

Computing Entropy Rate Of Symbol Sources & A Distribution-free Limit Theorem

Ishanu Chattopadhyay Hod Lipson
ic99@cornell.edu hod.lipson@cornell.edu

Abstract—Entropy rate of sequential data-streams naturally quantifies the complexity of the generative process. Thus entropy rate fluctuations could be used as a tool to recognize dynamical perturbations in signal sources, and could potentially be carried out without explicit background noise characterization. However, state of the art algorithms to estimate the entropy rate have markedly slow convergence; making such entropic approaches non-viable in practice. We present here a fundamentally new approach to estimate entropy rates, which is demonstrated to converge significantly faster in terms of input data lengths, and is shown to be effective in diverse applications ranging from the estimation of the entropy rate of English texts to the estimation of complexity of chaotic dynamical systems. Additionally, the convergence rate of entropy estimates do not follow from any standard limit theorem, and reported algorithms fail to provide any confidence bounds on the computed values. Exploiting a connection to the theory of probabilistic automata, we establish a convergence rate of $O(\log |s| / \sqrt[3]{|s|})$ as a function of the input length $|s|$, which then yields explicit uncertainty estimates, as well as required data lengths to satisfy pre-specified confidence bounds.

Index Terms—Entropy rate, Stochastic processes, Probabilistic automata, Symbolic dynamics

I. MOTIVATION, BACKGROUND & CONTRIBUTION

The entropy rate of a stationary and ergodic process converges in probability to the per-letter Kolmogorov complexity of a single sufficiently long sample path [1]. While Kolmogorov complexity is incomputable, entropy rates can, in principle, be estimated. Ability to quantify the complexity of a signal source, even in the average sense, can provide valuable insights into the driving dynamics; and can potentially be used as a tool to detect dynamical anomalies without explicit knowledge of background noise processes.

However, source entropy rate estimation from an observed sample path is computationally non-trivial. Even with the assumptions of ergodicity and stationarity, one cannot fruitfully apply the defining relation in Eq.(6) due to the exponential increase in the number of different words with the word-length. This is particularly important if there are long-range dependencies in the symbol stream. Such dependencies introduce additional long-range structure; decreasing the source entropy in the process. In such cases unacceptably long words or *blocks* must be considered, and pre-mature truncation of the computation would lead to large errors.

The best known algorithms that carry out a more efficient computation are based on Lempel-Ziv (LZ) source coding [2], [3], [4]. The LZ coding algorithms are asymptotically optimal, i.e. their compression rate approaches the source entropy rate for any ergodic stationary stochastic process. The key idea here is adaptive dictionary compression: parse the input string into distinct phrases, and represent them with codewords, making sure that short codewords are assigned to common phrases. Done optimally, one ends up with a compressed string, such that the ratio of the input and output lengths approach the source entropy rate. Different variations on this idea have been reported [5], [6]. Techniques distinct from LZ parsing are also known, e.g., Rissanen [7] reported a universal compression scheme, which instead of gathering parsed segments of the input along with their occurrence counts, collects the “contexts” in which each symbol of the input string occurs, together with conditional occurrence counts.

Importantly, a majority of the reported techniques do more than just compute the entropy rate; they are indeed full-scale data compression utilities, that produce a decodable representation of the input. Can we do better if we are only interested in the former? This paper provides an affirmative answer to this possibility.

Secondly, existing techniques lack convergence rate estimates; computation of error bars for reported approaches do not follow from any standard limit theorem. There is indeed no analytical

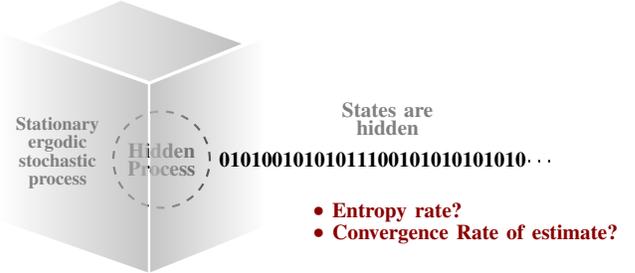


Fig. 1. Problem description. Given a quantized data stream, how do we compute the entropy rate of the hidden process? Even with the assumption of stationarity and ergodicity for the generator, reported algorithms converge very slowly. Additionally, these convergence rates are unknown for such approaches; implying that we cannot put uncertainty bounds on the computed values in practice. We show that a significantly faster computation of the entropy rate is possible; and derive a universal lower bound on how slowly this convergence might occur.

way to check for the internal consistency of the estimation or its accuracy. We may observe gradual convergence to a limiting value, and this is indeed guaranteed by theory; but are unable to provide uncertainty bounds on the computed estimate with finite inputs. Typically observed slow convergence in all non-trivial scenarios, for all reported algorithms, makes this a key issue. An empirical relationship, without proof or theoretical backing, has been suggested [8], which conjectures the $|s|$ -dependence ($|s|$ being the length of the input s) of the estimated entropy rate \tilde{H} to follow $\tilde{H} \simeq H_{\text{actual}} + c \frac{\log |s|}{|s|^\gamma}$, where c, γ are fit parameters. In this paper, we show that, at least with our algorithm, the convergence rate is given by $O(\log |s| / \sqrt[3]{|s|})$. This is a distribution-free result, in the sense that the asymptotic bound does not depend on the source characteristics. In consequence, we can derive explicit uncertainty estimates at specified confidence bounds on the estimated entropy rate for finite-length input data.

A. Key Insight

Our approach is based on modeling discrete and finite-valued stationary and ergodic sources as probabilistic automata. Our automata is distinct from that of Paz [9], and each model in our case is in fact an encoding of a measure defined on the space of strictly infinite strings over a finite alphabet. While the formalisms are completely different, some aspects of this approach has subtle parallels to that of Rissanen’s “context algorithm” [7]; his search for contexts which yield similar probabilities of generating future symbols is analogous to our search for a synchronizing string in the input stream - a finite sequence of symbols that, once executed on a probabilistic automaton, leads to a fixed state irrespective of the initial conditions. Of course we do not know anything about the hidden model a priori; but nevertheless we establish that such a string, at least in a well-defined approximate sense, always exists and is identifiable efficiently. Finally, we show that, given such an approximate synchronizing string, we can use results from non-parametric statistics to bound the probability of error as a function of the input length.

B. Entropy & Entropy Rate

Entropy $H(X)$ of a discrete random variable X , taking values in the alphabet Σ , is defined as:

$$H(X) = - \sum_{x \in \Sigma} p(x) \log p(x) \quad (1)$$

where $p(x)$ is the probability of occurrence of $x \in \Sigma$. The base of the logarithm is generally taken to be 2, and then the entropy is being expressed in *bits*. While the definition of entropy of a random variable may be obtained axiomatically, a perhaps more compelling approach is to show that it arises as the average length of the shortest description of a random variable [10].

The joint entropy of a set of random variables X_1, \dots, X_n , with X_i taking values in the alphabet Σ_i , is defined in the usual manner:

$$H(X_1, \dots, X_n) = - \sum_{x_i \in \Sigma_i} p(x_1, \dots, x_n) \log_2 p(x_1, \dots, x_n) \quad (2)$$

The chain rule for entropy calculations [10] follows from the definitions, and is of particular importance:

$$H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1) \quad (3)$$

The notion of entropy formalizes the Asymptotic Equipartition Property (AEP): If discrete random variables X_1, \dots, X_n are i.i.d. and have probability mass function $p(x)$, then we have:

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \xrightarrow{a.s.} H(X) = - \sum_{x \in \Sigma} p(x) \log p(x) \quad (4)$$

The AEP implies that $nH(X)$ bits suffice on average to describe n i.i.d. random variables. If the random variables are not independent, the entropy $H(X_1, \dots, X_n)$ still grows asymptotically linearly with n at a rate known as the entropy rate of the process. In particular, if the random variables define a stationary ergodic stochastic process $\mathcal{X} = \{X_i\}$, then the AEP still holds:

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \xrightarrow{a.s.} H(\mathcal{X}) \quad (5)$$

where $H(\mathcal{X})$ is the entropy rate of the process defined as:

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) \quad (6)$$

As in the case of the i.i.d. variables, *typical sequences* of length n may be represented using approximately $nH(\mathcal{X})$ bits. Thus the entropy rate quantifies the average description length of the process, and hence its expected complexity [1].

II. STOCHASTIC PROCESSES & PROBABILISTIC AUTOMATA

As mentioned earlier, our approach hinges upon effectively using probabilistic automata to model stationary, ergodic processes. Our automata models are distinct to those reported in the literature [9], [11]. The details of this formalism can be found in [12]; we include a brief overview here for the sake of completeness.

Notation 1. Σ denotes a finite alphabet of symbols. The set of all finite but possibly unbounded strings on Σ is denoted by Σ^* [13]. The set of finite strings over Σ form a concatenative monoid, with the empty word λ as identity. The set of strictly infinite strings on Σ is denoted as Σ^ω , where ω denotes the first transfinite cardinal. For a string x , $|x|$ denotes its length, and for a set A , $|A|$ denotes its cardinality. Also, $\Sigma^{\leq d} = \{x \in \Sigma^* \text{ s.t. } |x| \leq d\}$.

Definition 1 (QSP). A QSP \mathcal{H} is a discrete time Σ -valued strictly stationary, ergodic stochastic process, i.e.

$$\mathcal{H} = \{X_t : X_t \text{ is a } \Sigma\text{-valued random variable, } t \in \mathbb{N} \cup \{0\}\} \quad (7)$$

A process is ergodic if moments may be calculated from a sufficiently long realization, and strictly stationary if moments are time-invariant.

We next formalize the connection of QSPs to PFSA generators. We develop the theory assuming multiple realizations of the QSP \mathcal{H} , and fixed initial conditions. Using ergodicity, we will be then able to apply our construction to a single sufficiently long realization, where initial conditions cease to matter.

Definition 2 (σ -Algebra On Infinite Strings). For the set of infinite strings on Σ , we define \mathfrak{B} to be the smallest σ -algebra generated by the family of sets $\{x\Sigma^\omega : x \in \Sigma^*\}$.

Lemma 1. Every QSP induces a probability space $(\Sigma^\omega, \mathfrak{B}, \mu)$.

