be a sequence

$$\{s_{\ell N+1}, s_{\ell N+2}, \dots, s_{\ell N+N}\}$$

of N consecutive symbols starting at an index $k = \ell N + 1$, with $\ell \in \mathbb{N}_{\geq 0}$. An M -of- N encoder is an encoder for which every N -symbol codeword has M or fewer non-free symbols, i.e.,

$$\sum_{k=\ell N+1}^{\ell N+N} I_{s_k \neq 0} \leq M, \quad \forall \ell \in \mathbb{N}_{\geq 0}. \quad (7)$$

The total number of distinct N -symbol codewords available to an M -of- N encoder is thus given by

$$L(N, M, S) := \sum_{i=0}^{\lfloor M \rfloor} \binom{N}{i} S^i, \quad (8)$$

where the i th term in the summation counts the number of N -symbol codewords with exactly i non-free symbols.

Note that in keeping with the problem setup, the M -of- N encoders considered here each draw their symbols from the symbol library $\mathcal{A} := \{0, 1, \dots, S\}$ and transmit symbols with sampling period T .

An intuitive property of M -of- N encoders is that an M -of- N encoder has average cost per symbol not exceeding M/N .

To see this, suppose M and N are fixed and let $N_1, N_2 \in \mathbb{N}_{\geq 0}$ be arbitrary with $N_2 \geq N_1 \geq 0$. The width of the interval $[N_1, N_2]$ may not be an integer multiple of N , but nevertheless it is between two integer multiples of N : we have $N(k-2) \leq N_2 - N_1 \leq \hat{N}_2 - \hat{N}_1 = kN$ for some non-negative integer $k \in \mathbb{N}_{\geq 0}$. From the first inequality we obtain

$$k \leq \frac{N_2 - N_1 + 1}{N} + 2 - \frac{1}{N},$$

and from the second inequality we have

$$\sum_{k=N_1}^{N_2} I_{s_k \neq 0} \leq kM,$$

because in each of the k N -length intervals containing $[N_1, N_2]$ there are at most M non-free symbols. Combining these yields

$$\sum_{k=N_1}^{N_2} I_{s_k \neq 0} \leq \frac{M}{N}(N_2 - N_1 + 1) + N_0, \quad (9)$$

where N_0 is some constant. Dividing both sides by $(N_2 - N_1 + 1)$ and rearranging, (3) follows.

The fact that an M -of- N encoder refrains from sending “expensive” codewords effectively reduces its ability to transmit information. Indeed, since an M -of- N encoder takes NT time units to transmit one of $L(N, M, S)$ codewords, the encoder transmits merely $\frac{\log_2 L(N, M, S)}{NT}$ bits per time unit. For $M < N$, this is strictly less than the encoder’s average bit-rate r .

The first lemma, proved in the appendix, shows that the set of M -of- N encoders is “complete” in the sense that every encoder with average cost per symbol not exceeding a finite threshold γ_{\max} is actually an M -of- N encoder for N sufficiently large and $M \approx \gamma_{\max} N$.

Lemma 1: For every encoder/decoder pair with average bit-rate r and average cost per symbol not exceeding $\gamma_{\max} \in [0, 1]$, and every constant $\epsilon > 0$, there exist positive integers M and N with $M < N\gamma_{\max}(1 + \epsilon)$ such that the encoder is an M -of- N encoder. \square

The next lemma establishes a relationship between the number of codewords $L(N, M, S)$ available to an M -of- N encoder and the function f defined in (6).

