Wisconsin, Madison, WI 53706

Computational Mechanics: Pattern and Prediction, Structure and Simplicity

Cosma Rohilla Shalizi* and James P. Crutchfield Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501 $Electronic\ addresses:\ \{shalizi, chaos\} @santafe.edu$ (February 1, 2008)

Computational mechanics, an approach to structural complexity, defines a process's causal states and gives a procedure for finding them. We show that the causal-state representation—an ϵ -machine—is the minimal one consistent with accurate prediction. We establish several results on ϵ -machine optimality and uniqueness and on how ϵ -machines compare to alternative representations. Further results relate measures of randomness and structural complexity obtained from ϵ -machines to those from ergodic and information theories.			Capturing a Pattern Defined E The Lessons of History Old Country Lemma	
			Computational Mechanics A Causal States	9 9 9 10
	Santa Fe Institute Working Paper 99-07-044 Keywords : complexity, computation, entropy, information, pattern, statistical mechanics. <i>Running Head</i> : Computational Mechanics		Conditional on a Causal State	10 10 10 10 10
02.50.Wp, 05.45, 05.65+b, 89.70.+c			B Causal State-to-State Transitions	11 11 11
	Contents		C ϵ -Machines	11 11 11
Ι	Introduction 2		ϵ -Machines Are Deterministic	11 12
II	Patterns 3 A Algebraic Patterns			13
	A Algebraic Patterns	V	Causal States Are Maximally Prescient Causal States Are Sufficient Statistics . Prescient Rivals Defined Refinement Lemma Causal States Are Minimal	13 13 14 14 15
III	Paddling around Occam's Pool6A Hidden Processes6Processes Defined6		Causal States Are Unique	15 15
	Stationarity	VI	Excess Entropy	16 16 16 17
	D Patterns in Ensembles	VII	A Discussion	18 18 18

APP	ENDIXES	20			
A	Information-Theoretic Formulæ 2				
В	The Equivalence Relation that Induces Causal States 2				
\mathbf{C}	Time Reversal 2				
D	ϵ -Machines are Monoids				
\mathbf{E}	Alternate Proof of the Refinement Lemma 2				
F	Finite Entropy for the Semi-Infinite Future 2 The Finite-Control Theorem				
G	Relations to Other Fields 1 Time Series Modeling				
Refer	rences	2 5			
Gloss	ary of Notation	29			

I. INTRODUCTION

Organized matter is ubiquitous in the natural world, but the branch of physics which ought to handle itstatistical mechanics—lacks a coherent, principled way of describing, quantifying, and detecting the many different kinds of structure nature exhibits. Statistical mechanics has good measures of disorder in thermodynamic entropy and in related quantities, such as the free energies. When augmented with theories of critical phenomena [1] and pattern formation [2], it also has an extremely successful approach to analyzing patterns formed through symmetry breaking, both in equilibrium [3] and, more recently, outside it [4]. Unfortunately, these successes involve many ad hoc procedures—such as guessing relevant order parameters, identifying small parameters for perturbation expansion, and choosing appropriate function bases for spatial decomposition. It is far from clear that the present methods can be extended to handle all the many kinds of organization encountered in nature, especially those produced by biological processes.

Computational mechanics [5] is an approach that lets us directly address the issues of pattern, structure, and organization. While keeping concepts and mathematical tools already familiar from statistical mechanics, it is distinct from the latter and complementary to it. In essence, from either empirical data or from a probabilistic description of behavior, it shows how to infer a model of the hidden process that generated the observed behavior. This representation—the ϵ -machine—captures the patterns and regularities in the observations in a way that reflects the causal structure of the process. Usefully, with this model in hand, one can extrapolate beyond the original observational data to make predictions of future behavior. Moreover, in a well defined sense that is the subject of the following, the ϵ -machine is the unique maximally efficient model of the observed datagenerating process.

 ϵ -Machines themselves reveal, in a very direct way, how information is stored in the process, and how that stored information is transformed by new inputs and by the passage of time. This, and not using computers for simulations and numerical calculations, is what makes computational mechanics "computational", in the sense of "computation theoretic".

The basic ideas of computational mechanics were introduced a decade ago [6]. Since then they have been used to analyze dynamical systems [7,8], cellular automata [9], hidden Markov models [10], evolved spatial computation [11], stochastic resonance [12], globally coupled maps [13], and the dripping faucet experiment [14]. Despite this record of successful application, there has been some uncertainty about the mathematical foundations of the subject. In particular, while it seemed evident from construction that an ϵ -machine captured the patterns inherent in a process and did so in a minimal way, no explicit proof of this was published. Moreover, there was no proof that, if the ϵ -machine was optimal in this way, it was the *unique* optimal representation of a process. These little-needed gaps have now been filled. Subject to some (reasonable) restrictions on the statistical character of a process, we prove that the ϵ -machine is indeed the unique optimal causal model. The rigorous proof of these results is the main burden of this paper. We gave preliminary versions of the optimality results—but not the uniqueness theorem, which is new here—in Ref. [15].

The outline of the exposition is as follows. We begin by showing how computational mechanics relates to other approaches to pattern, randomness, and causality. The upshot of this is to focus our attention on patterns within a statistical ensemble and their possible representations. Using ideas from information theory, we state a quantitative version of Occam's Razor for such representations. At that point we define causal states [6], equivalence classes of behaviors, and the structure of transitions between causal states—the ϵ -machine. We then show that the causal states are ideal from the point of view of Occam's Razor, being the simplest way of attaining the maximum possible predictive power. Moreover, we show

that the causal states are uniquely optimal. This combination allows us to prove a number of other, related optimality results about ϵ -machines. We examine the assumptions made in deriving these optimality results, and we note that several of them can be lifted without unduly upsetting the theorems. We also establish bounds on a process's intrinsic computation as revealed by ϵ -machines and by quantities in information and ergodic theories. Finally, we close by reviewing what has been shown and what seem like promising directions for further work on the mathematical foundations of computational mechanics.

A series of appendices provide supplemental material on information theory, equivalence relations and classes, ϵ -machines for time-reversed processes, semi-group theory, and connections and distinctions between computational mechanics and other fields.

To set the stage for the mathematics to follow and to motivate the assumptions used there, we begin now by reviewing prior work on pattern, randomness, and causality. We urge the reader interested only in the mathematical development to skip directly to Sec. II F—a synopsis of the central assumptions of computational mechanics—and continue from there.

II. PATTERNS

To introduce our approach to—and even to argue that *some* approach is necessary for—discovering and describing patterns in nature we begin by quoting Jorge Luis Borges:

These ambiguities, redundancies, and deficiencies recall those attributed by Dr. Franz Kuhn to a certain Chinese encyclopedia entitled Celestial Emporium of Benevolent *Knowledge*. On those remote pages it is written that animals are divided into (a) those that belong to the Emperor, (b) embalmed ones, (c) those that are trained, (d) suckling pigs, (e) mermaids, (f) fabulous ones, (g) stray dogs, (h) those that are included in this classification, (i) those that tremble as if they were mad, (j) innumerable ones, (k) those drawn with a very fine camel's hair brush, (1) others, (m) those that have just broken a flower vase, (n) those that resemble flies from a distance.

—J. L. Borges, "The Analytical Language of John Wilkins", in Ref. [16, p. 103]; see also discussion in Ref. [17].

The passage illustrates the profound gulf between patterns, and classifications derived from patterns, that are appropriate to the world and help us to understand it and those patterns which, while perhaps just as legitimate as prosaic regularities, are not at all informative.

What makes the *Celestial Emporium's* scheme inherently unsatisfactory, and not just strange, is that it tells us nothing about animals. We want to find patterns in a process that "divide it at the joints, as nature directs, not breaking any limbs in half as a bad carver might" [18, Sec. 265D].

Computational mechanics is not directly concerned with pattern formation per se [4]; though we suspect it will ultimately be useful in that domain. Nor is it concerned with pattern recognition as a practical matter as found in, say, neuropsychology [19], psychophysics [20], cognitive ethology [21], computer engineering [22], and signal and image processing [23,24]. Instead, it is concerned with the questions of what patterns are and how patterns should be represented. One way to highlight the difference is to call this pattern discovery, rather than pattern recognition.

The bulk of the intellectual discourse on what patterns are has been philosophical. One distinct subset has been conducted under the broad rubric of mathematical logic. Within this there are approaches, on the one hand, that draw on (highly) abstract algebra and the theory of relations; on the other, that approach patterns via the theory of algorithms and effective procedures.

The general idea, in both approaches, is that some object \mathcal{O} has a pattern \mathcal{P} — \mathcal{O} has a pattern "represented", "described", "captured", and so on by \mathcal{P} —if and only if we can use \mathcal{P} to predict or compress \mathcal{O} . Note that the ability to predict implies the ability to compress, but not vice versa; here we stick to prediction. The algebraic and algorithmic strands differ mainly on how \mathcal{P} itself should be represented; that is, they differ in how it is expressed in the vocabulary of some formal scheme.

We should emphasize here that "pattern" in this sense implies a kind of regularity, structure, symmetry, organization, and so on. In contrast, ordinary usage sometimes accepts, for example, speaking about the "pattern" of pixels in a particular slice of between-channels video "snow"; but we prefer to speak of that as the *configuration* of pixels.

A. Algebraic Patterns

Although the problem of pattern discovery appears early, in Plato's Meno [25] for example, perhaps the first attempt to make the notion of "pattern" mathematically rigorous was that of Whitehead and Russell in Principia Mathematica. They viewed pattern as a property, not of sets, but of relations within or between sets, and accordingly they work out an elaborate relation-arithmetic [26, vol. II, part IV]; cf. [27, ch. 5–6]. This starts by defining the relation-number of a relation between two sets as the class of all the relations that are equivalent to it under one-to-one, onto mappings of the two sets. In this framework relations share a common pattern or structure if they have the same relation-number. For

instance, all square lattices have similar structure since their elements share the same neighborhood relation; as do all hexagonal lattices. Hexagonal and square lattices, however, exhibit different patterns since they have non-isomorphic neighborhood relations—i.e., since they have different relation-numbers. (See also recoding equivalence defined in Ref. [28].) Less work has been done on this than they—especially Russell [29]—had hoped. This may be due in part to a general lack of familiarity with Volume II of Ref. [26].

A more recent attempt at developing an algebraic approach to patterns builds on semi-group theory and its Krohn-Rhodes decomposition theorem. Ref. [30] discusses a range of applications of this approach to patterns. Along these lines, Rhodes and Nehaniv have tried to apply semi-group complexity theory to biological evolution [31]. They suggest that the complexity of a biological structure can be measured by the number of subgroups in the decomposition of an automaton that describes the structure.

Yet another algebraic approach has been developed by Grenander and co-workers, primarily for pattern recognition [32]. Essentially, this is a matter of trying to invent a minimal set of *generators* and *bonds* for the pattern in question. Generators can adjoin each other, in a suitable n-dimensional space, only if their bonds are compatible. Each pair of compatible bonds at once specifies a binary algebraic operation and an observable element of the configuration built out of the generators. (Our construction in App. D, linking an algebraic operation with concatenations of strings, is analogous in a rough way.) Probabilities can be attached to these bonds, leading in a natural way to a (Gibbsian) probability distribution over entire configurations. Grenander and his colleagues have used these methods to characterize, inter alia, several biological phenomena [33,34].

B. Turing Mechanics: Patterns and Effective Procedures

The other path to patterns follows the traditional exploration of the logical foundations of mathematics, as articulated by Frege and Hilbert and pioneered by Church, Gödel, Post, Russell, Turing, and Whitehead. A more recent and relatively more popular approach goes back to Kolmogorov and Chaitin, who were interested in the exact reproduction of an individual object [35–38]; in particular, their focus was discrete symbol systems, rather than (say) real numbers or other mathematical objects. The candidates for expressing the pattern \mathcal{P} were universal Turing machine (UTM) programs—specifically, the shortest UTM program that can exactly produce the object \mathcal{O} . This program's length is called \mathcal{O} 's Kolmogorov-Chaitin complexity. Note that any scheme—automaton, grammar, or what-not—that is Turing equivalent and for which a notion of "length" is well defined will do as a representational scheme. Since we can convert from one such device to another—say, from a Post tag system [39] to a Turing machine—with only a finite description of the first system, such constants are easily assimilated when measuring complexity in this approach.

In particular, consider the first n symbols \mathcal{O}_n of \mathcal{O} and the shortest program \mathcal{P}_n that produces them. We ask, What happens to the limit

$$\lim_{n \to \infty} \frac{|\mathcal{P}_n|}{n} \,, \tag{1}$$

where $|\mathcal{P}|$ is the length in bits of program \mathcal{P} ? On the one hand, if there is a fixed-length program \mathcal{P} that generates arbitrarily many digits of \mathcal{O} , then this limit vanishes. Most of our interesting numbers, rational or irrationalsuch as π , e, $\sqrt{2}$ —are of this sort. These numbers are eminently compressible: the program \mathcal{P} is the compressed description, and so it captures the pattern obeyed by the sequence describing \mathcal{O} . If the limit goes to 1, on the other hand, we have a completely incompressible description and conclude, following Kolmogorov, Chaitin, and others, that \mathcal{O} is random [35–38,40,41]. This conclusion is the desired one: the Kolmogorov-Chaitin framework establishes, formally at least, the randomness of an individual object without appeals to probabilistic descriptions or to ensembles of reproducible events. And it does so by referring to a deterministic, algorithmic representation the UTM.

There are many well-known difficulties with applying Kolmogorov complexity to natural processes. First, as a quantity, it is uncomputable in general, owing to the halting problem [38]. Second, it is maximal for random sequences; this can be construed either as desirable, as just noted, or as a failure to capture structure, depending on one's aims. Third, it only applies to a single sequence; again this is either good or bad. Fourth, it makes no allowance for noise or error, demanding exact reproduction. Finally, $\lim_{n\to\infty} |\mathcal{P}_n|/n$ can vanish, although the computational resources needed to run the program, such as time and storage, grow without bound.

None of these impediments have kept researchers from attempting to use Kolmogorov-Chaitin complexity for practical tasks—such as measuring the complexity of natural objects (e.g. Ref. [42]), as a basis for theories of inductive inference [43,44], and generally as a means of capturing patterns [45]. As Rissanen [46, p. 49] says, this is akin to "learn[ing] the properties [of a data set] by writing programs in the hope of finding short ones!"