Proof: Assuming stationarity, we can construct a probability measure $\mu : \mathfrak{B} \rightarrow [0, 1]$ by defining for any sequence $x \in \Sigma^* \setminus \{\lambda\}$, and

a sufficiently large number of realizations N_R (assuming ergodicity):

$$\mu(x\Sigma^\omega) = \lim_{N_R \rightarrow \infty} \frac{\# \text{ of initial occurrences of } x}{\# \text{ of initial occurrences of all sequences of length } |x|}$$

and extending the measure to elements of $\mathfrak{B} \setminus \mathfrak{B}$ via at most countable sums. Thus $\mu(\Sigma^\omega) = \sum_{x \in \Sigma^*} \mu(x\Sigma^\omega) = 1$, and for the null word $\mu(\lambda\Sigma^\omega) = \mu(\Sigma^\omega) = 1$. ■

Notation 2. For notational brevity, we denote $\mu(x\Sigma^\omega)$ as $\Pr(x)$.

Classically, automaton states are equivalence classes for the Nerode relation; two strings are equivalent if and only if any finite extension of the strings is either both in the language under consideration, or neither are [13]. We use a probabilistic extension [14].

Definition 3 (Probabilistic Nerode Equivalence Relation). $(\Sigma^\omega, \mathfrak{B}, \mu)$ induces an equivalence relation \sim_N on the set of finite strings Σ^* as:

$$\forall x, y \in \Sigma^*, x \sim_N y \iff \forall z \in \Sigma^* \left(\Pr(xz) = \Pr(yz) = 0 \right) \vee \left| \Pr(xz)/\Pr(x) - \Pr(yz)/\Pr(y) \right| = 0 \quad (8)$$

Notation 3. For $x \in \Sigma^*$, the equivalence class of x is $[x]$.

It is easy to see that \sim_N is right invariant, i.e.

$$x \sim_N y \Rightarrow \forall z \in \Sigma^*, xz \sim_N yz \quad (9)$$

A right-invariant equivalence on Σ^* always induces an automaton structure; and hence the probabilistic Nerode relation induces a probabilistic automaton: states are equivalence classes of \sim_N , and the transition structure arises as follows: For states q_i, q_j , and $x \in \Sigma^*$,

$$([x] = q) \wedge ([x\sigma] = q') \Rightarrow q \xrightarrow{\sigma} q' \quad (10)$$

Before formalizing the above construction, we introduce the notion of probabilistic automata with initial, but no final, states.

Definition 4 (Initial-Marked PFSA). An initial marked probabilistic finite state automaton (a Initial-Marked PFSA) is a quintuple $(Q, \Sigma, \delta, \tilde{\pi}, q_0)$, where Q is a finite state set, Σ is the alphabet, $\delta : Q \times \Sigma \rightarrow Q$ is the state transition function, $\tilde{\pi} : Q \times \Sigma \rightarrow [0, 1]$ specifies the conditional symbol-generation probabilities, and $q_0 \in Q$ is the initial state. δ and $\tilde{\pi}$ are recursively extended to arbitrary $y = \sigma x \in \Sigma^*$ as follows:

$$\forall q \in Q, \delta(q, \lambda) = q \quad (11)$$

$$\delta(q, \sigma x) = \delta(\delta(q, \sigma), x) \quad (12)$$

$$\forall q \in Q, \tilde{\pi}(q, \lambda) = 1 \quad (13)$$

$$\tilde{\pi}(q, \sigma x) = \tilde{\pi}(q, \sigma) \tilde{\pi}(\delta(q, \sigma), x) \quad (14)$$

Additionally, we impose that for distinct states $q_i, q_j \in Q$, there exists a string $x \in \Sigma^*$, such that $\delta(q_i, x) = q_j$, and $\tilde{\pi}(q_i, x) > 0$.

Note that the probability of the null word is unity from each state.

If the current state and the next symbol is specified, our next state is fixed; similar to Probabilistic Deterministic Automata [15]. However, unlike the latter, we lack final states in the model. Additionally, we assume our graphs to be strongly connected.

Later we will remove initial state dependence using ergodicity. Next we formalize how a PFSA arises from a QSP.

Lemma 2 (PFSA Generator). Every Initial-Marked PFSA $G = (Q, \Sigma, \delta, \tilde{\pi}, q_0)$ induces a unique probability measure μ_G on the measurable space $(\Sigma^\omega, \mathfrak{B})$.

Proof: Define set function μ_G on the measurable space $(\Sigma^\omega, \mathfrak{B})$:

$$\mu_G(\emptyset) \triangleq 0 \quad (15)$$

$$\forall x \in \Sigma^*, \mu_G(x\Sigma^\omega) \triangleq \tilde{\pi}(q_0, x) \quad (16)$$

$$\forall x, y \in \Sigma^*, \mu_G(\{x, y\}\Sigma^\omega) \triangleq \mu_G(x\Sigma^\omega) + \mu_G(y\Sigma^\omega) \quad (17)$$

Countable additivity of μ_G is immediate, and (See Definition 4):

$$\mu_G(\Sigma^\omega) = \mu_G(\lambda\Sigma^\omega) = \tilde{\pi}(q_0, \lambda) = 1 \quad (18)$$

implying that $(\Sigma^\omega, \mathfrak{B}, \mu_G)$ is a probability space. ■

We refer to $(\Sigma^\omega, \mathfrak{B}, \mu_G)$ as the probability space generated by the Initial-Marked PFSA G .

Lemma 3 (Probability Space To PFSA). If the probabilistic Nerode relation corresponding to a probability space $(\Sigma^\omega, \mathfrak{B}, \mu)$ has a finite index, then the latter has an initial-marked PFSA generator.

Proof: Let Q be the set of equivalence classes of the probabilistic Nerode relation (Definition 3), and define functions $\delta : Q \times \Sigma \rightarrow Q$, $\tilde{\pi} : Q \times \Sigma \rightarrow [0, 1]$ as:

$$\delta([x], \sigma) = [x\sigma] \quad (19)$$

$$\tilde{\pi}([x], \sigma) = \frac{\Pr(x'\sigma)}{\Pr(x')} \text{ for any choice of } x' \in [x] \quad (20)$$

where we extend $\delta, \tilde{\pi}$ recursively to $y = \sigma x \in \Sigma^*$ as

$$\delta(q, \sigma x) = \delta(\delta(q, \sigma), x) \quad (21)$$

$$\tilde{\pi}(q, \sigma x) = \tilde{\pi}(q, \sigma)\tilde{\pi}(\delta(q, \sigma), x) \quad (22)$$

For verifying the null-word probability, choose a $x \in \Sigma^*$ such that $[x] = q$ for some $q \in Q$. Then, from Eq. (20), we have:

$$\tilde{\pi}(q, \lambda) = \frac{\Pr(x'\lambda)}{\Pr(x')} \text{ for any } x' \in [x] \Rightarrow \tilde{\pi}(q, \lambda) = \frac{\Pr(x')}{\Pr(x')} = 1 \quad (23)$$

Finite index of \sim_N implies $|Q| < \infty$, and hence denoting $[\lambda]$ as q_0 , we conclude: $G = (Q, \Sigma, \delta, \tilde{\pi}, q_0)$ is an Initial-Marked PFSA. Lemma 2 implies that G generates $(\Sigma^\omega, \mathfrak{B}, \mu)$, which completes the proof. ■

The above construction yields a *minimal realization* for the Initial-Marked PFSA, unique up to state renaming.

Lemma 4 (QSP to PFSA). *Any QSP with a finite index Nerode equivalence is generated by an Initial-Marked PFSA.*

Proof: Follows immediately from Lemma 1 (QSP to Probability Space) and Lemma 3 (Probability Space to PFSA generator). ■

A. Canonical Representations

We have defined a QSP as both ergodic and stationary, whereas the Initial-Marked PFSA has a designated initial state. Next we introduce canonical representations to remove initial-state dependence. We use $\tilde{\Pi}$ to denote the matrix representation of $\tilde{\pi}$, i.e., $\tilde{\Pi}_{ij} = \tilde{\pi}(q_i, \sigma_j)$, $q_i \in Q, \sigma_j \in \Sigma$. We need the notion of transformation matrices Γ_σ .

Definition 5 (Transformation Matrices). *For an initial-marked PFSA $G = (Q, \Sigma, \delta, \tilde{\pi}, q_0)$, the symbol-specific transformation matrices $\Gamma_\sigma \in [0, 1]^{|Q| \times |Q|}$ are:*

$$\Gamma_\sigma|_{ij} = \begin{cases} \tilde{\pi}(q_i, \sigma), & \text{if } \delta(q_i, \sigma) = q_j \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

Transformation matrices have a single non-zero entry per row, reflecting our generation rule that given a state and a generated symbol, the next state is fixed.

First, we note that, given an initial-marked PFSA G , we can associate a probability distribution \wp_x over the states of G for each $x \in \Sigma^*$ in the following sense: if $x = \sigma_{r_1} \cdots \sigma_{r_m} \in \Sigma^*$, then we have:

$$\wp_x = \wp_{\sigma_{r_1} \cdots \sigma_{r_m}} = \frac{1}{\underbrace{\|\wp_\lambda \prod_{j=1}^m \Gamma_{\sigma_{r_j}}\|_1}_{\text{Normalizing factor}}} \wp_\lambda \prod_{j=1}^m \Gamma_{\sigma_{r_j}} \quad (25)$$

where \wp_λ is the stationary distribution over the states of G . Note that there may exist more than one string that leads to a distribution \wp_x , beginning from the stationary distribution \wp_λ . Thus, \wp_x is an equivalence class of strings, i.e., x is not unique.

Definition 6 (Canonical Representation). *An initial-marked PFSA $G = (Q, \Sigma, \delta, \tilde{\pi}, q_0)$ uniquely induces a canonical representation $(Q^c, \Sigma, \delta^c, \tilde{\pi}^c)$, where Q^c is a subset of the set of probability distributions over Q , and $\delta^c : Q^c \times \Sigma \rightarrow Q^c$, $\tilde{\pi}^c : Q^c \times \Sigma \rightarrow [0, 1]$ are constructed as follows:*

- 1) *Construct the stationary distribution on Q using the transition probabilities of the Markov Chain induced by G , and include this as the first element \wp_λ of Q^c . Note that the transition matrix for G is the row-stochastic matrix $M \in [0, 1]^{|Q| \times |Q|}$, with $M_{ij} = \sum_{\sigma: \delta(q_i, \sigma) = q_j} \tilde{\pi}(q_i, \sigma)$, and hence \wp_λ satisfies:*

$$\wp_\lambda M = \wp_\lambda \quad (26)$$

- 2) *Define δ^c and $\tilde{\pi}^c$ recursively:*

$$\delta^c(\wp_x, \sigma) = \frac{1}{\|\wp_x \Gamma_\sigma\|_1} \wp_x \Gamma_\sigma \triangleq \wp_{x\sigma} \quad (27)$$

$$\tilde{\pi}^c(\wp_x, \sigma) = \wp_x \tilde{\Pi} \quad (28)$$

For a QSP \mathcal{H} , the canonical representation is denoted as $\mathcal{C}_{\mathcal{H}}$.