Lemma 2: For any $N, S \in \mathbb{N}_{>0}$ and $\gamma \in [0, 1]$, the function L defined in (8) and the function f defined in (6) satisfy

$$\frac{\log_2 L(N, N\gamma, S)}{N} \leq \log_2(S+1)f(\gamma, S), \quad (10)$$

with equality holding only when $\gamma = 0$ or $\gamma = 1$. Moreover, we have asymptotic equality in the sense that

$$\lim_{N \rightarrow \infty} \frac{\log_2 L(N, N\gamma, S)}{N} = \log_2(S+1)f(\gamma, S). \quad (11)$$

\square

Proof of Lemma 2. Let N and S be arbitrary positive integers. First we prove (10) for $\gamma \in (0, \frac{S}{S+1}]$. Applying the

Binomial Theorem to the identity $1 = (\gamma + (1 - \gamma))^N$, we obtain

$$1 = \sum_{i=0}^N \binom{N}{i} \gamma^i (1 - \gamma)^{N-i}$$

Since each term in the summation is positive, keeping only the first $\lfloor N\gamma \rfloor$ terms yields the inequality

$$1 > \sum_{i=0}^{\lfloor N\gamma \rfloor} \binom{N}{i} \gamma^i (1 - \gamma)^{N-i} \quad (12)$$

Next, a calculation presented as Lemma 5 in the appendix reveals that

$$\gamma^i (1 - \gamma)^{N-i} \geq 2^{-N H(\gamma)} \frac{S^i}{S^{N\gamma}} \quad (13)$$

for all $N, S \in \mathbb{N}_{>0}$, $\gamma \in (0, S/(S+1)]$, and $i \in [0, N\gamma]$. Using this in (12) and simplifying yields

$$\frac{\log_2 L(N, N\gamma, S)}{N} < H(\gamma) + \gamma \log_2 S \quad (14)$$

for all $N, S \in \mathbb{N}_{>0}$ and $\gamma \in (0, S/(S+1)]$. Since $\log_2(S+1)f(\gamma, S) = H(\gamma) + \gamma \log_2 S$ when $\gamma \in [0, \frac{S}{S+1}]$, (14) proves (10) for $\gamma \in (0, \frac{S}{S+1}]$. Next, suppose $\gamma \in (\frac{S}{S+1}, 1)$ and observe from (8) that $L(N, M, S)$ is a sum of positive terms whose index reaches $\lfloor M \rfloor$, hence $L(N, N\gamma, S)$ is strictly less than $L(N, N, S)$ for any $\gamma < 1$. We conclude that

$$\frac{\log_2 L(N, N\gamma, S)}{N} < \frac{\log_2 L(N, N, S)}{N} \quad (15)$$

$$= \log_2(S+1) \quad (16)$$

Since $\log_2(S+1)f(\gamma, S) = \log_2(S+1)$ for $\gamma \in (\frac{S}{S+1}, 1)$, this concludes the proof of (10) for $\gamma \in (0, 1)$. The proof of (10) for $\gamma = 0$, follows merely from inspection of (10), and the $\gamma = 1$ case follows from the equality in (15). The proof of the asymptotic result (11) appears in [8]. This concludes the proof of Lemma 2. ■

The following lemma provides a necessary condition for an M -of- N encoder to bound process (1).

Lemma 3: Consider an M -of- N encoder/decoder pair using symbols $\{0, \dots, S\}$ with sampling period T . If the pair keeps the state of (1) bounded for every initial condition, then we must have

$$\frac{\ln L(N, M, S)}{NT} > \sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A]. \quad (17)$$

□

Proof of Lemma 3. Consider an encoder/decoder/controller triple that bounds the state of the process (1) and whose encoder is an M -of- N encoder using symbols $\{0, \dots, S\}$ and sampling period T . Since the encoder sends one of $L(M, N, S)$ codewords every NT time units, the average bit-rate of the encoder is $\log_2 L(M, N, S)/NT$. Theorem 1 of [2] proves that if the state of the process (1) is bounded by an encoder/decoder/controller triple under the communication constraints described in our problem setup, then the pair's

average bit-rate must not be less than $\frac{1}{\ln 2} \sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A]$. Hence, we have

$$\frac{\log_2 L(M, N, S)}{NT} \geq \frac{1}{\ln 2} \sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A],$$

from which (17) follows. This proves the lemma. ■

Now we are ready to prove Theorem 1.