Various of the difficulties just listed have been addressed by subsequent work. Bennett's $logical\ depth$ accounts for time resources [47]. (In fact, it is the time for the minimal-length program \mathcal{P} to produce \mathcal{O} .) Koppel's sophistication attempts to separate out the "regularity" portion of the program from the random or instance-specific input data [48,49]. Ultimately, these extensions and generalizations remain in the UTM, exact-reproduction setting and so inherit inherent uncomputability.

C. Patterns with Error

Motivated by these theoretical difficulties and practical concerns, an obvious next step is to allow our pattern \mathcal{P} some degree of approximation or error, in exchange for shorter descriptions. As a result, we lose perfect reproduction of the original configuration from the pattern. Given the ubiquity of noise in nature, this is a small price to pay. We might also say that sometimes we are willing to accept small deviations from a regularity, without really caring what the precise deviation is. As pointed out in Ref. [17]'s conclusion, this is certainly a prime motivation in thermodynamic descriptions, in which we explicitly throw away, and have no interest in, vast amounts of microscopic detail in order to find a workable description of macroscopic observations.

Some interesting philosophical work on patterns-witherror has been done by Dennett, with reference not just to questions about the nature of patterns and their emergence but also to psychology [50]. The intuition is that truly random processes can be modeled very simply—"to model coin-tossing, toss a coin." Any prediction scheme that is more accurate than assuming complete independence ipso facto captures a pattern in the data. There is thus a spectrum of potential pattern-capturers ranging from the assumption of pure noise to the exact reproduction of the data, if that is possible. Dennett notes that there is generally a trade-off between the simplicity of a predictor and its accuracy, and he plausibly describes emergent phenomena [51,52] as patterns that allow for a large reduction in complexity for only a small reduction in accuracy. Of course, Dennett was by no means the first to consider predictive schemes that tolerate error and noise; we discuss some of the earlier work in App. G. However, to our knowledge, he was the first to have made such predictors a central part of an explicit account of what patterns are. It must be noted that this account lacks the mathematical detail of the other approaches we have considered so far, and that it relies on the inexact prediction of a single configuration. In fact, it relies on exact predictors that are "fuzzed up" by noise. The introduction of noise, however, brings in probabilities, and their natural setting is in ensembles. It is in that setting that the ideas we share with Dennett can receive a proper quantitative treatment.

D. Randomness: The Anti-Pattern?

We should at this point say a bit about the relations between randomness, complexity, and structure, at least as we use those words. Ignoring some foundational issues, randomness is actually rather well understood and well handled by classical tools introduced by Boltzmann [53]; Fisher, Neyman, and Pearson [54]; Kolmogorov [35]; and Shannon [55], among others. One tradition in the study of complexity in fact identifies complexity with random-

ness and, as we have just seen, this is useful for some purposes. As these purposes are *not* those of analyzing patterns in processes and in real-world data, however, they are not ours. Randomness simply does not correspond to a notion of pattern or structure at all and, by implication, neither Kolmogorov-Chaitin complexity nor any of its spawn measure pattern.

Nonetheless, some approaches to complexity conflate "structure" with the opposite of randomness, as conventionally understood and measured in physics by thermodynamic entropy or a related quantity, such as Shannon entropy. In effect, structure is defined as "one minus disorder". In contrast, we see pattern—structure, organization, regularity, and so on—as describing a coordinate "orthogonal" to a process's degree of randomness. That is, complexity (in our sense) and randomness each capture a useful property necessary to describe how a process manipulates information. This complementarity is even codified by the complexity-entropy diagrams introduced in Ref. [6]. It should be clear now that when we use the word "complexity" we mean "degrees" of pattern, not degrees of randomness.

E. Causation

We want our representations of patterns in dynamical processes to be causal—to say how one state of affairs leads to or produces another. Although a key property, causality enters our development only in an extremely weak sense, the weakest one can use mathematically, which is Hume's [56]: one class of event causes another if the latter always follows the former; the effect invariably succeeds the cause. As good indeterminists, in the following we replace this invariant-succession notion of causality with a more probabilistic one, substituting a homogeneous distribution of successors for the solitary invariable successor. (A precise statement appears in Sec. IV A's definition of causal states.) This approach results in a purely phenomenological statement of causality, and so it is amenable to experimentation in ways that stronger notions of causality—e.g., that of Ref. [57]—are not. Ref. [58] independently reaches a concept of causality essentially the same ours via philosophical arguments.

F. Synopsis of Pattern

In line with these observations, the ideal, synthesizing approach to patterns would be at once:

- 1. Algebraic, giving us an explicit breakdown or decomposition of the pattern into its parts;
- 2. Computational, showing how the process stores and uses information;
- 3. Calculable, analytically or by systematic approximation;

- 4. Causal, telling us how instances of the pattern are actually produced; and
- 5. Naturally stochastic, not merely tolerant of noise but explicitly formulated in terms of ensembles.

This mix is precisely the brew we claim, in all modesty, to have on tap.

III. PATTERNS IN ENSEMBLES: PADDLING AROUND OCCAM'S POOL

Here a pattern \mathcal{P} is something knowledge of which lets us predict, at better than chance rates, if possible, the future of sequences drawn from an ensemble \mathcal{O} : \mathcal{P} has to be statistically accurate and confer some leverage or advantage as well. Let's fix some notation and state the assumptions that will later let us prove the basic results.

A. Hidden Processes

We restrict ourselves to discrete-valued, discrete-time stationary stochastic processes. (See Sec. VIIB for discussion of these assumptions.) Intuitively, such processes are sequences of random variables S_i , the values of which are drawn from a countable set \mathcal{A} . We let i range over all the integers, and so get a bi-infinite sequence

$$\stackrel{\leftrightarrow}{S} = \dots S_{-1} S_0 S_1 \dots \tag{2}$$

In fact, we define a process in terms of the distribution of such sequences; cf. Ref. [59].

Definition 1 (A Process) Let \mathcal{A} be a countable set. Let $\Omega = \mathcal{A}^{\mathbb{Z}}$ be the set of bi-infinite sequences composed from \mathcal{A} , $T_i: \Omega \mapsto \mathcal{A}$ be the function that returns the i^{th} element s_i of a bi-infinite sequence $\omega \in \Omega$, and \mathcal{F} the field of cylinder sets of Ω . Adding a probability measure P gives us a probability space (Ω, \mathcal{F}, P) , with an associated random variable S. A process is a sequence of random variables $S_i = T_i(S), i \in \mathbb{Z}$.

Here, and throughout, we follow the convention of using capital letters to denote random variables and lower-case letters their particular values.

It follows from Def. 1 that there are well defined probability distributions for sequences of every finite length. Let \vec{S}_t^L be the sequence of $S_t, S_{t+1}, \ldots, S_{t+L-1}$ of L random variables beginning at S_t . $\vec{S}_t \equiv \lambda$, the null sequence. Likewise, $\overset{\leftarrow}{S}_t$ denotes the sequence of L random variables going up to S_t , but not including it; $\overset{\leftarrow}{S}_t^L = \vec{S}_{t-L}^L$. Both $\vec{S}_t^L = \overset{\leftarrow}{S}_{t-L}^L$ and $\overset{\leftarrow}{S}_t^L$ take values from $s^L \in \mathcal{A}^L$. Similarly, \vec{S}_t

and $\overset{\leftarrow}{S}_t$ are the semi-infinite sequences starting from and stopping at t and taking values $\overset{\leftarrow}{s}$ and $\overset{\leftarrow}{s}$, respectively.

Intuitively, we can imagine starting with distributions for finite-length sequences and extending them gradually in both directions, until the infinite sequence is reached as a limit. While this can be a useful picture to have in mind, defining a process in this way raises some subtle measure-theoretic issues, such as how finite-dimensional distributions limit on an infinite-dimensional one [60, ch. 7]. To avoid these we start with the infinite-dimensional distribution.

Definition 2 (Stationarity) A process S_i is stationary if and only if

$$P(\overrightarrow{S}_t^L = s^L) = P(\overrightarrow{S}_0 = s^L) , \qquad (3)$$

for all $t \in \mathbb{Z}$, $L \in \mathbb{Z}^+$, and all $s^L \in \mathcal{A}^L$.

In other words, a stationary process is one that is time-translation invariant. Consequently, $P(\overrightarrow{S}_t = \overrightarrow{s}) = P(\overrightarrow{S}_0 = \overrightarrow{s})$ and $P(\overrightarrow{S}_t = \overleftarrow{s}) = P(\overleftarrow{S}_0 = \overleftarrow{s})$, and so we drop the subscripts from now on.

B. The Pool

Our goal is to predict all or part of \overrightarrow{S} using some function of some part of \overleftarrow{S} . We begin by taking the set \overleftarrow{S} of all pasts and partitioning it into mutually exclusive and jointly comprehensive subsets. That is, we make a class \mathcal{R} of subsets of pasts. (See Fig. 1 for a schematic example.) Each $\rho \in \mathcal{R}$ will be called a *state* or an *effective state*. When the current history \overleftarrow{s} is included in the set ρ , we will speak of the process being in state ρ . Thus, we define a function from histories to effective states:

$$\eta : \stackrel{\leftarrow}{\mathbf{S}} \mapsto \mathcal{R} .$$
(4)

A specific individual history $s \in \mathbf{S}$ maps to a specific state $\rho \in \mathbf{R}$; the random variable s for the past maps to the random variable s for the effective states. It makes little difference whether we think of s as being a function from a history to a subset of histories or a function from a history to the *label* of that subset. Each interpretation is convenient at different times, and we will use both.

Note that we could use *any* function defined on $\overleftarrow{\mathbf{S}}$ to partition that set, by assigning to the same ρ all the histories \overleftarrow{s} on which the function takes the same value. Similarly, any equivalence relation on $\overleftarrow{\mathbf{S}}$ partitions it. (See

 $^{^{1}}$ At several points our constructions require referring to sets of sets. To help mark the distinction, we call the set of sets of histories a *class*.

App. B for more on equivalence relations.) Due to the way we defined a process's distribution, each effective state has a well defined distribution of futures, though not necessarily a unique one.² Specifying the effective state thus amounts to making a prediction about the process's future. All the histories belonging to a given effective state are treated as equivalent for purposes of predicting the future. (In this way, the framework formally incorporates traditional methods of time-series analysis; see App. G 1.)

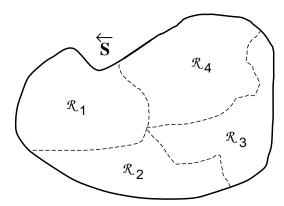


FIG. 1. A schematic picture of a partition of the set \mathbf{S} of all histories into some class of effective states: $\mathcal{R} = \{\mathcal{R}_i : i = 1, 2, 3, 4\}$. Note that the \mathcal{R}_i need not form compact sets; we simply draw them that way for clarity. One should have in mind Cantor sets or other more pathological structures.

We call the collection of all partitions \mathcal{R} of the set of histories $\stackrel{\leftarrow}{\mathbf{S}}$ Occam's pool.

C. A Little Information Theory

Since the bulk of the following development will be consumed with notions and results from information theory [55], we now review several highlights briefly, for the benefit of readers unfamiliar with the theory and to fix notation. Appendix A lists a number of useful information-theoretic formulæ, which get called upon in our proofs. Throughout, our notation and style of proof follow those in Ref. [62].

1. Entropy Defined

Given a random variable X taking values in a countable set \mathcal{A} , the entropy of X is

$$H[X] \equiv -\sum_{x \in A} P(X = x) \log_2 P(X = x) , \qquad (5)$$

taking $0\log 0=0$. Notice that H[X] is the expectation value of $-\log_2 P(X=x)$ and is measured in *bits* of information. Caveats of the form "when the sum converges to a finite value" are implicit in all statements about the entropies of infinite countable sets \mathcal{A} .

Shannon interpreted H[X] as the uncertainty in X. (Those leery of any subjective component in notions like "uncertainty" may read "effective variability" in its place.) He showed, for example, that H[X] is the mean number of yes-or-no questions needed to pick out the value of X on repeated trials, if the questions are chosen to minimize this average [55].

2. Joint and Conditional Entropies

We define the joint entropy H[X,Y] of two variables X (taking values in \mathcal{A}) and Y (taking values in \mathcal{B}) in the obvious way,

$$H[X,Y] \equiv$$

$$-\sum_{(x,y)\in\mathcal{A}\times\mathcal{B}} P(X=x, Y=y) \log_2 P(X=x, Y=y) .$$
(6)

We define the conditional entropy H[X|Y] of one random variable X with respect to another Y from their joint entropy:

$$H[X|Y] \equiv H[X,Y] - H[Y] . \tag{7}$$

This also follows naturally from the definition of conditional probability, since $P(X=x|Y=y) \equiv P(X=x,Y=y)/P(Y=y)$. H[X|Y] measures the mean uncertainty remaining in X once we know Y.

3. Mutual Information

The mutual information I[X;Y] between two variables is defined to be

$$I[X;Y] \equiv H[X] - H[X|Y] . \tag{8}$$

This is the average reduction in uncertainty about X produced by fixing Y. It is non-negative, like all entropies here, and symmetric in the two variables.

D. Patterns in Ensembles

It will be convenient to have a way of talking about the uncertainty of the future. Intuitively, this would just be $H[\vec{S}]$, but in general that quantity is infinite and awkward to manipulate. (The special case in which $H[\vec{S}]$

²This is not necessarily true if η is sufficiently pathological. To paraphrase Ref. [61], readers should assume that all our functions are sufficiently tame, measure-theoretically, that whatever induced distributions we invoke will exist.

is finite is dealt with in App. F.) Normally, we evade this by considering H[S], the uncertainty of the next L symbols, treated as a function of L. On occasion, we will refer to the entropy per symbol or entropy rate [55,62]:

$$h[\overrightarrow{S}] \equiv \lim_{L \to \infty} \frac{1}{L} H[\overrightarrow{S}^L] , \qquad (9)$$

and the *conditional entropy rate*,

$$h[\overrightarrow{S}|X] \equiv \lim_{L \to \infty} \frac{1}{L} H[\overrightarrow{S}^L|X] , \qquad (10)$$

where X is some random variable and the limits exist. For stationary stochastic processes, the limits always exist [62, Theorem 4.2.1, p. 64].

These entropy rates are also always bounded above by H[S]; which is a special case of Eq. (A3). Moreover, if $h[\overrightarrow{S}] = H[S]$, the process consists of independent variables—independent, identically distributed (IID) variables, in fact, since we are only concerned with stationary processes here.