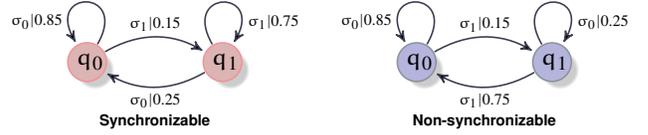


Fig. 2. **Synchronizable and non-synchronizable machines.** Identifying contexts is a key step in estimating the entropy rate of stochastic signals sources; and for PFSA generators, this translates to a state-synchronization problem. However, not all PFSA are synchronizable, e.g., while the top machine is synchronizable, the bottom one is not. Note that a history of just one symbol suffices to determine the current state in the synchronizable machine (top), while no finite history can do the same in the non-synchronizable machine (bottom). However, we show that a ϵ -synchronizable string always exists (Theorem 1).

Lemma 5 (Properties of Canonical Representation). *Given an initial-marked PFSA $G = (Q, \Sigma, \delta, \tilde{\pi}, q_0)$:*

- 1) *The canonical representation is independent of the initial state.*
- 2) *The canonical representation $(Q^c, \Sigma, \delta^c, \tilde{\pi}^c)$ contains a copy of G in the sense that there exists a set of states $Q' \subset Q^c$, such that there exists a one-to-one map $\zeta : Q \rightarrow Q'$, with:*

$$\forall q \in Q, \forall \sigma \in \Sigma, \begin{cases} \tilde{\pi}(q, \sigma) = \tilde{\pi}^c(\zeta(q), \sigma) \\ \delta(q, \sigma) = \delta^c(\zeta(q), \sigma) \end{cases} \quad (29)$$

- 3) *If during the construction (beginning with \wp_λ) we encounter $\wp_x = \zeta(q)$ for some $x \in \Sigma^*$, $q \in Q$ and any map ζ as defined in (2), then we stay within the graph of the copy of the initial-marked PFSA for all right extensions of x .*

Proof: (1) follows the ergodicity of QSPs, which makes \wp_λ independent of the initial state in the initial-marked PFSA.

(2) The canonical representation subsumes the initial-marked representation in the sense that the states of the latter may themselves be seen as degenerate distributions over Q , i.e., by letting

$$\mathcal{E} = \{e^i \in [0, 1]^{|Q|}, i = 1, \dots, |Q|\} \quad (30)$$

denote the set of distributions satisfying:

$$e^i|_{j} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad (31)$$

(3) follows from the strong connectivity of G . ■

Lemma 5 implies that initial states are unimportant; we may denote the initial-marked PFSA induced by a QSP \mathcal{H} , with the initial marking removed, as $\mathcal{P}_{\mathcal{H}}$, and refer to it simply as a “PFSA”. States in $\mathcal{P}_{\mathcal{H}}$ are representable as states in $\mathcal{C}_{\mathcal{H}}$ as elements of \mathcal{E} . Next we show that we always encounter a state arbitrarily close to some element in \mathcal{E} (See Eq. (30)) in the canonical construction starting from the stationary distribution \wp_λ on the states of $\mathcal{P}_{\mathcal{H}}$.

Next we introduce the notion of ϵ -synchronization of probabilistic automata (See Figure 2), which would be of fundamental importance to our entropy estimation algorithm in the next section. Synchronization of automata is fixing or determining the current state; thus it is analogous to contexts in Rissanen’s “context algorithm” [7]. We show that not all PFSA are synchronizable, but all are ϵ -synchronizable.

Theorem 1 (ϵ -Synchronization of Probabilistic Automata). *For any QSP \mathcal{H} over Σ , the PFSA $\mathcal{P}_{\mathcal{H}}$ satisfies:*

$$\forall \epsilon' > 0, \exists x \in \Sigma^*, \exists \theta \in \mathcal{E}, \|\wp_x - \theta\|_\infty \leq \epsilon' \quad (32)$$

Proof: We show that all PFSA are at least approximately synchronizable [16], [17], which is not true for deterministic automata. If the graph of $\mathcal{P}_{\mathcal{H}}$ (i.e., the deterministic automaton obtained by removing the arc probabilities) is synchronizable, then Eq. (32) trivially holds true for $\epsilon' = 0$ for any synchronizing string x . Thus, we assume the graph of $\mathcal{P}_{\mathcal{H}}$ to be non-synchronizable. From definition of non-synchronizability, it follows:

$$\forall q_i, q_j \in Q, \text{ with } q_i \neq q_j, \forall x \in \Sigma^*, \delta(q_i, x) \neq \delta(q_j, x) \quad (33)$$

If the PFSA has a single state, then every string satisfies the condition in Eq. (32). Hence, we assume that the PFSA has more than one state. Now if we have:

$$\forall x \in \Sigma^*, \frac{\Pr(x'x)}{\Pr(x')} = \frac{\Pr(x''x)}{\Pr(x'')} \text{ where } [x'] = q_i, [x''] = q_j \quad (34)$$

then, by the Definition 3, we have a contradiction $q_i = q_j$. Hence $\exists x_0$ such that

$$\frac{\Pr(x'x_0)}{\Pr(x')} \neq \frac{\Pr(x''x_0)}{\Pr(x'')} \text{ where } [x'] = q_i, [x''] = q_j \quad (35)$$

$$\text{Since : } \sum_{x \in \Sigma^*} \frac{\Pr(x'/x)}{\Pr(x')} = 1, \text{ for any } x' \text{ where } [x'] = q_i \quad (36)$$

we conclude without loss of generality $\forall q_i, q_j \in Q$, with $q_i \neq q_j$:

$$\exists x^{ij} \in \Sigma^*, \frac{\Pr(x'x^{ij})}{\Pr(x')} > \frac{\Pr(x''x^{ij})}{\Pr(x'')} \text{ where } [x'] = q_i, [x''] = q_j$$

It follows from induction that if we start with a distribution \wp on Q such that $\wp_i = \wp_j = 0.5$, then for any $\epsilon' > 0$ we can construct a finite string x_0^{ij} such that if $\delta(q_i, x_0^{ij}) = q_r, \delta(q_j, x_0^{ij}) = q_s$, then for the new distribution \wp' after execution of x_0^{ij} will satisfy $\wp'_s > 1 - \epsilon'$. Recalling that $\mathcal{P}_{\mathcal{J}C}$ is strongly connected, we note that, for any $q_t \in Q$, there exists a string $y \in \Sigma^*$, such that $\delta(q_s, y) = q_t$. Setting $x_1^{ij \rightarrow t} = x_0^{ij} y$, we can ensure that the distribution \wp'' obtained after execution of x_1^{ij} satisfies $\wp''_t > 1 - \epsilon'$ for any q_t of our choice. For arbitrary initial distributions \wp^A on Q , we must consider contributions arising from simultaneously executing $x_1^{ij \rightarrow t}$ from states other than just q_i and q_j . Nevertheless, it is easy to see that executing $x_1^{ij \rightarrow t}$ implies that in the new distribution $\wp^{A'}$, we have $\wp^{A'}_t > \wp^A_i + \wp^A_j - \epsilon'$. It follows that executing the string $x^{1,2 \rightarrow |Q|} x^{3,4 \rightarrow |Q|} \dots x^{n-1,n \rightarrow |Q|}$, where

$$n = \begin{cases} |Q| & \text{if } |Q| \text{ is even} \\ |Q| - 1 & \text{otherwise} \end{cases} \quad (37)$$

would result in a final distribution $\wp^{A''}$ which satisfies $\wp^{A''}_i > 1 - \frac{1}{2}\epsilon'$. Appropriate scaling of ϵ' then completes the proof. \blacksquare

Theorem 1 induces the notion of ϵ -synchronizing strings, and guarantees their existence for arbitrary PFSA.

Definition 7 (ϵ -synchronizing Strings). A string $x \in \Sigma^*$ is ϵ -synchronizing for a PFSA if:

$$\exists \vartheta \in \mathcal{E}, \|\wp_x - \vartheta\|_\infty \leq \epsilon \quad (38)$$

Theorem 1 is an existential result, and does not yield an algorithm for computing synchronizing strings (See Theorem 3). We may estimate an asymptotic upper bound on such a search.

Corollary 1 (To Theorem 1). At most $O(1/\epsilon)$ strings from the lexicographically ordered set of all strings over the given alphabet need to be analyzed to find an ϵ -synchronizing string.

Proof: Theorem 1 works by multiplying entries from the $\tilde{\Pi}$ matrix, which cannot be all identical (otherwise the states would collapse). Let the minimum difference between two unequal entries be η . Then, following the construction in Theorem 1, the length ℓ of the synchronizing string, up to linear scaling, satisfies: $\eta^\ell = O(\epsilon)$, implying $\ell = O(\log(1/\epsilon))$. Hence, the number of strings to be analyzed is at most all strings of length ℓ , where $|\Sigma|^\ell = |\Sigma|^{O(\log(1/\epsilon))} = O(1/\epsilon)$. \blacksquare

B. Symbolic Derivatives

Computation of ϵ -synchronizing strings requires the notion of symbolic derivatives. Note that, PFSA states are not observable; we observe symbols generated from hidden states. A symbolic derivative at a given string specifies the distribution of the next symbol over the alphabet.

Notation 4. We denote the set of probability distributions over a finite set of cardinality k as $\mathcal{D}(k)$.

Definition 8 (Symbolic Count Function). For a string s over Σ , the count function $\#^s : \Sigma^* \rightarrow \mathbb{N} \cup \{0\}$, counts the number of times a particular substring occurs in s . The count is overlapping, i.e., in a string $s = 0001$, we count the number of occurrences of 00s as $\underline{00}01$ and $0\underline{00}1$, implying $\#^s 00 = 2$.