Proof of Theorem 1. By Lemma 1, for any $\epsilon > 0$ there exist $M, N \in \mathbb{N}_{>0}$ with $M < N\gamma_{\max}(1 + \epsilon)$ for which the encoder/decoder is an M -of- N encoder. Since the state of the process is kept bounded, by Lemma 3 we have

$$\sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A] < \frac{\log_2 L(N, M, S)}{NT} \ln 2. \quad (18)$$

Since L is monotonically nondecreasing in its second argument and $M < N\gamma_{\max}(1 + \epsilon)$, we have

$$\frac{\log_2 L(N, M, S)}{NT} \leq \frac{\log_2 L(N, N\gamma_{\max}(1 + \epsilon), S)}{NT}. \quad (19)$$

Lemma 2 implies that

$$\frac{\log_2 L(N, N\gamma_{\max}(1 + \epsilon), S)}{NT} \leq rf(\gamma_{\max}(1 + \epsilon), S). \quad (20)$$

Combining these and letting $\epsilon \rightarrow 0$, we obtain (5). This completes the proof of Theorem 1. ■

IV. SUFFICIENT CONDITION FOR BOUNDEDNESS WITH LIMITED-COMMUNICATION ENCODERS

The previous section established a necessary condition (5) on the average bit-rate and average cost per symbol of an encoder/decoder pair in order to bound the state of (1). In this section, we show that that condition is also sufficient for a bounding encoder/decoder to exist. The proof is constructive in that we provide the encoder/decoder.

Theorem 2: Assume that A is diagonalizable. For every $S \in \mathbb{N}_{>0}$, $T > 0$, and $\gamma_{\max} \in [0, 1]$ satisfying

$$rf(\gamma_{\max}, S) \ln 2 > \sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A], \quad (21)$$

where r is defined in (2) and the function f is defined in (6), there exists a controller and an encoder/decoder pair using S nonfree symbols, sampling period T , and average cost per symbol not exceeding γ_{\max} , that keeps the state of the process (1) bounded for every initial condition $x_0 \in \mathcal{X}_0$. □

The proof of Theorem 2 relies on the following lemma, which provides a sufficient condition for the existence of an M -of- N encoder to bound the state of process (1).

Lemma 4: Assume that A is diagonalizable. For every $T > 0$ and $N, M, S \in \mathbb{N}_{>0}$ with $N \geq M \geq 0$ satisfying

$$\frac{\ln L(N, M, S)}{NT} > \sum_{i: \mathfrak{R}\lambda_i[A] \geq 0} \lambda_i[A] \quad (22)$$

there exists an M -of- N encoder on alphabet $\{0, \dots, S\}$ with sampling period T that keeps the state of the process (1) bounded for every initial condition. □

Proof of Lemma 4. This proof builds on Theorem 2 from [2], which provides sufficient conditions on an encoder's average bit-rate to ensure the existence of a stabilizing controller and encoder/decoder pair. The result states that if an encoder's average bit-rate r satisfies

$$r \geq \frac{1}{\ln 2} \sum_{i: \mathcal{R}\lambda_i[A] \geq 0} \lambda_i[A], \quad (23)$$

where A is the continuous-time process matrix, then there exists a controller and encoder/decoder pair that bound the state of the process. An outline of the proof of this result from [2] follows. The encoder works by placing a bounding rectangle around the volume where the state is known to lie, and partitions it into two pieces along its largest axis. Then the encoder transmits a 0 or 1, depending on which sub-rectangle the state lies in. Since the decoder can calculate the bounding rectangles with its knowledge of \mathcal{X}_0 and process dynamics in (1), it estimates the state to be the centroid of whichever sub-rectangle the received symbol corresponds to. The decoder transmits this state estimate to the controller, which can be any stabilizing linear state-feedback controller.

