# Data Smashing

Ishanu Chattopadhyay
ic99@cornell.edu

Hod Lipson
hod.lipson@cornell.edu

*Abstract*—Investigation of the underlying physics or biology from empirical data requires a quantifiable notion of similarity - when do two observed data sets indicate nearly identical generating processes, and when they do not. The discriminating characteristics to look for in data is often determined by heuristics designed by experts, *e.g.*, distinct shapes of "folded" lightcurves may be used as "features" to classify variable stars, while determination of pathological brain states might require a Fourier analysis of brainwave activity. Finding good features is non-trivial. Here, we propose a universal solution to this problem: we delineate a principle for quantifying similarity between sources of arbitrary data streams, without a priori knowledge, features or training. We uncover an algebraic structure on a space of symbolic models for quantized data, and show that such stochastic generators may be added and uniqely inverted; and that a model and its inverse always sum to the generator of flat white noise. Therefore, every data stream has an anti-stream: data generated by the inverse model. Similarity between two streams, then, is the degree to which one, when summed to the other's anti-stream, mutually annihilates all statistical structure to noise. We call this data smashing. We present diverse applications, including disambiguation of brainwaves pertaining to epileptic seizures, detection of anomalous cardiac rhythms, and classification of astronomical objects from raw photometry. In our examples, the data smashing principle, without access to any domain knowledge, meets or exceeds the performance of specialized algorithms tuned by domain experts.

*Index Terms*—feature-free classification, universal metric, probabilistic automata

## I. Motivation & Contribution

The term "data smashing" might conjure up images of erasing information or destroying hard drives. But just as smashing atoms can reveal their composition, "colliding" quantitative data streams can reveal their hidden structure.

We describe here a new principle, where quantitative data streams have corresponding anti-streams, which inspite of being non-unique, are tied to the stream's unique statistical structure. We then describe "data smashing", a process by which streams and anti-streams can be algorithmically collided to reveal differences that are difficult to detect using conventional techniques. We establish this principle formally, describe how we implemented it in practice, and report its performance on a number of real-world cases. The results show that without access to any domain knowledge, data smashing meets or exceeds the accuracy achieved by specialized algorithms and heuristics devised by domain experts.

Nearly all automated discovery systems today rely, at their core, on the ability to compare data: From automatic image recognition to discovering new astronomical objects, such systems must be able to compare and contrast data records in order to group them, classify them, or identify the odd-one-out. Despite rapid growth in the amount of data collected and the increasing rate at which it can be processed, analysis of quantitative data streams still relies heavily on knowing what to look for.

Any time a data mining algorithm searches beyond simple correlations, a human expert must help define a notion of similarity - by specifying important distinguishing "features" of the data to compare, or by training learning algorithms using copious amounts of examples. The data smashing principle removes the reliance on expert-defined features or examples, and in many cases, does so faster and with better accuracy than traditional methods.

This paper is organized as follows: Sections I-VI describe the key concepts, along with a brief but complete description of the approach. The mathematical details, including proffs of correctness, are presenetd in Sections VII-IX. Qunatization schemes are discussed in Section X. Comparisons with some standard notions of statistical dependencies is carried out in Section XI, and the paper is concluded in Secion XII.
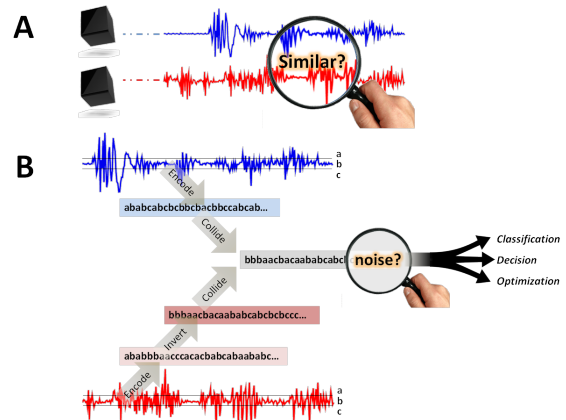


Fig. 1. **Data smashing:** (A) determining the similarity between two data streams is key to any data mining process, but relies heavily on human-prescribed criteria. (B) Data smashing first encodes each data stream, then collides one with the inverse of the other. The randomness of the resulting stream reflects the similarity of the original streams, leading to a cascade of downstream applications involving classification, decision and optimization.

## II. Anti-streams

The notion of data smashing applies only to data in the form of an ordered series of digits or symbols, such as acoustic waves from a microphone, light intensity over time from a telescope, traffic density along a road, or network activity from a router. The anti-stream contains the "opposite" information from the original data stream, and is produced by algorithmically inverting the statistical distribution of symbol sequences appearing in the original stream. For example, sequences of digits that were common in the original stream will be rare in the anti-stream, and vice versa. Streams and anti-streams can then be algorithmically "collided" in a way that systematically cancels any common statistical structure in the original streams, leaving only information relating to their statistically significant differences. We call this the principle of *Information Annihilation* (See Fig. 1).

Data smashing involves two data streams and proceeds in three steps: raw data streams are first quantized, by converting continuous value to a string of characters or symbols. The simplest example of such quantization is where all positive values are mapped to the symbol "1" and all negative values to "0", thus generating a string of bits. Next, we select one of the quantized input streams, and generate its anti-stream. Finally, we annihilate this anti-stream against the remaining quantized input stream and measure what information remains. The remaining information is estimated from the deviation of the resultant stream from flat white noise (FWN). Since a data stream is perfectly annihilated by a correct realization of its anti-stream, any deviation of the collision product from noise quantifies statistical dissimilarity. Using this causal similarity metric, we can cluster streams, classify them, or identify stream segments that are unusual or different. The algorithms are linear in input data, implying they can be applied efficiently to streams in near-real time. Importantly, data smashing can be applied without understanding where the streams were generated, how they are encoded, and what they represent.

Ultimately, from a collection of data streams and their pairwise similarities, it is possible to automatically "back out" the underlying metric embedding of the data, revealing its hidden structure for use with traditional machine learning methods.

Dependence across data streams is often quantified using mutual information (1). However, mutual information and data smashing are
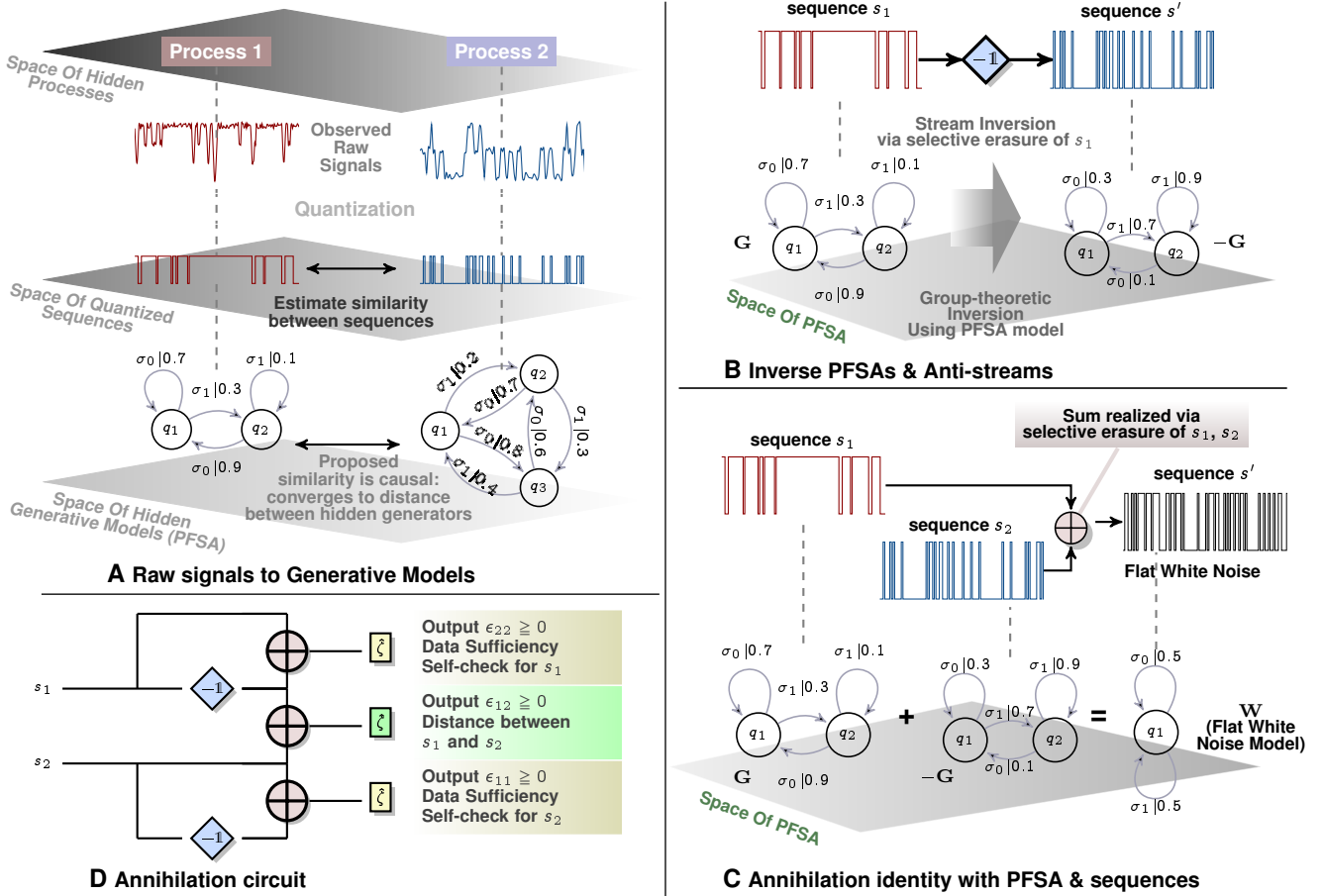
Fig. 2. **Calculation of causal similarity using information annihilation.** (A) We quantize raw signals to symbolic sequences over the chosen alphabet, and compute a causal similarity between such sequences. The underlying theory is established assuming the existence of generative probabilistic automata for these sequences, but our algorithms do not require explicit model construction, or a priori knowledge of their structures. (B) Concept of stream inversion; while we can find the group inverse of a given PFSA algebraically, we can also transform a generated sequence directly to one that represents the inverse model, without constructing the model itself. (C) Summing PFSAs $G$ and its inverse $-G$ yields the zero PFSA $W$. We can carry out this annihilation purely at the sequence level to get flat white noise. (D) Circuit that allows us to measure similarity distance between streams $s_1, s_2$ via computation of $\epsilon_{11}, \epsilon_{22}$ and $\epsilon_{12}$ (See Table 1). Given a threshold $\epsilon^\star > 0$, if $\epsilon_{kk} < \epsilon^\star$, then we have sufficient data for stream $s_k$ ($k = 1, 2$). Additionally if $\epsilon_{12} \leqq \epsilon^\star$, then we conclude that $s_1, s_2$ have the same stochastic source with high probability (which converges exponentially fast to 1 with length of input).

distinct concepts. The former measures dependence between streams; the latter computes a distance between the generative processes themselves. Two sequences of independent coin-flips necessarily have zero mutual information, but data smashing will identify the streams as similar; being generated by the same stochastic process. Moreover, smashing only works correctly if the streams are independent or nearly so (See Section XI-A).

Similarity computed via data smashing is clearly a function of the statistical information buried in the input streams. However, it might not be easy to find the right statistical tool, that reveals this hidden information, particularly without domain knowledge, or without first constructing a good system model (See Section XI-B for an example where smashing reveals non-trivial categories missed by simple statistical measures). We describe in detail the process of computing anti-streams, and the process of comparing information. In Section VII-IX we provide theoretical bounds on the confidence levels, minimal data lengths required for reliable analysis, and scalability of the process as function of the signal encodings.

We have limitations. Data smashing is not directly applicable to learning tasks that do not depend or require a notion of similarity, *e.g.*, identifying a specific time instant at which some event of interest transpired within a data set, or predicting the next step in a time series. Even with the problems to which smashing is applicable, we do not claim strictly superior quantitative performance to the state-of-art in any and all applications; carefully chosen approaches tuned to specific problems can certainly do as well, or better. Our claim is not that we uniformly outperform existing methods, but that we are on

par, as evidenced in multiple example applications; yet do so without requiring expert knowledge, or a training set. Additionally, technical reasons preclude applicability to data from strictly deterministic systems (See section on Limitations & Assumptions).
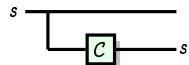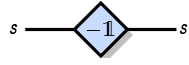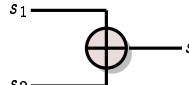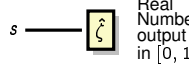
### III. THE HIDDEN MODELS

The notion of a universal comparison metric makes sense only in the context of a featureless approach, where one considers pairwise similarity (or dissimilarity) between individual measurement sets. However, while the advantage of considering the notion of similarity between data sets instead of between feature vectors has been recognized [2], [3], [4], the definition of similarity measures has remained intrinsically heuristic and application dependent, with the possibility of a universal metric been summarily rejected. We show that such universal comparison is indeed realizable, at least under some general assumptions on the nature of the generating process.

We consider sequential observations, *e.g.*, time series of sensor data. The first step is mapping the possibly continuous-valued sensory observations to discrete symbols via pre-specified quantization of the data range (See Section X and Fig. 11). Each symbol represents a slice of the data range, and the total number of slices define the symbol alphabet $\Sigma$ (where $|\Sigma|$ denotes the alphabet size). The coarsest quantization has a binary alphabet consisting of say 0 and 1 (it is not important what symbols we use, we can as well represent the letters of the alphabet with $a$ and $b$), but finer quantizations with larger alphabets are also possible. An observed data stream is thus mapped to a symbolic sequence over this pre-specified alphabet. We

## TABLE I
### ALGORITHMS FOR STREAM OPERATIONS

**(Procedures below are used to assemble the annihilation circuit shown in Fig. 2D, which carries out data smashing)**

| Stream Operation | Algorithmic Procedure (Pseudocode) |
|---|---|
| ■ **Independent Stream Copy**[†]<br><br>$s$ ──── $\mathcal{C}$ ──── $s'$<br><br>**Generate an independent sample path from the same hidden stochastic source.** | 1  Generate stream $\omega_0$ from FWN<br>2  Read current symbol $\sigma_1$ from $s_1$, and $\sigma_2$ from $\omega_0$<br>3  If $\sigma_1 = \sigma_2$, then write $\sigma_1$ to output $s'$<br>4  Move read positions one step to right, and go to step 1<br>*This operation is required internally in stream inversion.* |
| ■ **Stream Inversion**[†]<br><br>$s$ ──── $\langle -1 \rangle$ ──── $s'$<br><br>**Generate sample path from inverse model of hidden source.** | 1  Generate $\|\Sigma\| - 1$ independent copies of $s_1$: $s_1, \cdots, s_{\|\Sigma\|-1}$<br>2  Read current symbols $\sigma_i$ from $s_i$ ($i = 1, \cdots, \|\Sigma\| - 1$)<br>3  If $\sigma_i \neq \sigma_j$ for all distinct $i, j$, then write $\Sigma \setminus \bigcup_{i=1}^{\|\Sigma\|-1} \sigma_i$ to output $s'$<br>4  Move read positions one step to right, and go to step 1 |
| ■ **Stream Summation**[†]<br><br>$s_1$ ───┐<br>　　　　$\oplus$ ─── $s'$<br>$s_2$ ───┘<br><br>**Generating sample path from sum of hidden sources.** | 1  Read current symbols $\sigma_i$ from $s_i$ ($i = 1, 2$)<br>2  If $\sigma_1 = \sigma_2$, then write to output $s'$<br>3  Move read positions one step to right, and go to step 1 |
| ■ **Deviation from FWN**[‡]<br><br>$s$ ──── $\hat{\zeta}$ ──── Real Number output in $[0,1]$<br><br>**Estimating the deviation of a symbolic stream from FWN.**<br><br>(Symbolic derivatives (Definition 9) in Section VII formalizes $\phi^s(\cdot)$. If $s$ is generated by a FWN process, then $\phi^s(x) \to \mathcal{U}_\Sigma$ for any $x \in \Sigma^*$, and hence $\hat{\zeta}(s, \ell) \to 0$.) | $$\hat{\zeta}(s, \ell) = \frac{\|\Sigma\| - 1}{\|\Sigma\|} \sum_{x : \|x\| \leq \ell} \frac{\|\phi^s(x) - \mathcal{U}_\Sigma\|_\infty}{\|\Sigma\|^{2\|x\|}}, \text{ where}$$<br><br>• $\|\Sigma\|$ is the alphabet size, $\|x\|$ is the length of the string $x$<br>• $\ell$ is the maximum length of strings upto which the sum is evaluated. For a given $\epsilon^*$, we choose $\ell = \ln(1/\epsilon^*) / \ln(\|\Sigma\|)$ (See Proposition 14)<br>• $\mathcal{U}_\Sigma$ is the uniform probability vector of length $\|\Sigma\|$<br>• For $\sigma_i \in \Sigma$, $\phi^s(x)\big|_i = \dfrac{\text{\# of occurrences of } x\sigma_i \text{ in string } s}{\text{\# of occurrences of } x \text{ in string } s}$ |

[†]See Section IX for proof of correctness

[‡]See Definition 22 and Propositions 13 and 14 in Section IX

[§] Symbolic derivatives underlie the rigorous proofs. However, for the actual implementation, they are only needed in the final step to compute deviation from FWN

---

assume that the symbol alphabet and its interpretation is fixed for a particular task.

Quantization involves some information loss which can be reduced with finer alphabets at the expense of increased computational complexity (See Section X). We use quantization schemes (See Fig. 11) which require no domain expertise.

### A. Inverting and combining hidden models

Quantized Stochastic Processes (QSPs) which capture the statistical structure of symbolic streams can be modeled using probabilistic automata, provided the processes are ergodic and stationary [5], [6], [7]. For the purpose of computing our similarity metric, we require that the number of states in the automata be finite (*i.e.* we only assume the existence of a generative Probabilistic Finite State Automata (PFSA)); we do not attempt to construct explicit models or require knowledge of either the exact number of states or any explicit bound thereof (See Fig. 2).

A slightly restricted subset of the space of all PFSA over a fixed alphabet admits an Abelian group structure (See Section VIII); wherein the operations of commutative addition and inversion are well-defined. A trivial example of an Abelian group is the set of reals with the usual addition operation; addition of real numbers is commutative and each real number $a$ has a unique inverse $-a$, which when summed produce the unique identity $0$. We have previously discussed the Abelian group structure on PFSAs in the context of model selection [8]. Here, we show that key group operations, necessary for classification, can be carried out on the observed sequences alone, without any state synchronization or reference to the hidden generators of the sequences.

Existence of a group structure implies that given PFSAs $G$ and $H$, sums $G + H, G - H$, and unique inverses $-G$ and $-H$ are well-defined. Individual symbols have no notion of a "sign", and hence the models $G$ and $-G$ are *not* generators of sign-inverted sequences which would not make sense as our generated sequences are symbol streams. For example, the anti-stream of a sequence 10111 is not

$-1\ 0\ -1\ -1\ -1$, but a fragment that has inverted statistical properties in terms of the occurrence patterns of the symbols 0 and 1 (See Table I). For a PFSA $G$, the unique inverse $-G$ is the PFSA which when added to $G$ yields the group identity $W = G + (-G)$, *i.e.*, the zero model. Note, the zero model $W$ is characterized by the property that for any arbitrary PFSA $H$ in the group, we have $H + W = W + H = H$.
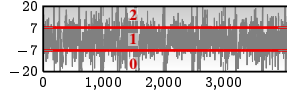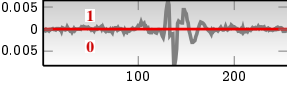
For any fixed alphabet size, the zero model is the unique single-state PFSA (up to minimal description [9]) that generates symbols as consecutive realizations of independent random variables with uniform distribution over the symbol alphabet. Thus $W$ generates flat white noise (FWN), and the entropy rate of FWN achieves the theoretical upper bound among the sequences generated by arbitrary PFSA in the model space. Two PFSAs $G, H$ are identical if and only if $G + (-H) = W$.