Definition 3 (Capturing a Pattern) \mathcal{R} captures a pattern if and only if there exists an L such that

$$H[\stackrel{\rightarrow}{S}^{L}|\mathcal{R}] < LH[S] \ . \tag{11}$$

This says that \mathcal{R} captures a pattern when it tells us something about how the distinguishable parts of a process affect each other: \mathcal{R} exhibits their dependence. (We also speak of η , the function associated with pasts, as capturing a pattern, since this is implied by \mathcal{R} capturing a pattern.) Supposing that these parts do not affect each other, then we have IID random variables, which is as close to the intuitive notion of "patternless" as one is likely to state mathematically. Note that, because of the independence bound on joint entropies (Eq. (A3)), if the inequality is satisfied for some L, it is also satisfied for every L' > L. Thus, we can consider the difference $H[S] - H[\vec{S}|\mathcal{R}]/L$, for the smallest L for which it is nonzero, as the strength of the pattern captured by \mathcal{R} . We will now mark an upper bound (Lemma 1) on the strength of patterns; later we will show how to attain this upper bound (Thm. 1).

E. The Lessons of History

We are now in a position to prove a result about patterns in ensembles that will be useful in connection with our later theorems about causal states.

Lemma 1 (Old Country Lemma) For all \mathcal{R} and for all $L \in \mathbb{Z}^+$,

$$H[\vec{S}^{L}|\mathcal{R}] \ge H[\vec{S}^{L}|\overset{\leftarrow}{S}]$$
 (12)

Proof. By construction (Eq. (4)), for all L,

$$H[\stackrel{\rightarrow}{S}^{L}|\mathcal{R}] = H[\stackrel{\rightarrow}{S}^{L}|\eta(\stackrel{\leftarrow}{S})]$$
 (13)

But

$$H[\vec{S}^{L}|\eta(\overset{\leftarrow}{S})] \ge H[\vec{S}^{L}|\overset{\leftarrow}{S}],$$
 (14)

since the entropy conditioned on a variable is never more than the entropy conditioned on a function of the variable (Eq. (A14)). QED.

Remark 1. That is, conditioning on the whole of the past reduces the uncertainty in the future to as small a value as possible. Carrying around the whole semi-infinite past is rather bulky and uncomfortable and is a somewhat dismaying prospect. Put a bit differently: we want to forget as much of the past as possible and so reduce its burden. It is the contrast between this desire and the result of Eq. (12) that leads us to call this the Old Country Lemma.

Remark 2. Lemma 1 establishes the promised upper bound on the strength of patterns: viz., the strength of the pattern is at most $H[S] - H[\vec{S} \mid \overset{\leftarrow}{S}]/L_{past}$, where L_{past} is the least value of L such that $H[\vec{S} \mid \overset{\leftarrow}{S}] < LH[S]$.

F. Minimality and Prediction

Let's invoke Occam's Razor: "It is vain to do with more what can be done with less" [63]. To use the razor, we need to fix what is to be "done" and what "more" and "less" mean. The job we want done is accurate prediction, i.e., reducing the conditional entropies $H[S] |\mathcal{R}|$ as far as possible, the goal being to attain the bound set by Lemma 1. But we want to do this as simply as possible, with as few resources as possible. On the road to meeting these two constraints—minimal uncertainty and minimal resources—we will need a measure of the second. Since P(S = s) is well defined, there is an induced measure on the η -states; i.e., $P(\mathcal{R} = \rho)$, the probability of being in any particular effective state, is well defined. Accordingly, we define the following measure of resources.

Definition 4 (Complexity of State Classes) The statistical complexity of a class \mathcal{R} of states is

$$C_{\mu}(\mathcal{R}) \equiv H[\mathcal{R}]$$

$$= -\sum_{\rho \in \mathcal{R}} P(\mathcal{R} = \rho) \log_2 P(\mathcal{R} = \rho) ,$$
(15)

when the sum converges to a finite value.

The μ in C_{μ} reminds us that it is a measure-theoretic property and depends ultimately on the distribution over the process's sequences, which induces a measure over states.

The statistical complexity of a state class is the average uncertainty (in bits) in the process's current state. This, in turn, is the same as the average amount of memory (in bits) that the process appears to retain about the past, given the chosen state class \mathcal{R} . (We will later, in Def. 12, see how to define the statistical complexity of a process itself.) The goal is to do with as little of this memory as possible. Restated then, we want to minimize statistical complexity, subject to the constraint of maximally accurate prediction.

The idea behind calling the collection of all partitions of \mathbf{S} Occam's pool should now be clear: One wants to find the shallowest point in the pool. This we now do.

IV. COMPUTATIONAL MECHANICS

Those who are good at archery learnt from the bow and not from Yi the Archer. Those who know how to manage boats learnt from the boats and not from Wo.

—Anonymous in Ref. [64].

The ultimate goal of computational mechanics is to discern the patterns intrinsic to a process. That is, as much as possible, the goal is to let the process describe itself, on its own terms, without appealing to a priori assumptions about the process's structure. Here we simply explore the consistency and well-definedness of these goals. Of course, practical constraints may keep us from doing more than approximating these ideals more or less grossly. Naturally, such problems, which always turn up in implementation, are much easier to address if we start from secure foundations.

A. Causal States

Definition 5 (A Process's Causal States) The causal states of a process are the members of the range of the function $\epsilon : \stackrel{\leftarrow}{\mathbf{S}} \mapsto 2^{\stackrel{\leftarrow}{\mathbf{S}}}$ —the power set of $\stackrel{\leftarrow}{\mathbf{S}}$:

$$\epsilon(\stackrel{\leftarrow}{s}) \equiv \{\stackrel{\leftarrow}{s}' | P(\stackrel{\rightarrow}{S} = \stackrel{\rightarrow}{s} | \stackrel{\leftarrow}{S} = \stackrel{\leftarrow}{s}) = P(\stackrel{\rightarrow}{S} = \stackrel{\rightarrow}{s} | \stackrel{\leftarrow}{S} = \stackrel{\leftarrow}{s}'),
\text{for all } \stackrel{\rightarrow}{s} \in \stackrel{\rightarrow}{S}, \stackrel{\leftarrow}{s}' \in \stackrel{\leftarrow}{S}\},$$
(16)

that maps from histories to sets of histories. We write the i^{th} causal state as S_i and the set of all causal states as S; the corresponding random variable is denoted S, and its realization σ .

The cardinality of S is unspecified. S can be finite, countably infinite, a continuum, a Cantor set, or something stranger still. Examples of these are given in Refs. [5] and [10]; see especially the examples for hidden Markov models given there.

Alternately and equivalently, we could define an equivalence relation \sim_{ϵ} such that two histories are equivalent if and only if they have the same conditional distribution of futures, and then define causal states as the equivalence classes generated by \sim_{ϵ} . (In fact, this was the original approach [6].) Either way, the divisions of this partition of \mathbf{S} are made between regions that leave us in different conditions of ignorance about the future.

This last statement suggests another, still equivalent, description of ϵ :

$$\epsilon(\stackrel{\leftarrow}{s}) = \{\stackrel{\leftarrow}{s}' | P(\stackrel{\rightarrow}{S}^{L} = \stackrel{\rightarrow}{s}^{L} | \stackrel{\leftarrow}{S} = \stackrel{\leftarrow}{s}) = P(\stackrel{\rightarrow}{S}^{L} = \stackrel{\rightarrow}{s}^{L} | \stackrel{\leftarrow}{S} = \stackrel{\leftarrow}{s}'), \\
\stackrel{\rightarrow}{s}^{L} \in \stackrel{\rightarrow}{S}^{L}, \quad \stackrel{\leftarrow}{s}' \in \stackrel{\leftarrow}{S}, L \in \mathbb{Z}^{+}\}. \tag{17}$$

Using this we can make the original definition, Eq. (16), more intuitive by picturing a sequence of partitions of the space \mathbf{S} of all histories in which each new partition, induced using L+1, is a refinement of the previous one induced using L. At the coarsest level, the first partition (L=1) groups together those histories that have the same distribution for the very next observable. These classes are then subdivided using the distribution of the next two observables, then the next three, four, and so on. The limit of this sequence of partitions—the point at which every member of each class has the same distribution of futures, of whatever length, as every other member of that class—is the partition of \mathbf{S} induced by \sim_{ϵ} . See App. B for a detailed discussion and review of the equivalence relation \sim_{ϵ} .

Although they will not be of direct concern in the following, due to the time-asymptotic limits taken, there are transient causal states in addition to those (recurrent) causal states defined above in Eq. (16). Roughly speaking, the transient causal states describe how a lengthening sequence (a history) of observations allows us to identify the recurrent causal states with increasing precision. See the developments in App. B and in Refs. [10] and [65] for more detail on transient causal states.

Causal states are a particular kind of effective state, and they have all the properties common to effective states (Sec. IIIB). In particular, each causal state S_i has several structures attached:

- 1. The index i—the state's "name".
- 2. The set of histories that have brought the process to S_i , which we denote $\{\stackrel{\leftarrow}{s} \in S_i\}$.
- 3. A conditional distribution over futures, denoted $P(\vec{S} \mid \mathcal{S}_i)$, and equal to $P(\vec{S} \mid \overleftarrow{s})$, $\overleftarrow{s} \in \mathcal{S}_i$. Since we refer to this type of distribution frequently and since it is the "shape of the future", we call it the state's *morph*.

Ideally, each of these should be denoted by a different symbol, and there should be distinct functions linking each of these structures to their causal state. To keep the growth of notation under control, however, we shall be strategically vague about these distinctions. Readers may variously picture ϵ as mapping histories to (i) simple indices, (ii) subsets of histories, or (iii) ordered triples of indices, subsets, and morphs; or one may even leave ϵ uninterpreted, as preferred, without interfering with the development that follows.

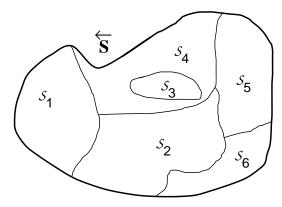


FIG. 2. A schematic representation of the partitioning of the set S of all histories into causal states $S_i \in S$. Within each causal state all the individual histories S have the same morph—the same conditional distribution $P(S \mid S)$ for future observables.

1. Morphs

Each causal state has a unique morph, i.e., no two causal states have the same conditional distribution of futures. This follows directly from Def. 5, and it is not true of effective states in general. Another immediate consequence of that definition is that

$$P(\overrightarrow{S} = \overrightarrow{s} \mid \mathcal{S} = \epsilon(\overleftarrow{s})) = P(\overrightarrow{S} = \overrightarrow{s} \mid \overleftarrow{S} = \overleftarrow{s}). \tag{18}$$

(Again, this is not generally true of effective states.) This observation lets us prove a useful lemma about the conditional independence of the past $\stackrel{\leftarrow}{S}$ and the future $\stackrel{\rightarrow}{S}$.

Lemma 2 The past and the future are independent, conditioning on the causal states.

Proof. Recall that two random variables X and Z are conditionally independent if and only if there is a third variable Y such that

$$P(X = x, Y = y, Z = z)$$

= $P(X = x|Y = y)P(Z = z|Y = y)P(Y = y)$. (19)

That is, all of the dependence of Z on X is mediated by Y. For convenience below we note that, re-factoring the conditional probabilities, this is equivalent to the requirement that:

$$P(X = x, Y = y, Z = z)$$

= $P(Z = z|Y = y)P(Y = y|X = x)P(X = x)$. (20)

Let us consider $P(\stackrel{\leftarrow}{S} = \stackrel{\leftarrow}{s}, S = \sigma, \stackrel{\rightarrow}{S} = \stackrel{\rightarrow}{s})$.

P (
$$\overleftarrow{S} = \overleftarrow{s}, \mathcal{S} = \sigma, \overrightarrow{S} = \overrightarrow{s}$$
)
= $P(\overrightarrow{S} = \overrightarrow{s} | \mathcal{S} = \sigma, \overleftarrow{S} = \overleftarrow{s}) P(\mathcal{S} = \sigma, \overleftarrow{S} = \overleftarrow{s})$ (21)
= $P(\overrightarrow{S} = \overrightarrow{s} | \mathcal{S} = \sigma, \overleftarrow{S} = \overleftarrow{s}) P(\mathcal{S} = \sigma | \overleftarrow{S} = \overleftarrow{s}) P(\overleftarrow{S} = \overleftarrow{s})$.

Now, $P(S = \sigma | \overleftarrow{S} = \overleftarrow{s}) = 0$, unless $\sigma = \epsilon(\overleftarrow{s})$, which case $P(S = \sigma | \overleftarrow{S} = \overleftarrow{s}) = 1$. Either way, the first two factors in the last line of Eq. (21) can be written, by Eq. (18),

P (
$$\overrightarrow{S} = \overrightarrow{s} \mid \mathcal{S} = \sigma$$
, $\overleftarrow{S} = \overleftarrow{s}$)P($\mathcal{S} = \sigma \mid \overleftarrow{S} = \overleftarrow{s}$)
$$= P(\overrightarrow{S} = \overrightarrow{s} \mid \mathcal{S} = \sigma)P(\mathcal{S} = \sigma \mid \overleftarrow{S} = \overleftarrow{s}), \qquad (22)$$

so that, substituting Eq. (22) into Eq. (21),

P (
$$\overleftarrow{S} = \overleftarrow{s}, \mathcal{S} = \sigma, \overrightarrow{S} = \overrightarrow{s}$$
)
= $P(\overrightarrow{S} = \overrightarrow{s} | \mathcal{S} = \sigma)P(\mathcal{S} = \sigma | \overleftarrow{S} = \overleftarrow{s})P(\overleftarrow{S} = \overleftarrow{s})$. (23)

QED.

2. Homogeneity

Following Ref. [58], we introduce two new definitions and a lemma which are required later on, especially in the proof of Lemma 7 and the theorems depending on that lemma.

Definition 6 (Strict Homogeneity) A set X is strictly homogeneous with respect to a certain random variable Y when the conditional distribution P(Y|X) for Y is the same for all subsets of X.

Definition 7 (Weak Homogeneity) A set X is weakly homogeneous with respect to Y if X is not strictly homogeneous with respect to Y, but $X \setminus X_0$ (X with X_0 removed) is, where X_0 is a subset of X of measure X_0 .

Lemma 3 (Strict Homogeneity of Causal States) A process's causal states are the largest subsets of histories that are all strictly homogeneous with respect to futures of all lengths.

Proof. We must show that, first, the causal states are strictly homogeneous with respect to futures of all lengths and, second, that no larger strictly homogeneous subsets of histories could be made. The first point, the strict homogeneity of the causal states, is evident from Eq. (17): By construction, all elements of a causal state have the same morph, so any part of a causal state will have the same morph as the whole state. The second point likewise follows from Eq. (17), since the causal state by construction contains all the histories with a given morph.