We now provide a summary of the proof of Lemma 4, and refer the reader to [8] for the detailed proof. Suppose T, N, M, S satisfy the constraints in the statement of the lemma as well as (22). Theorem 2 from [2] guarantees the existence of a stabilizing encoder/decoder/controller triple with average bit-rate $\log_2 L(N, M, S)/NT$. All that remains is to adapt this encoder/decoder pair into our framework, i.e., an M -of- N encoder that uses $S + 1$ symbols with sampling period T . We do this by building an encoder which runs a copy of the encoder from [2] internally. The outer encoder feeds samples of the state to the inner encoder, from which the inner encoder generates a string of 1's and 0's. Due to the bit-rate being $\log_2 L(N, M, S)/NT$, after NT time units the inner encoder has generated a string of length $\log_2 L(N, M, S)$. There are $L(N, M, S)$ possible such strings, and so the outer encoder can map this string uniquely to one of the $L(N, M, S)$ possible N -length codewords with M or fewer non-free symbols from the alphabet $\{0, \dots, S\}$. The outer encoder transmits this codeword across the channel to the decoder, which performs the inverse mapping to recover the string of $\log_2 L(N, M, S)$ 1's and 0's, which it delivers to its inner decoder. From this bit-string and knowledge of the process dynamics, the inner decoder computes a state estimate, which it delivers to the stabilizing linear state-feedback controller. This process repeats each NT time units. Since the inner encoder/decoder pair bound the process, the outer M -of- N encoder described here does as well. This proves the lemma. ■

Now we are ready to prove Theorem 2.

Proof of Theorem 2. Assume that S, T , and γ_{\max} satisfy (21), so that

$$\epsilon := rf(\gamma_{\max}, S) \ln 2 - \sum_{i: \mathcal{R}\lambda_i[A] \geq 0} \lambda_i[A] > 0. \quad (24)$$

Equation (11) establishes that $\frac{\ln L(N, N\gamma_{\max}, S)}{NT}$ gets arbitrarily close to $rf(\gamma_{\max}, S)$ as we increase N , so we pick N sufficiently large to satisfy

$$rf(\gamma_{\max}, S) \ln 2 - \frac{\ln L(N, N\gamma_{\max}, S)}{NT} < \epsilon/2. \quad (25)$$

By (8), we have $L(N, \lfloor N\gamma_{\max} \rfloor, S) = L(N, N\gamma_{\max}, S)$ for every N, γ_{\max} , and S . Setting $M := \lfloor N\gamma_{\max} \rfloor$, then by (24) and (25) we have found N and M satisfying $\frac{\ln L(N, M, S)}{NT} > \sum_{i: \mathcal{R}\lambda_i[A] \geq 0} \lambda_i[A]$. Hence by Lemma 4, there exists an M -of- N encoder/decoder which bounds the state of (1). Since all M -of- N encoder/decoders have average cost per symbol not exceeding M/N , and this encoder/decoder satisfies $M/N \leq \gamma_{\max}$, we conclude that this encoder has an average cost per symbol not exceeding γ_{\max} . This concludes the proof, since we have found the desired encoder. ■

An unexpected consequence of Theorems 1 and 2 is that when it is possible to keep the state of a process bounded with a given bit-rate $r := \log_2(S + 1)/T$, one can always find M -of- N encoders that bound it for (essentially) the same bit-rate and average cost per symbol not exceeding $S/(S + 1)$, i.e., approximately a fraction $1/(S + 1)$ of the symbols will not consume communication resources. In the most advantageous case, the encoder/decoder use the alphabet $\{0, 1\}$ and the encoder's symbol stream consumes no more than 50% of the communication resources. The price paid for using an encoder/decoder with average cost per symbol near $S/(S + 1)$ is that it may require prohibitively long codewords (large N) as compared to an encoder with higher average cost per symbol. This is due to the fact that $\ln L(N, N\gamma_{\max}, S)/NT$ is monotonically nondecreasing in γ_{\max} , so that a smaller γ_{\max} generally requires a larger N to satisfy (25).