### B. Metric Structure on Model Space

In addition to the Abelian group, the PFSA space admits a metric structure (See Section VII). The distance between two models thus can be interpreted as the deviation of their group-theoretic difference from a FWN process. Information annihilation exploits the possibility of estimating causal similari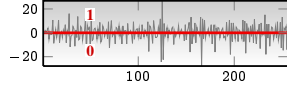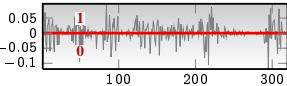ty between observed data streams by estimating this distance from the observed sequences alone without requiring the models themselves.

We can estimate the distance of the hidden generative model from FWN given only an observed stream $s$. This is achieved by the function $\hat{\zeta}$ (See Table I, row 4). Intuitively, given an observed sequence fragment $x$, we first compute the deviation of the distribution of the next symbol from the uniform distribution over the alphabet. $\hat{\zeta}(s, \ell)$ is the sum of these deviations for all historical fragments $x$ with length up to $\ell$, weighted by $1/\|\Sigma\|^{2\|x\|}$. The weighted sum ensures that deviation of the distributions for longer $x$ have smaller contribution to $\hat{\zeta}(s, \ell)$, which addresses the issue that the occurrence frequencies of longer sequences are more variable.

TABLE II
APPLICATION PROBLEMS, & RESULTS[‡]

| System | Input Description | Classification Performance |
|---|---|---|

**1. Identify epileptic pathology (10)**

- 495 EEG excerpts, each 23.6$s$ sampled at 173.61hz
- Signal derivative as input
- Quantization[⋆] (3 letter):

| IA accuracy | 98.9% |
|---|---|
| State of art | NA |

No comparable result is available in the literature. However, IA reveals a 1D manifold structure in the dataset, while (10) with additional assumptions on the nature of hidden processes, fails to yield such insight.

**2. Identify heart murmur (11)**

- 65 .wav files sampled at 44.1kHz ($\sim$ 10$s$ each)
- Quantization[⋆] (2 letter):

| IA precision (murmur) | 75.2% |
|---|---|
| State of art | 67% |

State of the art (11) achieved in supervised learning with task-specific features

**3. Classify variable stars (Cepheid variable vs RR Lyrae) from photometry (OGLE II) (12)**

- 10699 photometric series
- Differentiated folded/raw photometry used as input
- Quantization[⋆] (3 letter):

| IA accuracy | 99.8% | **Folded** |
|---|---|---|
| State of art | 99.6% | **Photometry** |

State of the art (12) achieved with task-specific features and multiple hand-optimized classification steps

| IA accuracy | 94.3% | **Unfolded Photometry** |
|---|---|---|
| State of art | NA | (*This capability is beyond the state of art*) |

**4. EEG based Biometric Authentication (13) with visually evoked potentials (VEP)**

- 122 subjects, multi-variate data from 61 standard electrodes.
- 256 data points for each trial for each electrode.
- Total # of data series: 5477 (each with 61 variables).
- Quantization[⋆] (2 letter):

| | kNN | SVM |
|---|---|---|
| IA accuracy | 97.96% | 99.65% |
| State of art | 95.6% | 98.96 % |

State of the art (14) achieved with task-specific features, and after eliminating 2 subjects from consideration

**5. Text-independent speaker identification using ELSDSR database (15)**

- 23 speakers (9 female, 14 male), 16kHZ recording
- $\sim$ 100$s$ recording/speaker
- 2$s$ snippets used as time series excerpts
- Total # of time series: 1270
- Quantization[⋆] (2 letter):

| IA accuracy | 80.2% |
|---|---|
| State of art | 73.73% |

State of the art (16) achieved with task-specific features and multiple hand-optimized classification steps

[⋆] See Section X for details on choosing quantization schemes

## IV. KEY INSIGHT: THE INFORMATION ANNIHILATION PRINCIPLE

Our key insight is the following: two sets of sequential observations have the same generative process if the *inverted* copy of one can *annihilate* the statistical information contained in the other. We claim, that given two symbol streams $s_1$ and $s_2$, we can check if the underlying PFSAs (say $G_1, G_2$) satisfy the *annihilation equality*: $G_1 + (-G_2) = W$ without explicitly knowing or constructing the models themselves.

Data smashing is predicated on being able to invert and sum streams, and to compare streams to noise. Inversion generates a stream $s'$ given a stream $s$, such that if PFSA $G$ is the source for $s$, then $-G$ is the source for $s'$. Summation collides two streams: Given streams $s_1$ and $s_2$, generate a new stream $s'$ which is a realization of FWN if and only if the hidden models $G_1, G_2$ satisfy $G_1 + G_2 = W$. Finally, deviation of a stream $s$ from that generated by a FWN process can be calculated directly.

Importantly, for a stream $s$ (with generator $G$), the inverted stream $s'$ is not unique. Any symbol stream generated from the inverse model $-G$ qualifies as an inverse for $s$; thus anti-streams are non-unique. What is indeed unique is the generating inverse PFSA model. Since, our technique compares the hidden stochastic processes and not their possibly non-unique realizations, the non-uniqueness of anti-streams is not problematic.

Despite the possibility of mis-synchronization between hidden model states, applicability of the algorithms shown in Table I for disambiguation of hidden dynamics is valid. We show in Section IX that the algorithms evaluate distinct models to be distinct, and nearly identical hidden models to be nearly identical.

Estimating the deviation of a stream from FWN is straightforward (as specified by $\hat{\zeta}(s, \ell)$ in Table I, row 4). All subsequences of a given length must necessarily occur with the same frequency for a FWN process; and we simply estimate the deviation from this behavior in the observed sequence. The other two tasks are carried out via selective erasure of symbols from the input stream(s) (See Table I, rows 1-3). For example, summation of streams is realized as follows: given two streams $s_1, s_2$, we read a symbol from each stream, and if they match then we copy it to our output, and ignore the symbols read when they do not match.

Thus, data smashing allows us to manipulate streams via selective erasure, to estimate a distance between the hidden stochastic sources. Specifically, we estimate the degree to which the sum of a stream and its anti-stream brings the entropy rate of the resultant stream close to its theoretical upper bound.

### A. Contrast with Feature-based State of Art

Contemporary research in machine learning is dominated by the search for good "features" (17), which are typically understood to be heuristically chosen discriminative attributes characterizing objects or phenomena of interest. Finding such attributes is not easy [18], [19]. Moreover, the number of characterizing features *i.e.* the size of the feature set, needs to be relatively small to avoid intractability of the subsequent learning algorithms. Additionally, their heuristic definition precludes any notion of optimality; it is impossible to quantify the quality of a given feature set in any absolute terms; we can only compare how it performs in the context of a specific task against a few selected variations.

In addition to the heuristic nature of feature selection, machine learning algorithms typically necessitate the choice of a distance metric in the feature space. For example, the classic "nearest neighbor" k-NN classifier (20) requires definition of proximity,

## A  Covergence of Self-Annihilation Error

## B  Computation time for annihilation



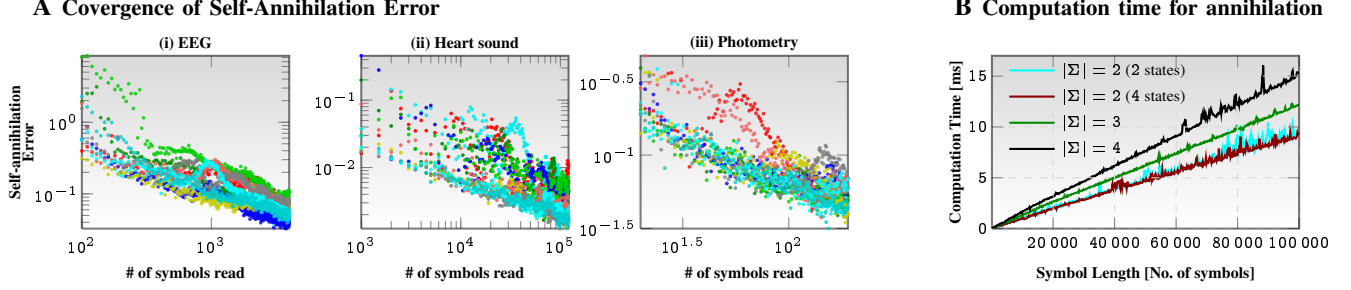**(i) EEG**  **(ii) Heart sound**  **(iii) Photometry**

Fig. 3.  **Computational complexity and convergence rates for information annihilation.** (A) Illustrates exponential convergence of the self-annihilation error for a small set of data series for different applications (plate (i) for EEG data, plate (ii) for heart sound recordings, and plate (iii) for photometry). (B) Computation times for carrying out annihilation using the circuit shown in Fig. 2D as a function of the length of input streams for different alphabet sizes (and for diffrent number of states in the hidden models). Note that the asymptotic time complexity of obtaining the similarity distances scales as $O(|\Sigma|n)$, where $n$ is the length of the shorter of the two input streams.

and the k-means algorithm (21) depends on pairwise distances in the feature space for clustering. To side-step the heuristic metric problem, recent approaches often learn appropriate metrics directly from data, attempting to "back out" a metric from side information or labeled constraints (22). Unsupervised approaches use dimensionality reduction and embedding strategies to uncover the geometric structure of geodesics in the feature space (*e.g.* see manifold learning (23), (24), (25)). However, automatically inferred data geometry in the feature space is, again, strongly dependent on the initial choice of features. Since Euclidean distances between feature vectors are often misleading (23), heuristic features make it impossible to conceive of a task-independent universal metric.

In contrast, smashing is based on an application-independent notion of similarity between quantized sample paths observed from hidden stochastic processes. Our universal metric quantifies the degree to which the summation of the inverted copy of any one stream to the other annihilates the existing statistical dependencies, leaving behind flat white noise. We circumvent the need for features altogether (See Fig. 1B) and do not require training.

Despite the fact that the estimation of similarities between two data streams is performed in absence of the knowledge of the underlying source structure or its parameters, we establish that this universal metric is causal, *i.e.*, with sufficient data it converges to a well-defined distance between the hidden stochastic sources themselves, without ever knowing them explicitly.

### B. Self-annihilation Test for Data-sufficiency Check

Statistical process characteristics dictate the amount of data required for estimation of the proposed distance. With no access to the hidden models, we cannot estimate the required data length a priori; however it is possible to check for data-sufficiency for a specified error threshold via self-annihilation. Since the proposed metric is causal, the distance between two independent samples from the same source always converges to zero. We estimate the degree of self-annihilation achieved in order to determine data sufficiency; *i.e.*, a stream is sufficiently long if it can sufficiently annihilate an inverted self-copy to FWN.

The self-annihilation based data-sufficiency test consists of two steps: given an observed symbolic sequence $s$, we first generate an independent copy (say $s'$). This is the independent stream copy operation (See Table I, row 1), which can be carried out via selective symbol erasure without any knowledge of the source itself. Once we have $s$ and $s'$, we check if the inverted version of one annihilates the other to a pre-specified degree. In particular, we generate $s''$ from $s$ via stream inversion, and use stream summation of $s'$ and $s''$ to produce the final output stream $s'''$, and check if $\hat{\zeta}(s''', \ell)$ is less than some specified threshold $\epsilon^\star > 0$. We show that considering only histories up to a length $\ell = \frac{\ln(1/\epsilon^\star)}{\ln(|\Sigma|)}$ in the computation of $\hat{\zeta}(s''', \ell)$ is sufficient (See Section IX).

The self-annihilation error is also useful to rank the effectiveness of different quantization schemes. Better quantization schemes (*e.g.* ternary instead of binary) will be able to produce better self-annihilation while maintaining the ability to discriminate different streams (See Section X).

## V. Feature-free Classification and Clustering

Given $n$ data streams $s_1, \cdots, s_n$, we construct a matrix $E$, such that $E_{ij}$ represents the estimated distance between the streams $s_i, s_j$. Thus, the diagonal elements of $E$ are the self-annihilation errors, while the off-diagonal elements represent inter-stream similarity estimates (See Fig. 2D for the basic annihilation circuit). Given a positive threshold $\epsilon^\star > 0$, the self-annihilation tests are passed if $\epsilon_{kk} \leqq \epsilon^\star$ ($k = i, j$), and for sufficient data the streams $s_i, s_j$ have identical sources with high probability if and only if $\epsilon_{ij} \leqq \epsilon^\star$. Once $E$ is constructed, we can determine clusters by rearranging $E$ into prominent diagonal blocks. Any standard technique (26) can be used for such clustering; information annihilation is only used to find the causal distances between observed data streams, and the resultant distance matrix can then used as input to state-of-the-art clustering methodologies, or finding geometric structures (such as lower dimensional embedding manifolds (23)) induced by the similarity metric on the data sources.

The matrix $H$, obtained from $E$ by setting the diagonal entries to zero, estimates a distance matrix. An Euclidean embedding (27) of $H$ then leads to deeper insight into the geometry of the space of the hidden generators, *e.g.*, in the case of the EEG data, the time series' describe a one-dimensional manifold (a curve), with data from similar phenomena clustered together along the curve (See Fig. 4A(ii)).

### A. Computational Complexity & Data Requirements

The asymptotic time complexity of carrying out the stream operations scales linearly with input length, and the granularity of the alphabet (See Section IX and Fig. 3B for illustration of the linear time complexity of estimating inter-stream similarity).

To pass the self-annihilation test, a data stream must be sufficiently long; and the required length $|s|$ of the input $s$ with a specified threshold $\epsilon^\star$ is dictated by the characteristics of the generating process. Selective erasure in annihilation (See Table I) implies that the output tested for being FWN is shorter compared to the input stream, and the expected shortening ratio $\beta$ can be explicitly computed (See Section IX). We refer to $\beta$ as the *annihilation efficiency*, since the convergence rate of the self-annihilation error scales as $1/\sqrt{\beta|s|}$. In other words, the required length $|s|$ of the data stream to achieve a self-annihilation error of $\epsilon^\star$ scales as $1/\beta(\epsilon^\star)^2$. It is important to note that our analysis shows that the annihilation efficiency is independent of the descriptional complexity of the process, *e.g.*, in Fig. 10 the self-annihilation error for a simpler two state process converges faster to a four state process. However the convergence rate always scales as $O(1/\sqrt{|s|})$ as dictated by the the Central Limit Theorem (CLT) (28).

### B. Limitations & Assumptions

Data smashing is not useful in problems which do not require a notion of similarity, *e.g.*, predicting the future course of a time series, or analyzing a data set to pinpoint the occurrence time of an event of interest.

For problems to which smashing is applicable, we implicitly assume the *existence* of PFSA generators; although we never find

**A** EEG (epileptic pathology detection)

**B** Heart mumur detection from digital stethoscope

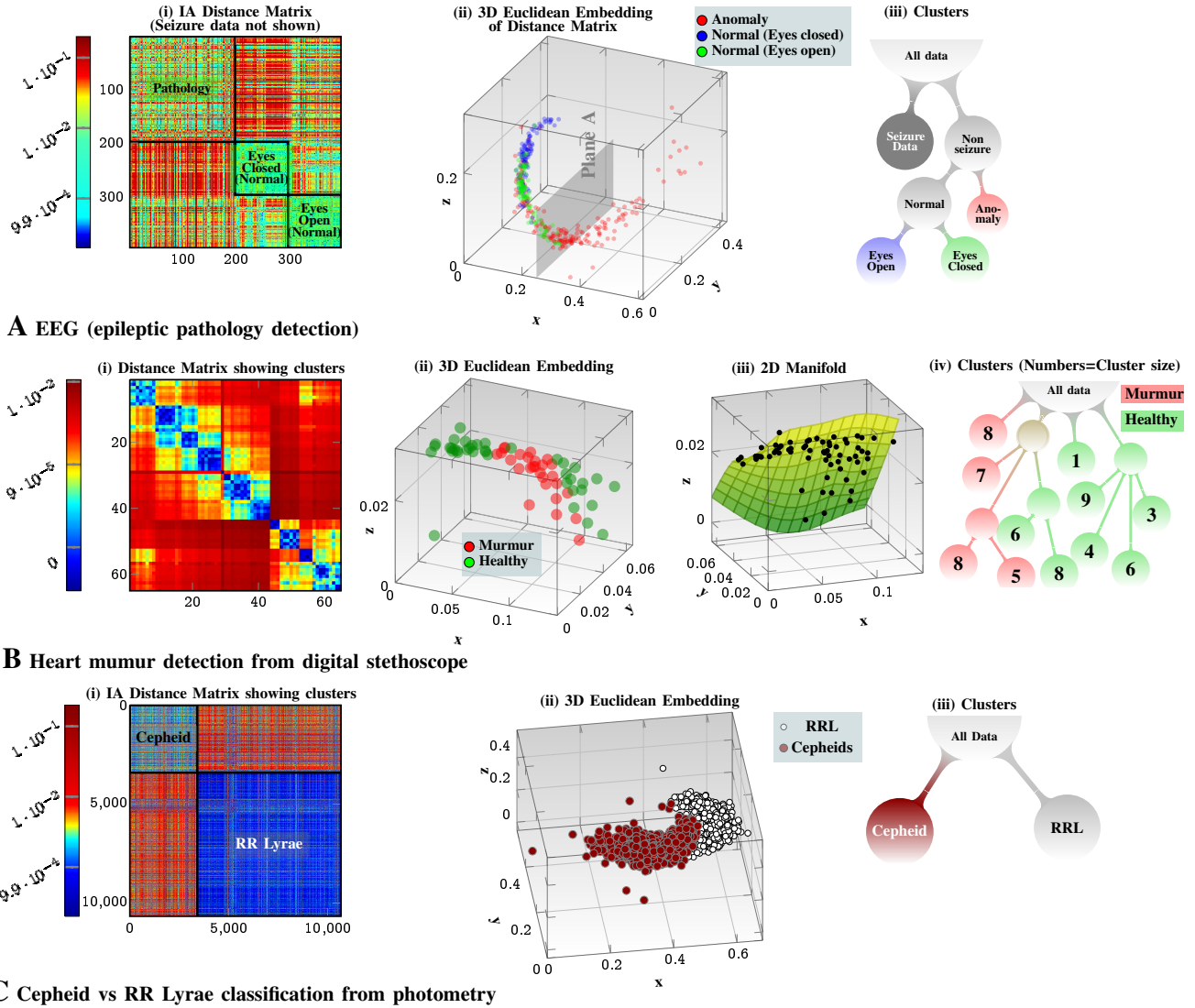**C** Cepheid vs RR Lyrae classification from photometry

Fig. 4. **Data smashing applications.** Pairwise distance matrices, identified clusters and 3D projections of Euclidean embeddings for epileptic pathology identification (shown in (A)), identification of heart murmur (shown in (B)), and classification of variable stars from photometry (shown in (C)). In these applications, the relevant clusters are found unsupervised.

these models explicitly. It follows that what we actually assume is not any particular modeling framework, but that the systems of interest satisfy the properties of ergodicity, stationarity, and have a finite (but not a priori bounded) number of states (See Section VII). In practice, our technique performs well even if these properties are only approximately satisfied (*e.g.* quasi-stationarity instead of stationarity, see example in Section XI-B). The algebraic structure of the space of PFSAs (in particular, existence of unique group inverses) is key to the information annihilation principle; however we argue that any quantized ergodic stationary stochastic process is indeed representable as a probabilistic automata (See Section VII).