Any other set strictly homogeneous with respect to futures must be smaller than a causal state, and any set that includes a causal state as a proper subset cannot be *strictly* homogeneous. QED.

Remark. The statistical explanation literature would say that causal states are the "statistical-relevance basis for causal explanations". The elements of such a basis are, precisely, the largest classes of combinations of independent variables with homogeneous distributions for the dependent variables. See Ref. [58] for further discussion along these lines.

B. Causal State-to-State Transitions

The causal state at any given time and the next value of the observed process together determine a new causal state; this is proved shortly in Lemma 5. Thus, there is a natural relation of succession among the causal states; recall the discussion of causality in Sec. II E. Moreover, given the current causal state, all the possible next values have well defined conditional probabilities. In fact, by construction the entire semi-infinite future does. Thus, there is a well defined probability $T_{ij}^{(s)}$ of the process generating the value $s \in \mathcal{A}$ and going to causal state \mathcal{S}_j , if it is in state \mathcal{S}_i .

Definition 8 (Causal Transitions) The labeled transition probability $T_{ij}^{(s)}$ is the probability of making the transition from state S_i to state S_j while emitting the symbol $s \in A$:

$$T_{ij}^{(s)} \equiv P(\mathcal{S}' = \mathcal{S}_j, \stackrel{\rightarrow}{S}^1 = s | \mathcal{S} = \mathcal{S}_i) ,$$
 (24)

where \mathcal{S} is the current causal state and \mathcal{S}' its successor on emitting s. We denote the set $\{T_{ij}^{(s)}:s\in\mathcal{A}\}$ by \mathbf{T} .

Lemma 4 (Transition Probabilities) $T_{ij}^{(s)}$ is given by

$$T_{ij}^{(s)} = P(\overleftarrow{s}s \in \mathcal{S}_j | \overleftarrow{s} \in \mathcal{S}_i)$$
 (25)

$$= \frac{P(\overleftarrow{s} \in \mathcal{S}_i, \overleftarrow{s} s \in \mathcal{S}_j)}{P(\overleftarrow{s} \in \mathcal{S}_i)}, \qquad (26)$$

where $\overleftarrow{s}s$ is read as the semi-infinite sequence obtained by concatenating $s \in \mathcal{A}$ onto the end of \overleftarrow{s} .

Proof.

$$T_{ij}^{(s)} = P(\mathcal{S}' = \mathcal{S}_j, \stackrel{\rightarrow}{S}^1 = s | \mathcal{S} = \mathcal{S}_i)$$
 (27)

$$= \frac{P(S' = S_j, \overrightarrow{S} = s, S = S_i)}{P(S = S_i)}.$$
 (28)

Now $S = S_i$ if and only if $\overleftarrow{s} \in S_i$, and $S' = S_j$ if and only $\overleftarrow{s}' \in S_j$, where by \overleftarrow{s}' we mean the history that is the immediate successor to \overleftarrow{s} ; for consistency, $\overleftarrow{s}' = \overleftarrow{s}s$. So we can rewrite Eq. (28) as

$$T_{ij}^{(s)} = \frac{P(\stackrel{\leftarrow}{s} \in \mathcal{S}_i, \stackrel{\rightarrow}{S}^1 = s, \stackrel{\leftarrow}{s}' \in \mathcal{S}_j)}{P(\mathcal{S} = \mathcal{S}_i)}$$
(29)

$$= \frac{P(\stackrel{\leftarrow}{s} \in \mathcal{S}_i, \stackrel{\rightarrow}{S}^1 = s, \stackrel{\leftarrow}{s} s \in \mathcal{S}_j)}{P(\mathcal{S} = \mathcal{S}_i)}$$
(30)

$$= \frac{P(\overleftarrow{s} \in \mathcal{S}_i, \overleftarrow{s}s \in \mathcal{S}_j)}{P(\mathcal{S} = \mathcal{S}_i)}.$$
 (31)

In the third line we used the fact that $\overset{\leftarrow}{S} = \overset{\leftarrow}{s}$ and $\overset{\leftarrow}{S} = \overset{\leftarrow}{s}s$ jointly imply $\vec{S} = s$, making that condition redundant. QED.

Notice that $T_{ij}^{(\lambda)} = \delta_{ij}$; that is, the transition labeled by the null symbol λ is the identity.

C. ϵ -Machines

The combination of the function ϵ from histories to causal states with the labeled transition probabilities $T_{ij}^{(s)}$ is called the ϵ -machine of the process [5,6].

Definition 9 (An ϵ -Machine Defined)

The ϵ -machine of a process is the ordered pair $\{\epsilon, \mathbf{T}\}$, where ϵ is the causal state function and \mathbf{T} is set of the transition matrices for the states defined by ϵ .

Equivalently, we may denote an ϵ -machine by $\{S, T\}$. To satisfy the algebraic requirement outlined in Sec. IIF, we make explicit the connection with semi-group theory.

Proposition 1 (ϵ -Machines Are Monoids) The algebra generated by the ϵ -machine $\{\epsilon, \mathbf{T}\}$ is a semi-group with an identity element, i.e., it is a monoid.

Proof. See App. D.

Remark. Due to this, ϵ -machines can be interpreted as capturing a process's generalized symmetries. Any subgroups of an ϵ -machine's semi-group are, in fact, symmetries in the more familiar sense.

Lemma 5 (ϵ -Machines Are Deterministic) For each S_i and $s \in A$, $T_{ij}^{(s)} > 0$ only for that S_j for which $\epsilon(s) = S_j$ if and only if $\epsilon(s) = S_i$, for all pasts s.

Proof. The lemma is equivalent to asserting that for all $s \in \mathcal{A}$ and $\overleftarrow{s}, \overleftarrow{s}' \in \mathbf{S}$, if $\epsilon(\overleftarrow{s}) = \epsilon(\overleftarrow{s}')$, then $\epsilon(\overleftarrow{s}s) = \epsilon(\overleftarrow{s}s)$. ($\overleftarrow{s}s$ is just another history and belongs to one or another causal state.)

Suppose this were not true. Then there would have to exist at least one future \overrightarrow{s} such that

$$P(\overrightarrow{S} = \overrightarrow{s} \mid \overleftarrow{S} = \overleftarrow{s}s) \neq P(\overrightarrow{S} = \overrightarrow{s} \mid \overleftarrow{S} = \overleftarrow{s}s),$$
 (32)

when nonetheless $\epsilon(s) = \epsilon(s')$. Equivalently, we would have

$$\frac{P(\stackrel{\hookrightarrow}{S} = \stackrel{\longleftarrow}{ss} \stackrel{\rightarrow}{s})}{P(\stackrel{\longleftarrow}{S} = \stackrel{\longleftarrow}{ss})} \neq \frac{P(\stackrel{\hookrightarrow}{S} = \stackrel{\longleftarrow}{s} \stackrel{\rightarrow}{s} \stackrel{\rightarrow}{s})}{P(\stackrel{\longleftarrow}{S} = \stackrel{\longleftarrow}{s} \stackrel{\rightarrow}{s})}, \tag{33}$$

where we read s $\stackrel{\rightarrow}{s}$ as the semi-infinite string that begins s and continues $\stackrel{\rightarrow}{s}$. (Remember, the point at which we break the stochastic process into a past and a future is arbitrary.) However, the probabilities in the denominators are equal to $P(\stackrel{\rightarrow}{S}=s|\stackrel{\leftarrow}{S}=\stackrel{\leftarrow}{s})P(\stackrel{\leftarrow}{S}=\stackrel{\leftarrow}{s})$ and $P(\stackrel{\rightarrow}{S}=s|\stackrel{\leftarrow}{S}=\stackrel{\leftarrow}{s}')P(\stackrel{\leftarrow}{S}=\stackrel{\leftarrow}{s}')$, respectively, and by assumption $P(\stackrel{\rightarrow}{S}=s|\stackrel{\rightarrow}{S}=\stackrel{\leftarrow}{s}')=P(\stackrel{\rightarrow}{S}=s|\stackrel{\leftarrow}{S}=\stackrel{\leftarrow}{s})$, since $\epsilon(\stackrel{\leftarrow}{s}')=\epsilon(\stackrel{\leftarrow}{s})$. Therefore, we would need

$$\frac{P(\stackrel{\hookrightarrow}{S} = \stackrel{\leftarrow}{s} s \stackrel{\rightarrow}{s})}{P(\stackrel{\hookrightarrow}{S} = \stackrel{\leftarrow}{s})} \neq \frac{P(\stackrel{\hookrightarrow}{S} = \stackrel{\leftarrow}{s} \stackrel{\rightarrow}{s} \stackrel{\rightarrow}{s})}{P(\stackrel{\hookrightarrow}{S} = \stackrel{\leftarrow}{s} \stackrel{\rightarrow}{s})}.$$
 (34)

This is the same, though, as

$$P(\overrightarrow{S} = s \overrightarrow{s} \mid \overleftarrow{S} = \overleftarrow{s}) \neq P(\overrightarrow{S} = s \overrightarrow{s} \mid \overleftarrow{S} = \overleftarrow{s}')$$
. (35)

This is to say that there is a future $s\overrightarrow{s}$ that has different probabilities depending on whether we conditioned on \overleftarrow{s} or on \overleftarrow{s}' . But this contradicts the assumption that the two histories belong to the same causal state. Therefore, there is no such future \overrightarrow{s} , and the alternative statement of the lemma is true. QED.

Remark 1. In automata theory [66], a set of states and transitions is said to be deterministic if the current state and the next input—here, the next result from the original stochastic process—together fix the next state. This use of the word "deterministic" is often confusing, since many stochastic processes (e.g., simple Markov chains) are deterministic in this sense.

Remark 2. Starting from a fixed state, a given symbol always leads to at most one single state. But there can be several transitions from one state to another, each labeled with a different symbol.

Remark 3. Clearly, if $T_{ij}^{(s)} > 0$, then $T_{ij}^{(s)} = P(\vec{S} = s | \mathcal{S} = \mathcal{S}_i)$. In automata theory the "disallowed" transitions $(T_{ij}^{(s)} = 0)$ are sometimes explicitly represented and lead to a "reject" state indicating that the particular history does not occur.

Lemma 6 (Causal States Are Independent) The probability distributions over causal states at different times are conditionally independent.

Proof. What we wish to show is that, writing \mathcal{S} , \mathcal{S}' , \mathcal{S}'' for the sequence of causal states at three successive times, \mathcal{S} and \mathcal{S}'' are conditionally independent, given \mathcal{S}' . We can do this directly:

P (
$$S = \sigma, S' = \sigma', S'' = \sigma''$$
)
= $P(S'' = \sigma'' | S = \sigma, S' = \sigma') P(S = \sigma, S' = \sigma')$
= $P(\overrightarrow{S} \in a | S = \sigma, S' = \sigma') P(S = \sigma, S' = \sigma')$, (36)

where a is the subset of all symbols that lead from σ' to σ'' . This is a well defined subset, in virtue of Lemma 5 immediately preceding, which also guarantees the equality of conditional probabilities we have used. Likewise,

$$P(\mathcal{S}'' = \sigma'' | \mathcal{S}' = \sigma') = P(\overrightarrow{S} \in a | \mathcal{S}' = \sigma'). \tag{37}$$

But, by construction,

$$P(\overrightarrow{S} \in a | \mathcal{S} = \sigma, \mathcal{S}' = \sigma') = P(\overrightarrow{S} \in a | \mathcal{S}' = \sigma'), \quad (38)$$

and hence

$$P(S'' = \sigma'' | S' = \sigma') = P(S'' = \sigma'' | S = \sigma, S' = \sigma').$$
(39)

So, to resume,

P (
$$S = \sigma$$
, $S' = \sigma'$, $S'' = \sigma''$)
= P($S'' = \sigma'' | S' = \sigma'$)P($S = \sigma$, $S' = \sigma'$)
= P($S'' = \sigma'' | S' = \sigma'$)P($S' = \sigma' | S = \sigma$)P($S = \sigma$). (40)

The last line follows from the definition of conditional probability and is equivalent to the more easily interpreted expression given by

$$P(S''|S')P(S|S')P(S')$$
. (41)

Thus, applying mathematical induction to Eq. (41), causal states at different times are independent, conditioning on the intermediate causal states. QED.

Remark 1. This lemma strengthens the claim that the causal states are, in fact, the causally efficacious states: given knowledge of the present state, what has gone before makes no difference. (Again, recall the philosophical preliminaries of Sec. II E.)

Remark 2. This result indicates that the causal states, considered as a process, define a kind of Markov chain. Thus, causal states can be roughly considered to be a generalization of Markovian states. We say "kind of" since the class of ϵ -machines is substantially richer [5,10] than what one normally associates with Markov chains [67,68].

Definition 10 (ϵ -Machine Reconstruction)

 ϵ -Machine reconstruction is any procedure that given a process P(S), or an approximation of P(S), produces the process's ϵ -machine $\{S, T\}$.

Given a mathematical description of a process, one can often calculate analytically its ϵ -machine. (For example, see the computational mechanics analysis of spin systems in Ref. [65].) There is also a wide range of algorithms which reconstruct ϵ -machines from empirical estimates of P(S). Some, such as those used in Refs. [5–7,69], operate in "batch" mode, taking the raw data as a whole and producing the ϵ -machine. Others could operate incrementally, in "on-line" mode, taking in individual measurements and re-estimating the set of causal states and their transition probabilities.

V. OPTIMALITIES AND UNIQUENESS

We now show that: causal states are maximally accurate predictors of minimal statistical complexity; they are unique in sharing both properties; and their state-to-state transitions are minimally stochastic. In other words, they satisfy both of the constraints borrowed from Occam, and they are the only representations that do so. The overarching moral here is that causal states and ϵ -machines are the goals in any learning or modeling scheme. The argument is made by the time-honored means of proving optimality theorems. We address, in our concluding remarks (Sec. VII), the practicalities involved in attaining these goals.

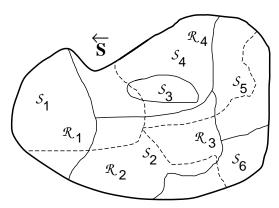


FIG. 3. An alternative class \mathcal{R} of states (delineated by dashed lines) that partition \mathbf{S} overlaid on the causal states \mathcal{S} (outlined by solid lines). Here, for example, \mathcal{S}_2 contains parts of \mathcal{R}_1 , \mathcal{R}_2 , \mathcal{R}_3 and \mathcal{R}_4 . The collection of all such alternative partitions form *Occam's pool*. Note again that the \mathcal{R}_i need not be compact nor simply connected, as drawn.