Corollary 1: If the process (1) can be bounded with an encoder/decoder pair with symbol alphabet $\{0, 1, \dots, S\}$ and sampling period T , then for any $\epsilon > 0$ with $T > \epsilon$, there exists an M -of- N encoder with alphabet $\{0, 1, \dots, S\}$, sampling period $T - \epsilon$, and average cost per symbol not exceeding $S/(S + 1)$ that bounds its state. □

Proof of Corollary 1. Let ϵ be arbitrary in $(0, T)$. Since the controller and encoder/decoder pair bound the system and the average cost per symbol never exceeds 1, by Theorem 1 we have

$$rf(1, S) \ln 2 \geq \sum_{i: \mathcal{R}\lambda_i[A] \geq 0} \lambda_i[A], \quad (26)$$

where $r := \log_2(S + 1)/T$. Since $f(S/(S + 1), S) = f(1, S)$ for all $S \in \mathbb{N}_{>0}$, we have

$$\begin{aligned} \frac{\log_2(S + 1)}{T - \epsilon} f\left(\frac{S}{S + 1}, S\right) &> rf\left(\frac{S}{S + 1}, S\right) \\ &= rf(1, S) \ln 2 \geq \sum_{i: \mathcal{R}\lambda_i[A] \geq 0} \lambda_i[A]. \end{aligned} \quad (27)$$

By Theorem 2, there exists an encoder/decoder pair with symbol alphabet $\{0, \dots, S\}$, sampling period $T - \epsilon$, and average cost per symbol not exceeding $S/(S + 1)$ which

bounds the state of (1). By Lemma 1, this encoder is an M -of- N encoder for appropriately chosen M and N . ■

V. CONCLUSION AND FUTURE WORK

In this paper we considered the problem of bounding the state of a continuous-time linear process under communication constraints. We considered constraints on both the average bit-rate and the average cost per symbol of an encoding scheme. Our main contribution was a necessary and sufficient condition on the process and constraints for which a bounding encoder/decoder/controller exists. In the absence of a limit on the average cost per symbol, the conditions recovered previous work. A surprising corollary to our main result was the observation that one may impose a constraint on the average cost per symbol without necessarily needing to loosen the bit-rate constraint. Specifically, we proved that if a process may be bounded with a particular bit-rate, then there exists a (possibly very complex) encoder/decoder that can bound it using no more than 50% non-free symbols on average, yet obeying the *same* bit-rate. This was surprising because one would expect the prohibition of some codewords would require that the encoder always compensate by transmitting at a higher bit-rate.

We observed in Remark 1 that smaller alphabets incur a smaller penalty on the conditions for boundedness in Theorems 1 and 2. This suggests that encoding schemes with small alphabets may be able to bound the state of the process with bit-rates and average costs not far above the minimum theoretical bounds as established in Theorems 1 and 2. *Event-based* control strategies comprise one such class of encoders; they use a small number of non-free symbols to notify the decoder/controller about certain state-dependent events. We have preliminary results showing that event-based encoders can indeed be used to produce encoder/decoder pairs that almost achieve the minimum achievable average cost per symbol that appears in Theorems 1 and 2.

Finally, our problem setup considered merely whether there exists a bounding encoder/decoder/controller triple. It seems natural to extend this setup to finding stabilizing triples.