Data smashing is not applicable to data from strictly deterministic systems. Such systems are representable by probabilistic automata; however transitions occur with probabilities which are either 0 or 1. PFSAs with zero-probability transitions are non-invertible, which invalidates the underlying theoretical guarantees (See Section VIII). Similarly, data streams in which some alphabet symbol is exceedingly rare would be difficult to invert (See Section IX for the notion of annihilation efficiency).

Symbolization invariably introduces quantization error. This can be made small by using larger alphabets. However, larger alphabet sizes demand longer observed sequences (See Section IX, Fig. 9), implying that the length of observation limits the quantization granularity,

and in the process limits the degree to which the quantization error can be mitigated. Importantly, with coarse quantizations distinct processes may evaluate to be similar. However, identical processes will still evaluate to be identical (or nearly so), provided the streams pass the self-annihilation test. The self-annihilation test thus offers an application-independent way to compare and rank quantization schemes (See Section X).

The algorithmic steps (See Table I) require no synchronization (we can start reading the streams anywhere), implying that non-equal length of time-series, and phase mismatches are of no consequence.

## VI. APPLICATION EXAMPLES

Data smashing begins with quantizing streams to symbolic sequences, followed by the use of the annihilation circuit (Fig. 2D) to compute pairwise causal similarities. Details of the quantization schemes, computed distance matrices, and identified clusters and Euclidean embeddings are summarized in Table II and Fig. 4.

Our first application is classification of brain electrical activity from different physiological and pathological brain states (10). We used sets of electroencephalographic (EEG) data series consisting of surface EEG recordings from healthy volunteers with eyes closed and open, and intracranial recordings from epilepsy patients during

seizure free intervals from within and from outside the seizure generating area, as well as intracranial recordings of seizures.

Starting with the data series of electric potentials, we generated sequences of relative changes between consecutive values before quantization. This step allows a common alphabet for sequences with wide variability in the sequence mean values.

The distance matrix from pairwise smashing yielded clear clusters corresponding to seizure, normal eyes open (EO), normal eyes closed (EC) and epileptic pathology in non-seizure conditions. (See Fig. 4A, seizures not shown due to large differences from the rest).

Embedding the distance matrix (See Fig. 4A, plate (i)) yields a one-dimensional manifold (a curve), with contiguous segments corresponding to different brain states, *e.g.*, right hand side of plane A correspond to epileptic pathology. This provides a particularly insightful picture, which eludes complex non-linear modeling(10).

Next we classify cardiac rhythms from noisy heat-sound data recorded using a digital stethoscope (11). We analyzed 65 data series (ignoring the labels) corresponding to healthy rhythms and murmur, to verify if we could identify clusters without supervision that correspond to the expert-assigned labels.

We found 11 clusters in the distance matrix (See Fig. 4B), 4 of which consisted of mainly data with murmur (as determined by the expert labels), and the rest consisting of mainly healthy rhythms (See Fig. 4B, plate (iv)). Classification precision for murmur is noted in Table 2 (75.2%). Embedding of the distance matrix revealed a two dimensional manifold (See Fig. 4B, plate (iii)).

Our next problem is the classification of variable stars using light intensity series (photometry) from the Optical Gravitational Lensing Experiment (OGLE) survey (12). Supervised classification of photometry proceeds by first "folding" each light-curve to its known period to correct phase mismatches. In our first analysis, we started with folded light-curves; and generated data series of the relative changes between consecutive brightness values in the curves before quantization, which allows for the use of a common alphabet for light curves with wide variability in the mean brightness values. Using data for Cepheids and RRLs (3426 Cepheids, 7273 RRL), we obtained a classification accuracy of 99.8% which marginally outperforms the state of art (See Table II). Clear clusters (obtained unsupervised) corresponding to the two classes can be seen in the computed distance matrix (See Fig. 4C, plate (i)), and the 3D projection of its Euclidean embedding (See Fig. 4C, plate (ii)). The 3D embedding was very nearly constrained within a 2D manifold (See Fig. 4C plate (ii)).

Additionally, in our second analysis, we asked if data smashing can work without knowledge of the period of the variable star; skipping the folding step. Smashing raw photometry data yielded a classification accuracy of 94.3% for the two classes (See Table II). This direct approach is beyond state of the art techniques.

Our fourth application is biometric authentication using visually evoked EEG potentials (VEP). The public database used (13). considered 122 subjects, each of whom was exposed to pictures of objects chosen from the standardized Snodgrass set (29).

Note that while this application is supervised (since we are not attempting to find clusters unsupervised), no actual training is involved; we merely mark the randomly chosen subject-specific set of data series as the library set representing each individual subject. If "unknown" test data series is smashed against each element of each of the libraries corresponding to the individual subjects, we expected that the data series from the same subject will annihilate each other correctly, while those from different subjects will fail to do so to the same extent. We outperformed the state of art for both kNN and SVM based approaches (See Table II).

Our fifth application is text independent speaker identification using the ELSDSR database (15), which includes recording from 23 speakers (9 female, and 14 male, with possibly non-native accents). As before, training involved specifying the library series for each speaker. We computed the distance matrix by smashing the library data series against each other, and trained a simple kNN on the Euclidean embedding of the distance matrix. The test data then yielded a classification accuracy of 80.2%, which beat the state of art figure of 73.73% for 2$s$ snippets of recording data (16) (See Table II).

In the suceeding sections, we develop the mathematical details of the information annihilatin principle, and establish the correctness

of the data smashing algorithm. Section VII presents the theory of probabilistic automata as a modeling framework for ergodic stationary quantized stochastic processes. Section VIII describes the relevant algebraic structures, including that of an Abelian group, definable on the space of probabilistic automata. This is central to the notion of anti-streams. Section IX then establishes that the stream operations delineated in Table I are indeed correct. Section X discusses quantization schemes; specifically describing how to choose the granularity of the quantization. Section XI expounds the differences between the data smashing approach and some specific standard notions often used to quantify statistical dependencies, *e.g.* mutual information between data streams. We also discuss a specific example to illustrate that simple statistical features may miss important dynamical artifacts in data, which is easily revealed via data smashing. The paper is summarized and concluded in Section XII.

## VII. Stochastic Processes & Probabilistic Automata

To establish the correctness of the data smashing algorithm, we first establish the possibility of using probabilistic automata to model stationary, ergodic processes. Our automata models (5) are distinct to those reported in the literature (30), [31]. We include a brief overview here for the sake of completeness.

**Notation 1.** $\Sigma$ *denotes a finite alphabet of symbols. The set of all finite but possibly unbounded strings on $\Sigma$ is denoted by $\Sigma^\star$ (32). The set of finite strings over $\Sigma$ form a concatenative monoid, with the empty word $\lambda$ as identity. The set of strictly infinite strings on $\Sigma$ is denoted as $\Sigma^\omega$, where $\omega$ denotes the first transfinite cardinal. For a string $x$, $|x|$ denotes its length, and for a set $A$, $|A|$ denotes its cardinality.*

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

$$\mathcal{H} = \{ X_t : X_t \text{ is a } \Sigma\text{-valued random variable}, t \in \mathbb{N} \cup \{0\} \} \quad (1)$$

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

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

*Proof:* Assuming stationarity, we can construct a probability measure $\mu : \mathfrak{B} \to [0, 1]$ by defining for any sequence $x \in \Sigma^\star \setminus \{\lambda\}$, and a sufficiently large number of realizations $N_R$ (assuming ergodicity):

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

and extending the measure to elements of $\mathfrak{B} \setminus B$ via at most countable sums. Thus $\mu(\Sigma^\omega) = \sum_{x \in \Sigma^\star} \mu(x\Sigma^\omega) = 1$, 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 (32). We use a probabilistic extension (9).

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

$$\forall x, y \in \Sigma^\star, x \sim_N y \iff \forall z \in \Sigma^\star \left( \Big( Pr(xz) = Pr(yz) = 0 \Big) \right.$$

$$\left. \bigvee \Big| Pr(xz)/Pr(x) - Pr(yz)/Pr(y) \Big| = 0 \right) \quad (2)$$

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

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

$$x \sim_N y \Rightarrow \forall z \in \Sigma^\star, xz \sim_N yz \qquad (3)$$

A right-invariant equivalence on $\Sigma^\star$ 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 $x \in \Sigma^\star$,

$$([x] = q) \wedge ([x\sigma] = q') \Rightarrow q \xrightarrow{\sigma} q' \qquad (4)$$

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 $(Q, \Sigma, \delta, \widetilde{\pi}, q_0)$, where $Q$ is a finite state set, $\Sigma$ is the alphabet, $\delta : Q \times \Sigma \to Q$ is the state transition function, $\widetilde{\pi} : Q \times \Sigma \to [0, 1]$ specifies the conditional symbol-generation probabilities, and $q_0 \in Q$ is the initial state. $\delta$ and $\widetilde{\pi}$ are recursively extended to arbitrary $y = \sigma x \in \Sigma^\star$ as follows:*

$$\forall q \in Q, \delta(q, \lambda) = q \qquad (5)$$

$$\delta(q, \sigma x) = \delta(\delta(q, \sigma), x) \qquad (6)$$

$$\forall q \in Q, \widetilde{\pi}(q, \lambda) = 1 \qquad (7)$$

$$\widetilde{\pi}(q, \sigma x) = \widetilde{\pi}(q, \sigma)\widetilde{\pi}(\delta(q, \sigma), x) \qquad (8)$$

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

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 (33). 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 $G = (Q, \Sigma, \delta, \widetilde{\pi}, q_0)$ induces a unique probability measure $\mu_G$ on the measurable space $(\Sigma^\omega, \mathfrak{B})$.*

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

$$\mu_G(\varnothing) \triangleq 0 \qquad (9)$$

$$\forall x \in \Sigma^\star, \mu_G(x\Sigma^\omega) \triangleq \delta(q_0, x) \qquad (10)$$

$$\forall x, y \in \Sigma^\star, \mu_G(\{x, y\}\Sigma^\omega) \triangleq \mu_G(x\Sigma^\omega) + \mu_G(y\Sigma^\omega) \qquad (11)$$

Countable additivity of $\mu_G$ is immediate, and (See Definition 4):

$$\mu_G(\Sigma^\omega) = \mu_G(\lambda\Sigma^\omega) = \delta(q_0, \lambda) = 1 \qquad (12)$$

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 $\delta : Q \times \Sigma \to Q$, $\widetilde{\pi} : Q \times \Sigma \to [0, 1]$ as:

$$\delta([x], \sigma) = [x\sigma] \qquad (13)$$

$$\widetilde{\pi}([x], \sigma) = \frac{Pr(x'\sigma)}{Pr(x')} \text{ for any choice of } x' \in [x] \qquad (14)$$

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

$$\delta(q, \sigma x) = \delta(\delta(q, \sigma), x) \qquad (15)$$

$$\widetilde{\pi}(q, \sigma x) = \widetilde{\pi}(q, \sigma)\widetilde{\pi}(\delta(q, \sigma), x) \qquad (16)$$

For verifying the null-word probability, choose a $x \in \Sigma^\star$ such that $[x] = q$ for some $q \in Q$. Then, from Eq. (14), we have:

$$\widetilde{\pi}(q, \lambda) = \frac{Pr(x'\lambda)}{Pr(x')} \text{ for any } x' \in [x] \Rightarrow \widetilde{\pi}(q, \lambda) = \frac{Pr(x')}{Pr(x')} = 1 \qquad (17)$$

Finite index of $\sim_N$ implies $|Q| < \infty$, and hence denoting $[\lambda]$ as $q_0$, we conclude: $G = (Q, \Sigma, \delta, \widetilde{\pi}, q_0)$ 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 PFSAs have a designated initial state. Next we introduce canonical representations to remove initial-state dependence. We use $\widetilde{\Pi}$ to denote the matrix representation of $\widetilde{\pi}$, *i.e.*, $\widetilde{\Pi}_{ij} = \widetilde{\pi}(q_i, \sigma_j)$, $q_i \in Q, \sigma_j \in \Sigma$. We need the notion of transformation matrices $\Gamma_\sigma$.

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

$$\Gamma_\sigma \big|_{ij} = \begin{cases} \widetilde{\pi}(q_i, \sigma), & \text{if } \delta(q_i, \sigma) = q_j \\ 0, & \text{otherwise} \end{cases} \qquad (18)$$

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 $x \in \Sigma^\star$ in the following sense: if $x = \sigma_{r_1} \cdots \sigma_{r_m} \in \Sigma^\star$, then we have:

$$\wp_x = \wp_{\sigma_{r_1} \cdots \sigma_{r_m}} = \underbrace{\frac{1}{\|\wp_\lambda \prod_{j=1}^m \Gamma_{\sigma_{r_j}}\|_1}}_{\text{Normalizing factor}} \wp_\lambda \prod_{j=1}^m \Gamma_{\sigma_{r_j}} \qquad (19)$$

where $\wp_\lambda$ 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 $\wp_\lambda$. Thus, $\wp_x$ is an equivalence class of strings, *i.e.*, $x$ is not unique.

**Definition 6** (Canonical Representation). *An initial-marked PFSA $G = (Q, \Sigma, \delta, \widetilde{\pi}, q_0)$ uniquely induces a canonical representation $(Q^C, \Sigma, \delta^C, \widetilde{\pi}^C)$, where $Q^C$ is a subset of the set of probability distributions over $Q$, and $\delta^C : Q^C \times \Sigma \to Q^C$, $\widetilde{\pi}^C : Q^C \times \Sigma \to [0, 1]$ 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 $\wp_\lambda$ of $Q^C$. Note that the transition matrix for $G$ is the row-stochastic matrix $M \in [0, 1]^{|Q| \times |Q|}$, with $M_{ij} = \sum_{\sigma : \delta(q_i, \sigma) = q_j} \widetilde{\pi}(q_i, \sigma)$, and hence $\wp_\lambda$ satisfies:*

$$\wp_\lambda M = \wp_\lambda \qquad (20)$$

*2) Define $\delta^C$ and $\widetilde{\pi}^C$ recursively:*

$$\delta^C(\wp_x, \sigma) = \frac{1}{\|\wp_x \Gamma_\sigma\|_1} \wp_x \Gamma_\sigma \triangleq \wp_{x\sigma} \qquad (21)$$

$$\widetilde{\pi}^C(\wp_x, \sigma) = \wp_x \widetilde{\Pi} \qquad (22)$$

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

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

*1) The canonical representation is independent of the initial state.*

*2) The canonical representation $(Q^C, \Sigma, \delta^C, \widetilde{\pi}^C)$ 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 $\zeta : Q \to Q'$, with:*

$$\forall q \in Q, \forall \sigma \in \Sigma, \begin{cases} \widetilde{\pi}(q, \sigma) = \widetilde{\pi}^C(\zeta(q), \sigma) \\ \delta(q, \sigma) = \delta^C(\zeta(q), \sigma) \end{cases} \qquad (23)$$

*3) If during the construction (beginning with $\wp_\lambda$) we encounter $\wp_x = \zeta(q)$ for some $x \in \Sigma^\star$, $q \in Q$ and any map $\zeta$ as defined in (2), then we stay within the graph of the copy of the initial-marked PFSA for all right extensions of $x$.*
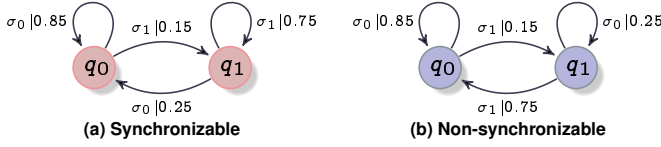
**Fig. 5. Synchronizable and non-synchronizable machines.** Synchronization is determination of the current state from observed past symbols. Not all PFSAs are synchronizable, *e.g.*, while the top machine is synchronizable, the bottom one is not. Note that a history of just one symbol suffices to determine the current state in the synchronizable machine (top), while no finite history can do the same in the non-synchronizable machine (bottom). A $\epsilon$-synchronizing string always exists (5) for a PFSA, which is not true for deterministic automata (34),(35).

*Proof:* (1) follows the ergodicity of QSPs, which makes $\wp_\lambda$ 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

$$\mathcal{E} = \left\{ e^i \in [0\ 1]^{|Q|}, i = 1, \cdots, |Q| \right\} \tag{24}$$

denote the set of distributions satisfying:

$$e^i|_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \tag{25}$$

(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. (24)) in the canonical construction starting from the stationary distribution $\wp_\lambda$ on the states of $\mathcal{P}_\mathcal{H}$.

Next we introduce the notion of $\epsilon$-synchronization of probabilistic automata (See Figure 5). Synchronization of automata is fixing or determining the current state. Not all PFSAs are synchronizable, but all are $\epsilon$-synchronizable (5).

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

$$\exists \vartheta \in \mathcal{E}, ||\wp_x - \vartheta||_\infty \leqq \epsilon \tag{26}$$

We next introduce 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 $\mathscr{D}(k)$.*

**Definition 8** (Symbolic Count Function). *For a string $s$ over $\Sigma$, the count function $\#^s : \Sigma^\star \to \mathbb{N} \cup \{0\}$, 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 $\Sigma$, the symbolic derivative $\phi^s : \Sigma^\star \to \mathscr{D}(|\Sigma| - 1)$ is defined:*

$$\phi^s(x)\big|_i = \frac{\#^s x \sigma_i}{\sum_{\sigma_i \in \Sigma} \#^s x \sigma_i} \tag{27}$$

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

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

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

**Proposition 1** ($\epsilon$-Convergence). *If $x \in \Sigma^\star$ is $\epsilon$-synchronizing, then:*

$$\forall \epsilon > 0, \lim_{|s| \to \infty} ||\phi^s(x) - \widetilde{\pi}([x], \cdot)||_\infty \leqq_{a.s} \epsilon \tag{28}$$

*Proof:* We use the Glivenko-Cantelli theorem (36) on uniform convergence of empirical distributions. Since $x$ is $\epsilon$-synchronizing:

$$\forall \epsilon > 0, \exists \vartheta \in \mathcal{E}, ||\wp_x - \vartheta||_\infty \leqq \epsilon \tag{29}$$

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

$$e^i|_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \tag{30}$$

Let $x$ $\epsilon$-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:

$$||\wp_x - \vartheta||_\infty \leqq \epsilon \Rightarrow \wp_x = \alpha \vartheta + (1 - \alpha) u \tag{31}$$

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

$$\lim_{|s| \to \infty} ||\phi^s(x) - \widetilde{\pi}(q, \cdot)||_\infty = \lim_{|s| \to \infty} ||\phi^s(x) - A_\alpha + A_\alpha - \widetilde{\pi}(q, \cdot)||_\infty$$

$$\leqq \underbrace{\lim_{|s| \to \infty} ||\phi^s(x) - A_\alpha||_\infty}_{\text{a.s. } \mathbf{0} \text{ by Glivenko-Cantelli}} + \lim_{|s| \to \infty} ||A_\alpha - \widetilde{\pi}(q, \cdot)||_\infty$$

$$\leqq_{a.s} (1 - \alpha) \left( ||\widetilde{\pi}(q, \cdot) - u||_\infty \right) \leqq_{a.s} \epsilon$$

This completes the proof. ∎

The notion of canonical representations, along with that of the symbolic derivatives will be used to establish the correctness of the stream operations in Section IX. Note that the canonical representation is free from the notion of initial states; intuitively this translates to our ability to carry out the stream operations (Table 1, main text) without knowledge of the initial states of the hidden models. The notion of the symbolic derivatives, along with Proposition 1 establishes that if the derivatives computed from two sufficiently long observed sequences $s_1, s_2$ match up closely, then the underlying generative PFSAs are also close. The detailed formulation in (5) proves that we can conclude that the distance between these underlying models is small with a high probability (in the PAC sense).