Definition 9 (Symbolic Derivative). For a string s generated by a QSP over Σ , the symbolic derivative $\phi^s : \Sigma^* \rightarrow \mathcal{D}(|\Sigma| - 1)$ is defined:

$$\phi^s(x)|_i = \frac{\#^s x \sigma_i}{\sum_{\sigma_i \in \Sigma} \#^s x \sigma_i} \quad (39)$$

Thus, $\forall x \in \Sigma^*$, $\phi^s(x)$ is a probability distribution over Σ . $\phi^s(x)$ is referred to as the symbolic derivative at x .

Note that $\forall q_i \in Q$, $\tilde{\pi}$ induces a probability distribution over Σ as $[\tilde{\pi}(q_i, \sigma_1), \dots, \tilde{\pi}(q_i, \sigma_{|\Sigma|})]$. We denote this as $\tilde{\pi}(q_i, \cdot)$.

We next show that the symbolic derivative at x can be used to estimate this distribution for $q_i = [x]$, provided x is ϵ -synchronizing.

Theorem 2 (ϵ -Convergence). If $x \in \Sigma^*$ is ϵ -synchronizing, then:

$$\forall \epsilon > 0, \lim_{|s| \rightarrow \infty} \|\phi^s(x) - \tilde{\pi}([x], \cdot)\|_\infty \leq_{a.s} \epsilon \quad (40)$$

Proof: We use the Glivenko-Cantelli theorem [18] on uniform convergence of empirical distributions. Since x is ϵ -synchronizing:

$$\forall \epsilon > 0, \exists \vartheta \in \mathcal{E}, \|\wp_x - \vartheta\|_\infty \leq \epsilon \quad (41)$$

Recall that $\mathcal{E} = \{e^i \in [0, 1]^{|Q|}, i = 1, \dots, |Q|\}$ denotes the set of distributions over Q satisfying:

$$e^i|_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad (42)$$

Let x ϵ -synchronize to $q \in Q$. Thus, when we encounter x while reading s , we are guaranteed to be distributed over Q as \wp_x , where:

$$\|\wp_x - \vartheta\|_\infty \leq \epsilon \Rightarrow \wp_x = \alpha \vartheta + (1 - \alpha)u \quad (43)$$

where $\alpha \in [0, 1]$, $\alpha \geq 1 - \epsilon$, and u is an unknown distribution over Q . Defining $A_\alpha = \alpha \tilde{\pi}(q, \cdot) + (1 - \alpha) \sum_{j=1}^{|Q|} u_j \tilde{\pi}(q_j, \cdot)$, we note that $\phi^s(x)$ is an empirical distribution for A_α , implying:

$$\begin{aligned} \lim_{|s| \rightarrow \infty} \|\phi^s(x) - \tilde{\pi}(q, \cdot)\|_\infty &= \lim_{|s| \rightarrow \infty} \|\phi^s(x) - A_\alpha + A_\alpha - \tilde{\pi}(q, \cdot)\|_\infty \\ &\stackrel{\text{a.s. 0 by Glivenko-Cantelli}}{\leq} \overbrace{\lim_{|s| \rightarrow \infty} \|\phi^s(x) - A_\alpha\|_\infty}^{\leq \epsilon} + \lim_{|s| \rightarrow \infty} \|A_\alpha - \tilde{\pi}(q, \cdot)\|_\infty \\ &\leq_{a.s} (1 - \alpha) (\|\tilde{\pi}(q, \cdot) - u\|_\infty) \leq_{a.s} \epsilon \end{aligned}$$

This completes the proof. \blacksquare

C. Computation of ϵ -synchronizing Strings

Next we describe identification of ϵ -synchronizing strings given a sufficiently long observed string (i.e. a sample path) s . Theorem 1 guarantees existence, and Corollary 1 establishes that $O(1/\epsilon)$ substrings need to be analyzed till we encounter an ϵ -synchronizing string. These do not provide an executable algorithm, which arises from an inspection of the geometric structure of the set of probability vectors over Σ , obtained by constructing $\phi^s(x)$ for different choices of the candidate string x .

Definition 10 (Derivative Heap). Given a string s generated by a QSP, a derivative heap $\mathcal{D}^s : 2^{\Sigma^*} \rightarrow \mathcal{D}(|\Sigma| - 1)$ is the set of probability distributions over Σ calculated for a subset of strings $L \subset \Sigma^*$ as:

$$\mathcal{D}^s(L) = \{\phi^s(x) : x \in L \subset \Sigma^*\} \quad (44)$$

Lemma 6 (Limiting Geometry). Let us define:

$$\mathcal{D}_\infty = \lim_{|s| \rightarrow \infty} \lim_{L \rightarrow \Sigma^*} \mathcal{D}^s(L) \quad (45)$$

If \mathcal{U}_∞ is the convex hull of \mathcal{D}_∞ , and u is a vertex of \mathcal{U}_∞ , then

$$\exists q \in Q, \text{ such that } u = \tilde{\pi}(q, \cdot) \quad (46)$$

Proof: Recalling Theorem 2, the result follows from noting that any element of \mathcal{D}_∞ is a convex combination of elements from the set $\{\tilde{\pi}(q_1, \cdot), \dots, \tilde{\pi}(q_{|Q|}, \cdot)\}$. \blacksquare

Lemma 6 does not claim that the number of vertices of the convex hull of \mathcal{D}_∞ equals the number of states, but that every vertex corresponds to a state. We cannot generate \mathcal{D}_∞ since we have a finite observed string s , and we can calculate $\phi^s(x)$ for a finite number of x . Instead, we show that choosing a string corresponding to the vertex of the convex hull of the heap, constructed by considering $O(1/\epsilon)$ strings, gives us an ϵ -synchronizing string with high probability.

Theorem 3 (Derivative Heap Approx.). For s generated by a QSP, let $\mathcal{D}^s(L)$ be computed with $L = \Sigma^{O(\log(1/\epsilon))}$. If for $x_0 \in \Sigma^{O(\log(1/\epsilon))}$, $\phi^s(x_0)$ is a vertex of the convex hull of $\mathcal{D}^s(L)$, then

$$\text{Prob}(x_0 \text{ is not } \epsilon\text{-synchronizing}) \leq e^{-|s| \epsilon p_0} \quad (47)$$

where p_0 is the probability of encountering x_0 in s .

Proof: The result follows from Sanov's Theorem [19] for convex set of probability distributions. If $|s| \rightarrow \infty$, then x_0 is guaranteed to be ϵ -synchronizing (Theorem 1, and Corollary 1). Denoting the number of times we encounter x_0 in s as $n(|s|)$, and since \mathcal{D}_∞ is a convex set of distributions (allowing us to drop the polynomial factor in Sanov's bound), we apply Sanov's Theorem to the case of finite s :

$$\text{Prob}\left(\text{KL}(\phi^s(x_0) \parallel \wp_{x_0} \tilde{\Pi}) > \epsilon\right) \leq e^{-n(|s|) \epsilon} \quad (48)$$

where $\text{KL}(\cdot|\cdot)$ is the Kullback-Leibler divergence [20]. From the bound [21]:

$$\frac{1}{4} \|\phi^s(x_0) - \wp_{x_0} \tilde{\Pi}\|_\infty^2 \leq \text{KL}(\phi^s(x_0) \|\wp_{x_0} \tilde{\Pi}) \quad (49)$$

and $n(|s|) \rightarrow |s|p_0$, where $p_0 > 0$ is the stationary probability of encountering x_0 in s , we conclude:

$$\text{Prob}(\|\phi^s(x_0) - \wp_{x_0} \tilde{\Pi}\|_\infty > \epsilon) \leq 2e^{-\frac{1}{2}|s|\epsilon p_0} \quad (50)$$

which completes the proof. ■

III. ENTROPY RATE FOR PFSA-GENERATED PROCESSES

Given the PFSA model, the entropy rate is easily computable.

Theorem 4 (Entropy Rate For PFSA). *The entropy rate $H(G)$, in bits, for the QSP generated by a PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$ is given by:*

$$H(G) = \sum_{i=1}^{|\mathcal{Q}|} \wp_{\lambda|_i} \sum_{\sigma_j \in \Sigma} \tilde{\pi}(q_i, \sigma_j) \log \tilde{\pi}(q_i, \sigma_j) \quad (51)$$

where the base of the logarithms is 2.

Proof: Denote the QSP generated by G as $\mathcal{X} = \{X_i\}$. Using the chain rule (See Eq. (3)), we have:

$$H(G) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1) \quad (52)$$

Since G is always at some state $q \in Q$, we conclude that for any i :

$$H(X_i | X_{i-1}, \dots, X_1) \in \left\{ \sum_{\sigma_j \in \Sigma} \tilde{\pi}(q, \sigma_j) \log \tilde{\pi}(q, \sigma_j) : q \in Q \right\}$$

Furthermore, since G is strongly connected, and therefore has a unique stationary distribution \wp_λ [22], the number of times state q_i occurs approaches $n\wp_{\lambda|_i}$ as $n \rightarrow \infty$. This completes the proof. ■

If the underlying PFSA model is not available, and we have only a symbolic stream generated by a QSP, then Eq. (51) cannot be directly employed to estimate the entropy rate. In that case, one possibility is to first infer the hidden PFSA using the algorithm reported in [12], and then estimate the entropy rate from Eq. (51). However, if we are only interested in the latter, then we do not need to infer the complete generative model; and there exists a more parsimonious approach to estimate the entropy rate directly.

First, we need a lemma which bounds the deviation in entropy for deviations in the probability distribution in the discrete case.

Lemma 7 (Bound on Entropy Deviation). *For probability distributions p, q on a finite set Σ , we have for all $\epsilon \in (0, 1)$,*

$$\|p - q\|_\infty \leq \epsilon \Rightarrow$$

$$|H(p) - H(q)| < \epsilon' \log \frac{|\Sigma| - 1}{\epsilon'} + (1 - \epsilon') \log \frac{1}{1 - \epsilon'} \quad (53)$$

$$\text{where } \epsilon' = \begin{cases} \epsilon & \text{if } \epsilon \leq 1/2 \\ 1 - \epsilon & \text{otherwise} \end{cases}$$

where $H(p), H(q)$ are entropies for distributions p, q respectively.