As part of our strategy, though, we also prove several results that are not optimality results; we call these lemmas to indicate their subordinate status. All of our theorems, and some of our lemmas, will be established by comparing causal states, generated by ϵ , with other rival sets of states, generated by other functions η . In short, none of the rival states—none of the other patterns—can out-perform the causal states.

It is convenient to fix some additional notation. Let S be the random variable for the current causal state,

 $\vec{S} \in \mathcal{A}$ the next "observable" we get from the original stochastic process, \mathcal{S}' the next causal state, \mathcal{R} the current state according to η , and \mathcal{R}' the next η -state. σ will stand for a particular value (causal state) of \mathcal{S} and ρ a particular value of \mathcal{R} . When we quantify over alternatives to the causal states, we quantify over \mathcal{R} .

Theorem 1 (Causal States are Maximally Prescient) [15]

For all \mathcal{R} and all $L \in \mathbb{Z}^+$,

$$H[\stackrel{\rightarrow}{S}^L|\mathcal{R}] \ge H[\stackrel{\rightarrow}{S}^L|\mathcal{S}] \ . \tag{42}$$

Proof. We have already seen that $H[\overrightarrow{S}^L | \mathcal{R}] \geq H[\overrightarrow{S}^L | \overrightarrow{S}]$ (Lemma 1). But by construction (Def. 5),

$$P(\overrightarrow{S}^{L} = \overrightarrow{s}^{L} | \overleftarrow{S} = \overleftarrow{s}) = P(\overrightarrow{S}^{L} = \overrightarrow{s}^{L} | \mathcal{S} = \epsilon(\overleftarrow{s})) . \quad (43)$$

Since entropies depend only on the probability distribution, $H[\overrightarrow{S}^L | \mathcal{S}] = H[\overrightarrow{S}^L | \overleftarrow{S}]$ for every L. Thus, $H[\overrightarrow{S}^L | \mathcal{R}] \geq H[\overrightarrow{S}^L | \mathcal{S}]$, for all L. QED.

Remark. That is to say, causal states are as good at predicting the future—are as prescient—as complete histories. In this, they satisfy the first requirement borrowed from Occam. Since the causal states are well defined and since they can be systematically approximated, we have shown that the upper bound on the strength of patterns (Def. 3 and Lemma 1, Remark) can in fact be reached. Intuitively, the causal states achieve this because, unlike effective states in general, they do not throw away any information about the future which might be contained in S. Even more colloquially, to paraphrase the definition of information in Ref. [70], the causal states record every difference (about the past) that makes a difference (to the future). We can actually make this intuition quite precise, in an easy corollary to the theorem.

Corollary 1 (Causal States Are Sufficient Statistics) The causal states S of a process are sufficient statistics for predicting it.

Proof. It follows from Thm. 1 and Eq. (8) that, for all $L \in \mathbb{Z}^+$,

$$I[\overrightarrow{S}^{L}; \mathcal{S}] = I[\overrightarrow{S}^{L}; \overleftarrow{S}],$$
 (44)

where I was defined in Eq. (8). Consequently, the causal state is a *sufficient statistic*—see Refs. [62, p. 37] and [71, sec. 2.4-2.5]—for predicting futures of any length. QED.

All subsequent results concern rival states that are as prescient as the causal states. We call these *prescient rivals* and denote a class of them $\hat{\mathcal{R}}$.

Definition 11 (Prescient Rivals) Prescient rivals $\hat{\mathcal{R}}$ are states that are as predictive as the causal states; viz., for all $L \in \mathbb{Z}^+$,

$$H[\overrightarrow{S}^{L}|\hat{\mathcal{R}}] = H[\overrightarrow{S}^{L}|\mathcal{S}] .$$
 (45)

Remark. Prescient rivals are also sufficient statistics.

Lemma 7 (Refinement Lemma) For all prescient rivals $\hat{\mathcal{R}}$ and for each $\hat{\rho} \in \hat{\mathcal{R}}$, there is a $\sigma \in \mathcal{S}$ and a measure-0 subset $\hat{\rho}_0 \subset \hat{\rho}$, possibly empty, such that $\hat{\rho} \setminus \hat{\rho}_0 \subseteq \sigma$, where \setminus is set subtraction.

Proof. We invoke a straightforward extension of Thm. 2.7.3 of Ref. [62]: If X_1, X_2, \ldots, X_n are random variables over the same set \mathcal{A} , each with distinct probability distributions, Θ a random variable over the integers from 1 to n such that $P(\Theta = i) = \lambda_i$, and Z a random variable over \mathcal{A} such that $Z = X_{\Theta}$, then

$$H[Z] = H\left[\sum_{i=1}^{n} \lambda_i X_i\right]$$

$$\geq \sum_{i=1}^{n} \lambda_i H[X_i] . \tag{46}$$

In words, the entropy of a mixture of distributions is at least the mean of the entropies of those distributions. This follows since H is strictly concave, which in turn follows from $x \log x$ being strictly convex for $x \geq 0$. We obtain equality in Eq. (46) if and only if all the λ_i are either 0 or 1, i.e., if and only if Z is at least weakly homogeneous (Def. 7).

The conditional distribution of futures for each rival state ρ can be written as a weighted mixture of the morphs of one or more causal states. (Cf. Fig. 3.) Thus, by Eq. (46), unless every ρ is at least weakly homogeneous with respect to S (for each L), the entropy of S conditioned on S will be higher than the minimum, the entropy conditioned on S. So, in the case of the maximally predictive \hat{R} , every $\hat{\rho} \in \hat{R}$ must be at least weakly homogeneous with respect to all S. But the causal states are the largest classes that are strictly homogeneous with respect to all S. Thus, the strictly homogeneous part of each $\hat{\rho} \in \hat{R}$ must be a subclass, possibly improper, of some causal state $\sigma \in S$. QED.

Remark 1. An alternative proof appears in App. E. Remark 2. The content of the lemma can be made quite intuitive, if we ignore for a moment the measure-0 set $\hat{\rho}_0$ of histories mentioned in its statement. It then asserts that any alternative partition $\hat{\mathcal{R}}$ that is as prescient as the causal states must be a refinement of the causal-state partition. That is, each $\hat{\mathcal{R}}_i$ must be a (possibly

improper) subset of some S_j . Otherwise, at least one $\hat{\mathcal{R}}_i$ would have to contain parts of at least two causal states. And so, using this $\hat{\mathcal{R}}_i$ to predict the future observables would lead to more uncertainty about \overrightarrow{S} than using the causal states. This is illustrated by Fig. 4, which should be contrasted with Fig. 3.

Adding the measure-0 set $\hat{\rho}_0$ of histories to this picture does not change its heuristic content much. Precisely because these histories have zero probability, treating them in an "inappropriate" way makes no discernible difference to predictions, morphs, and so on. There is a problem of terminology, however, since there seems to be no standard name for the relationship between the partitions $\hat{\mathcal{R}}$ and \mathcal{S} . We propose to say that the former is a refinement of the latter almost everywhere or, simply, a refinement a.e.

Remark 3. One cannot work the proof the other way around to show that the causal states have to be a refinement of the equally prescient $\hat{\mathcal{R}}$ -states. This is precluded because applying the theorem borrowed from Ref. [62], Eq. (46), hinges on being able to reduce uncertainty by specifying from which distribution one chooses. Since the causal states are constructed so as to be strictly homogeneous with respect to futures, this is not the case. Lemma 3 and Thm. 1 together protect us.

Remark 4. Because almost all of each prescient rival state is wholly contained within a single causal state, we can construct a function $g: \hat{\mathcal{R}} \mapsto \mathcal{S}$, such that, if $\eta(s) = \hat{\rho}$, then $\epsilon(s) = g(\hat{\rho})$ almost always. We can even say that $\mathcal{S} = g(\hat{\mathcal{R}})$ almost always, with the understanding that this means that, for each $\hat{\rho}$, $P(\mathcal{S} = \sigma | \hat{\mathcal{R}} = \hat{\rho}) > 0$ if and only if $\sigma = g(\hat{\rho})$.

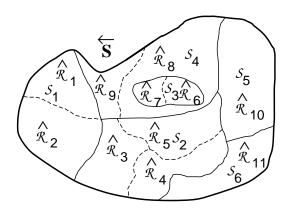


FIG. 4. A prescient rival partition $\hat{\mathcal{R}}$ must be a refinement of the causal-state partition almost everywhere. That is, almost all of each $\hat{\mathcal{R}}_i$ must contained within some \mathcal{S}_j ; the exceptions, if any, are a set of histories of measure 0. Here for instance \mathcal{S}_2 contains the positive-measure parts of $\hat{\mathcal{R}}_3$, $\hat{\mathcal{R}}_4$, and $\hat{\mathcal{R}}_5$. One of these rival states, say $\hat{\mathcal{R}}_3$, could have member-histories in any or all of the other causal states, provided the total measure of such exceptional histories is zero. Cf. Fig. 3.

Theorem 2 (Causal States Are Minimal) [15] For all prescient rivals $\hat{\mathcal{R}}$,

$$C_{\mu}(\hat{\mathcal{R}}) \ge C_{\mu}(\mathcal{S}) \ . \tag{47}$$

Proof. By Lemma 7, Remark 4, there is a function g such that $S = g(\hat{R})$ almost always. But $H[f(X)] \leq H[X]$ (Eq. (A11)) and so

$$H[S] = H[g(\hat{R})] \le H[\hat{R}]. \tag{48}$$

but $C_{\mu}(\hat{\mathcal{R}}) = H[\hat{\mathcal{R}}]$ (Def. 4). QED.

Remark 1. We have just established that no rival pattern, which is as good at predicting the observations as the causal states, is any simpler, in the sense given by Def. 4, than the causal states. (This is the theorem of Ref. [6].) Occam therefore tells us that there is no reason not to use the causal states. The next theorem shows that causal states are uniquely optimal, and so that Occam's Razor all but forces us to use them.

Remark 2. Here it becomes important that we are trying to predict the whole of $\overset{\rightharpoonup}{S}$ and not just some piece, $\overset{\rightharpoonup}{S}$. Suppose two histories $\overset{\rightharpoonup}{s}$ and $\overset{\rightharpoonup}{s}'$ have the same conditional distribution for futures of lengths up to L, but differing ones after that. They would then belong to different causal states. An η -state that merged those two causal states, however, would have just as much ability to predict $\overset{\rightharpoonup}{S}$ as the causal states. More, these \mathcal{R} -states would be simpler, in the sense that the uncertainty in the current state would be lower. We conclude that causal states are optimal, but for the hardest job—that of predicting futures of all lengths.

Remark 3. We have already seen (Thm. 1, Remark 2) that causal states are sufficient statistics for predicting futures of all lengths; so are all prescient rivals. A minimal sufficient statistic is one that is a function of all other sufficient statistics [62, p. 38]. Since, in the course of the proof of Thm. 2, we have shown that there is a function g from any $\hat{\mathcal{R}}$ to \mathcal{S} , we have also shown that causal states are minimal sufficient statistics.

We may now, as promised, define the *statistical complexity of a process* [5,6].

Definition 12 (Statistical Complexity of a Process) The statistical complexity " $C_{\mu}(\mathcal{O})$ " of a process \mathcal{O} is that of its causal states: $C_{\mu}(\mathcal{O}) \equiv C_{\mu}(\mathcal{S})$.

Due to the minimality of causal states we see that the statistical complexity measures the average amount of historical memory stored in the process. Without the minimality theorem, this interpretation would not be possible, since we could trivially elaborate internal states, while still generating the same observed process. C_{μ} for those states would grow without bound and so be arbitrary and not a characteristic property of the process [17].

Theorem 3 (Causal States Are Unique) For all prescient rivals $\hat{\mathcal{R}}$, if $C_{\mu}(\hat{\mathcal{R}}) = C_{\mu}(\mathcal{S})$, then there exists an invertible function between $\hat{\mathcal{R}}$ and \mathcal{S} that almost always preserves equivalence of state: $\hat{\mathcal{R}}$ and η are the same as \mathcal{S} and ϵ , respectively, except on a set of histories of measure 0.

Proof. From Lemma 7, we know that $S = g(\hat{R})$ almost always. We now show that there is a function f such that $\hat{R} = f(S)$ almost always, implying that $g = f^{-1}$ and that f is the desired relation between the two sets of states. To do this, by Eq. (A12) it is sufficient to show that $H[\hat{R}|S] = 0$. Now, it follows from an information-theoretic identity (Eq. (A8)) that

$$H[\mathcal{S}] - H[\mathcal{S}|\hat{\mathcal{R}}] = H[\hat{\mathcal{R}}] - H[\hat{\mathcal{R}}|\mathcal{S}]. \tag{49}$$

Since, by Lemma 7 $H[S|\hat{R}] = 0$, both sides of Eq. (49) are equal to H[S]. But, by hypothesis, $H[\hat{R}] = H[S]$. Thus, $H[\hat{R}|S] = 0$ and so there exists an f such that $\hat{R} = f(S)$ almost always. We have then that $f(g(\hat{R})) = \hat{R}$ and g(f(S)) = S, so $g = f^{-1}$. This implies that f preserves equivalence of states almost always: for almost all $s, s' \in S$, $\eta(s) = \eta(s')$ if and only if $\epsilon(s) = \epsilon(s')$. QED.

Remark. As in the case of the Refinement Lemma 7, on which the theorem is based, the measure-0 caveats seem unavoidable. A rival that is as predictive and as simple (in the sense of Def. 4) as the causal states, can assign a measure-0 set of histories to different states than the ϵ -machine does, but no more. This makes sense: such a measure-0 set makes no difference, since its members are never observed, by definition. By the same token, however, nothing prevents a minimal, prescient rival from disagreeing with the ϵ -machine on those histories.

Theorem 4 (ϵ -Machines Are Minimally Stochastic) [15] For all prescient rivals $\hat{\mathcal{R}}$,

$$H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}] \ge H[\mathcal{S}'|\mathcal{S}] ,$$
 (50)

where S' and \hat{R}' are the next causal state of the process and the next η -state, respectively.