We also need to briefly describe the concept of a metric on the space of probabilistic automata established in (5).

**Proposition 2** (Metric For Probabilistic Automata). *For two strongly connected PFSAs $G_1, G_2$, denote the symbolic derivative at $x \in \Sigma^\star$ as $\phi^s_{G_1}(x)$ and $\phi^s_{G_2}(x)$ respectively. Then,*

$$\Theta(G_1, G_2) = \frac{|\Sigma| - 1}{|\Sigma|} \lim_{\substack{|s_1| \to \infty, \\ |s_2| \to \infty}} \sum_{x \in \Sigma^\star} \left\{ \frac{||\phi^{s_1}_{G_1}(x) - \phi^{s_2}_{G_2}(x)||_\infty}{|\Sigma|^{2|x|}} \right\}$$

*defines a metric on the space of probabilistic automata on $\Sigma$.*

*Proof:* The above metric is slightly different from the one introduced in (5). However, the proof of the metric properties follows almost identically. ∎

The following result is immediate, and justifies the expression given in Table 1 of main text (Row 4).

**Corollary 1** (For Proposition 2). *For any two PFSA $G_1, G_2$:*

$$0 \leqq \Theta(G_1, G_2) \leqq 1 \tag{32}$$

*Proof:* The lower bound is immediate by setting $G_1 = G_2$. For the upper bound, we note:

$$\Theta(G_1, G_2) = \frac{|\Sigma| - 1}{|\Sigma|} \lim_{\substack{|s_1| \to \infty, \\ |s_2| \to \infty}} \sum_{x \in \Sigma^\star} \left\{ \frac{||\phi^{s_1}_{G_1}(x) - \phi^{s_2}_{G_2}(x)||_\infty}{|\Sigma|^{2|x|}} \right\}$$

$$\leqq \frac{|\Sigma| - 1}{|\Sigma|} \lim_{\substack{|s_1| \to \infty, \\ |s_2| \to \infty}} \sum_{x \in \Sigma^\star} \left\{ \frac{\max \left( ||\phi^{s_1}_{G_1}(x) - \phi^{s_2}_{G_2}(x)||_\infty \right)}{|\Sigma|^{2|x|}} \right\}$$

$$= \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \frac{1}{|\Sigma|^{2|x|}} = \frac{|\Sigma| - 1}{|\Sigma|} \sum_{k=0}^\infty \frac{|\Sigma|^k}{|\Sigma|^{2k}} = \frac{|\Sigma| - 1}{|\Sigma|} \sum_{k=0}^\infty \frac{1}{|\Sigma|^k}$$

where the last two steps follow from the fact that there are $|\Sigma|^{|x|}$ strings of length $|x|$, which allows us to replace the sum over $x \in \Sigma^\star$ to over $k = |x|$. Finally, noting that $\frac{|\Sigma| - 1}{|\Sigma|} \sum_{k=0}^\infty \frac{1}{|\Sigma|^k} = 1$, completes the proof. ∎

Next, we elucidate the relevant algebraic structures on the space of PFSA.

## VIII. Algebraic Structures On PFSA Space

The material presented in this section is reproduced from the first author's previous work (8), and is included here for the sake of completeness.

The formulation in Section VII indicates that a symbolic dynamical process has a probabilistic finite state description if and only if the corresponding Nerode equivalence has a finite index.

**Definition 10** (Space of PFSA). *The space of all PFSA over a given symbol alphabet is denoted by $\mathscr{A}$ and the space of all probability measures $p$ inducing a finite-index probabilistic Nerode equivalence on the corresponding measure space $(\Sigma^\omega, \mathfrak{B}_\Sigma, p)$ is denoted by $\mathcal{P}$.*

As expected, there is a close relationship between $\mathscr{A}$ and $\mathcal{P}$, which is made explicit in the sequel.

**Definition 11** (PFSA Map $\mathbb{H}$). *Let $p \in \mathscr{P}$ and $G = (Q, \Sigma, \delta, q_0, \widetilde{\Pi}) \in \mathscr{A}$. The map $\mathbb{H} : \mathscr{A} \to \mathscr{P}$ is defined as $\mathbb{H}(G) = p$ such that the following condition is satisfied:*

$$\forall x = \sigma_1 \cdots \sigma_r \in \Sigma^\star, \tag{33}$$

$$p(x) = \widetilde{\Pi}(q_0, \sigma_1) \prod_{k=1}^{r-1} \widetilde{\Pi}(\delta^\star(q_0, \sigma_1 \cdots \sigma_k), \sigma_{k+1}) \tag{34}$$

*where $r \in \mathbb{N}$, the set of positive integers.*

**Definition 12** (Right Inverse $\mathbb{H}_{-1}$). *The right inverse of the map $\mathbb{H}$ is denoted by $\mathbb{H}_{-1} : \mathscr{P} \to \mathscr{A}$ such that*

$$\forall p \in \mathcal{P}, \ \mathbb{H}(\mathbb{H}_{-1}(p)) = p \tag{35}$$

An explicit construction of $\mathbb{H}_{-1}$ is reported in (9); we only require that such a map exists.

**Definition 13** (Perfect Encoding). *Given an alphabet $\Sigma$, a PFSA $G = (Q, \Sigma, \delta, q_0, \widetilde{\Pi})$ is said to be a perfect encoding of the measure space $(\Sigma^\omega, \mathscr{B}_\Sigma, p)$ if $p = \mathbb{H}(G)$.*

There are possibly many PFSA realizations that encode the same probability measure on $\mathscr{B}_\Sigma$ due to existence of non-minimal realizations and state relabeling; neither of them affect the underlying encoded measure. From this perspective, a notion of PFSA equivalence is introduced as follows:

**Definition 14** (PFSA Equivalence). *Two PFSA $G_1$ and $G_2$ are defined to be equivalent if $\mathbb{H}(G_1) = \mathbb{H}(G_2)$. In this case, we say $G_1 = G_2$.*

In the sequel, a PFSA $G$ implies the equivalence class of $G$, *i.e.*, $\{P \in \mathscr{A} : \mathbb{H}(P) = \mathbb{H}(G)\}$.

**Definition 15** (Structural Equivalence). *Two PFSA $G_i = (Q_i, \Sigma, \delta_i, q_0^i, \widetilde{\Pi}_i) \in \mathscr{A}$, $i = 1, 2$, are defined to have the equivalent (or identical) structure if $Q_1 = Q_2, q_0^1 = q_0^2$ and $\delta_1(q, \sigma) = \delta_2(q, \sigma), \forall q \in Q_1 \ \forall \sigma \in \Sigma$.*

**Definition 16** (Synchronous Composition of PFSA). *The binary operation of synchronous composition of two PFSA $G_i = (Q_i, \Sigma, \delta, q_0^{(i)}, \widetilde{\Pi}_i) \in \mathscr{A}$ where $i = 1, 2$, denoted by $\otimes : \mathscr{A} \times \mathscr{A} \to \mathscr{A}$ is defined as*

$$G_1 \otimes G_2 = \left(Q_1 \times Q_2, \Sigma, \delta', (q_0^{(1)}, q_0^{(2)}), \widetilde{\Pi}'\right) \tag{36}$$

*where $\delta'$ and $\widetilde{\Pi}'$ is computed as follows:*

$$\forall q_i \in Q_1, q_j \in Q_2, \sigma \in \Sigma, \begin{cases} \delta'((q_i, q_j), \sigma) = (\delta_1(q_i, \sigma), \delta_2(q_j, \sigma)) \\ \widetilde{\Pi}'((q_i, q_j), \sigma) = \widetilde{\Pi}_1(q_i, \sigma) \end{cases} \tag{37}$$

In general, $\otimes$ *i.e.* synchronous composition is non-commutative.

**Proposition 3** (Synchronous Composition of PFSA). *Let $G_1, G_2 \in \mathscr{A}$. Then, $\mathbb{H}(G_1) = \mathbb{H}(G_1 \otimes G_2)$ and therefore $G_1 = G_1 \otimes G_2$ in the sense of Definition 14.*

*Proof:* See Theorem 4.5 in (9). ∎

Synchronous composition of PFSA allows transformation of PFSA with disparate structures to non-minimal descriptions that have the same underlying graphs. This assertion is crucial for the development in the sequel, since any binary operation defined for two PFSA with an identical structure can be extended to the general case on account of Definition 16 and Proposition 3.

Next we show that a restricted PFSA subspace can be assigned the algebraic structure of an Abelian group. We first construct the Abelian group on a subspace of probability measures, and then induce the group structure on this subspace of PFSA via the isomorphism between the two spaces.

**Definition 17** (Restricted PFSA Space). *Let $\mathscr{A}^+ = \{G = (Q, \Sigma, \delta, q_0, \widetilde{\Pi}) : \widetilde{\Pi}(q, \sigma) > 0 \ \forall q \in Q \ \forall \sigma \in \Sigma\}$ that is a proper subset of $\mathscr{A}$. It follows that the transition map of any PFSA in the subset $\mathscr{A}^+$ is a total function. We restrict the map $\mathbb{H} : \mathscr{A} \to \mathscr{P}$ on a smaller domain $\mathscr{A}^+$, that is, $\mathbb{H}^+ : \mathscr{A}^+ \to \mathscr{P}^+$, i.e., $\mathbb{H}^+ = \mathbb{H}|_{\mathscr{A}^+}$.*

**Definition 18** (Restricted Probability Measure). *Let $\mathscr{P}^+ \triangleq \{p \in \mathscr{P} : p(x) \neq 0, \forall x \in \Sigma^\star\}$ that is a proper subset of $\mathscr{P}$. Each element of $\mathscr{P}^+$ is a probability measure that assigns a non-zero probability to each string on $\mathfrak{B}_\Sigma$. Similar to Definition 17, we restrict $\mathbb{H}_{-1}$ on $\mathscr{P}^+$, i.e., $\mathbb{H}_{-1}^+ = \mathbb{H}_{-1}|_{\mathscr{P}^+}$.*

Since we do not distinguish PFSA in the same equivalence class (See Definition 14), we have the following result.

**Proposition 4** (Isomorphism of $\mathbb{H}^+$). *The map $\mathbb{H}^+$ is an isomorphism between the spaces $\mathscr{A}^+$ and $\mathscr{P}^+$, and its inverse is $\mathbb{H}_{-1}^+$.*

*Proof:* Immediate from preceding discussion. ∎

**Definition 19** (Abelian Operation on $\mathscr{P}^+$). *The addition operation $\oplus : \mathscr{P}^+ \times \mathscr{P}^+ \to \mathscr{P}^+$ is defined by $p_3 \triangleq p_1 \oplus p_2, \forall p_1, p_2 \in \mathscr{P}^+$ such that*

1) $p_3(\epsilon) = 1$.
2) $\forall x \in \Sigma^\star$ and $\tau \in \Sigma$, $\dfrac{p_3(x\tau)}{p_3(x)} = \dfrac{p_1(x\tau)p_2(x\tau)}{\sum_{\alpha \in \Sigma} p_1(x\alpha)p_2(x\alpha)}$

$p_3$ is a well-defined probability measure on $\mathscr{P}^+$, since $\forall x \in \Sigma^\star$:

$$\Sigma_{\tau \in \Sigma} p_3(x\tau) = \Sigma_{\tau \in \Sigma} \frac{p_1(x\tau)p_2(x\tau)}{\sum_{\alpha \in \Sigma} p_1(x\alpha)p_2(x\alpha)} p_3(x) = p_3(x) \tag{38}$$

**Proposition 5** (abelian Group of PFSA). *The algebra $(\mathscr{P}^+, \oplus)$ forms an Abelian group.*

*Proof:* Closure property and commutativity of $(\mathscr{P}^+, \oplus)$ are obvious. The associativity, existence of identity and existence of inverse element are established next.

*(1) Associativity i.e.* $(p_1 \oplus p_2) \oplus p_3 = p_1 \oplus (p_2 \oplus p_3)$. Now, $\forall x \in \Sigma^\star, \tau \in \Sigma$, we have:

$$\frac{((p_1 \oplus p_2) \oplus p_3)(x\tau)}{((p_1 \oplus p_2) \oplus p_3)(x)} = \frac{(p_1 \oplus p_2)(x\tau)p_3(x\tau)}{\sum_{\beta \in \Sigma}(p_1 \oplus p_2)(x\beta)p_3(x\beta)}$$

$$= \frac{p_1(x\tau)(p_2 \oplus p_3)(x\tau)}{\sum_{\beta \in \Sigma} p_1(x\beta)(p_2 \oplus p_3)(x\beta)} = \frac{(p_1 \oplus (p_2 \oplus p_3))(x\tau)}{(p_1 \oplus (p_2 \oplus p_3))(x)} \tag{39}$$

*(2) Existence of identity*: Let us introduce a probability measure $\mathbf{i}_\circ$ of symbol strings such that:

$$\forall x \in \Sigma^\star, \ \mathbf{i}_\circ(x) = \left(\frac{1}{|\Sigma|}\right)^{|x|} \tag{40}$$

where $|x|$ denotes the length of the string $x$. Then, $\forall \tau \in \Sigma$ that $\frac{\mathbf{i}_\circ(x\tau)}{\mathbf{i}_\circ(x)} = \frac{1}{|\Sigma|}$. For a measure $p \in \mathscr{P}^+$ and $\forall \tau \in \Sigma$,

$$\frac{(p \oplus \mathbf{i}_\circ)(x\tau)}{(p \oplus \mathbf{i}_\circ)(x)} = \frac{p(x\tau)\mathbf{i}_\circ(x\tau)}{\sum_{\alpha \in \Sigma} p(x\alpha)\mathbf{i}_\circ(x\alpha)} = \frac{p(x\tau)}{p(x)}$$

This implies that $p \oplus \mathbf{i}_\circ = \mathbf{i}_\circ \oplus p = p$ by Definition 19 and by commutativity. Therefore, $\mathbf{i}_\circ$ is the identity of the monoid $(\mathscr{P}^+, \oplus)$.
*(3) Existence of inverse*: $\forall p \in \mathscr{P}^+, \forall x \in \Sigma^\star$ and $\forall \tau \in \Sigma$, let $-p$ be defined by the following relations:

$$(-p)(\epsilon) = 1 \tag{41}$$

$$\frac{(-p)(x\tau)}{(-p)(x)} = \frac{p^{-1}(x\tau)}{\sum_{\alpha \in \Sigma} p^{-1}(x\alpha)} \tag{42}$$

**(a) Summing arbitray PFSA models via Non-minimal realizations**



**(b) Annihilation identity with PFSA defined on binary alphabet**

Fig. 6. Addition of arbitrary PFSAs with the same alphabet, using non-minimal realizations to equate structures (via synchronous composition)



**(a) Zero PFSA for binary alphabet**   **(b) Zero PFSA for trinary alphabet**

Fig. 7. Zero PFSAs for different alphabet sizes
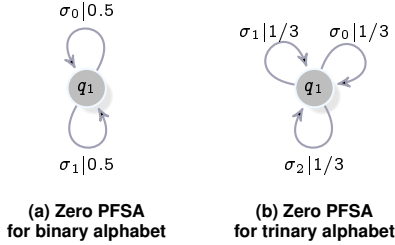
Then, we have:

$$\frac{(p \oplus (-p))(x\tau)}{(p \oplus (-p))(x)} = \frac{p(x\tau)(-p)(x\tau)}{\sum_{\alpha \in \Sigma} p(x\alpha)(-p)(x\alpha)} = \frac{1}{|\Sigma|} \quad (43)$$

This gives $p \oplus (-p) = \mathbf{i}_\circ$ which completes the proof. ∎

We denote the zero-element $\mathbf{i}_\circ$ of the Abelian group $(\mathscr{P}^+, \oplus)$ as *flat white noise (FWN)*.

### A. Explicit Computation of the Abelian Operation $\oplus$

The isomorphism between $\mathscr{P}^+$ and $\mathscr{A}^+$ (See Proposition 4) induces the following Abelian operation on $\mathscr{A}^+$.

**Definition 20** (Addition Operation on PFSA). *Given any $G_1, G_2 \in \mathscr{P}^+$, the addition operation $+ : \mathscr{A}^+ \times \mathscr{A}^+ \to \mathscr{A}^+$ is defined as:*

$$G_1 + G_2 = \mathbb{H}_{-1}^+(\mathbb{H}^+(G_1) \oplus \mathbb{H}^+(G_2))$$

If the summand PFSA have identical structure (i.e., their underlying graphs are identical), then the explicit computation of this sum is stated as follows.

**Proposition 6** (PFSA Addition). *If two PFSA $G_1, G_2 \in \mathscr{A}^+$ are of the same structure, i.e., $G_i = (Q, \Sigma, \delta, q_0, \widetilde{\Pi}_i), i = \{1, 2\}$, then we have $G_1 + G_2 = (Q, \Sigma, \delta, q_0, \widetilde{\Pi})$ where*

$$\widetilde{\Pi}(q, \sigma) = \frac{\widetilde{\Pi}_1(q, \sigma)\widetilde{\Pi}_2(q, \sigma)}{\sum_{\alpha \in \Sigma} \widetilde{\Pi}_1(q, \alpha)\widetilde{\Pi}_2(q, \alpha)} \quad (44)$$

*Proof:* Let $p_i = \mathbb{H}^+(G_i)$, $i = \{1, 2\}$ and since $G_1, G_2$ have the same structure, we have from Eq. (33):

$$\forall \sigma \in \Sigma, \forall x \text{ s.t. } \delta^\star(q_0, x) = q \in Q,$$

$$\frac{p_i(x\sigma)}{p_i(x)} = \widetilde{\Pi}_i(\delta^\star(q_0, x), \sigma) = \widetilde{\Pi}_i(q, \sigma) \quad (45)$$

Now, by Definition 19 and Definition 11,

$$\widetilde{\Pi}(q, \sigma) = \frac{(p_1 \oplus p_2)(x\sigma)}{(p_1 \oplus p_2)(x)} = \frac{p_1(x\sigma)p_2(x\sigma)}{\sum_{\alpha \in \Sigma} p_1(x\alpha)p_2(x\alpha)}$$

$$= \frac{\frac{p_1(x\sigma)p_2(x\sigma)}{p_1(x)p_2(x)}}{\sum_{\alpha \in \Sigma} \frac{p_1(x\alpha)p_2(x\alpha)}{p_1(x)p_2(x)}} = \frac{\widetilde{\Pi}_1(q, \sigma)\widetilde{\Pi}_2(q, \sigma)}{\sum_{\alpha \in \Sigma} \widetilde{\Pi}_1(q, \alpha)\widetilde{\Pi}_2(q, \alpha)}$$

∎

The extension to the general case is achieved by using synchronous composition of probabilistic machines.