Proof: We have from definition:

$$H(p) - H(q) = \sum_i p_i \log \frac{1}{p_i} - \sum_i q_i \log \frac{1}{q_i}$$

We note that the function $f(x) = x \log \frac{1}{x}$ satisfies:

$$\delta f = \left(\log \frac{1}{x} - \frac{1}{\ln 2} \right) \delta x \quad (53)$$

implying that perturbations of x cause maximum change in f , when x is in the neighborhood of 0, which in turn implies that deviation in entropy for a perturbed distribution p is the maximized when:

$$p \rightarrow p^* = (0 \quad \dots \quad 0 \quad 1) \text{ upto permutations} \quad (54)$$

Since, $\|p - q\|_\infty \leq \epsilon$, the perturbed distribution q from $p = p^*$ is non-unique. We claim (Claim A), that the perturbed distribution resulting in maximum entropy deviation, is given by:

$$q^* = \left(\frac{\epsilon'}{|\Sigma| - 1} \quad \dots \quad \frac{\epsilon'}{|\Sigma| - 1} \quad 1 - \epsilon' \right) \text{ upto permutations} \quad (55)$$

$$\text{where } \epsilon' = \begin{cases} \epsilon & \text{if } \epsilon \leq 1/2 \\ 1 - \epsilon & \text{otherwise} \end{cases} \quad (56)$$

To establish this claim, we first note that q^* satisfies the constraints:

$$\forall i \ q_i^* > 0, \sum_i q_i^* = 1, \|p^* - q^*\|_\infty = \epsilon \quad (57)$$

Let q' be a perturbation of q^* , defined as:

$$q'_i = q_i^* + \alpha_i, \text{ with } \sum_i \alpha_i = 0 \quad (58)$$

satisfying the constraint:

$$\|q' - p^*\|_\infty \leq \epsilon \quad (59)$$

Note that the above constraint, and the definition of q^* implies that:

$$\alpha_{|\Sigma|} \geq 0 \quad (60)$$

Then we claim that for small perturbations,

$$H(q') < H(q^*) \quad (61)$$

We find differential perturbations in contribution to the entropy from perturbation of each entry in q^* . For terms $i \in \{1, \dots, |\Sigma| - 1\}$, we note that the perturbed term is of the form:

$$g(x) = \frac{\epsilon' + x}{|\Sigma| - 1} \log \frac{|\Sigma| - 1}{\epsilon' + x} \quad (62)$$

$$\Rightarrow \delta g(0) = \frac{1}{|\Sigma| - 1} \left(\log \frac{|\Sigma| - 1}{\epsilon'} - \ln 2 \right) \alpha_i \quad (63)$$

if α_i is small. And the $|\Sigma|$ -th term is of the form:

$$f(x) = (1 - \epsilon' + x) \log \frac{1}{1 - \epsilon' + x} \quad (64)$$

$$\Rightarrow \delta f(0) = \left(\log \frac{1}{1 - \epsilon'} - \ln 2 \right) \alpha_{|\Sigma|} \quad (65)$$

if $\alpha_{|\Sigma|}$ is small. This implies that the perturbation of entropy, for small perturbations in the distribution q^* , is given by:

$$\begin{aligned} \delta H(q^*) &= \frac{1}{|\Sigma| - 1} \left(\log \frac{|\Sigma| - 1}{\epsilon'} - \ln 2 \right) \sum_{i=1}^{|\Sigma| - 1} \alpha_i \\ &\quad + \left(\log \frac{1}{1 - \epsilon'} - \ln 2 \right) \alpha_{|\Sigma|} \end{aligned} \quad (66)$$

Noting that $\sum_{i=1}^{|\Sigma| - 1} \alpha_i = -\alpha_{|\Sigma|}$, and setting $b = |\Sigma| - 1$, we have:

$$\begin{aligned} \delta H(q^*) &= \left(-\frac{1}{b} \left(\log \frac{b}{\epsilon'} - \ln 2 \right) + \left(\log \frac{1}{1 - \epsilon'} - \ln 2 \right) \right) \alpha_{|\Sigma|} \\ &= \left(\underbrace{\left(\frac{1}{b} - 1 \right) \ln 2}_{t_1} + \underbrace{\log \frac{1}{1 - \epsilon'} - \frac{1}{b} \log \frac{b}{\epsilon'}}_{t_2} \right) \alpha_{|\Sigma|} \end{aligned} \quad (67)$$

We note that since $|\Sigma| \geq 2$, $t_1 \leq 0$. Then, since $\epsilon' \leq 1/2$, we have:

$$\log \frac{1}{1 - \epsilon'} \leq 1 \text{ with equality for } \epsilon' = 1/2 \quad (68)$$

And we note that $\frac{1}{b} \log \frac{b}{\epsilon'}$ attains its minimum value of $1/b + 1/b \log b$ at $\epsilon' = 1/2$, implying:

$$\delta H(q^*) \leq \left(\frac{1}{b} - 1 \right) (\ln 2 - 1) \leq 0 \quad (69)$$

This establishes that within the set of admissible perturbed distributions q' , from q^* , all infinitesimally small perturbations necessarily reduce the entropy, i.e., $H(q^*)$ attains a locally maximum value.

We note that for all arbitrary admissible perturbations q' from q^* , $\|q' - p^*\|_\infty \leq \epsilon$ and definition of ϵ' implies that each entry in q' is either always in $[0, 1/2]$, or in $[1/2, 1]$, and not both. Noting that each summand in the calculation of entropy is of the form $x \log x$, which is monotonic in both intervals, we conclude that $H(q^*)$ is indeed the globally maximum entropy within all admissible perturbations q' . It follows that any perturbation of q^* , satisfying the constraint of Eq. (59), leads to a smaller difference of entropy from p^* , which establishes claim A. Noting that:

$$|H(p^*) - H(q^*)| = \epsilon' \log \frac{|\Sigma| - 1}{\epsilon'} + (1 - \epsilon') \log \frac{1}{1 - \epsilon'}$$

completes the proof. ■

This bound on entropy deviation for ∞ -norm bounded deviations in distribution will be important in the sequel. We denote this as the generalized binary entropy function $\mathbb{B}(\epsilon, |\Sigma|)$.

Definition 11 (Generalized Binary Entropy Function).

$$\mathbb{B}(\epsilon, |\Sigma|) = \epsilon' \log \frac{|\Sigma| - 1}{\epsilon'} + (1 - \epsilon') \log \frac{1}{1 - \epsilon'} \quad (70)$$

$$\text{where } \epsilon' = \begin{cases} \epsilon & \text{if } \epsilon \leq 1/2 \\ 1 - \epsilon & \text{otherwise} \end{cases}$$

Corollary 2 (To Lemma 7). *Given a symbol stream generated by a*

PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$, and an ϵ -synchronizing string x_0 , we have:

$$\left| \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0x)) - H(G) \right| < \mathbb{B}(\epsilon, |\Sigma|)$$

Proof: We first establish the following claim (Claim A): x_0 is ϵ -synchronizing implies that any right extension x_0x is also ϵ -synchronizing (where $x \in \Sigma^*$). To see this, note that x_0 is ϵ -synchronizing implies $\exists \vartheta \in \mathcal{E}$ with:

$$\varrho_{x_0} = \alpha\vartheta + (1 - \alpha)u, \text{ with } \alpha \in [0, 1], \alpha \geq 1 - \epsilon \quad (71)$$

where u is an unknown distribution over Q . It follows that: $\forall \sigma \in \Sigma$,

$$\varrho_{x_0\sigma} = \frac{1}{\|\varrho_{x_0}\Gamma_\sigma\|_1} (\alpha\vartheta\Gamma_\sigma + (1 - \alpha)u\Gamma_\sigma) \quad (72)$$

Now, $\forall \sigma \in \Sigma$, there is a unique $\vartheta' \in \mathcal{E}$, such that $\vartheta' = \vartheta\Gamma_\sigma$, and since $\|\varrho_{x_0}\Gamma_\sigma\|_1 \leq 1$, it follows that: $\forall \sigma \in \Sigma, \exists \vartheta' \in \mathcal{E}$, such that:

$$\varrho_{x_0\sigma} = \alpha'\vartheta' + \text{additional terms}, \text{ with } \alpha' \in [0, 1], \alpha' \geq 1 - \epsilon$$

By straightforward induction, we conclude that:

$$\forall x \in \Sigma^*, \exists \vartheta(x) \in \mathcal{E}, \text{ such that } \|\varrho_{x_0x} - \vartheta(x)\|_\infty \leq \epsilon \quad (73)$$

which establishes Claim A.

Next we claim (Claim B) that $H(G)$ can be written as:

$$H(G) = \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} H(\tilde{\pi}([x], \cdot)) \quad (74)$$

To see this, note that the PFSA G is strongly connected with a unique stationary distribution ϱ_λ , and Theorem 4 implies:

$$H(G) = \sum_{i=1}^{|\mathcal{Q}|} \varrho_{\lambda|_i} H(\tilde{\pi}(q_i, \cdot)) \quad (75)$$

Set the initial state of G to be $q \in Q$, where x_0 ϵ -synchronizes to q . For any n , and each $x \in \Sigma^n$, $[x]$ is the equivalence class corresponding to some $q_i \in Q$. Let the number of times $[x]$ corresponds to q_i , for $x \in \Sigma^n$, be n_i . Then, uniqueness of ϱ_λ implies that $\lim_{n \rightarrow \infty} n_i/n = \varrho_{\lambda|_i}$, which implies:

$$\sum_{i=1}^{|\mathcal{Q}|} \varrho_{\lambda|_i} H(\tilde{\pi}(q_i, \cdot)) = \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} H(\tilde{\pi}([x], \cdot)) \quad (76)$$

establishing Claim B. Thus, we can write:

$$\left| \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0x)) - H(G) \right| \quad (77)$$

$$= \left| \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \left\{ \lim_{|s| \rightarrow \infty} H(\phi^s(x_0x)) - H(\tilde{\pi}([x], \cdot)) \right\} \right| \quad (78)$$

$$\leq \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \left| \lim_{|s| \rightarrow \infty} H(\phi^s(x_0x)) - H(\tilde{\pi}([x], \cdot)) \right| \quad (79)$$

We note that Claim A implies that x_0x ϵ -synchronizes to $[x]$ in G , which then implies from Theorem 2, and Lemma 7:

$$\lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \left| \lim_{|s| \rightarrow \infty} H(\phi^s(x_0x)) - H(\tilde{\pi}([x], \cdot)) \right| \quad (80)$$

$$< \lim_{n \rightarrow \infty} \frac{1}{|\Sigma^n|} \sum_{x \in \Sigma^n} \mathbb{B}(\epsilon, |\Sigma|) \quad (81)$$

which completes the proof. \blacksquare

Next we modify the Dvoretzky-Kiefer-Wolfowitz inequality, to be applicable to the case where the number of samples drawn is itself a random variable.