Proof. From Lemma 5, S' is fixed by S and $\stackrel{\rightarrow}{S}^1$ together, thus $H[S'|S,\stackrel{\rightarrow}{S}^1]=0$ by Eq. (A12). Therefore, from the chain rule for entropies Eq. (A6),

$$H[\overrightarrow{S}^{1}|\mathcal{S}] = H[\mathcal{S}', \overrightarrow{S}^{1}|\mathcal{S}]. \tag{51}$$

We have no result like the Determinism Lemma 5 for the rival states $\hat{\mathcal{R}}$, but entropies are always nonnegative: $H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}, \vec{S}^1] \geq 0$. Since for all L, $H[\vec{S}^L|\hat{\mathcal{R}}] = H[\vec{S}^L|\mathcal{S}]$ by the definition, Def. (11), of prescient rivals, $H[\vec{S}^1|\hat{\mathcal{R}}] = H[\vec{S}^1|\mathcal{S}]$. Now we apply the chain rule again,

$$H[\hat{\mathcal{R}}', \vec{S}^1 | \hat{\mathcal{R}}] = H[\vec{S}^1 | \hat{\mathcal{R}}] + H[\hat{\mathcal{R}}' | \vec{S}^1, \hat{\mathcal{R}}]$$
 (52)

$$\geq H[\vec{S}^1|\hat{\mathcal{R}}] \tag{53}$$

$$=H[\stackrel{\rightarrow}{S}^{1}|\mathcal{S}] \tag{54}$$

$$=H[\mathcal{S}', \overrightarrow{S}^{1}|\mathcal{S}] \tag{55}$$

$$= H[\mathcal{S}'|\mathcal{S}] + H[\overrightarrow{S}^{1}|\mathcal{S}',\mathcal{S}]. \tag{56}$$

In going from Eq. (54) to Eq. (55) we have used Eq. (51), and in the last step we have used the chain rule once more.

Using the chain rule one last time, we have

$$H[\hat{R}', \vec{S}^{1}|\hat{R}] = H[\hat{R}'|\hat{R}] + H[\vec{S}^{1}|\hat{R}', \hat{R}].$$
 (57)

Putting these expansions, Eqs. (56) and (57), together we get

$$H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}] + H[\overrightarrow{S}^{1}|\hat{\mathcal{R}}',\hat{\mathcal{R}}] \ge H[\mathcal{S}'|\mathcal{S}] + H[\overrightarrow{S}^{1}|\mathcal{S}',\mathcal{S}] \quad (58)$$

$$H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}] - H[\mathcal{S}'|\mathcal{S}] \ge H[\overrightarrow{S}^{1}|\mathcal{S}',\mathcal{S}] - H[\overrightarrow{S}^{1}|\hat{\mathcal{R}}',\hat{\mathcal{R}}] .$$

From Lemma 7, we know that $S = g(\hat{R})$, so there is another function g' from ordered pairs of η -states to ordered pairs of causal states: $(S', S) = g'(\hat{R}', \hat{R})$. Therefore, Eq. (A14) implies

$$H[\stackrel{\rightarrow}{S}^1 | \mathcal{S}', \mathcal{S}] \ge H[\stackrel{\rightarrow}{S}^1 | \hat{\mathcal{R}}', \hat{\mathcal{R}}] .$$
 (59)

And so, we have that

$$H[\overrightarrow{S}^{1}|\mathcal{S}',\mathcal{S}] - H[\overrightarrow{S}^{1}|\hat{\mathcal{R}}',\hat{\mathcal{R}}] \ge 0$$

$$H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}] - H[\mathcal{S}'|\mathcal{S}] \ge 0$$

$$H[\hat{\mathcal{R}}'|\hat{\mathcal{R}}] \ge H[\mathcal{S}'|\mathcal{S}] . \tag{60}$$

QED.

Remark. What this theorem says is that there is no more uncertainty in transitions between causal states, than there is in the transitions between any other kind of prescient effective states. In other words, the causal states approach as closely to perfect determinism—in the usual physical, non-computation-theoretic sense—as any rival that is as good at predicting the future. This sort of internal determinism has long been held to be a desideratum of scientific models [72].

VI. BOUNDS

In this section we develop bounds between measures of structural complexity and entropy derived from ϵ -machines and those from ergodic and information theories, which are perhaps more familiar.

Definition 13 (Excess Entropy) The excess entropy **E** of a process is the mutual information between its semi-infinite past and its semi-infinite future:

$$\mathbf{E} \equiv I[\overrightarrow{S}; \overleftarrow{S}] \ . \tag{61}$$

The excess entropy is a frequently-used measure of the complexity of stochastic processes and appears under a variety of names; e.g., "predictive information", "stored information", "effective measure complexity", and so on [73–79]. **E** measures the amount of apparent information stored in the observed behavior about the past. As we now establish, **E** is not, in general, the amount of memory that the process stores internally about its past; a quantity measured by C_{μ} .

Theorem 5 (The Bounds of Excess) The statistical complexity C_{μ} bounds the excess entropy **E**:

$$\mathbf{E} \le C_{\mu} , \qquad (62)$$

with equality if and only if $H[S|\overrightarrow{S}] = 0$.

Proof. $\mathbf{E} = I[\overrightarrow{S}; \overleftarrow{S}] = H[\overrightarrow{S}] - H[\overrightarrow{S} \mid \overleftarrow{S}]$ and, by the construction of causal states, $H[\overrightarrow{S} \mid \overleftarrow{S}] = H[\overrightarrow{S} \mid \mathcal{S}]$, so

$$\mathbf{E} = H[\overrightarrow{S}] - H[\overrightarrow{S} | \mathcal{S}] = I[\overrightarrow{S}; \mathcal{S}] . \tag{63}$$

Thus, since the mutual information between two variables is never larger than the self-information of either one of them (Eq. (A9)), $\mathbf{E} \leq H[S] = C_{\mu}$, with equality if and only if H[S|S] = 0. QED.

Remark 1. Note that we have invoked $H[\overrightarrow{S}]$, not $H[\overrightarrow{S}]$, but only while subtracting off quantities like $H[\overrightarrow{S} \mid \overleftarrow{S}]$. We need not worry, therefore, about the existence of a finite $L \to \infty$ limit for $H[\overrightarrow{S}]$, just that of a finite $L \to \infty$ limit for $I[\overrightarrow{S} \mid \overleftarrow{S}]$ and $I[\overrightarrow{S} \mid \overleftarrow{S}]$. There are many elementary cases (e.g., the fair coin process) where the latter limits exist while the former do not.

Remark 2. At first glance, it is tempting to see ${\bf E}$ as the amount of information stored in a process. As Thm. 5 shows, this temptation should be resisted. ${\bf E}$ is only a lower bound on the true amount of information the process stores about its history, namely C_{μ} . We can, however, say that ${\bf E}$ measures the apparent information in the process, since it is defined directly in terms of observed sequences and not in terms of hidden, intrinsic states, as C_{μ} is.

Remark \hat{g} . Perhaps another way to describe what \mathbf{E} measures is to note that, by its implicit assumption of block-Markovian structure, it takes sequence-blocks as states. But even for the class of block-Markovian sources, for which such an assumption is appropriate, excess entropy and statistical complexity measure different kinds

of information storage. Refs. [65] and [80] showed that in the case of one-dimensional range-R spin systems, or any other block-Markovian source where block configurations are isomorphic to causal states:

$$C_{\mu} = \mathbf{E} + Rh_{\mu} , \qquad (64)$$

for finite R. Only for zero-entropy-rate block-Markovian sources will the excess entropy, a quantity estimated directly from sequence blocks, equal the statistical complexity, the amount of memory stored in the process. Examples of such sources include periodic processes, for which we have $C_{\mu} = \mathbf{E} = \log_2 p$, where p is the period.

Corollary 2 For all prescient rivals $\hat{\mathcal{R}}$,

$$\mathbf{E} \le H[\hat{\mathcal{R}}] \ . \tag{65}$$

Proof. This follows directly from Thm. 2, since $H[\hat{\mathcal{R}}] \geq C_{u}$. QED.

Lemma 8 (Conditioning Does Not Affect Entropy Rate) For all prescient rivals $\hat{\mathcal{R}}$,

$$h[\vec{S}] = h[\vec{S} \mid \hat{\mathcal{R}}] , \qquad (66)$$

where the entropy rate $h[\vec{S}]$ and the conditional entropy rate $h[\vec{S} \mid \hat{\mathcal{R}}]$ were defined in Eq. (9) and Eq. (10), respectively.

Proof. From Thm. 5 and its Corollary 2, we have

$$\lim_{L \to \infty} \left(H[\vec{S}^{\perp}] - H[\vec{S}^{\perp}|\hat{\mathcal{R}}] \right) \le \lim_{L \to \infty} H[\hat{\mathcal{R}}] , \qquad (67)$$

or,

$$\lim_{L \to \infty} \frac{H[\overrightarrow{S}^L] - H[\overrightarrow{S}^L | \hat{\mathcal{R}}]}{L} \le \lim_{L \to \infty} \frac{H[\hat{\mathcal{R}}]}{L} . \tag{68}$$

Since, by Eq. (A4), $H[\vec{S}^L] - H[\vec{S}^L|\hat{\mathcal{R}}] \ge 0$, we have

$$h[\overrightarrow{S}] - h[\overrightarrow{S} \mid \hat{\mathcal{R}}] = 0 . \tag{69}$$

QED.

Remark. Forcing the process into a certain state $\hat{\mathcal{R}} = \hat{\rho}$ is akin to applying a controller, once. But in the infinite-entropy case, $H[\overrightarrow{S}^{\ }] \to_{L \to \infty} \infty$, with which we are concerned, the future could contain (or consist of) an infinite sequence of disturbances. In the face of this "grand disturbance", the effects of the finite control are simply washed out.

Another way of viewing this is to reflect on the fact that $h[\overrightarrow{S}]$ accounts for the effects of all the dependencies between all the parts of the entire semi-infinite future.

This, owing to the time-translation invariance of stationarity, is equivalent to taking account of all the dependencies in the entire process, including those between past and future. But these are what is captured by $h[\vec{S} \mid \hat{\mathcal{R}}]$. It is not that conditioning on \mathcal{R} fails to reduce our uncertainty about the future; it does so, for all finite times, and conditioning on \mathcal{S} achieves the maximum possible reduction in uncertainty. Rather, the lemma asserts that such conditioning cannot effect the asymptotic rate at which such uncertainty grows with time.

Theorem 6 (Control Theorem) Given a class $\hat{\mathcal{R}}$ of prescient rivals,

$$H[S] - h[\overrightarrow{S} \mid \hat{\mathcal{R}}] \le C_{\mu} , \qquad (70)$$

where H[S] is the entropy of a single symbol from the observable stochastic process.

Proof. As is well known (Ref. [62, Thm. 4.2.1, p. 64]), for any stationary stochastic process,

$$\lim_{L \to \infty} \frac{H[\overrightarrow{S}^L]}{L} = \lim_{L \to \infty} H[S_L|\overrightarrow{S}^{L-1}] . \tag{71}$$

Moreover, the limits always exist. Up to this point, we have defined $h[\vec{S}]$ in the manner of the left-hand side; recall Eq. (9). It will be convenient in the following to use that of the right-hand side.

From the definition of conditional entropy, we have

$$H[\overset{\leftarrow}{S}^{L}] = H[\overset{\leftarrow}{S}^{1} | \overset{\leftarrow}{S}^{L-1}] + H[\overset{\leftarrow}{S}^{L-1}]$$

$$= H[\overset{\leftarrow}{S}^{L-1} | \overset{\leftarrow}{S}^{1}] + H[\overset{\leftarrow}{S}^{1}]. \tag{72}$$

So we can express the entropy of the last observable the process generated before the present as

$$H[\overset{\leftarrow}{S}^{1}] = H[\overset{\leftarrow}{S}^{L}] - H[\overset{\leftarrow}{S}^{L-1}|\overset{\leftarrow}{S}^{1}]$$
 (73)

$$= H[\overset{\leftarrow}{S}^{1}|\overset{\leftarrow}{S}^{L-1}] + H[\overset{\leftarrow}{S}^{L-1}] - H[\overset{\leftarrow}{S}^{L-1}|\overset{\leftarrow}{S}^{1}] \quad (74)$$

$$= H[\overset{\leftarrow}{S} | \overset{\leftarrow}{S} \overset{L-1}{S}] + I[\overset{\leftarrow}{S} ; \overset{\leftarrow}{S}] . \tag{75}$$

We go from Eq. (73) to Eq. (74) by substituting the first RHS of Eq. (72) for $H[\stackrel{\leftarrow}{S}^L]$.

Taking the $L \to \infty$ limit has no effect on the LHS,

$$H[\overset{\leftarrow}{S}^{1}] = \lim_{L \to \infty} \left(H[\overset{\leftarrow}{S}^{1}|\overset{\leftarrow}{S}^{L-1}] + I[\overset{\leftarrow}{S}^{L-1};\overset{\leftarrow}{S}^{1}] \right) . \quad (76)$$

Since the process is stationary, we can move the first term in the limit forward to $H[S_L|\overrightarrow{S}^{L-1}]$. This limit is $h[\overrightarrow{S}]$, by Eq. (71). Furthermore, because of stationarity, $H[\overrightarrow{S}] = H[\overrightarrow{S}] = H[S]$. Shifting the entropy rate $h[\overrightarrow{S}]$ to the LHS of Eq. (76) and appealing to time-translation once again, we have

$$H[S] - h[\overrightarrow{S}] = \lim_{L \to \infty} I[\overrightarrow{S}^{L-1}; \overleftarrow{S}]$$
 (77)

$$=I[\overset{\leftarrow}{S};\vec{S}^{1}] \tag{78}$$

$$= H[\overrightarrow{S}] - H[\overrightarrow{S}] \stackrel{\leftarrow}{S}$$

$$= H[\overrightarrow{S}] - H[\overrightarrow{S}] \stackrel{\leftarrow}{S}$$

$$= H[\overrightarrow{S}] - H[\overrightarrow{S}] \stackrel{\rightarrow}{S}$$

$$= I[\overrightarrow{S}] : \mathcal{S}$$
(80)
$$= I[\overrightarrow{S}] : \mathcal{S}$$
(81)

$$=H[\overrightarrow{S}^{1}]-H[\overrightarrow{S}^{1}|\mathcal{S}] \tag{80}$$

$$=I[\overrightarrow{S}^{1};\mathcal{S}] \tag{81}$$

$$\leq H[\mathcal{S}] = C_{\mu} , \qquad (82)$$

where the last inequality comes from Eq. (A9). QED.

Remark 1. The Control Theorem is inspired by, and is a version of, Ashby's law of requisite variety [81, ch. 11]. This states that applying a controller can reduce the uncertainty in the controlled variable by at most the entropy of the control variable. (This result has recently been rediscovered in Ref. [82].) Thinking of the controlling variable as the causal state, we have here a limitation on the controller's ability to reduce the entropy rate.

Remark 2. This is the only result so far where the difference between the finite-L and the infinite-L cases is important. For the analogous result in the finite case, see App. F, Thm. 7.