**Proposition 7** (PFSA Addition (General case)). *Given two PFSA $G_1, G_2 \in \mathscr{A}^+$, the sum $G_1 + G_2$ is computed via Proposition 6 and Definition 16 as follows:*

$$G_1 + G_2 = (G_1 \otimes G_2) + (G_2 \otimes G_1) \quad (46)$$

*Proof:* Noting that $G_1 \otimes G_2$ and $G_2 \otimes G_1$ have the same structure up to state relabeling, it follows from Proposition 3:

$$\mathbb{H}^+(G_1 + G_2) = \mathbb{H}^+(G_1) \oplus \mathbb{H}^+(G_2) \quad (\text{See Definition 20})$$

$$= \mathbb{H}^+(G_1 \otimes G_2) \oplus \mathbb{H}^+(G_2 \otimes G_1)$$

$$= \mathbb{H}^+ \left( (G_1 \otimes G_2) + (G_2 \otimes G_1) \right)$$

which completes the proof. ∎

**Example 1.** *Let $G_1$ and $G_2$ be two PFSA with identical structures, such that the probability morph matrices are:*
$$\widetilde{\Pi}_1 = \begin{pmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \end{pmatrix} \text{ and } \widetilde{\Pi}_2 = \begin{pmatrix} 0.1 & 0.9 \\ 0.6 & 0.4 \end{pmatrix} \qquad (47)$$

*Then the $\widetilde{\Pi}$-matrix for the sum $G_1 + G_2$, denoted by $\widetilde{\Pi}_{12}$, is*
$$\widetilde{\Pi}_{12} = \begin{pmatrix} 0.1 \times 0.2 & 0.9 \times 0.8 \\ 0.6 \times 0.4 & 0.4 \times 0.6 \end{pmatrix} \xrightarrow[rows]{Normalize} \begin{pmatrix} 0.027 & 0.973 \\ 0.5 & 0.5 \end{pmatrix}$$

### IX. Correctness Of Stream Operations

In this section (Section IX), we prove that the stream operations described in Table I of main text are indeed correct.

#### A. Independent Stream Copy

We show that the "Independent Stream Copy" operation produces an independent realization from a pseudo-copy of the PFSA model generating the input stream. First, we formalize the notion of pseudo-copies.

**Definition 21.** *Given a PFSA $G = (Q, \Sigma, \delta, \widetilde{\pi})$ in the canonical representation, a pseudo-copy is a canonical PFSA $\mathbb{P}_\gamma(G) = (Q, \Sigma, \delta, \mathbb{P}(\widetilde{\pi}))$, where we have:*
$$\mathbb{P}_\gamma(\Pi) = \gamma[\mathbb{I} - (1-\gamma)\Pi]^{-1}\Pi \qquad (48)$$
*for some scalar $\gamma \in (0,1)$.*

We note that that while the row-stochastic matrix $\Pi$ may not be invertible, and $[\mathbb{I} - \Pi]$ is definitely singular (since $\Pi$ has a eigenvalue at 1); the matrix $\gamma[\mathbb{I} - (1-\gamma)\Pi]^{-1}$ is always well-defined for $\gamma \in (0,1)$, and additionally is a non-negative row-stochastic matrix (37).

We use the following notation:

**Notation 5.** *For a given string $s$, the underlying PFSA generator is denoted as $G \leftarrow s$, and for a given PFSA $G$, $G \rightarrow s$ is a realization generated by $G$. Note that $G \leftarrow s$ automatically implies that we are referring to the PFSA generator as $|s| \rightarrow \infty$, since one cannot have a unique generator for bounded strings.*

**Proposition 8** (Independent Stream Copy). *Given a symbol stream $s$ with a hidden PFSA generator $G$, let stream $s'$ be generated via:*

*1 Generate stream $\omega_0$ from FWN*
*2 Read current symbol $\sigma_1$ from $s$, and $\sigma'$ from $\omega_0$*
*3 If $\sigma = \sigma'$, then write to output $s'$*
*4 Move read positions one step to right, and go to step 1*

*Then, we have:*

*1) Well-defined convergence of underlying models:*
$$\lim_{|s'| \rightarrow \infty} (G' \leftarrow s') = \lim_{|s| \rightarrow \infty} \mathbb{P}_{\frac{1}{2}}(G \leftarrow s) \qquad (49)$$

*2) If $s', s''$ are generated from the above algorithm from the same input stream $s$, then, in the limit of infinite length, $s', s''$ are independent realizations of $\mathbb{P}_{\frac{1}{2}}(G \leftarrow s)$.*

*3) If $s'_1, s'_2$ are generated by the algorithm for input streams $s_1, s_2$ respectively, then we have:*
$$\forall \epsilon > 0, \Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \leqq \epsilon$$
$$\Rightarrow \Theta(G'_1 \leftarrow s'_1, G'_2 \leftarrow s'_2) \leqq |\Sigma|\epsilon \qquad (50)$$

$$\forall \epsilon > 0, \Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \geqq \epsilon$$
$$\Rightarrow \Theta(G'_1 \leftarrow s'_1, G'_2 \leftarrow s'_2) \geqq \frac{|\Sigma|}{(2|\Sigma| - 1)^2}\epsilon \qquad (51)$$

*Proof:* (1) Let $\gamma$ be the probability that the first symbol in the input stream $s$ is recorded in the output. Since, the stream $\omega_0$ is FWN, we conclude that $\gamma = \frac{1}{2}$, and that $\gamma$ is also the constant probability that any symbol in $s$ is recorded. Thus, assuming that the symbolic derivatives computed are exact (*i.e.* the input stream is infinite), the transition matrix $M$ of a realization for the PFSA $G' \leftarrow s'$ can be expressed as a function of the transition matrix $\Pi$ for $G \leftarrow s$ as:

$$M = \gamma\Pi + (1-\gamma)\gamma\Pi^2 + (1-\gamma)^2\gamma\Pi^3 + \cdots \qquad (52)$$
$$\Rightarrow M = \gamma[\mathbb{I} - (1-\gamma)\Pi]^{-1}\Pi = \mathbb{P}_\gamma(\Pi) \qquad (53)$$

Since, the transformation $[\mathbb{I} - (1-\gamma)\Pi]^{-1}$ is invertible, the rank of $M$ is the same as $\Pi$, which implies that states in $G \leftarrow s$ cannot collapse when we pass to $\mathbb{P}_\gamma(G)$, implying in turn that $G' \leftarrow s'$ is has indeed the same minimal structure as $G \leftarrow s$, which establishes claim (1).

(2) Claim (1) implies:
$$\lim_{|s'| \rightarrow \infty} G' \leftarrow s' = \lim_{|s''| \rightarrow \infty} G'' \leftarrow s'' \qquad (54)$$

The independence claim then follows immediately from noting that random erasure, as executed by the stated algorithm, eliminates any possibility of synchronization between the states of the same underlying model $G' \leftarrow s'$ in the limit of infinite string lengths .

(3) Consider the PFSAs $G_1 \leftarrow s_1, G_2 \leftarrow s_2$ as $|s_1|, |s_2| \rightarrow \infty$. Let us bring them to the same structure, via the transformations $(G_1 \leftarrow s_1) \otimes (G_2 \leftarrow s_2)$ and $(G_2 \leftarrow s_2) \otimes (G_1 \leftarrow s_1)$ respectively (See [REF]). Let us denote the transition matrices of the PFSAs in their transformed representations as $\Pi_1, \Pi_2$ respectively. Then denoting $\Delta = \Pi_1 - \Pi_2$, and $\Delta' = \mathbb{P}_\gamma(\Pi_1) - \mathbb{P}_\gamma(\Pi_2)$, we claim:
$$\frac{\gamma}{(2-\gamma)^2}||\Delta||_\infty \leqq ||\Delta'||_\infty \leqq \frac{1}{\gamma}||\Delta||_\infty \qquad (55)$$

To establish this claim, we first note that for any stochastic matrix $A$, we have:
$$\gamma[\mathbb{I} - (1-\gamma)A]^{-1}A = \frac{1}{1-\gamma}\gamma[\mathbb{I} - (1-\gamma)A]^{-1} - \frac{\gamma}{1-\gamma}\mathbb{I} \quad (56)$$

which implies the upper bound in Eq. (105) as follows:
$$(1-\gamma)\left(\mathbb{P}_\gamma(\Pi_1) - \mathbb{P}_\gamma(\Pi_2)\right)$$
$$= \gamma[\mathbb{I} - (1-\gamma)\Pi_1]^{-1} - \gamma[\mathbb{I} - (1-\gamma)\Pi_2]^{-1}$$
$$= \gamma(1-\gamma)[\mathbb{I} - (1-\gamma)\Pi_2]^{-1}(\Pi_2 - \Pi_1)[\mathbb{I} - (1-\gamma)\Pi_1]^{-1}$$
This implies:
$$\Rightarrow \Delta' = -\gamma[\mathbb{I} - (1-\gamma)\Pi_2]^{-1}\Delta[\mathbb{I} - (1-\gamma)\Pi_1]^{-1}$$
$$\Rightarrow ||\Delta'||_\infty \leqq \gamma||[\mathbb{I} - (1-\gamma)\Pi_2]^{-1}||_\infty$$
$$\times ||[\mathbb{I} - (1-\gamma)\Pi_1]^{-1}||_\infty ||\Delta||_\infty$$
$$\Rightarrow ||\Delta'||_\infty \leqq \gamma \times \frac{1}{\gamma} \times \frac{1}{\gamma} \times ||\Delta||_\infty = \frac{1}{\gamma}||\Delta||_\infty$$

And the lower bound follows from noting:
$$\Rightarrow \Delta' = -\gamma[\mathbb{I} - (1-\gamma)\Pi_2]^{-1}\Delta[\mathbb{I} - (1-\gamma)\Pi_1]^{-1}$$
$$\Rightarrow -\gamma\Delta = [\mathbb{I} - (1-\gamma)\Pi_2]\Delta'[\mathbb{I} - (1-\gamma)\Pi_1]$$
$$\Rightarrow \gamma||\Delta||_\infty \leqq ||[\mathbb{I} - (1-\gamma)\Pi_2]||_\infty ||[\mathbb{I} - (1-\gamma)\Pi_1]||_\infty ||\Delta'||_\infty$$
$$\Rightarrow ||\Delta'||_\infty \geqq \frac{\gamma}{||[\mathbb{I} - (1-\gamma)\Pi_2]||_\infty ||[\mathbb{I} - (1-\gamma)\Pi_1]||_\infty}||\Delta||_\infty$$
$$\geqq \frac{\gamma}{(2-\gamma)^2}||\Delta||_\infty \qquad (57)$$

Next, we compute bounds on the probability morph matrices.

We denote the probability morph matrices for the relevant PFSAs as follows (PFSAs on left, morph matrices on right):

| | |
|---|---|
| $(G_1 \leftarrow s_1) \otimes (G_2 \leftarrow s_2)$ | $\widetilde{\Pi}_1$ |
| $(G_2 \leftarrow s_2) \otimes (G_1 \leftarrow s_1)$ | $\widetilde{\Pi}_2$ |
| $\mathbb{P}_\gamma((G_1 \leftarrow s_1) \otimes (G_2 \leftarrow s_2))$ | $\mathbb{P}(\widetilde{\Pi}_1)$ |
| $\mathbb{P}_\gamma((G_2 \leftarrow s_2) \otimes (G_1 \leftarrow s_1))$ | $\mathbb{P}(\widetilde{\Pi}_2)$ |

And, additionally, we use the notation:
$$\widetilde{\Delta} = \widetilde{\Pi}_1 - \widetilde{\Pi}_2 \qquad (58)$$
$$\widetilde{\Delta}' = \mathbb{P}(\widetilde{\Pi}_1) - \mathbb{P}(\widetilde{\Pi}_2) \qquad (59)$$

Without loss of generality, we assume that for each PFSA, given a state, each symbol leads to a distinct state. This can be arranged via state splitting if necessary, and implies:
$$||\Delta||_\infty = ||\widetilde{\Pi}_1 - \widetilde{\Pi}_2||_\infty \qquad (60)$$

$$||\Delta'||_\infty = ||\mathbb{P}(\widetilde{\Pi}_1) - \mathbb{P}(\widetilde{\Pi}_2)||_\infty \qquad (61)$$

which therefore leads to the bounds:

$$\frac{\gamma}{(2-\gamma)^2}||\widetilde{\Pi}_1 - \widetilde{\Pi}_2||_\infty \leqq ||\mathbb{P}(\widetilde{\Pi}_1) - \mathbb{P}(\widetilde{\Pi}_2)||_\infty$$
$$\leqq \frac{1}{\gamma}||\widetilde{\Pi}_1 - \widetilde{\Pi}_2||_\infty \qquad (62)$$

Recall the definition of the PFSA metric (See Proposition 2):

$$\Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2)$$
$$= \frac{|\Sigma|-1}{|\Sigma|} \lim_{|s_1| \to \infty, |s_2| \to \infty} \sum_{x \in \Sigma^\star} \left\{ \frac{||\phi_{G_1}^{s_1}(x) - \phi_{G_2}^{s_2}(x)||_\infty}{|\Sigma|^{2|x|}} \right\} \qquad (63)$$

and, note that:

$$\Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \leqq \epsilon$$
$$\Rightarrow \forall x \in \Sigma^\star, ||\phi_{G_1}^{s_1}(x) - \phi_{G_2}^{s_2}(x)||_\infty \leqq \epsilon \qquad (64)$$

$$\Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \geqq \epsilon$$
$$\Rightarrow \forall x \in \Sigma^\star, ||\phi_{G_1}^{s_1}(x) - \phi_{G_2}^{s_2}(x)||_\infty \geqq \epsilon \qquad (65)$$

Since the bounds established in Eq. (62) is applicable to any non-minimal realization of the PFSAs $(G_1 \leftarrow s_1) \otimes (G_2 \leftarrow s_2)$ and $(G_2 \leftarrow s_2) \otimes (G_1 \leftarrow s_1)$, considering the full $\Sigma$-ary tree as the limiting "unfolded" realization, we conclude from Eq. (62) that:

$$\Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \leqq \epsilon$$
$$\Rightarrow \Theta(G_1' \leftarrow s_1', G_2' \leftarrow s_2') \leqq \frac{1}{\gamma}\epsilon \qquad (66)$$

and also:

$$\Theta(G_1 \leftarrow s_1, G_2 \leftarrow s_2) \geqq \epsilon$$
$$\Rightarrow \Theta(G_1' \leftarrow s_1', G_2' \leftarrow s_2') \geqq \frac{\gamma}{(2-\gamma)^2}\epsilon \qquad (67)$$

The desired bounds then follow from noting that in the stated algorithm, we have $\gamma = \frac{1}{|\Sigma|}$. This completes the proof. ∎

**Remark 1.** *Proposition 8 establishes that the stream $s'$ obtained from an input stream $s$ may not be a realization from the hidden generator for the latter (i.e. stream $s$), but is a realization from a PFSA which is a pseudo-copy of the generator for $s$.*

*Also, note that it is not true in general that a pseudo-copy is close to the original PFSA, in the sense of our metric.*

*Nevertheless, Proposition 8 shows that if the distance between two PFSAs is small, then so is the distance between their pseudo-copies; and if the distance between two PFSAs is large, then so is the distance between the pseudo-copies.*

*Thus, if we determine the distance between pseudo-copies, we have a good estimate of the distance between the original machines.*

### B. Stream Summation

**Proposition 9** (Stream Summation). *Given a symbol streams $s_1, s_2$ with hidden PFSA generators $G_1, G_2$, let stream $s'$ be generated via:*
1 *Read current symbols $\sigma_i$ from $s_i$ $(i = 1, 2)$*
2 *If $\sigma_1 = \sigma_2$, then write to output $s'$*
3 *Move read positions one step to right, and go to step 1*

*Then, denoting the FWN generator as $W$, we have:*
1) *If $(G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ has a single state in its minimal realization, then we have:*
$$(G' \leftarrow s') = (G_1 \leftarrow s_1) + (G_2 \leftarrow s_2) \qquad (68)$$
*i.e., then $s'$ is an exact realization of $(G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ in the limit $|s_1|, |s_2| \to \infty$.*
2) *If $(G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ has $N > 1$ states in its minimal realization, we have the lower bound:*
$$\Theta((G_1 \leftarrow s_1) + (G_2 \leftarrow s_2), W) \geqq \epsilon$$
$$\Rightarrow \Theta((G' \leftarrow s'), W) \geqq \frac{\epsilon}{|\Sigma|^N + |\Sigma|^{N-1}} \qquad (69)$$
3) *We have the upper bound:*

$$\Theta((G_1 \leftarrow s_1) + (G_2 \leftarrow s_2), W) \leqq \epsilon$$
$$\Rightarrow \Theta\left((G' \leftarrow s'), W\right) \leqq \epsilon \qquad (70)$$

*Proof:* (1) If $(G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ has a single causal state:
$$\forall x \in \Sigma^\star, \lim_{|s'| \to \infty} \phi^{s'}(x) = v, \text{ where } \sum_i v_i = 1, v_i \geqq 0 \qquad (71)$$

We assume without loss of generality that $G_1 \leftarrow s_1, G_2 \leftarrow s_2$ in their canonical representations (See Definition 6). Thus, the streams $s_1, s_2$ can be assumed to start at states of the canonical representation which maps to the corresponding stationary distribution over the states of the corresponding initial-marked PFSAs (See Definition 6, and the discussion immediately after). Also, note that we can delete arbitrary prefixes from $s_1$ and $s_2$, and still assume that they start at these states. Thus, we delete prefixes of $s_1, s_2$ up to the point where the next symbols match, and we see our first output symbol.

Since we see a symbol in the output if there is a match in both $s_1$ and $s_2$, it follows that the probability of seeing the first symbol in $s'$ as $\sigma_i$ is given by:
$$\lim_{|s_1|, |s_2| \to \infty} \frac{1}{\sum_i \phi^{s_1}(\lambda)|_i \phi^{s_2}(\lambda)|_i} \phi^{s_1}(\lambda)|_i \phi^{s_2}(\lambda)|_i = v_i \qquad (72)$$

Also, since $G_1 \leftarrow s_1$ and $G_2 \leftarrow s_2$ can be assumed to have the same graph without loss of generality (via considering non-minimal realizations if necessary), it follows that the next hidden states $q, q'$ in $s_1, s_2$ after seeing the first output symbol, are still synchronized. Thus, we conclude, that if the first observed symbol in $s'$ is $\sigma$, then the probability that the next symbol is $\sigma_j$, is also given by:
$$\lim_{|s_1|, |s_2| \to \infty} \frac{1}{\sum_i \phi^{s_1}(\lambda)|_j \phi^{s_2}(\lambda)|_j} \phi^{s_1}(\lambda)|_i \phi^{s_2}(\lambda)|_i = v_j \qquad (73)$$

It follows from straightforward induction, that at any point in $s'$, the distribution of the next symbol is given by $v$, *i.e*, in the limit $s_1 \to \infty, s_2 \to \infty$, $s'$ is an exact realization from the sum $(G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ (if the latter has a single causal state).