Lemma 8 (DKW-bound for symbolic derivatives). *For a string s generated by a PFSA, and a given ϵ -synchronizing string x_0 :*

$\forall x \in \Sigma^*$ such that x_0x occurs in s with probability $\zeta > 0$,

$$\Pr \left(\left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty > \epsilon \right) < 8 \left(1 + \frac{1}{e} \right) e^{-|\zeta| \frac{\epsilon^2}{1 + \epsilon^2}}$$

Proof: We note that $\phi^s(x_0x)$ is an empirical distribution with the limiting distribution given by $\lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \triangleq \phi^*$. Using the DKW inequality [23], and denoting the number of occurrences of x_0x in s with the random variable N_{x_0x} , we have:

$$\Pr \left(\left\{ \left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty > \epsilon \right\} \wedge \{N_{x_0x} = n'\} \right)$$

$$\begin{aligned} & \leq 2e^{-2\epsilon^2 n'} \Pr(\{N_{x_0x} = n'\}) \\ \Rightarrow \Pr \left(\left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty > \epsilon \right) & \leq \sum_{n' \in \mathbb{N}} 2e^{-2\epsilon^2 n'} \Pr(\{N_{x_0x} = n'\}) \quad (82) \end{aligned}$$

We partition \mathbb{N} into disjoint sets U_r and $V_r = \mathbb{N} \setminus U_r$, parametrized by $r > 0$, where:

$$U_r = \left[\left[|s|\zeta(1-r) \right], \left[|s|\zeta(1+r) \right] \right] \quad (83)$$

Using Chernoff bounds for the probability of $n' \in V_r$, we have:

$$\begin{aligned} \Pr \left(\left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty > \epsilon \right) & \leq \sum_{n' \in U_r} 2e^{-2\epsilon^2 n'} \Pr(\{N_{x_0x} = n'\}) \\ & \quad + \sum_{n' \in V_r} 2e^{-2\epsilon^2 n'} \Pr(\{N_{x_0x} = n'\}) \\ & \leq \left(\lceil 2|s|\zeta r \rceil \times 2e^{-2\epsilon^2 |s|\zeta(1-r)} \times 1 \right) + \left(1 \times 2e^{-\frac{r^2 |s|\zeta}{2+r}} \right) \\ & \leq 4 \lceil |s|\zeta r \rceil e^{-2\epsilon^2 |s|\zeta(1-r)} + 2e^{-\frac{r^2 |s|\zeta}{2+r}} \quad (84) \end{aligned}$$

Denoting $|s|\zeta$ as t , we have the bound:

$$\forall r > 0, f(r) = 4 \lceil rt \rceil e^{-2\epsilon^2 t(1-r)} + 2e^{-\frac{r^2 t}{2+r}} \quad (85)$$

We note that the two terms are equal if:

$$2\epsilon^2 t(1-r) = \frac{r^2 t}{2+r} + \ln(2 \lceil rt \rceil) \quad (86)$$

It follows that if we solve for r in terms of ϵ after dropping the non-negative log-term, then the first term would be bigger or equal compared to the second. Solving the resulting quadratic, we get:

$$r < \frac{\epsilon^2}{1 + \epsilon^2} \quad (87)$$

A larger value of r makes the first term larger, and the second term smaller; hence we use $r = \frac{\epsilon^2}{1 + \epsilon^2}$, leading to the non-tight bound:

$$f(r) < 8 \left[\frac{\epsilon^2}{1 + \epsilon^2} t \right] e^{-2 \frac{\epsilon^2}{1 + \epsilon^2} t} < 8 \left(1 + \frac{\epsilon^2}{1 + \epsilon^2} t \right) e^{-2 \frac{\epsilon^2}{1 + \epsilon^2} t} \quad (88)$$

Using the fact that $\forall y \in \mathbb{R}, 1 - y \leq e^{-y}$, we have:

$$f(r) < 8e^{-2 \frac{\epsilon^2}{1 + \epsilon^2} t} + 8e^{-\frac{\epsilon^2}{1 + \epsilon^2} t - 1} < 8 \left(1 + \frac{1}{e} \right) e^{-\frac{\epsilon^2}{1 + \epsilon^2} t} \quad (89)$$

which completes the proof. \blacksquare

Corollary 3 (To Lemma 8). *For s generated by a PFSA, and an ϵ -synchronizing x_0 , we have for any $x \in \Sigma^*$:*

$$\begin{aligned} \Pr \left(\left| H(\phi^s(x_0x)) - H \left(\lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right) \right| > \mathbb{B}(\epsilon, |\Sigma|) \right) \\ < 8 \left(1 + \frac{1}{e} \right) e^{-|\zeta| \frac{\epsilon^2}{1 + \epsilon^2}} \end{aligned}$$

Proof: It follows from Lemma 7 and continuity of entropy that

$$\begin{aligned} \left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty \leq \epsilon \\ \Rightarrow \left| H(\phi^s(x_0x)) - H \left(\lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right) \right| \leq \mathbb{B}(\epsilon, |\Sigma|) \end{aligned}$$

Using Lemma 8, we have:

$$\begin{aligned} \Pr \left(\left| \phi^s(x_0x) - \lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right|_\infty \leq \epsilon \right) & \geq 1 - 8 \left(1 + \frac{1}{e} \right) e^{-|\zeta| \frac{\epsilon^2}{1 + \epsilon^2}} \\ \Rightarrow \Pr \left(\left| H(\phi^s(x_0x)) - H \left(\lim_{|s'| \rightarrow \infty} \phi^{s'}(x_0x) \right) \right| \leq \mathbb{B}(\epsilon, |\Sigma|) \right) \\ & \geq 1 - 8 \left(1 + \frac{1}{e} \right) e^{-|\zeta| \frac{\epsilon^2}{1 + \epsilon^2}} \end{aligned}$$

which completes the proof. \blacksquare

Theorem 5 (Bound on Entropy Calculation with Finite Samples). *For any string x generated by a PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$, and a given ϵ -synchronizing string x_0 , there exist C_0, C_1 depending only on the size of the alphabet $|\Sigma|$, such that, for any independently chosen set*

of strings $\mathcal{N} \subseteq \Sigma^*$:

$$\Pr \left(\left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - \lim_{n \rightarrow \infty} \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| > \mathbb{B}(\epsilon, |\Sigma|) + \epsilon \right) \leq C_0 \frac{1 + \epsilon^2}{|s|e^3} + 2e^{-C_1 |\mathcal{N}| \epsilon^2}$$

Proof: We note that:

$$\begin{aligned} A &= \left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - \lim_{n \rightarrow \infty} \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| \\ &\leq \left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| \\ &\quad + \left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) - \lim_{n \rightarrow \infty} \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| \end{aligned}$$

We denote the two RHS terms as B and C, and note:

$$\begin{aligned} B &\triangleq \left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| \\ &= \left| \sum_{i=1}^{|\mathcal{Q}|} \tilde{\rho}_i \left(H(\phi^s(x_0 x'_i)) - \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x'_i)) \right) \right| \quad (90) \end{aligned}$$

where x_0 ϵ -synchronizes to $q_0 \in \mathcal{Q}$, $\delta(q_0, x'_i) = q_i$, and $\tilde{\rho}$ is the empirical estimate of the stationary distribution. Using the bound from Corollary 3:

$$\begin{aligned} \Pr(B > \mathbb{B}(\epsilon, |\Sigma|)) &< 8 \left(1 + \frac{1}{e}\right) \sum_{i=1}^{|\mathcal{Q}|} \tilde{\rho}_i e^{-|s| \tilde{\rho}_i \frac{\epsilon^2}{1 + \epsilon^2}} \\ &< 8 \left(1 + \frac{1}{e}\right) \frac{(1 + \epsilon^2) |\mathcal{Q}|}{e |s| \epsilon^2} \quad (91) \end{aligned}$$

For the second RHS term:

$$\begin{aligned} C &\triangleq \left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) - \lim_{n \rightarrow \infty} \frac{1}{|\Sigma_n^+|} \sum_{x \in \Sigma_n^+} \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x)) \right| \\ &= \left| \sum_{i=1}^{|\mathcal{Q}|} (\tilde{\rho}_i - \rho_{\lambda_i}) \lim_{|s| \rightarrow \infty} H(\phi^s(x_0 x'_i)) \right| \leq \|\tilde{\rho} - \rho_{\lambda}\|_1 \log |\Sigma| \quad (92) \end{aligned}$$

where x_0 ϵ -synchronizes to $q_0 \in \mathcal{Q}$ and $\delta(q_0, x'_i) = q_i$. Using DKW:

$$\begin{aligned} \Pr(\|\tilde{\rho} - \rho_{\lambda}\|_{\infty} > \epsilon) &\leq 2e^{-2|\mathcal{N}| \epsilon^2} \\ \Rightarrow \Pr(\|\tilde{\rho} - \rho_{\lambda}\|_1 \log |\Sigma| \leq \epsilon) &> 1 - 2e^{-\frac{2}{\log^2 |\Sigma|} |\mathcal{N}| \epsilon^2} \quad (93) \end{aligned}$$

Using the bounds in Eq. (91), and (93), we get:

$$\begin{aligned} E &\triangleq \Pr(B + C \leq \mathbb{B}(\epsilon, |\Sigma|) + \epsilon) \\ &> \left(1 - 8 \left(1 + \frac{1}{e}\right) \frac{(1 + \epsilon^2) |\mathcal{Q}|}{e |s| \epsilon^2}\right) \times \left(1 - 2e^{-\frac{2}{\log^2 |\Sigma|} |\mathcal{N}| \epsilon^2}\right) \end{aligned}$$

Since we are using ϵ -synchronization, it follows that the number of states $|\mathcal{Q}|$ is upper bounded by $(|\Sigma| - 1)/\epsilon$ which then yields:

$$E > \left(1 - C_0 \frac{1 + \epsilon^2}{|s| e^3}\right) \left(1 - 2e^{-C_1 |\mathcal{N}| \epsilon^2}\right) \quad (94)$$

with $C_0 = (8/e + 8/e^2)(|\Sigma| - 1)$, and $C_1 = \frac{2}{\log^2 |\Sigma|}$

$$\Rightarrow \Pr(A \leq \mathbb{B}(\epsilon, |\Sigma|) + \epsilon) > 1 - C_0 \frac{1 + \epsilon^2}{|s| e^3} - 2e^{-C_1 |\mathcal{N}| \epsilon^2}$$

which completes the proof. ■

Theorem 6 (Main Theorem). *Given a finite string s generated by a PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$, and a string $x_0 \in \Sigma^*$ satisfying the pre-conditions described in Theorem 3, we have for any independently chosen set of strings $\mathcal{N} \subseteq \Sigma^*$:*

$$\begin{aligned} \Pr \left(\left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - H(G) \right| > \epsilon + 2\mathbb{B}(\epsilon, |\Sigma|) \right) \\ \leq C_0 \frac{1 + \epsilon^2}{|s| e^3} + 2e^{-C_1 |\mathcal{N}| \epsilon^2} + e^{-\epsilon p_0 |s|} \quad (95) \end{aligned}$$

where $C_0 = (8/e + 8/e^2)(|\Sigma| - 1)$, $C_1 = 2/\log^2 |\Sigma|$ and p_0 is the non-zero occurrence probability of x_0 in s .