Remark 3. By applying Thm. 2 and Lemma 8, we could go from the theorem as it stands to $H[S] - h[\vec{S}]$ $|\hat{\mathcal{R}}| \leq H[\hat{\mathcal{R}}]$. This has a pleasing appearance of symmetry to it, but is actually a weaker limit on the strength of the pattern or, equivalently, on the amount of control that fixing the causal state (or one of its rivals) can exert.

VII. CONCLUDING REMARKS

A. Discussion

Let's review, informally, what we have shown. We began with questions about the nature of patterns and about pattern discovery. Our examination of these issues lead us to want a way of describing patterns that was at once algebraic, computational, intrinsically probabilistic, and causal. We then defined patterns in ensembles, in a very general and abstract sense, as equivalence classes of histories, or sets of hidden states, used for prediction. We defined the strength of such patterns (by their forecasting ability or prescience) and their statistical complexity (by the entropy of the states or the amount of information retained by the process about its history). We showed that there was a limit on how strong such patterns could get for each particular process, given by the predictive ability of the entire past. In this way, we narrowed our goal to finding a predictor of maximum strength and minimum complexity.

Optimal prediction led us to the equivalence relation \sim_{ϵ} and the function ϵ and so to representing patterns by causal states and their transitions—the ϵ -machine. Our first theorem showed that the causal states are maximally

prescient; our second, that they are the simplest way of representing the pattern of maximum strength; our third theorem, that they are unique in having this double optimality. Further results showed that ϵ -machines are the least stochastic way of capturing maximum-strength patterns and emphasized the need to employ the efficacious but hidden states of the process, rather than just its gross observables, such as sequence blocks.

Why are ϵ -machine states causal? First, ϵ -machine architecture (say, as given by its semi-group algebra) delineates the dependency between the morphs $P(\overrightarrow{S} \mid \overleftarrow{S})$, considered as events in which each new symbol determines the succeeding morph. Thus, if state B follows state A then A is a cause of B and B is an effect of A. Second, ϵ -machine minimality guarantees that there are no other events that intervene to render A and B independent [17].

The ϵ -machine is thus a causal representation of all the patterns in the process. It is maximally predictive and minimally complex. It is at once computational, since it shows how the process stores information (in the causal states) and transforms that information (in the state-tostate transitions), and algebraic (for details on which see App. D). It can be analytically calculated from given distributions and systematically approached from empirical data. It satisfies the basic constraints laid out in Sec. IIF.

These comments suggest that computational mechanics and ϵ -machines are related or may be of interest to a number of fields. Time series analysis, decision theory, machine learning, and universal coding theory explicitly or implicitly require models of observed processes. The theories of stochastic processes, formal languages and computation, and of measures of physical complexity are all concerned with representations of processes—concerns which also arise in the design of novel forms of computing devices. Very often the motivations of these fields are far removed from computational mechanics. But it is useful, if only by way of contrast, to touch briefly on these areas and highlight one or several connections with computational mechanics, and we do so in App. G.

B. Limitations of the Current Results

Let's catalogue the restrictive assumptions we made at the beginning and that were used by our development.

- 1. We know exact joint probabilities over sequence blocks of all lengths for a process.
- 2. The observed process takes on discrete values.
- 3. The process is discrete in time.
- 4. The process is a pure time series; e.g., without spatial extent.
- 5. The observed process is stationary.

6. Prediction can only be based on the process's past, not on any outside source of information.

The question arises, Can any be relaxed without much trouble?

One way to lift the first limitation is to develop a statistical error theory for ϵ -machine inference that indicates, say, how much data is required to attain a given level of confidence in an ϵ -machine with a given number of causal states. This program is underway and, given its initial progress, we describe several issues in more detail in the next section.

The second limitation probably can be addressed, but with a corresponding increase in mathematical sophistication. The information-theoretic quantities we have used are also defined for continuous random variables. It is likely that many of the results carry over to the continuous setting.

The third limitation also looks similarly solvable, since continuous-time stochastic process theory is moderately well developed. This may involve sophisticated probability theory or functional analysis.

As for the fourth limitation, there already exist tricks to make spatially extended systems look like time series. Essentially, one looks at all the paths through spacetime, treating each one as if it were a time series. While this works well for data compression [83], it is not yet clear whether it will be entirely satisfactory for capturing structure [84]. More work needs to be done on this subject.

It is unclear at this time how to relax the assumption of stationarity. One can formally extend most of the results in this paper to non-stationary processes without much trouble. It is, however, unclear how much substantive content these extensions have and, in any case, a systematic classification of non-stationary processes is (at best) in its infant stages.

Finally, one might say that the last restriction is a positive feature when it comes to thinking about patterns and the intrinsic structure of a process. "Pattern" is a vague word, of course, but even in ordinary usage it is only supposed to involve things inside the process, not the rest of the universe. Given two copies of a document, the contents of one copy can be predicted with an enviable degree of accuracy by looking at the other copy. This tells us that they share a common structure, but says absolutely nothing about what that pattern is, since it is just as true of well-written and tightly-argued scientific papers (which presumably are highly organized) as it is of monkey-at-keyboard pieces of gibberish (which definitely are not).

C. Conclusions and Directions for Future Work

Computational mechanics aims to understand the nature of patterns and pattern discovery. We hope that the foregoing development has convinced the reader that

we are neither being rash when we say that we have laid a foundation for those projects, nor that we are being flippant when we say that patterns are what ϵ -machines represent and that we discover them by ϵ -machine reconstruction. We would like to close by marking out two broad avenues for future work.

First, consider the mathematics of ϵ -machines them-We have just mentioned possible extensions selves. in the form of lifting assumptions made in this development, but there are many other ways to go. A number of measure-theoretic issues relating to the definition of causal states (omitted here for brevity) deserve careful treatment, along the lines of Ref. [10]. It would be helpful to have a good understanding of the measurement-resolution scaling properties of ϵ -machines for continuous-state processes, and of their relation to such ideas in automata theory as the Krohn-Rhodes decomposition [30]. Anyone who manages to absorb Volume II of Ref. [26] would probably be in a position to answer interesting questions about the structures that processes preserve, perhaps even to give a purely relationtheoretic account of ϵ -machines. We have alluded in a number of places to the trade-off between prescience and complexity. For a given process there is presumably a sequence of optimal machines connecting the one-state, zero-complexity machine with minimal prescience to the ϵ -machine. Each member of the path is the minimal machine for a certain degree of prescience; it would be very interesting to know what, if anything, we can say in general about the shape of this "prediction frontier".

Second, there is ϵ -machine reconstruction, an activity about which we have said next to nothing. As we mentioned above (p. 12), there are already several algorithms for reconstructing machines from data, even "on-line" ones. It is fairly evident that these algorithms will find the true machine in the limit of infinite time and infinite data. What is needed is an understanding of the error statistics [85] of different reconstruction procedures of the kinds of mistakes these procedures make and the probabilities with which they make them. Ideally, we want to find "confidence regions" for the products of reconstruction. The aim is to calculate (i) the probabilities of different degrees of reconstruction error for a given volume of data, (ii) the amount of data needed to be confident of a fixed bound on the error, or (iii) the rates at which different reconstruction procedures converge on the ϵ -machine. So far, an analytical theory has been developed that predicts the average number of estimated causal states as a function of the amount of data used when reconstructing certain kinds of processes [86]. Once we possess a more complete theory of statistical inference for ϵ -machines, analogous perhaps to what already exists in computational learning theory, we will be in a position to begin analyzing, sensibly and rigorously, the multitude of intriguing patterns and information-processing structures the natural world presents.

ACKNOWLEDGMENTS

We thank Dave Albers, Dave Feldman, Jon Fetter, Rob Haslinger, Wim Hordijk, Amihan Huesmann, Cris Moore, Mitch Porter, Erik van Nimwegen, and Karl Young for advice on the manuscript; and the participants in the 1998 SFI Complex Systems Summer School, the Madison probability seminar, the Madison Physics Department's graduate student mini-colloquium, and the Ann Arbor Complex Systems seminar for numerous helpful comments on earlier versions of these results. This work was supported at the Santa Fe Institute under the Computation, Dynamics, and Inference Program via ONR grant N00014-95-1-0975, NSF grant PHY-9970158, and DARPA contract F30602-00-2-0583.

APPENDIX A: INFORMATION-THEORETIC FORMULÆ

The following formulæ prove useful in the development. They are relatively intuitive, given our interpretation, and they can all be proved with little more than straight algebra; see Ref. [62, ch. 2]. Below, f is a function.

$$H[X,Y] = H[X] + H[Y|X] \tag{A1}$$

$$H[X,Y] \ge H[X] \tag{A2}$$

$$H[X,Y] \le H[X] + H[Y] \tag{A3}$$

$$H[X|Y] \le H[X] \tag{A4}$$

$$H[X|Y] = H[X]$$
 iff X is independent of Y (A5)

$$H[X, Y|Z] = H[X|Z] + H[Y|X, Z]$$
 (A6)

$$H[X, Y|Z] \ge H[X|Z] \tag{A7}$$

$$H[X] - H[X|Y] = H[Y] - H[Y|X]$$
 (A8)

$$I[X;Y] \le H[X] \tag{A9}$$

$$I[X;Y] = H[X] \text{ iff } H[X|Y] = 0$$
 (A10)

$$H[f(X)] \le H[X] \tag{A11}$$

$$H[X|Y] = 0 \text{ iff } X = f(Y) \tag{A12}$$

$$H[f(X)|Y] \le H[X|Y] \tag{A13}$$

$$\Pi[J(X)|Y] \leq \Pi[X|Y] \tag{1119}$$

$$H[X|f(Y)] \ge H[X|Y] \tag{A14}$$

Eqs. (A1) and (A6) are called the *chain rules* for entropies. Strictly speaking, the right hand side of Eq. (A12) should read "for each y, P(X=x|Y=y)>0 for one and only one x".

APPENDIX B: THE EQUIVALENCE RELATION THAT INDUCES CAUSAL STATES

Any relation that is reflexive, symmetric, and transitive is an *equivalence relation*.

Consider the set $\overleftarrow{\mathbf{S}}$ of all past sequences, of any length:

$$\stackrel{\leftarrow}{\mathbf{S}} = \{\stackrel{\leftarrow}{s}^L = s_{L-1} \cdots s_{-1} : s_i \in \mathcal{A}, L \in \mathbb{Z}^+\} . \tag{B1}$$

Recall that $\overset{\leftarrow}{s}^0 = \lambda$, the empty string. We define the relation \sim_{ϵ} over $\overset{\leftarrow}{\mathbf{S}}$ by

$$\stackrel{\leftarrow}{s_i}^K \sim_{\epsilon} \stackrel{\leftarrow}{s_j}^L \Leftrightarrow P(\overrightarrow{S} \mid \stackrel{\leftarrow}{s_i}^K) = P(\overrightarrow{S} \mid \stackrel{\leftarrow}{s_j}^L), \quad (B2)$$

for all semi-infinite $\overrightarrow{S} = s_0 s_1 s_2 \cdots$, where $K, L \in \mathbb{Z}^+$. Here we show that \sim_{ϵ} is an equivalence relation by reviewing the basic properties of relations, equivalence classes, and partitions. (The proof details are straightforward and are not included. See Ref. [87].) We will drop the length variables K and L and denote by $\overleftarrow{s}, \overleftarrow{s}', \overleftarrow{s}'' \in \overleftarrow{\mathbf{S}}$ members of any length in the set $\overleftarrow{\mathbf{S}}$ of Eq. (B1).

First, \sim_{ϵ} is a *relation* on **S** since we can represent it as a subset of the Cartesian product

$$\stackrel{\leftarrow}{\mathbf{S}} \times \stackrel{\leftarrow}{\mathbf{S}} = \{ (\stackrel{\leftarrow}{s}, \stackrel{\leftarrow}{s}') : \stackrel{\leftarrow}{s}, \stackrel{\leftarrow}{s}' \in \stackrel{\leftarrow}{\mathbf{S}} \} . \tag{B3}$$

Second, the relation \sim_{ϵ} is an equivalence relation on \mathbf{S} since it is

- 1. reflexive: $\overleftarrow{s} \sim_{\epsilon} \overleftarrow{s}$, for all $\overleftarrow{s} \in \overleftarrow{\mathbf{S}}$;
- 2. symmetric: $\overleftarrow{s} \sim_{\epsilon} \overleftarrow{s}' \Rightarrow \overleftarrow{s}' \sim_{\epsilon} \overleftarrow{s};$ and
- 3. transitive: $\overleftarrow{s} \sim_{\epsilon} \overleftarrow{s}'$ and $\overleftarrow{s}' \sim_{\epsilon} \overleftarrow{s}'' \Rightarrow \overleftarrow{s} \sim_{\epsilon} \overleftarrow{s}''$.

Third, if $\overleftarrow{s} \in \mathbf{S}$, the equivalence class of \overleftarrow{s} is

$$[\overleftarrow{s}] = \{ \overleftarrow{s}' \in \overleftarrow{\mathbf{S}} : \overleftarrow{s}' \sim_{\epsilon} \overleftarrow{s} \} .$$
 (B4)

The set of all equivalence classes in \mathbf{S} is denoted $\mathbf{S}/\sim_{\epsilon}$ and is called the *factor set* of \mathbf{S} with respect to \sim_{ϵ} . In Sec. IV A we called the individual equivalence classes causal states S_i and denoted the set of causal states $S = \{S_i : i = 0, 1, ..., k - 1\}$. That is, $S = \mathbf{S}/\sim_{\epsilon}$. (We noted in the main development that the cardinality k = |S| of causal states may or may not be finite.)

Finally, we list several basic properties of the causalstate equivalence classes.

1.
$$\bigcup_{\stackrel{\leftarrow}{s}\in \stackrel{\leftarrow}{\mathbf{S}}} [\stackrel{\leftarrow}{s}] = \stackrel{\leftarrow}{\mathbf{S}}$$
.

2.
$$\bigcup_{i=0}^{k-1} S_i = \mathbf{S}$$
.

3.
$$[\stackrel{\leftarrow}{s}] = [\stackrel{\leftarrow}{s}'] \Leftrightarrow \stackrel{\leftarrow}{s} \sim_{\epsilon} \stackrel{\leftarrow}{s}'$$
.

4. If
$$\overleftarrow{s}, \overleftarrow{s}' \in \overleftarrow{\mathbf{S}}$$
, either

(a)
$$[\stackrel{\leftarrow}{s}] \cap [\stackrel{\leftarrow}{s}'] = \emptyset$$
 or

(b)
$$[\stackrel{\leftarrow}{s}] = [\stackrel{\leftarrow}{s}']$$
.