(2) Let $H = (G_1 \leftarrow s_1) + (G_2 \leftarrow s_2)$ and assume $\Theta(H, W) \geqq \epsilon$. Also, let the set of states in the minimal realization of the canonical representation for $H$ be $Q$, and additionally let $\text{CARD}(Q) = N > 1$. It follows from the definition of our metric (and the fact that $\epsilon'$-synchronizing strings must occur for all $\epsilon' > 0$), that:
$$\forall q \in Q_H, ||\widetilde{\pi}(q, \cdot) - \mathcal{U}_\Sigma||_\infty \geqq \epsilon \qquad (74)$$
where, as before, $\mathcal{U}_\Sigma = \left( \frac{1}{|\Sigma|} \cdots \frac{1}{|\Sigma|} \right)$. We observe that the stream summation algorithm can be thought to be producing the output sequence by traversing the arcs in the canonical representation for $H$ augmented with "jumps" (See Fig. 8), where there are unreported and unlabeled transitions back to the state corresponding to the equivalence class $[\lambda]$ from each state: whenever we have a mismatch we jump back to $[\lambda]$. The probabilities of these back transitions can be easily computed, but not important here. However, this implies that if $q'$ is a state in the canonical representation for $G' \leftarrow s'$, then we have (assuming $\widetilde{\pi}'$ is the morph function for $G' \leftarrow s'$):
$$\widetilde{\pi}'(q', \cdot) = p(q')\widetilde{\pi}([\lambda], \cdot) + (1 - p(q'))\widetilde{\pi}(q, \cdot),$$
$$\text{for some } q \in Q, \text{ where } p(q') \in [0, 1] \qquad (75)$$
Since, $\widetilde{\pi}'(q', \cdot)$ is a weighted average of $\widetilde{\pi}([\lambda], \cdot)$, and $\widetilde{\pi}(q, \cdot)$ (both of which satisfy Eq. (82)), it is possible that:
$$||\widetilde{\pi}'(q', \cdot) - \mathcal{U}_\Sigma||_\infty \leqq \epsilon \qquad (76)$$
Assume, if possible, that ($Q'$ is the state set in the minimal realization of the canonical representation for $G' \leftarrow s'$):
$$\forall q' \in Q', ||\widetilde{\pi}'(q', \cdot) - \mathcal{U}_\Sigma||_\infty \leqq \epsilon \qquad (77)$$
We note that the same argument in claim (1) implies that:
$$\widetilde{\pi}'([\lambda], \cdot) = \widetilde{\pi}([\lambda], \cdot) \geqq \epsilon \qquad (78)$$
where the inequality follows from Eq. (82). Since, $\widetilde{\pi}'([\lambda], \cdot)$ is some weighted average of all $\widetilde{\pi}'(q', \cdot)$ vectors, and by assumption Eq. (77), the norm of the difference of each of these vectors from $\mathcal{U}_\Sigma|$ is bounded above by $\epsilon$, it follows that:
$$||\widetilde{\pi}'([\lambda], \cdot) - \mathcal{U}_\Sigma||_\infty \leqq \epsilon \qquad (79)$$
which is a contradiction. Thus, there exists at least one state $q_\star \in Q'$,
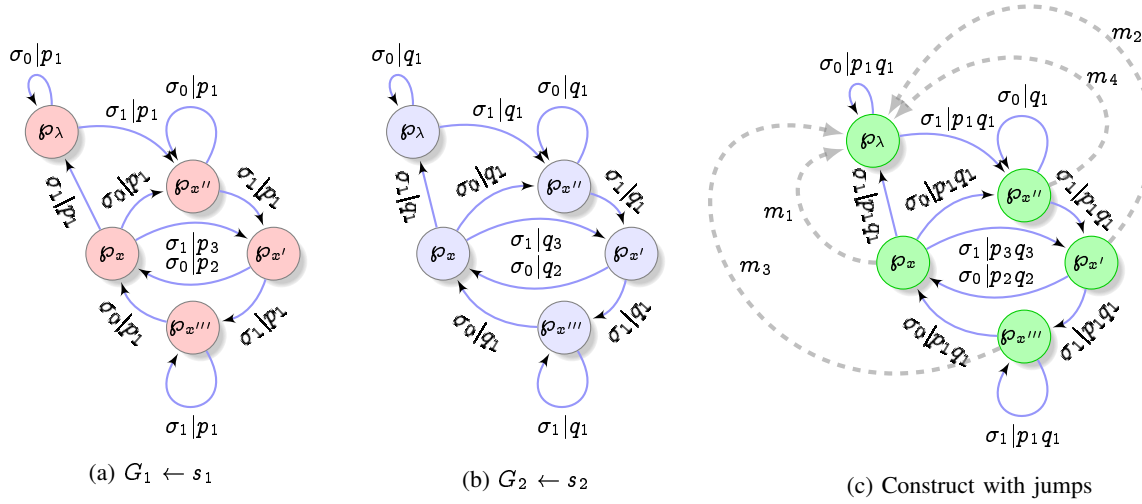
Fig. 8. Illustration for Proposition 9. Note that using the same structure for $G_1 \leftarrow s_1$ and $G_2 \leftarrow s_2$ causes no loss of generality, since arbitrary PFSAs over the same alphabet can be brought to the same structure via possibly non-minimal realizations

such that

$$||\widetilde{\pi}'(q_\star, \cdot) - \mathcal{U}_\Sigma||_\infty \geqq \epsilon \tag{80}$$

We also note that $G' \leftarrow s'$ has at most $N$ states, since, if none of the $\widetilde{\pi}'$ rows are equal, then we can represent $G' \leftarrow s'$ using the same graph as $H$, with the rows of $\widetilde{\pi}$ replaced with those from $\widetilde{\pi}'$.

Now, we compute $\Theta(G' \leftarrow s', W)$. We note that since $G' \leftarrow s'$ has at most $N$ states, $\lim_{|s'| \to \infty} \phi^{s'}(x)$ equals $\widetilde{\pi}'(q_\star, \cdot)$ at least once every $N$ levels, which implies:

$$\Theta(G' \leftarrow s', W) \geqq \frac{|\Sigma| - 1}{|\Sigma|} \sum_{i=0}^{\infty} \frac{1}{|\Sigma|^{2i+N}} \epsilon = \frac{\epsilon}{|\Sigma|^N + |\Sigma|^{N-1}} \tag{81}$$

(3) Assume $\Theta(H, W) \leqq \epsilon$. It follows from the definition of our metric (and the fact that $\epsilon'$-synchronizing strings must occur for all $\epsilon' > 0$), that:

$$\forall q \in Q_H, ||\widetilde{\pi}(q, \cdot) - \mathcal{U}_\Sigma||_\infty \leqq \epsilon \tag{82}$$

where, as before, $\mathcal{U}_\Sigma = \left( \frac{1}{|\Sigma|} \quad \cdots \quad \frac{1}{|\Sigma|} \right)$. Since, (using the notation used for claim (2) above) $\widetilde{\pi}'(q', \cdot)$ is a weighted average of $\widetilde{\pi}([\lambda], \cdot)$, it follows immediately that:

$$\forall q' \in Q', ||\widetilde{\pi}'(q', \cdot) - \mathcal{U}_\Sigma||_\infty \leqq \epsilon \tag{83}$$

and also, that the norm of the difference between any weighted average (with the weights being positive and summing to unity), of the rows of the $\widetilde{\pi}$ matrix from $\mathcal{U}_\Sigma$, is bounded above by $\epsilon$, This implies that each term in $\Theta(G' \leftarrow s', W)$ is bounded above by $\epsilon$, which establishes the desired bound:

$$\Theta(G' \leftarrow s', W) \leqq \epsilon \tag{84}$$

This completes the proof. ∎

**Remark 2.** *Note that the lower bound established in claim (2) is obviously not tight; since if $N = 1$, then we have exact summation, whereas the bound is off by a factor of $|\Sigma|$.*

**Remark 3.** *Proposition 9 establishes that the stream summation algorithm works perfectly if the summands sum to a one state machine, which includes FWN. For arbitrary inputs, the deviation of the realized sum from FWN is small if the deviation of the sum of the original models is small; and conversely the deviation of the realized sum from FWN is large if the deviation of the sum of the original models is large.*

**Corollary 2** (Contrapositives to Proposition 9). *Using the notation of Proposition 9, we have:*

*Lower Bound:* $\Theta(G' \leftarrow s', W) < \epsilon$
$$\Rightarrow \Theta((G_1 \leftarrow s_1) + (G_2 \leftarrow s_2), W) < \left( |\Sigma|^N + |\Sigma|^{N-1} \right) \epsilon \tag{85}$$

*Equality:* $\Theta(G' \leftarrow s', W) = 0$

$$\Rightarrow \Theta((G_1 \leftarrow s_1) + (G_2 \leftarrow s_2), W) = 0 \tag{86}$$

*Upper Bound:* $\Theta(G' \leftarrow s', W) > \epsilon$
$$\Rightarrow \Theta((G_1 \leftarrow s_1) + (G_2 \leftarrow s_2), W) > \epsilon \tag{87}$$

*Proof:* (Equality:) Note that $G' \leftarrow s' = W$, and the fact that every state in $G'$ is a convex combination of $[\lambda]$ and some state $q \in Q_H$, implies:

$$\forall x \in \Sigma^\star, \lim_{|s'| \to \infty} \phi^{s'}(x) = p\widetilde{\pi}([\lambda], \cdot) + (1 - p)\widetilde{\pi}(q, \cdot) = \mathcal{U}_\Sigma \tag{88}$$

Also, since $\widetilde{\pi}([\lambda], ) = \lim_{|s'| \to \infty} \phi^{s'}(\lambda)$ for arbitrary input streams, it follows that:

$$\forall q \in Q_H, p\mathcal{U}_\Sigma + (1 - p)\widetilde{\pi}(q, \cdot) = \mathcal{U}_\Sigma \Rightarrow \forall q \in Q_H, \widetilde{\pi}(q, \cdot) = \mathcal{U}_\Sigma \tag{89}$$

which establishes Eq. (86).

The other bounds follow by taking the contrapositive of the inequalities established in Proposition 9. ∎

**Remark 4.** *Corollary 2 implies that if the output sequence from the stream summation algorithm is FWN, then the summands are exact inverses of each other.*

### C. Stream Inversion

**Lemma 6** (Stream Summation to FWN). *Let streams $s_i, i = 1, \cdots, |\Sigma|$ be $|\Sigma|$ independent realizations from a PFSA $G$ defined over the alphabet $\Sigma$. And let $s'$ be generated as follows:*

*1 Read current symbols $\sigma_i$ from $s_i$ ($i = 1, \cdots, |\Sigma|$)*
*2 If $\sigma_i \neq \sigma_j$ for all distinct $i, j$, then write $\sigma_1$ to output $s'$*
*3 Move read positions one step to right, and go to step 1*

*Then, we have:*

$$\Theta(G' \leftarrow s', W) = 0 \tag{90}$$

*where $W$ is the FWN generator for the alphabet size $|\Sigma|$.*

*Proof:* Let the set of states for the minimal realization of the canonical representation for $G$ be $Q$, and the corresponding morph function be $\widetilde{\pi}$. Similarly, let the state set for $G' \leftarrow s'$ be $Q'$, and the associated probability morph function be $\widetilde{\pi}'$.

Let $s'_i$ be the sequence obtained by copying the current symbol from the input stream $s_i$ in Step (2) of the above scheme. (Thus, $s'_1 = s'$.) It is obvious from the symmetry of the scheme that:

$$\forall i, j \in \{1, \cdots, |\Sigma|\}, (G'_i \leftarrow s'_i) = (G'_j \leftarrow s'_j) \tag{91}$$

It follows that, if $f_i^j$ is the frequency of the $j^{th}$ symbol from the alphabet in the stream $s'_i$, then we have:

$$\forall i, \lim_{|s'_i| \to \infty} \frac{f_i^j}{|s'_i|} = \widetilde{\pi}'([\lambda], \cdot)\big|_j \tag{92}$$

$$\Rightarrow \lim_{|s'_i| \to \infty} \frac{\sum_i f_i^j}{|s'_i|} = \lim_{|s'_i| \to \infty} \frac{|\Sigma| f_i^j}{|s'_i|} = |\Sigma| \tilde{\pi}'([\lambda], \cdot)\big|_j \quad (93)$$

where we also used the fact that $|s'_i| = |s'_i|$. Next, noting that we have an output symbol in each $s'_i$ only when each new symbol is distinct, we conclude:

$$\forall j, k, \sum_i f_i^j = \sum_i f_i^k \quad (94)$$

which in turn implies for the $j^{th}$ and $k^{th}$ alphabet symbols:

$$\lim_{|s'_i| \to \infty} \frac{\sum_i f_i^j}{|s'_i|} = \lim_{|s'_i| \to \infty} \frac{|\Sigma| f_i^j}{|s'_i|} = |\Sigma| \tilde{\pi}'([\lambda], \cdot)\big|_j$$
$$= \lim_{|s'_i| \to \infty} \frac{\sum_i f_i^k}{|s'_i|} = \lim_{|s'_i| \to \infty} \frac{|\Sigma| f_i^k}{|s'_i|} = |\Sigma| \tilde{\pi}'([\lambda], \cdot)\big|_k \quad (95)$$

and hence we conclude:

$$\tilde{\pi}'([\lambda], \cdot) = \lim_{|s'| \to \infty} \phi^{s'}(\lambda) = \mathcal{U}_\Sigma \quad (96)$$

Next, denote the $r^{th}$ symbol in the stream $s'$ as $s'(r)$. Then, if we assume that the streams $s_i$ were all synchronized to the same state of $G$ just prior to the generation of $s'(r)$, we have:

$$\forall \sigma_i, \sigma_k \in \Sigma, Prob(s'(r) = \sigma_i)$$
$$= c \prod_{\sigma_j \in \Sigma} \tilde{\pi}(\wp_\lambda, \sigma_j) = Prob(s'(r) = \sigma_k) \quad (97)$$

which implies:

$$\forall \sigma_k \in \Sigma, Prob(s'(r) = \sigma_k) = \frac{1}{|\Sigma|} \quad (98)$$

Next, we consider the following construction: Consider the PFSA $G$, with the streams $s_i$ traversing the transitions via the symbol-labeled arcs, with each $s_i$ initialized to the state $[\lambda]$. Note that we have a new symbol in the output $s'$ if all current symbols in the $|\Sigma|$ input streams are distinct; which can occur in two possible ways:

1) all $s_i$ streams are synchronized to some state $q \in Q$, and a distinct symbol is generated for each $s_i$
2) no such synchronization; but the symbols generated are distinct

In the second case, we assume that a re-initialization occurs; $i.e.$, all the streams jump back to state $[\lambda]$ before the distinct symbols are generated causing the new output symbol.

Note that this construction causes no loss of generality, as we are simply defining a path, with jumps, for the given input streams through the PFSA $G$.

Denote the probability distribution of the output symbol, when the streams are synchronized at some state $q \in Q$ be $v^q$, $i.e.$ $v_i^q$ is the probability of seeing the $i^{th}$ symbol, given that we indeed have a new output symbol.

Next we observe that the streams $s_i$ can be assumed to be synchronized at $[\lambda]$ when the first symbol appears in the output $s'$ (since deletion of arbitrary leading prefixes has no effect in the limit of infinite data). Thus, we have (from Eq. (96)):

$$\lim_{|s'| \to \infty} \phi^{s'}(\lambda) = \mathcal{U}_\Sigma \quad (99)$$

Note that the next symbol may be produced after a "silent" jump to some state $q \in Q$. Additionally, the probability that the jump occurs to a specific state $q$ is an explicit function of the parameters (morph probabilities, and transition structure) of the PFSA $G$. However we do not need to compute these probabilities; we simply conclude that:

$$\forall x \in \Sigma^\star, \lim_{|s'| \to \infty} \phi^{s'}(x) = \sum_{q \in Q} p(q, x) v^q \quad (100)$$

where $p(q, x) \in [0, 1], \sum_q p(q, x) = 1 \quad (101)$

Noting that Eq. (98) establishes that $\forall q \in Q, v^q = \mathcal{U}_\Sigma$, we conclude:

$$\forall x \in \Sigma^\star, \lim_{|s'| \to \infty} \phi^{s'}(x) = \mathcal{U}_\Sigma \quad (102)$$

which establishes that $G' \leftarrow s'$ is the FWN generator. This completes the proof. ∎
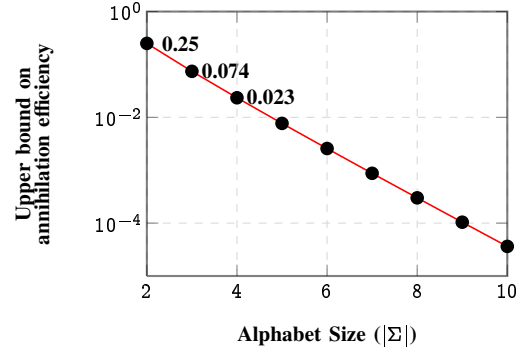


Fig. 9. Upper bound on annihilation efficiency $\frac{(|\Sigma|-1)!}{|\Sigma|^{|\Sigma|}}$ vs alphabet size $|\Sigma|$. This illustrates why having a fine quantization would necessitate large amounts of data to pass the self-annihilation test.

**Proposition 10** (Stream Inversion). *Given a stream $s$ which is generated by some hidden PFSA $G$, let stream $s'$ be generated via:*
*1 Generate $|\Sigma| - 1$ independent copies of $s_1$: $s_1, \cdots, s_{|\Sigma|-1}$*
*2 Read current symbols $\sigma_i$ from $s_i$ ($i = 1, \cdots, |\Sigma| - 1$)*
*3 If $\sigma_i \neq \sigma_j$ for all distinct $i, j$, then write $\Sigma \setminus \bigcup_{i=1}^{|\Sigma|-1} \sigma_i$ to output $s'$*
*4 Move read positions one step to right, and go to step 1*
*Then, we have:*

$$\Theta(-G, G' \leftarrow s') = 0 \quad (103)$$

*Proof:* Follows immediately from Lemma 6. ∎

**Proposition 11** (Asymptotic Complexity). *The asymptotic time complexity of carrying out the stream operations is $O(|s||\Sigma|$*

*Proof:* The algorithmic steps in each of the operations of stream copy and stream summation proceed in a symbol-by-symbol fashion, with no memory of previous symbols. Also, each step involves a constant number of integer comparisons. Assuming that each new symbol from the FWN processes involved can be generated with constant complexity, we conclude that the asymptotic time complexity of both stream summation and stream copy is $O(|s|)$. The stream inversion operation needs to generate $|\Sigma| - 1$ stream copies, implying that its asymptotic time complexity is $O(|s||\Sigma|)$. ∎

### D. Annihilation Efficiency

To pass the self-annihilation test, a data stream must be sufficiently long; and the required length $|s|$ of the input $s$ with a specified threshold $\epsilon^\star$ is dictated by the characteristics of the generating process. Thus the rate of convergence of the self-annihilation error as a function of $|s|$ quantifies the sample complexity of information annihilation. Let $s'$ be obtained from $s$ via stream inversion, and $s''$ be obtained via stream summation of $s$ and $s'$. Then, it follows that $s''$ is always a realization of the FWN process, which has an uniform probability of generating any symbol at any point. Thus, for any $x \in \Sigma^\star$, the vectors $\phi^{s''}(x)$ (See Table I, row 4) are empirical distributions which converge to the flat distribution as $|s''| \to \infty$. Additionally, the Central Limit Theorem (CLT) (28) dictates the convergence rate to scale as $1/\sqrt{|s''|}$ irrespective of the generating process for the input $s$. However, selective erasure in annihilation (See Table I) implies that $|s''| < |s|$, and the expected shortening ratio $\beta = \mathbf{E}(|s''|/|s|)$ does indeed depend on the generating process. We refer to $\beta$ as the *annihilation efficiency*, since the convergence rate of the self-annihilation error scales as $1/\sqrt{\beta|s|}$.