Algorithm 1: Detailed pseudocode for entropy rate estimation

Input: Data sequence s over alphabet Σ , ϵ , Confidence level α
Output: Entropy rate \mathbf{h} , Uncertainty \mathbf{E} at specified confidence level

- 1 Initialize $\mathbf{h} = 0$
- 2 Initialize $\text{COUNT}_{\text{total}} = 0$
- 3 Initialize $\text{COUNT}_{\text{map}} = \emptyset$ /* hashtable with keys as probability distributions, and values as doubles */
- 4 Set $C_0 = (8/e + 8/e^2)(|\Sigma| - 1)$, $C_1 = 2/\log^2 |\Sigma|$
- 5 Set $N_{\text{min}} = 10$ /* Any small integer suffices (See Section IV) */

/* I. ϵ -synchronization String Identification */

- 7 **foreach** $x \in \Sigma_n^+$ **do**
- 8 $D[x] \leftarrow \phi^s(x)$
- 9 $A \leftarrow \{x' : D[x'] \text{ is on the convex hull of the set of values in hashtable } D\}$
- 10 $x_0 \leftarrow \text{argmax}_{x \in A} \#^s x$ /* ϵ -synchronization string */
- 11 $p_0 \leftarrow (\#^s x_0) |s|$ /* Occurrence prob. of ϵ -synchronization string */

/* II. Entropy Rate Estimation */

- 13 Select $\mathcal{N} \subset \Sigma^*$ with length ℓ strings drawn with probability $\frac{1}{|\Sigma| \ell}$
- /* $|\mathcal{N}| \sim 10^7 \log^2 |\Sigma|$ sufficient for negligible uncertainty contribution */
- 14 **foreach** $x \in \mathcal{N}$ **do**
- 15 **if** $\#^s x_0 x > N_{\text{min}}$ **then**
- 16 Compute $\mathbf{u} \leftarrow \phi^s(x_0 x)$ /* Symbolic derivative at $x_0 x$ */
- 17 **if** \exists key $\mathbf{v} \in \text{COUNT}_{\text{map}}$ s.t. $\|\mathbf{u} - \mathbf{v}\|_{\infty} \leq \epsilon$ **then**
- 18 $\text{COUNT}_{\text{map}}[\mathbf{v}] \leftarrow \text{COUNT}_{\text{map}}[\mathbf{v}] + 1$
- 19 **else**
- 20 Set $\text{COUNT}_{\text{map}}[\mathbf{u}] = 1$
- 21 $\text{COUNT}_{\text{total}} \leftarrow \text{COUNT}_{\text{total}} + 1$
- 22 **else**
- 23 Delete x from \mathcal{N}

- 24 **foreach** key $\mathbf{v} \in \text{COUNT}_{\text{map}}$ **do**
- 25 $\mathbf{h} \leftarrow \mathbf{h} + \left(\frac{\text{COUNT}_{\text{map}}[\mathbf{v}]}{\text{COUNT}_{\text{total}}} \times H(\mathbf{v}) \right)$ /* $H(\mathbf{v})$: entropy of \mathbf{v} */

/* III. Uncertainty Estimation */

- 27 $\epsilon_* \leftarrow \min \epsilon_0$ satisfying: $\alpha + C_0 \frac{1 + \epsilon_0^2}{|s| e^3} + 2e^{-C_1 |\mathcal{N}| \epsilon_0^2} + e^{-\epsilon_0 p_0 |s|} \leq 1$
- 28 $\mathbf{E} \leftarrow \epsilon_* + 2\mathbb{B}(\epsilon_*, |\Sigma|)$
- 29 **return** \mathbf{h}, \mathbf{E}

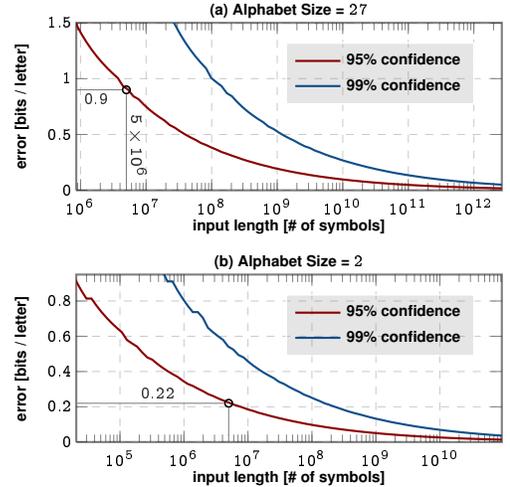


Fig. 3. Uncertainty bounds for different alphabet sizes. Note for a data length of 5×10^6 , we have an uncertainty of 0.9 bits at 95% confidence for a 27 letter alphabet (plate (a)); the corresponding uncertainty for a binary alphabet is 0.22 bits (plate(b)).

Proof: It follows from Corollary 7, and Theorem 5, that:

$$\begin{aligned} T &\triangleq \Pr \left(\left| \frac{1}{|\mathcal{N}|} \sum_{x \in \mathcal{N}} H(\phi^s(x_0 x)) - H(G) \right| \leq \epsilon + 2\mathbb{B}(\epsilon, |\Sigma|) \right) \\ &> \left(1 - C_0 \frac{1 + \epsilon^2}{|s| e^3} - 2e^{-C_1 |\mathcal{N}| \epsilon^2}\right) \times \Pr(x_0 \text{ is } \epsilon\text{-synchronizing}) \end{aligned}$$

Assuming x_0 satisfies the pre-conditions described in Theorem 3:

$$T > \left(1 - C_0 \frac{1 + \epsilon^2}{|s| e^3} - 2e^{-C_1 |\mathcal{N}| \epsilon^2}\right) \times (1 - e^{-\epsilon p_0 |s|})$$

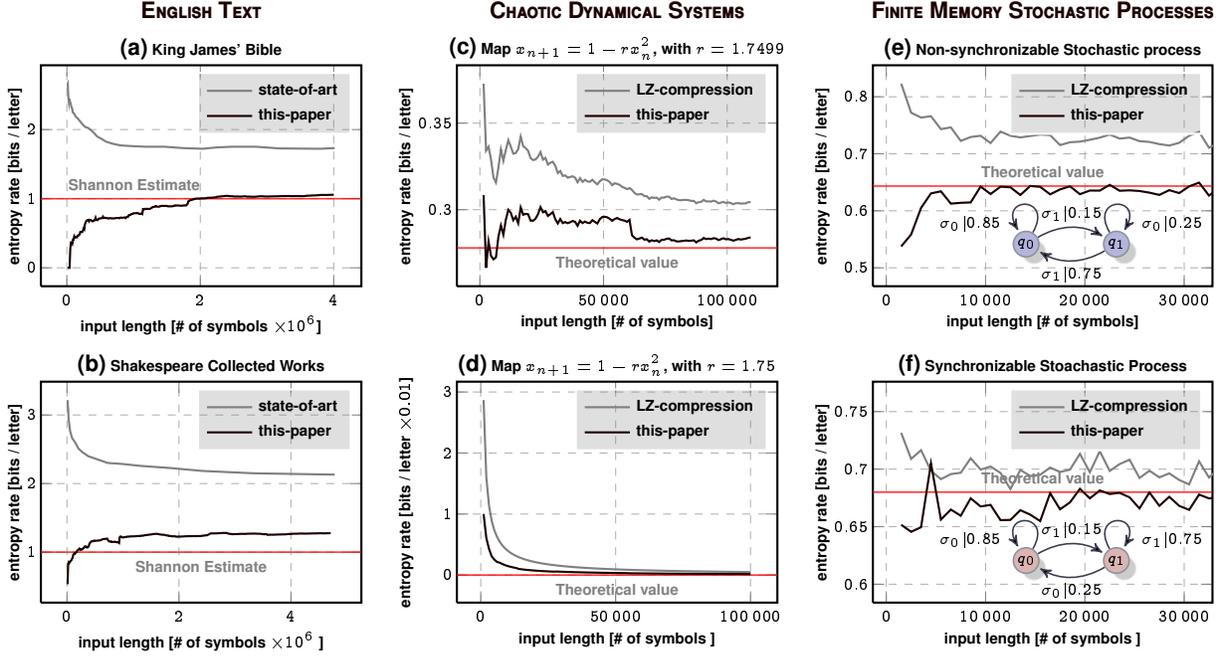


Fig. 4. **Applications.** Plates (a-b): Entropy rate of English text. Shannon’s experiment using human subjects puts the estimate around 1 bit per letter. We achieve very close estimates. The state-of-the-art plots are replicated from [8]. Plates (c-d): Entropy rate of sequences generated by a chaotic dynamical system and a binary generating partition. Plates (e-f): Entropy of symbol streams generated by probabilistic automata. Note that even with two states, and a binary alphabet, a non-synchronizable generating process leads to significantly larger errors with the LZ-based approaches.

$$> 1 - C_0 \frac{1 + \epsilon^2}{|s|e^3} - 2e^{-C_1|\mathcal{N}|e^2} - e^{\epsilon p_0|s|} \quad (96)$$

which completes the proof. ■

Remark 1. We note the following:

- For binary alphabets, we have: $C_0 \simeq 4.03, C_1 = 2$.
- Each term on the RHS of Eq. (95) reflects a specific contribution:

$$C_0 \underbrace{\frac{1 + \epsilon^2}{|s|e^3}}_{\text{Data-length Dependence}} + 2e^{-C_1|\mathcal{N}|e^2} + \underbrace{e^{-\epsilon p_0|s|}}_{\text{Synchronization Error Dependence}} \quad (97)$$

Eq. (95) bounds the maximum uncertainty at a given confidence level, which depends on the alphabet size. The uncertainty relationships for two alphabet sizes (2, 27) are shown in Figure 3.