APPENDIX

Proof of Lemma 1. By the definition of average cost per symbol not exceeding γ_{\max} in (3) there exists an integer $N_0 \in \mathbb{N}_{>0}$ such that for any symbol sequence $\{s_k\}$ that the encoder generates, we have

$$\sum_{i=1}^N I_{s_i \neq 0} < N_0 + N\gamma_{\max}, \quad \forall N \in \mathbb{N}_{>0}. \quad (28)$$

Pick $N \in \mathbb{N}_{>0}$ large enough to satisfy $N_0 + 2 < \epsilon N\gamma_{\max}$ and pick $M := \lfloor N_0 + 2 + N\gamma_{\max} \rfloor$. Combining these with (28), we obtain

$$\begin{aligned} \sum_{i=1}^N I_{s_i \neq 0} &< N_0 + N\gamma_{\max} < M \\ &\leq N_0 + 2 + N\gamma_{\max} < N\gamma_{\max}(1 + \epsilon), \end{aligned}$$

which establishes that $M < N\gamma_{\max}(1 + \epsilon)$. This completes the proof. ■

Lemma 5: The following inequality holds for all $N, S \in \mathbb{N}_{>0}$, $q \in (0, S/(S+1)]$, and $i \in [0, Nq]$:

$$q^i(1-q)^{N-i} \geq 2^{-NH(q)} \frac{S^i}{S^{Nq}} \quad (29)$$

where $H(q) := -q \log_2 q - (1-q) \log_2(1-q)$ is the base-2 entropy of a Bernoulli random variable with parameter q .

Proof of Lemma 5. Let N, S, q , and i take arbitrary values from the sets described in the lemma's statement. Since \log_2 is a monotone increasing function, $\log_2(q/(1-q))$ for $q > 0$ is maximized at the right endpoint value, $q = S/(S+1)$, where it equals $\log_2 S$. This leads to

$$\log_2 q - \log_2(1-q) \leq \log_2 S \quad (30)$$

for all $S \in \mathbb{N}_{>0}$ and $q \in (0, S/(S+1)]$. Next, since $i \in [0, Nq]$ by assumption, then $i - Nq \leq 0$. Multiplying (30) by $i - Nq$, standard algebraic manipulations yield

$$\begin{aligned} i \log_2 q + (N-i) \log_2(1-q) \\ \geq -NH(q) + (i - Nq) \log_2 S, \end{aligned}$$

from which (29) follows. ■

REFERENCES

- [1] R. Brockett and D. Liberzon, "Quantized feedback stabilization of linear systems," *Automatic Control, IEEE Transactions on*, vol. 45, no. 7, pp. 1279–1289, Jul 2000.
- [2] J. P. Hespanha, A. Ortega, and L. Vasudevan, "Towards the control of linear systems with minimum bit-rate," in *Proc. of the Int. Symp. on the Mathematical Theory of Networks and Syst.*, Aug. 2002.
- [3] G. Nair and R. Evans, "Communication-limited stabilization of linear systems," in *Decision and Control, 2000. Proceedings of the 39th IEEE Conference on*, vol. 1, 2000, pp. 1005–1010 vol.1.
- [4] S. Tatikonda and S. Mitter, "Control under communication constraints," *Automatic Control, IEEE Transactions on*, vol. 49, no. 7, pp. 1056–1068, July 2004.
- [5] G. N. Nair and R. J. Evans, "Exponential stabilisability of finite-dimensional linear systems with limited data rates," *Automatica*, vol. 39, no. 4, pp. 585–593, 2003.
- [6] A. Sahai and S. Mitter, "The necessity and sufficiency of anytime capacity for stabilization of a linear system over a noisy communication link — part i: Scalar systems," *Information Theory, IEEE Transactions on*, vol. 52, no. 8, pp. 3369–3395, Aug 2006.
- [7] A. Matveev and A. Savkin, "Multirate stabilization of linear multiple sensor systems via limited capacity communication channels," *SIAM Journal on Control and Optimization*, vol. 44, no. 2, pp. 584–617, 2005.
- [8] J. Pearson, J. P. Hespanha, and D. Liberzon, "Control with minimum communication cost per symbol," University of California, Santa Barbara, Tech. Rep., Mar. 2014. [Online]. Available: <http://www.ece.ucsb.edu/~hespanha/techrep.html>