Next, we compute $\beta$ in terms of the symbol frequencies:

**Proposition 12.** *Given an input stream $s$, let stream $s'$ be produced via stream inversion from $s$, and let $s''$ be produced via stream summation of $s$ and $s'$. Let $p_i$ be the probability of observing symbol $\sigma_i \in \Sigma$. Then, we have*

$$\beta = \mathbf{E}(|s''|/|s|) = (|\Sigma| - 1)! \prod_i p_i \quad (104)$$

*Proof:* To generate $s'$ from $s$, we first need to generate $|\Sigma| - 1$ independent stream copies of $s$. It is clear from the stated algorithm for independent stream copy (See Table I in main text and Proposition 8) that the expected length of each of these copies is $\frac{1}{|\Sigma|}|s|$. The probability of obtaining a symbol in the output by comparing these $|\Sigma| - 1$ streams (to get $s'$) is simply the probability of seeing a different symbol in each of the copied streams (as stated in the algorithm for stream inversion in Table I of main text). Denoting this probability as $\alpha$, we have:

$$\alpha = (|\Sigma| - 1)! \sum_{j=1}^{|\Sigma|} \prod_{i \neq j} p_i$$

$$= (|\Sigma| - 1)! \left\{ \prod_i p_i \right\} \left\{ \sum_i \frac{1}{p_i} \right\} = \frac{1}{H} |\Sigma|! \prod_i p_i \quad (105)$$

where $H$ is the harmonic mean of the probability vector $p$. Thus, the expected length of $s'$ is $\alpha|s|$. The final step is stream summation of $s$ and $s'$ to obtain $s''$. We note that the probability of seeing symbol $\sigma_i$ in the inverted stream $s'$ is $k/p_i$, where $k = \frac{1}{\sum_i \frac{1}{p_i}} = \frac{H}{|\Sigma|}$. It follows that stream summation of $s$ and $s'$, would result in an expected length of $\frac{H}{|\Sigma|}|s'|$, which when combined with Eq. (105) completes the proof. ∎

**Corollary 3.** *Annihilation efficiency satisfies:*

$$(|\Sigma| - 1)! \eta^{|\Sigma|} \left\{ \frac{1}{\eta} + 1 - |\Sigma| \right\} \leqq \beta \leqq \frac{(|\Sigma| - 1)!}{|\Sigma|^{|\Sigma|}} \quad (106)$$

*where $\eta$ is the probability of occurrence of the rarest symbol in the input stream $s$.*

*Proof:* The upper bound follows from noting that the product $\prod_i p_i$ is maximized when $\forall i, p_i = 1/|\Sigma|$. The lower bound is obtained by assuming that $\eta = \min_i p_i$, upon which the minimum value of the product $\prod_i p_i$ is given by:

$$\eta^{|\Sigma|-1}(1 - \eta(|\Sigma| - 1)) = \eta^{|\Sigma|} \left\{ \frac{1}{\eta} + 1 - |\Sigma| \right\} \quad (107)$$
∎

**Remark 5.** *The upper bound for the annihilation efficiency is realized if the input $s$ is FWN.*

### E. Distance Between Hidden Generators

**Definition 22** (FWN Deviation Estimators). *For a string $s$, the complete white noise deviation estimator $\zeta(s)$ is defined as:*

$$\zeta(s) = \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^s(x) - \mathcal{U}_\Sigma||_\infty \right\} \quad (108)$$

*And the partial white noise deviation estimator $\hat{\zeta}(s, \ell)$ is defined:*

$$\hat{\zeta}(s, \ell) = \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\ell} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^s(x) - \mathcal{U}_\Sigma||_\infty \right\} \quad (109)$$

*which only carries out the summation for all strings up to length $\ell$.*

**Proposition 13** (Causality Claim 1). *Given a stream $s$, and denoting the hidden generator for $s$ as $G_s$, and the zero model as $W$, we have:*

$$\lim_{|s| \to \infty} |\Theta(G_s, W) - \zeta(s)| = 0 \quad (110)$$

*Proof:* For a string $s'$ generated by the model $G(s)$, we denote $\lim_{|s'| \to \infty} \phi^{s'}(x)$ as $\phi^\star(x)$. Then:

$$\Theta(G(s), W) = \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \frac{1}{|\Sigma|^{2|x|}} ||\phi^\star(x) - \mathcal{U}_\Sigma||_\infty$$

$$\leqq \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^\star(x) - \phi^s(x)||_\infty \right\}$$

$$+ \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^s(x) - \mathcal{U}_\Sigma||_\infty \right\}$$
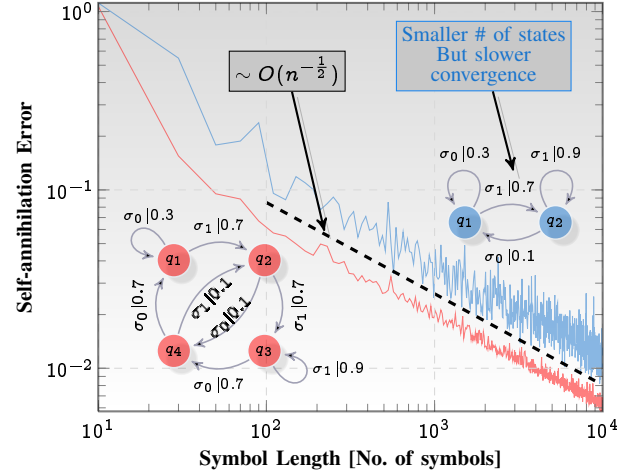


Fig. 10. **Convergence rate of the self-annihilation error** Shown to scale as $O(1/\sqrt{n})$ as dictated by the Central Limit Theorem. The convergence rates do not depend directly on the descriptional complexity of the generating processes; note that the data from the two state process has a slower convergence rate compared to that from the four state process. As discussed in the section on computational complexity, the convergence rate scales as $O(\sqrt{\beta n})$ where $\beta$ is the expected shortening of the input stream due to the selective erasure of the symbols in the different steps of the annihilation process. We establish in Proposition 12 that if $p_i$ is the occurrence probability of the symbol $\sigma_i$ in the stream $s''$, then we have: $\beta = (|\Sigma| - 1)! \prod_i p_i$.

$$\leqq \underbrace{\frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^\star(x) - \phi^s(x)||_\infty \right\}}_{B} + \zeta(s)$$

Starting with $\zeta(s)$ on the RHS, we end up with:

$$\zeta(s) \leqq \Theta(G(s), W)$$

$$+ \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^\star(x) - \phi^s(x)||_\infty \right\} \quad (111)$$

which then implies:

$$|\Theta(G_s, W) - \zeta(s)| \leqq \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^\star(x) - \phi^s(x)||_\infty \right\}$$

We note that $\phi^s(x)$ is an empirical estimate of $\phi^\star(x)$, which then implies via the Glivenko-Cantelli theorem (36) that

$$\forall x \in \Sigma^\star, ||\phi^s(x) - \phi^\star(x)||_\infty) \xrightarrow{a.s.} 0 \quad (112)$$

which completes the proof. ∎

Finally we establish our causality claim: while the deviation from FWN is estimated by function $\hat{\zeta}(s, \ell)$ from a finite observed string $s$ and consideration of finite histories of length bounded by $\ell$, it converges to the deviation of the underlying process from FWN in the limit of infinite data (See next proposition). It thus follows that the distance $\epsilon$ calculated by annihilating a stream $s$ against a second stream $s'$ converges to the absolute deviation of $G_s - G_{s'}$ from the FWN generator $W$.

**Proposition 14** (Causality Claim 2). *Given a stream $s$, and denoting the hidden generator for $s$ as $G_s$, and the zero model as $W$, we have:*

$$\lim_{|s| \to \infty} |\Theta(G_s, W) - \hat{\zeta}(s, \ell)| \leqq \epsilon \quad (113)$$

*where the partial estimator is evaluated upto length $\ell = \dfrac{\ln \frac{1}{\epsilon}}{\ln |\Sigma|}$*

*Proof:* Noting that we have:

$$|\Theta(G_s, W) - \hat{\zeta}(s, \ell)| \leqq |\Theta(G_s, W) - \zeta(s)|$$

$$+ \left| \frac{|\Sigma| - 1}{|\Sigma|} \sum_{x \in \Sigma^\star \setminus \Sigma^\ell} \left\{ \frac{1}{|\Sigma|^{2|x|}} ||\phi^s(x) - \mathcal{U}_\Sigma||_\infty \right\} \right|$$

$$\leq |\Theta(G_s, W) - \zeta(s)| + \frac{1}{|\Sigma|^\ell}$$

the result follows using Proposition 13. ∎

## X. Quantization Techniques

Information annihilation operates on symbolic sequences. Thus, we need to specify a quantization scheme to map possibly continuous-valued data streams to symbolic sequences. This is accomplished by the choice of a symbol alphabet, where each letter in the alphabet denotes a slice of the data range. Given a particular quantization scheme, we map each continuous-valued observation to the symbol representing the slice of the data range to which the observation belongs. Thus any chosen quantization scheme incurs error, which can be made small by using a fine quantization, *i.e.*, a large alphabet.

However, the length of the observed data limits the size of the alphabet that we can use. This is a direct consequence of the fact that the *annihilation efficiency* falls rapidly with the alphabet size (See Proposition 12, and Fig. 9). Thus, if $s$ is the input stream, $s'$ is obtained via stream inversion from $s$, and $s''$ is the output from stream summation of $s$ and $s'$, then the expected ratio of the lengths $|s''|/|s|$ falls rapidly as the alphabet size $|\Sigma|$ is increased, making the estimation of the deviation of $s''$ from FWN more and more difficult. Since the convergence rate of the self-annihilation error scales as $1/\sqrt{\beta|s|}$, it follows that the self-annihilation error increases rapidly with finer quantization (See Fig. 11 for illustration on the EEG dataset).

### A. Desired Properties of Quantization Schemes

It follows that a good quantization scheme is defined by the following properties:

1) The frequency of the rarest symbol in the quantized data streams are too small. This is to ensure that symbols are represented faithfully according to its generation probability from each state in the hidden model; too few occurrences of a particular symbol may represent statistical fluctuations rather than the generation probabilities.

2) The average self-annihilation error for the observed data streams is small, *i.e.*, if $\epsilon_{ii}$ is the self-annihilation error for the observed data stream $s_i$, then we require that $\frac{1}{T} \sum_{i=1}^{T} \epsilon_{ii}$ is small (where $T$ is the total number of observed data streams).

3) The average discrimination between data-streams is high, *i.e.*, if for two streams $s_i, s_j$, the similarity computed by information annihilation is $\epsilon_{ij}$, then we require that $\frac{1}{T(T-1)} \sum_{i=1}^{T} \sum_{\substack{j=1 \\ i \neq j}}^{T} \epsilon_{ij}$

is large (where $T$ is the total number of observed data streams).

One approach to choosing a quantization that satisfies the stated properties is the following: We restrict ourselves to maximum-entropy quantizations, *i.e.*, schemes in which each symbol occurs with the same frequency in the data set. In Fig. 11 plates (a)-(c), we show three such maximum-entropy schemes for the EEG-dataset. The alphabet size is increased from 2 to 4, and we choose the slices of the data-range such that each slice contains an approximately equal number of data points. For example, in plate (c) of Fig. 11, each of the four slices contains approximately 25% of the total number of observations in the data set. Such maximum-entropy schemes guarantee that property 1 (See above) is satisfied. For the remaining properties, we plot the mean self-annihilation error and the mean discrimination, for each alphabet size. As expected, we see that finer alphabets lead to high average discrimination, while at the same time incur high average self-annihilation errors (See Fig. 11(d)). The ratio of the two quantities is more useful, and in Fig. 11(e), we note that the trinary maximum-entropy quantization minimizes this ratio; implying high discrimination and low self-annihilation error.
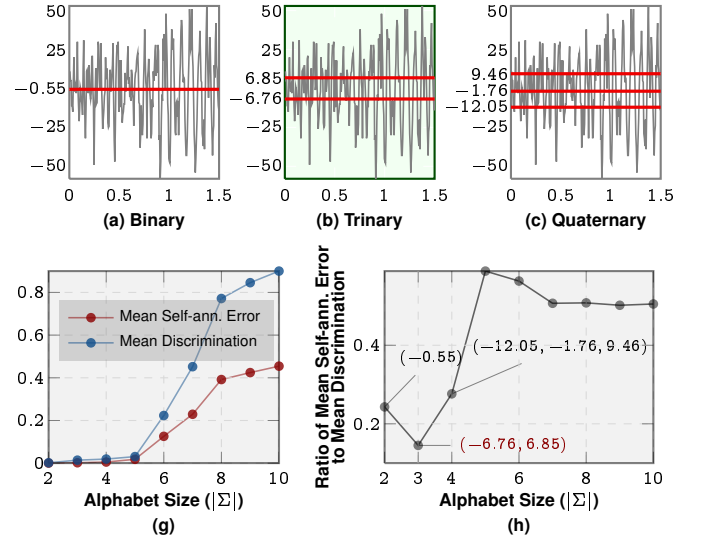


Fig. 11. Plates(a)-(c): Maximum-entropy quantization schemes for the EEG dataset with alphabet sizes 2, 3, and 4 respectively (1.5s of a single time series shown for clarity). As explained in the text, "maximum-entropy" in this context refers to the fact that each data slice contains approximately an equal number of observations. Plate (d) shows how the average self-annihilation error, as well as the average discrimination between different data streams, increases exponentially with alphabet size. Plate (e) shows that the ratio of the average self-annihilation error to average discrimination has a minimum at an alphabet size 3.
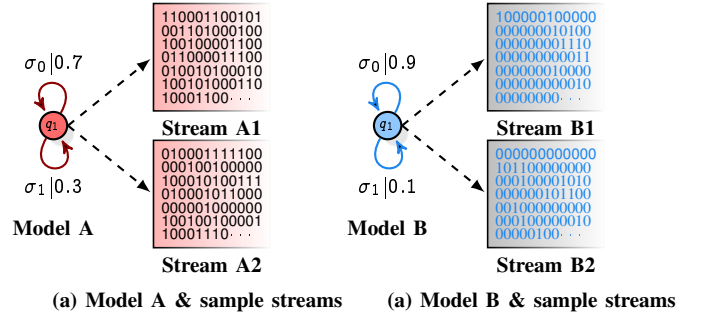


Fig. 12. Two distinct PFSA models (alphabet $\Sigma = \{\sigma_0, \sigma_1\}$) and initial sections of generated sample streams ($\sigma_0$ shown as 0, and $\sigma_1$ shown as 1)): While streams generated in independent runs of the same model have near zero mutual information, they are correctly evaluated as having similar generators via data smashing (See Tables III and IV). Also, runs from the different models also have near-zero mutual information, while smashing them correctly reveals a significant difference in the generators.

We note that if our chosen quantization is too coarse, then distinct processes may evaluate to be similar. However, too coarse an alphabet produces errors in only one direction; identical processes will still evaluate to be identical (or nearly so), provided the streams pass the self-annihilation test.

## XI. Comparison Against Simple Statistical Approches to Similarity

### A. Smashing & Mutual Information

Smashing two finite quantized data streams manipulates the statistical information contained in them. Notions of information-theoretic interdependence of sequential data have been investigated in the literature extensively; one such concept is *mutual information* between streams. For discrete random variables, mutual information quantifies the amount of information one random variable contains about another.

Formally, let $X, Y$ be discrete random variables with alphabets $\Sigma_X, \Sigma_Y$ and probability mass function $p(x) = Prob\{X = x : x \in \Sigma_X\}, p(y) = Prob\{Y = y : y \in \Sigma_Y\}$. Also, considering $(X, Y)$ as a single vector vector-valued random variable, we have

## TABLE III
DISTANCE MATRIX OBTAINED BY SMASHING STREAMS FROM MODELS A AND B. *(Note clear clusters corresponding to runs from the same model)*

| DATA SMASHING | A1 | A2 | B1 | B2 |
|---|---|---|---|---|
| A1 | **0.005** | **0.019** | 0.264 | 0.269 |
| A2 | **0.021** | **0.006** | 0.246 | 0.253 |
| B1 | 0.262 | 0.251 | **0.005** | **0.009** |
| B2 | 0.264 | 0.254 | **0.011** | **0.006** |

## TABLE IV
PAIRWISE MUTUAL INFORMATION OF STREAMS FROM MODELS A AND B. *(No indication of generative difference)*

| Mutual Inf. | A1 | A2 | B1 | B2 |
|---|---|---|---|---|
| A1 | 0.89 | 0.00000476 | 0.00017155 | 0.00002713 |
| A2 | 0.0000047 | 0.87 | 0.00001186 | 0.00003927 |
| B1 | 0.00017155 | 0.00001186 | 0.48 | 0.00000996 |
| B2 | 0.00002713 | 0.00003927 | 0.00000996 | 0.47 |



**(a)** $r = 0.209$

**(b)** $r = 0.257$

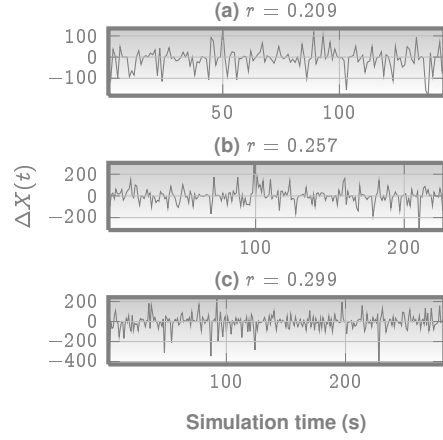**(c)** $r = 0.299$

**Simulation time (s)**

Fig. 13. **Model system.** The Lotka-Volterra system of reactions is a stochastic model that captures a simple two-species predator-prey dynamical system. We parameterize the system with $r$ (the propensity of the predation reaction), and generate time series data using Gillespie's stochastic simulation algorithm. Plates (a0-(c) show prey numbers varying with time in three different runs with different $r$ values.

the mass function $p(x, y)$. Then, mutual information between the discrete random variables $X, Y$ is defined as:

$$I(X, Y) = \sum_{y \in \Sigma_Y} \sum_{x \in \Sigma_X} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) \quad (114)$$

Mutual information is related to the notion of entropy: the entropy of a random variable is a measure of the amount of information required on the average to describe the random variable; while mutual information is the amount of information one variable contains about the other; or, more precisely, the degree to which the uncertainty in one can be reduced by knowing about the other.

Needless to say, if two data streams $X, Y$ are generated independently from the same underlying generator, then we have:

$$I(X, Y) = \sum_{y \in \Sigma_Y} \sum_{x \in \Sigma_X} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right)$$

$$= \sum_{y \in \Sigma_Y} \sum_{x \in \Sigma_X} p(x)p(y) \log \left( \frac{p(x)p(y)}{p(x)p(y)} \right) = 0 \quad (115)$$

Thus, sharing a common generative process does not imply a high mutual information; and conversely, high mutual information is indicative of some sort of statistical synchronization between the generative processes; which may be very different themselves.

Thus, the concept of mutual information and data smashing is "orthogonal" in the sense that while we measure statistical dependence for computing the former, the streams need to be statistically independent (or very nearly so) for the latter to work. Note that in the computation of the anti-stream, we generated streams that approximate independent copies of the input stream, which are then manipulated to yield the inverse. The algorithm requires this independence; in absence of which Proposition 10 falls apart.