Corollary 4 (To Theorem 6). As a function of the length of the observed data string s , the upper and lower confidence bands for the estimated entropy rate, with any fixed confidence level, converge at a rate $O\left(\frac{\log |s|}{|s|^{1/3}}\right)$.

Proof: For a given confidence level $k = k_{\text{data}} + k_{\text{depth}}$, where k_{data} captures the dependence on the data length through the first RHS term in Eq. (95), we get:

$$\epsilon^3 = \frac{C_0(1 + \epsilon^2)}{|s|k_{\text{data}}} \Rightarrow \epsilon < \left(\frac{2C_0}{|s|k_{\text{data}}}\right)^{1/3} \quad (98)$$

The distance between the confidence bands is given by:

$$B = 2\epsilon + 4B(\epsilon, |\Sigma|) \quad (99)$$

and using Eq. (98), along with the definition of the generalized binary entropy function (Definition 11), completes the proof. ■

IV. ALGORITHMIC IMPLEMENTATION

The algorithmic steps for the proposed entropy rate estimation technique is enumerated in Algorithm 1. The inputs to the algorithm is the data stream s , ϵ , and the confidence level α at which the error estimate is desired. Importantly, the size of the set of sampled string \mathcal{N} is not required to be an input; if computational effort is not a concern, then the uncertainty contribution from the term involving $|\mathcal{N}|$ (See Eq. (95)) can be reduced to negligible levels by using a sample set with $|\mathcal{N}| \simeq \frac{K}{2\epsilon^2} \log^2 |\Sigma|$, which would result in uncertainty contribution

of $\sim e^{-K}$. Using $|\mathcal{N}| \simeq 10^7 \log^2 |\Sigma|$ is generally sufficient to make this factor negligible; smaller sets may be used under computational constraints, which would lead to increased uncertainty in the entropy estimate.

Particularly rare strings may accumulate errors, which is prevented in the implementation by ignoring strings that occur too infrequently (Note N_{\min} in step 5 and step 15 of Algorithm 1).

A. Application to English text, Chaotic systems & Random walks

We demonstrate Algorithm 1 in three different applications. Our first application is the estimation of the entropy rate of English text. Shannon’s experimental approach with human subjects [24] suggests that English has an entropy of around one bit per letter. However, the large alphabet size (26 letters + space = 27), makes it computationally hard to verify this value. We apply our algorithm to relatively small corpora: the King James Bible (KJB) (which has a length $\sim 4 \times 10^6$ letters), and the collected works of Shakespeare (SHK, length $\sim 4.8 \times 10^6$ letters). These particular examples allow direct comparison against the results reported in [8]. We obtain entropy rates which are significantly closer to the Shannon estimate (See Figure 4): 1.05 **bits/letter** for KJB, and 1.25 **bits/letter** for SHK, while Schürmann *et al.* obtain the corresponding estimates to be 1.73 **bits/letter** and 2.13 **bits/letter**. The authors in [8] were able to improve the SHK estimate to 1.7 **bits/letter** using the “ansatz” mentioned before; Algorithm 1 yields an improved estimate without any such assumptions.

Our second application is entropy estimation of sequences produced by chaotic dynamical systems. We use the same iteration map used in [8]: namely $x_{n+1} = 1 - rx_n^2$, and use a binary generating partition at $x = 0$. We analyze the cases $r = 1.7499$ (Figure 4(b)) where it is very strongly intermittent, and $r = 1.75$ which is the Pomeau-Manneville intermittency point (Figure 4(c)). As before, we converge faster in the non-trivial case, and gets very close to the theoretical entropy given by the positive Lyapunov exponent due to Pesin’s identity [7].

Our third application analyzes sequences generated by finite memory ergodic stationary stochastic processes, modeled directly via probabilistic automata (Figure 4(e-f)). Thus, we are looking at generalized random walks, In spite of being somewhat more contrived compared to the first two applications, we can gain important insights

from this example. Even with two states, and with a binary alphabet, LZ-based approaches may perform significantly worse, particularly for short streams with long range dependencies. We note that the PFSA generator used in Figure 4(e) is non-synchronizable, i.e., no finite length of observed history tells us definitively what the current state is. Nevertheless, as we showed in Theorem 1, the machine is ϵ -synchronizable; and Algorithm 1 performs quite well, converging to the theoretical value with just under 10^4 symbols. In contrast, the LZ-compression based algorithm has an error of about 17% even after 3×10^4 symbols. This discrepancy in performance disappears if the generating process is synchronizable, e.g., if a finite history tells us precisely what the current state is. Indeed with a synchronizable PFSA in Figure 4(f) (here, the last symbol is sufficient to fix the current state), the algorithms have comparable performances.

V. SUMMARY & CONCLUSION

We delineate a new algorithm for estimating entropy rates of symbol streams, generated by hidden ergodic stationary processes. We establish the correctness of the algorithm by exploiting a connection with the theory of probabilistic automata, and that of finite measures on infinite strings. Importantly, we establish a distribution-free limit theorem. Using established results from non-parametric statistics, we show that entropy estimate converges at the rate $O(\log |s| / \sqrt[3]{|s|})$ as a function of the input data length $|s|$. In consequence, we are able to derive confidence bounds on the estimate, and dictate the worst-case data length required to guarantee a specified error bound at a given confidence level. Finally, we demonstrate that, in terms of data requirements, the proposed algorithm has superior performance to competing approaches, at least in the case of the chosen applications.

REFERENCES

- [1] "A note on kolmogorov complexity and entropy," *Applied Mathematics Letters*, vol. 16, no. 7, pp. 1129 – 1130, 2003.
- [2] J. Ziv and A. Lempel, "A universal algorithm for sequential data compression," *IEEE Tran on Information Theory*, vol. 23, no. 3, pp. 337–343, 1977.
- [3] —, "Compression of individual sequences via variable-rate coding," *IEEE Tran on Inf. Theory*, vol. 24, no. 5, pp. 530–536, 1978.
- [4] A. D. Wyner and J. Ziv, "The sliding-window Lempel-Ziv algorithm is asymptotically optimal," *Proceedings of the IEEE*, vol. 82, no. 6, pp. 872–877, 1994.
- [5] J. Langdon, G.G., "A note on the ziv - lempel model for compressing individual sequences (corresp.)," *Information Theory, IEEE Transactions on*, vol. 29, no. 2, pp. 284–287, 1983.
- [6] P. Grassberger, "Estimating the information content of symbol sequences and efficient codes," *IEEE Trans. on Inf. Theory*, vol. 35, no. 3, pp. 669–675, 1989.
- [7] J. Rissanen, "A universal data compression system," *IEEE Trans. on Inf. Theory*, vol. 29, no. 5, pp. 656–664, 1983.
- [8] T. Schürmann and P. Grassberger, "Entropy estimation of symbol sequences," *CHAOS*, vol. 6, no. 3, pp. 414–427, 1996.
- [9] A. Paz, *Introduction to probabilistic automata (Computer science and applied mathematics)*. Orlando, FL, USA: Academic Press, Inc., 1971.
- [10] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York, NY, USA: Wiley-Interscience, 1991.
- [11] E. Vidal, F. Thollard, C. de la Higuera, F. Casacuberta, and R. Carrasco, "Probabilistic finite-state machines - part i," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 27, no. 7, pp. 1013–1025, July 2005.
- [12] I. Chattopadhyay and H. Lipson, "Abductive learning of quantized stochastic processes with probabilistic finite automata," *Philos Trans A*, vol. 371, no. 1984, p. 20110543, Feb 2013.
- [13] J. E. Hopcroft, R. Motwani, and J. D. Ullman, *Introduction to Automata Theory, Languages, and Computation, 2nd ed.* Addison-Wesley, 2001.
- [14] I. Chattopadhyay and A. Ray, "Structural transformations of probabilistic finite state machines," *International Journal of Control*, vol. 81, no. 5, pp. 820–835, May 2008.
- [15] R. Gavaldà, P. W. Keller, J. Pineau, and D. Precup, "Pac-learning of markov models with hidden state," in *ECML*, ser. Lecture Notes in Computer Science, J. Fürnkranz, T. Scheffer, and M. Spiliopoulou, Eds., vol. 4212. Springer, 2006, pp. 150–161.
- [16] S. Bogdanovic, B. Imreh, M. Ciric, and T. Petkovic, "Directable automata and their generalizations - a survey," *Novi Sad Journal of Mathematics*, vol. 29, no. 2, pp. 31–74, 1999.
- [17] M. Ito and J. Duske, "On cofinal and definite automata," *Acta Cybern.*, vol. 6, pp. 181–189, 1984.
- [18] F. Topse, "On the glivenko-cantelli theorem," *Probability Theory and Related Fields*, vol. 14, pp. 239–250.
- [19] I. Csiszár, "Sanov property, generalized I-projection and a conditional limit theorem," *Ann. Probab.*, vol. 12, pp. 768–793, 1984.
- [20] E. L. Lehmann and J. P. Romano, *Testing statistical hypotheses*, 3rd ed., ser. Springer Texts in Statistics. New York: Springer, 2005.
- [21] A. Tsybakov, *Introduction to nonparametric estimation*, ser. Springer series in statistics.
- [22] W. Stewart, *Numerical methods for computing stationary distribution of finite irreducible Markov chains*. New York: Springer, 1999.
- [23] P. Massart, "The tight constant in the dkw inequality," *The Annals of Probability*, vol. 18, no. 3, pp. pp. 1269–1283, 1990.
- [24] C. E. Shannon, "Prediction and entropy of printed english," *Bell System Technical Journal*, vol. 30, pp. 50–64, Jan. 1951.