We can illustrate these points by a simple example (See Fig. 12). We consider two simple one-state PFSAs (A and B), with different event probabilities, and generated 10000 bit streams $A1, A2$ and $B1, B2$, Note that simply "running" a given PFSA twice, *i.e.* choosing a start state randomly, and generating symbols in accordance with the event probabilities, implies that the generations are independent. We smash the streams $A1, A2, B1, B2$ against each other, and compute the pairwise distance matrix shown in Table III. Note that streams $(A1, A2)$ annihilate nearly perfectly, as do the stream $(B1, B2)$ ; while streams $(A1, B1)$, $(A1, B2)$, $(A2, B1)$ and $(A2, B2)$ fail to do so. This results in clearly clustered values in Table III, which correctly indicate that streams $(A1, A2)$ and $(B1, B2)$ have identical generators, which differ significantly from each other.

Pairwise computation of mutual information between the streams $A1, A2, B1, B2$ is not expected to reveal this generative structure. Since the streams are generated independently, the mutual information between any two distinct streams would be zero (or nearly so

for finitely generated streams). This is illustrated in Table IV. Note that while the diagonal terms (which represent the self-information or entropy) are high; all off-diagonal terms are very nearly zero, and no clusters are discernible.

Thus, we can summarize:

- Mutual information measures the degree of statistical dependence between data streams; data smashing computes the distance between the generative processes, provided the data streams are independent or nearly so.
- We proved that maximizing entropy of a single stream maximizes the annihilation efficiency (See Proposition 12, and its corollary)
- Thus, data smashing is conceptually orthogonal to the notion of mutual information

### B. Smashing Vs Simple Statistical Measures

The pairwise distances computed via data smashing is clearly a function of the statistical information buried in the streams. However, it might not be easy to find the right statistical tool to mine this information for a particular problem. In this section we provide an example of a dynamical system, in which data smashing is able to recover meaningful nontrivial structure, which is missed by simple statistical measures.

We consider the Lotka-Volterra system of stochastic "reactions", modeling a simple closed eco-system of two species, one of which preys on the other. While deterministic differential equation models for this system do exist (and is widely studied), a more realistic model is this set of three simple reactions (See Fig. 13, plate A), primarily due to its ability to model the stochastic component. The generally accepted method to solve such systems, to produce the time traces of population numbers (See Fig. 13, plate B), is via the Gillespie's stochastic simulation algorithm. (Note: While the preceding theoretical development assumes ergodicity and stationarity, the theoretical considerations fall apart gracefully as we deviate from these idealizations).

In our simple model, as shown below:

$$X \xrightarrow{1.0} 2X \quad (R1)$$

$$X + Y \xrightarrow{0.005} 2Y \quad (R2)$$

$$Y \xrightarrow{r} \varnothing \quad (R3)$$

we consider the propensity of one of the reactions to be parameterized by $r$, which ranges between 0.2 to 0.3 in steps of 0.001. For each set of reaction parameters, we simulate the system 1000 times for a maximum of $200s$ using Gillespie's algorithm. In each simulation run, we initialize the system with $X = 128, Y = 256$. We

**A** Smashing distance vs difference of means vs difference of variance with $\triangle X$-time series



**B** Smashing distance vs difference of means vs difference of variance with $\triangle X$-time series after normalizing each series to have zero mean



**C** Smashing distance vs difference of means vs difference of variance with $\triangle X$-time series after normalizing each series to have zero mean and unit variance
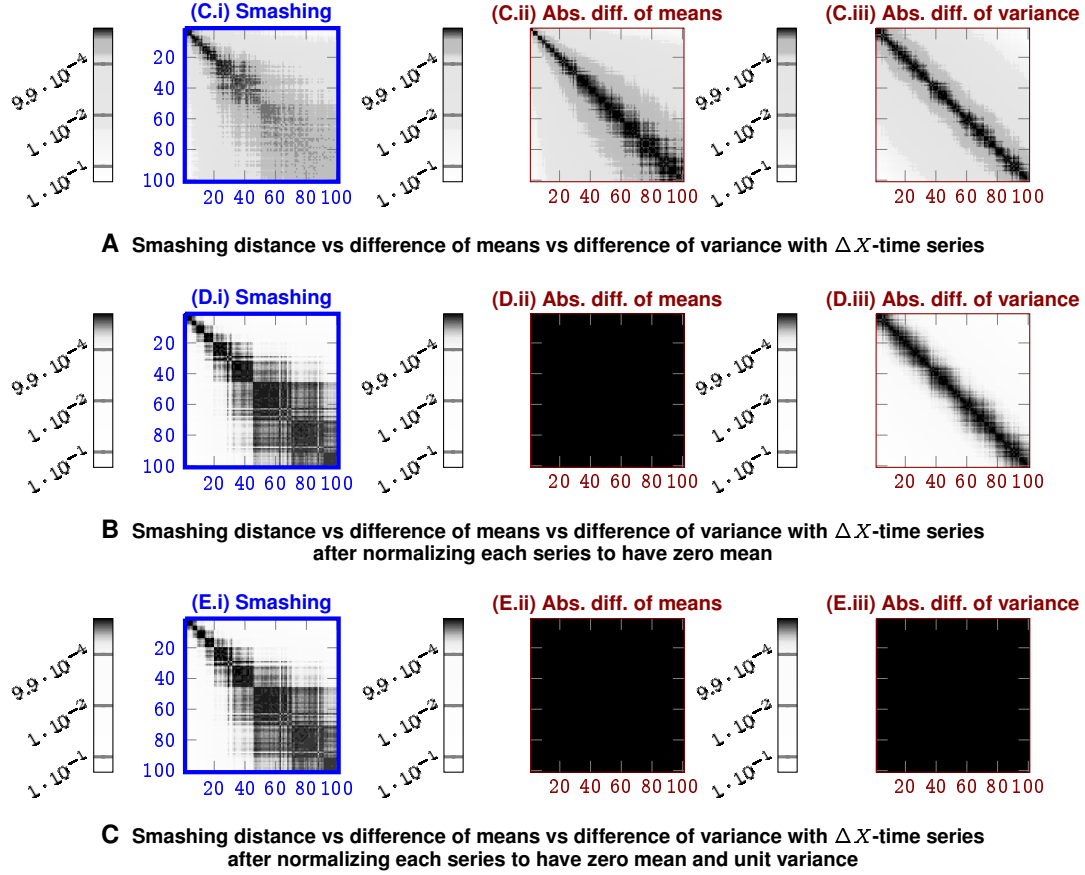
Fig. 14. We compute distances between the simulated time series for population numbers of the species $X$. First column is obtained via pairwise smashing. The second column is obtained as the absolute difference of means, and the third column is obtained as the absolute difference of variances. The second row illustrates the effect of normalizing the data to make each time series have zero mean. The third row illustrates the effect of normalization to make each series have zero mean and unit variance. Data smashing yields clear clusters, which are preserved through the normalization; whereas the simple statistical measures do not.

assume that we can make observations every $0.1s$ from the simulated dynamics. A few sample paths for the change in the number of prey with time are shown in Fig. 14.

Note that the probability of each reaction at any point in time is proportional to the number of combinatorial ways that particular reaction can transpire, as well as the propensity of the reaction itself. This combinatorial number is a function of the current population count of each species; and hence the reaction probabilities are strongly dependent on the current state vector. Since the simulation terminates when any one species becomes extinct, we cannot assume stationary behavior. Also, the initial state at least partially dictates the length of the time the ecosystem survives, implying non-ergodicity. Note that removing the restriction of a strictly positive and integer-valued population count, might result in a more well-behaved system.

Given our parameterization, we actually have 101 distinct systems with different sets of parameters, and for the system with index $i$, we have set the propensity of the third reaction as:

$$r = 0.2 + 0.001i \qquad (116)$$

Since the third reaction models predator death, we expect that increasing its propensity will make the predator degradation more probable. Thus, we can clearly expect smaller number of predators and larger number of prey on average, as $r$ is increased. However, a truly interesting structure would be uncovered if the behaviors exhibit some sort of clustering; as opposed to simply a monotonic dependence on $r$. We aanalyzed our set of $r$-parameterized dynamical systems as follows:

1) We concatenated all 1000 time series for species $X$ generated by simulating each system. Thus, the $i^{th}$ system generates the concatenated series $s_i$ for species $X$.
2) Next we generated $s_i^l$ from $s_i$ by taking one step differences

from $s_i$, i.e. $s_i^l$ is the time series of relative updates for the population of species $X$.
3) We mapped each sequential data series $s_i^l$, to a symbol stream using a binary partition function, which maps negative entries in the series to symbol 0, and positive entries to symbol 1.
4) We collided the symbolic streams pairwise, and compute the smashing distance matrix $H$. Thus, the $ij^{th}$ entry in $H$ is the deviation of the sum of $s_i^l$ and an inverted copy of $s_j^l$ from flat white noise. The result is shown in plate A(i) of Fig. 14.
5) We also generated the pairwise absolute differences of means, and variances. In each of these cases, the $ij^{th}$ entry of the corresponding matrix is the absolute difference between the corresponding statistical measure between the data series $s_i^l$ and $s_j^l$ (See plates A(ii-iii) in Fig. 14).

*Notably, the smashing matrix in plate A(i) of Fig. 14 shows clear clusters, whereas the matrices corresponding to mean and variance show trivial monotonic dependence on the parameter $r$.* To ascertain if the clustering obtained via data smashing is dependent on the mean or variance of the input data streams, we redid the analysis, after:

1) Normalization to zero mean signals prior to symbolization
2) Normalization to zero mean and unit variance signals prior to symbolization

In the first case, zeroing the mean makes the clusters appear more prominently (See plate B(i) of Fig. 14), while additionally normalizing the variance has little effect (See plate C(i) of Fig. 14). None of these changes allow the simple statistical measures to recover the clear clusters obtained via data smashing. The Lotka-Volterra system has a rich set of dynamical regimes, and it would not be surprising if such measures fail to capture this complexity. To that effect, we plotted the minimum number of predators after $100s$ of simulation

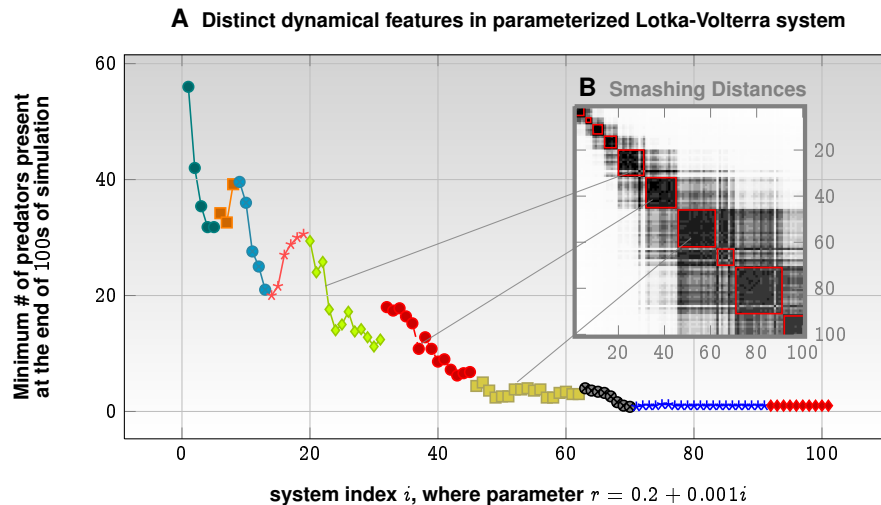## A  Distinct dynamical features in parameterized Lotka-Volterra system



**Fig. 15.** **Mapping clusters recovered via data smashing to a meaningful dynamical feature.** Plotting the minimum number of predators remaining in the system after a fixed simulation time of $100s$ (minimization carried out over the 1000 simulation runs of each system), illustrates that the clusters almost perfect correspond to monotonic domains of this function.

(minimum calculated over the 1000 simulation runs carried out for each parameter set, as discussed before). The result is shown in Fig. 15.

***The clusters identified via data smashing is now seen to correspond almost perfectly for each monotonic domain of this function.*** This illustrates that data smashing finds meaningful categorization, which simple statistical tools may miss. The differences discovered via smashing is obviously a function of the statistical structure of the observed data. However, the preceding example illustrates that it may not be easy to find the right statistical tool. Data smashing approach alleviates this challenge to a considerable degree.

### XII. Conclusion

We introduced data smashing to measure causal similarity between series of sequential observations. We demonstrated that our insight allows feature-less model-free classification in diverse applications, without the need for training, or expert tuned heuristics. Non-equal length of time-series, missing data, and possible phase mismatches are of no consequence.

While better classification algorithms may exist for specific problem domains, such algorithms are difficult to tune. The strength of data smashing lies in its ability to circumvent both the need for expert-defined heuristic features and expensive training; eliminating key bottlenecks in contemporary big data challenges.

### References

[1] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley, New York, 1991.

[2] R. P. W. Duin, D. d. Ridder, and D. M. J. Tax, "Experiments with a featureless approach to pattern recognition," *Pattern Recognition Letters*, pp. 1159–1166, 1997.

[3] V. Mottl, S. Dvoenko, O. Seredin, C. Kulikowski, and I. Muchnik, "Featureless pattern recognition in an imaginary hilbert space and its application to protein fold classification," *Machine Learning and Data Mining in Pattern Recognition*, pp. 322–336, 2001.

[4] E. Pekalska and R. P. W. Duin, "Dissimilarity representations allow for building good classifiers," *Pattern Recognition Letters*, vol. 23, no. 8, pp. 943–956, 2002.

[5] I. Chattopadhyay and H. Lipson, "Abductive learning of quantized stochastic processes using probabilistic finite automata," *Phil. Trans. of the Roy. Soc. A*, 2012, in press.

[6] J. P. Crutchfield and B. S. McNamara, "Equations of motion from a data series," *Complex systems*, vol. 1, no. 3, pp. 417–452, 1987.

[7] J. Crutchfield, "The calculi of emergence: computation, dynamics and induction," *Physica D: Nonlinear Phenomena*, vol. 75, no. 1, pp. 11–54, 1994.

[8] I. Chattopadhyay, Y. Wen, and A. Ray, "Pattern classification in symbolic streams via semantic annihilation of information," in *American Control Conference (ACC), 2010*, 30 2010-july 2 2010, pp. 492 –497.

[9] I. Chattopadhyay and A. Ray, "Structural transformations of probabilistic finite state machines," *International Journal of Control*, vol. 81, no. 5, pp. 820–835, 2008.

[10] R. G. Andrzejak, K. Lehnertz, F. Mormann, C. Rieke, P. David, and C. E. Elger, "Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: dependence on recording region and brain state," *Phys Rev E Stat Nonlin Soft Matter Phys*, vol. 64, no. 6 Pt 1, p. 061907, Dec 2001.

[11] P. Bentley, G. Nordehn, M. Coimbra, and S. Mannor, "The PASCAL Classifying Heart Sounds Challenge 2011 (CHSC2011) Results," http://www.peterjbentley.com/heartchallenge/index.html.

[12] M. K. Szymanski, "The optical gravitational lensing experiment. internet access to the ogle photometry data set: ogle-II bvi maps and i-band data," *Acta Astron.*, vol. 55, pp. 43–57, 2005.

[13] H. Begleiter, "EEG database data set," 1995, neurodynamics Laboratory, State University of New York Health Center Brooklyn, New York. [Online]. Available: http://archive.ics.uci.edu/ml/datasets/EEG+Database

[14] K. Brigham and B. Kumar, "Subject identification from electroencephalogram (eeg) signals during imagined speech," in *Biometrics: Theory Applications and Systems (BTAS), 2010 Fourth IEEE International Conference on*, 2010, pp. 1–8.

[15] "English language speech database for speaker recognition," Department of Informatics and mathematical modelling, Technical University of Denmark, 2005, department of Informatics and mathematical modelling, Technical University of Denmark. [Online]. Available: http://www2.imm.dtu.dk/$\sim$lfen/elsdsr/

[16] L. Feng and L. K. Hansen, "A new database for speaker recognition," Informatics and Mathematical Modelling, Technical University of Denmark, DTU, Richard Petersens Plads, Building 321, DK-2800 Kgs. Lyngby, Tech. Rep., 2005. [Online]. Available: http://www2.imm.dtu.dk/pubdb/p.php?3662

[17] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification (2nd Edition)*. Wiley-Interscience, 2000.

[18] G. Brumfiel, "High-energy physics: Down the petabyte highway," *Nature*, vol. 469, no. 7330, pp. 282–283, Jan 2011.

[19] R. G. Baraniuk, "More is less: Signal processing and the data deluge," *Science*, vol. 331, no. 6018, pp. 717–719, 2011.

[20] T. Cover and P. Hart, "Nearest neighbor pattern classification," *Information Theory, IEEE Transactions on*, vol. 13, no. 1, pp. 21 –27, january 1967.

[21] J. B. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.

[22] L. Yang, "An overview of distance metric learning," *Proc. Computer Vision and Pattern recognition, October*, vol. 7, 2007.

[23] J. Tenenbaum, V. de Silva, and J. Langford, "A global geometric framework for nonlinear dimensionality reduction." *Science (New York, N.Y.)*, vol. 290, no. 5500, pp. 2319–2323, 2000.

[24] S. Roweis and L. Saul, "Nonlinear dimensionality reduction by locally linear embedding." *Science (New York, N.Y.)*, vol. 290, no. 5500, pp. 2323–2326, 2000.

[25] H. Seung and D. Lee, "Cognition. the manifold ways of perception." *Science (New York, N.Y.)*, vol. 290, no. 5500, pp. 2268–2269, 2000.

[26] J. H. Ward, "Hierarchical grouping to optimize an objective function," *Journal of the American Statistical Association*, vol. 58, no. 301, pp. 236–244, 1963.

[27] M. J. Sippl and H. A. Scheraga, "Solution of the embedding problem and decomposition of symmetric matrices," *Proc. Natl. Acad. Sci. U.S.A.*, vol. 82, no. 8, pp. 2197–2201, Apr 1985.

[28] W. Feller, *The Fundamental Limit Theorems in Probability*. MacMillan, 1945.

[29] J. G. Snodgrass and M. V, "A standardized set of 260 pictures: Norms for name agreement, image agreement, familiarity, and visual complexity," *Journal of Experimental Psychology: Human Learning and Memory*, 1980.

[30] A. Paz, *Introduction to probabilistic automata (Computer science and applied mathematics)*. Orlando, FL, USA: Academic Press, Inc., 1971.

[31] 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.

[32] H. Hopcroft, R. Motwani, and J. Ullman, *Introduction to Automata Theory, Languages, and Computation, 2nd ed.* Addison Wesley, Boston, 2001.

[33] R. Gavaldà, P. W. Keller, J. Pineau, and D. Precup, "Pac-learning of markov models with hidden state," in *ECML*, 2006, pp. 150–161.

[34] 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.

[35] M. Ito and J. Duske, "On cofinal and definite automata," *Acta Cybern.*, vol. 6, pp. 181–189, 1984.

[36] F. Topsøe, "On the glivenko-cantelli theorem," *Probability Theory and Related Fields*, vol. 14, pp. 239–250, 1970, 10.1007/BF01111419. [Online]. Available: http://dx.doi.org/10.1007/BF01111419

[37] I. Chattopadhyay and A. Ray, "Language-measure-theoretic optimal control of probabilistic finite-state systems," *Int. J. Control*, vol. 80, no. 8, pp. 1271–1290, 2007.