- 5. The causal states \mathcal{S} are a partition of $\stackrel{\leftarrow}{\mathbf{S}}$. That is,
 - (a) $S_i \neq \emptyset$ for each i,

(b)
$$\bigcup_{i=0}^{k-1} \mathcal{S}_i = \stackrel{\leftarrow}{\mathbf{S}}$$
, and

(c)
$$S_i \cap S_j = \emptyset$$
 for all $i \neq j$.

We denote the start state with S_0 . The start state is the causal state associated with $s = \lambda$. That is, $S_0 = [\lambda]$.

APPENDIX C: TIME REVERSAL

The definitions and properties of the causal states obtained by scanning sequences in the opposite direction, i.e., the causal states $\overrightarrow{\mathbf{S}}/\sim_{\epsilon}$, follow similarly to those derived just above in App. B. In general, $\overleftarrow{\mathbf{S}}/\sim_{\epsilon} \neq \overrightarrow{\mathbf{S}}/\sim_{\epsilon}$. That is, past causal states are not necessarily the same as future causal states; past and future morphs can differ; unlike entropy rate [15], past and future statistical complexities need not be equal: $\overrightarrow{C}_{\mu} \neq \overrightarrow{C}_{\mu}$; and so on. The presence or lack of this type of time-reversal symmetry, as reflected in these inequalities, is a fundamental property of a process.

APPENDIX D: ϵ -MACHINES ARE MONOIDS

A semi-group is a set of elements closed under an associative binary operator, but without a guarantee that every, or indeed any, element has an inverse [88]. A monoid is a semi-group with an identity element. Thus, semi-groups and monoids are generalizations of groups. Just as the algebraic structure of a group is generally interpreted as a symmetry, we propose to interpret the algebraic structure of a semi-group as a generalized symmetry. The distinction between monoids and other semi-groups becomes important here: only semi-groups with an identity element—i.e., monoids—can contain subsets that are groups and so represent conventional symmetries.

We claim that the transformations that concatenate strings of symbols from \mathcal{A} onto other such strings form a semi-group G, the generators of which are the transformations that concatenate the elements of \mathcal{A} . The identity element is to be provided by concatenating the null symbol λ . The concatenation of string t onto the string t is forbidden if and only if strings of the form t have probability zero in a process. All such concatenations are to be realized by a single semi-group element denoted t. Since if t if t is t if t

Recall that, from our definition of the labeled transition probabilities, $T_{ij}^{(\lambda)} = \delta_{ij}$. Thus, $\mathbf{T}^{(\lambda)}$ is an identity element. This suggests using the labeled transition matrices to form a matrix representation of the semi-group. Accordingly, first define $U_{ij}^{(s)}$ by setting $U_{ij}^{(s)} = 0$ when $T_{ij}^{(s)} = 0$ and $U_{ij}^{(s)} = 1$ otherwise, to

remove probabilities. Then define the set of matrices $\mathbf{U} = \{\mathbf{T}^{(\lambda)}\} \bigcup \{\mathbf{U}^{(\mathbf{s})}, \mathbf{s} \in \mathcal{A}\}$. Finally, define G as the set of all matrices generated from the set \mathbf{U} by recursive multiplication. That is, an element g of G is

$$g^{(ab\dots cd)} = \mathbf{U}^{(d)}\mathbf{U}^{(c)}\dots\mathbf{U}^{(b)}\mathbf{U}^{(a)}, \qquad (D1)$$

where $a, b, \ldots c, d \in \mathcal{A}$. Clearly, G constitutes a semigroup under matrix multiplication. Moreover, $g^{(a \dots bc)} =$ $\mathbf{0}$ (the all-zero matrix) if and only if, having emitted the symbols $a \dots b$ in order, we must arrive in a state from which it is impossible to emit the symbol c. That is, the zero-matrix $\mathbf{0}$ is generated if and only if the concatenation of c onto $a \dots b$ is forbidden. The element \emptyset is thus the all-zero matrix $\mathbf{0}$, which clearly satisfies the necessary constraints. This completes the proof of Proposition 1.

We call the matrix representation—Eq. (D1) taken over all words in \mathcal{A}^k —of G the semi-group machine of the ϵ -machine $\{\mathcal{S}, \mathbf{T}\}$. See Ref. [89].

APPENDIX E: ALTERNATE PROOF OF THE REFINEMENT LEMMA

The proof of Lemma 7 carries through verbally, but we do not wish to leave loop-holes. Unfortunately, this means introducing two new bits of mathematics.

First of all, we need the largest classes that are strictly homogeneous (Def. 6) with respect to \vec{S}^L for fixed L; these are, so to speak, truncations of the causal states. Accordingly, we will talk about \mathcal{S}^L and σ^L , which are analogous to \mathcal{S} and σ . We will also need to define the function $\phi^L_{\sigma\rho} \equiv \mathrm{P}(\mathcal{S}^L = \sigma^L | \mathcal{R} = \rho)$.

Putting these together, for every L we have

$$H[\vec{S}^{L}|\mathcal{R} = \rho] = H[\sum_{\sigma^{L}} \phi_{\sigma\rho}^{L} P(\vec{S}^{L}|\mathcal{S}^{L} = \sigma^{L})]$$
 (E1)

$$\geq \sum_{\sigma^L} \phi^L_{\sigma\rho} H[\vec{S}^L | \mathcal{S}^L = \sigma^L] \ . \tag{E2}$$

Thus,

$$H[\overrightarrow{S}^{L} \mid \mathcal{R}] = \sum_{\alpha} P(\mathcal{R} = \rho) H[\overrightarrow{S}^{L} \mid \mathcal{R} = \rho]$$
 (E3)

$$\geq \sum_{\rho} P(\mathcal{R} = \rho) \sum_{\sigma^L} \phi_{\sigma\rho}^L H[\overrightarrow{S}^L | \mathcal{S}^L = \sigma^L]$$
 (E4)

$$= \sum_{\sigma^L, \rho} P(\mathcal{R} = \rho) \phi_{\sigma\rho}^L H[\overrightarrow{S}^L | \mathcal{S}^L = \sigma^L]$$
 (E5)

$$= \sum_{\sigma^L, \rho} P(\mathcal{S}^L = \sigma^L, \mathcal{R} = \rho) H[\vec{S}^L | \mathcal{S}^L = \sigma^L] \quad (E6)$$

$$= \sum_{\sigma^L} P(S^L = \sigma^L) H[\overrightarrow{S} \mid S^L = \sigma^L]$$
 (E7)

$$=H[\overrightarrow{S}^L|\mathcal{S}^L]. \tag{E8}$$

That is to say,

$$H[\vec{S}^{L}|\mathcal{R}] \ge H[\vec{S}^{L}|\mathcal{S}^{L}] ,$$
 (E9)

with equality if and only if every $\phi^L_{\sigma\rho}$ is either 0 or 1. Thus, if $H[\stackrel{\rightarrow}{S}^L|\hat{\mathcal{R}}]=H[\stackrel{\rightarrow}{S}|\mathcal{S}^L]$, every $\hat{\rho}$ is entirely contained within some σ^L ; except for possible subsets of measure 0. But if this is true for every L—which, in the case of a prescient rival $\hat{\mathcal{R}}$, it is—then every $\hat{\rho}$ is at least weakly homogeneous (Def. 7) with respect to all $\stackrel{\rightarrow}{S}$. Thus, by Lemma 3, all its members, except for that same subset of measure 0, belong to the same causal state. QED.

APPENDIX F: FINITE ENTROPY FOR THE SEMI-INFINITE FUTURE

While cases where H[S] is finite—more exactly, where $\lim_{L\to\infty} H[S]$ exists and is finite—may be uninteresting for information-theorists, they are of great interest to physicists, since they correspond, among other things, to periodic and limit-cycle behaviors. There are, however, only two substantial differences between what is true of the infinite-entropy processes considered in the main body of the development and the finite-entropy case.

First, we can simply replace statements of the form "for all L, $H[\overrightarrow{S}]$..." with $H[\overrightarrow{S}]$. For example, the optimal prediction theorem (Thm. 1) for finite-entropy processes becomes for all \mathcal{R} , $H[\overrightarrow{S}|\mathcal{R}] \geq H[\overrightarrow{S}|\mathcal{S}]$. The details of the proofs are, however, entirely analogous.

Second, we can prove a substantially stronger version of the control theorem (Thm. 6).

Theorem 7 (The Finite-Control Theorem) For all prescient rivals $\hat{\mathcal{R}}$,

$$H[\overrightarrow{S}] - H[\overrightarrow{S} \mid \hat{\mathcal{R}}] \le C_{\mu} .$$
 (F1)

Proof. By a direct application of Eq. (A9) and the definition of mutual information Eq. (8), we have that

$$H[\overrightarrow{S}] - H[\overrightarrow{S} \mid \mathcal{S}] \le H[\mathcal{S}] .$$
 (F2)

But, by the definition of prescient rivals (Def. 11), $H[S] = H[S] + \hat{S} = H[S]$, and, by definition, $C_{\mu} = H[S]$. Substituting equals for equals gives us the theorem. QED.

APPENDIX G: RELATIONS TO OTHER FIELDS

1. Time Series Modeling

The goal of time series modeling is to predict the future of a measurement series on the basis of its past. Broadly speaking, this can be divided into two parts: identify equivalent pasts and then produce a prediction for each class of equivalent pasts. That is, we first pick a function $\eta: \stackrel{\leftarrow}{\mathbf{S}} \mapsto \mathcal{R}$ and then pick another function $p: \mathcal{R} \mapsto \overrightarrow{S}$. Of course, we can choose for the range of p futures of some finite length (length 1 is popular) or even choose distributions over these. While practical applications often demand a single definite prediction— "You will meet a tall dark stranger", there are obvious advantages to predicting a distribution—"You have a .95 chance of meeting a tall dark stranger and a .05 chance of meeting a tall familiar albino." Clearly, the best choice for p is the actual conditional distribution of futures for each $\rho \in \mathcal{R}$. Given this, the question becomes what the best \mathcal{R} is; i.e., What is the best η ? At least in the case of trying to understand the whole of the underlying process, we have shown that the best η is, unambiguously, ϵ . Thus, our discussion has implicitly subsumed that of traditional time series modeling.

Computational mechanics—in its focus on letting the process speak for itself through (possibly impoverished) measurements—follows the spirit that motivated one approach to experimentally testing dynamical systems theory. Specifically, it follows in spirit the methods of reconstructing "geometry from a time series" introduced by Refs. [90] and [91]. A closer parallel is found, however, in later work on estimating minimal equations of motion from data series [92].

2. Decision-Theoretic Problems

The classic focus of decision theory is "rules of inductive behavior" [93–95]. The problem is to chose functions from observed data to courses of action that possess desirable properties. This task has obvious affinities to considering the properties of ϵ and its rivals η . We can go further and say that what we have done is consider a decision problem, in which the available actions consist of predictions about the future of the process. The calculation of the optimum rule of behavior in general faces formidable technicalities, such as providing an estimate of the utility of every different course of action under every different hypothesis about the relevant aspects of the world. On the one hand, it is not hard to concoct time-series tasks where the optimal rule of behavior does not use ϵ at all. On the other hand, if we simply aim to predict the process indefinitely far into the future, then because the causal states are minimal sufficient statistics for the distribution of futures (Thm. 2, Remark 4), the optimal rule of behavior will use ϵ .

3. Stochastic Processes

Clearly, the computational mechanics approach to patterns and pattern discovery involves stochastic processes in an intimate and inextricable way. Probabilists have, of course, long been interested in using information-theoretic tools to analyze stochastic processes, particularly their ergodic behavior [59,96–98]. There has also been considerable work in the hidden Markov model and optimal prediction literatures on inferring models of processes from data or from given distributions [10,99–102]. To the best of our knowledge, however, these two approaches have not been previously combined.

Perhaps the closest approach to the spirit of computational mechanics in the stochastic process literature is, surprisingly, the now-classical theory of optimal prediction and filtering for stationary processes, developed by Wiener and Kolmogorov [103–106]. The two theories share the use of information-theoretic notions, the unification of prediction and structure, and the conviction that "the statistical mechanics of time series" is a "field in which conditions are very remote from those of the statistical mechanics of heat engines and which is thus very well suited to serve as a model of what happens in the living organism" [106, p. 59]. So far as we have been able to learn, however, no one has ever used this theory to explicitly identify causal states and causal structure. leaving these implicit in the mathematical form of the prediction and filtering operators. Moreover, the Wiener-Kolmogorov framework forces us to sharply separate the linear and nonlinear aspects of prediction and filtering, because it has a great deal of trouble calculating nonlinear operators [105]. Computational mechanics is completely indifferent to this issue, since it packs all of the process's structure into the ϵ -machine, which is equally calculable in linear or strongly nonlinear situations.

4. Formal Language Theory and Grammatical Inference

A formal language is a set of symbol strings ("words" or "allowed words") drawn from a finite alphabet. Every formal language may be described either by a set of rules (a "grammar") for creating all and only the allowed words, by an abstract automaton which also generates the allowed words and rejects all "forbidden" words. Our ϵ -machines, stripped of probabilities, correspond to such automata—generative in the simple case or classificatory, if we add a reject state and move to it when none of the allowed symbols are encountered.

Since Chomsky [107,108], it has been known that formal languages can be classified into a hierarchy, the higher levels of which have strictly greater expressive power. The hierarchy is defined by restricting the form of the grammatical rules or, equivalently, by limiting the

amount and kind of memory available to the automata. The lowest level of the hierarchy is that of regular languages, which may be familiar to Unix-using readers as regular expressions. These correspond to finite-state machines and to hidden Markov models of finite dimension. In such cases, relatives of our minimality and uniqueness theorems are well known [66], and the construction of causal states is analogous to the "Nerode equivalence classing" procedure [66,109]. Our theorems, however, are not restricted to this low-memory, non-stochastic setting.

The problem of learning a language from observational data has been extensively studied by linguists, and by computer scientists interested in natural-language processing. Unfortunately, well developed learning techniques exist only for the two lowest classes in the Chomsky hierarchy, the regular and the context-free languages. (For a good account of these procedures see Ref. [110].) Adapting and extending this work to the reconstruction of ϵ -machines should form a useful area of future research, a point to which we alluded in the concluding remarks.

5. Computational and Statistical Learning Theory

The goal of computational learning theory [111,112] is to identify algorithms that quickly, reliably, and simply lead to good representations of a target "concept". The latter is typically defined to be a binary dichotomy of a certain feature or input space. Particular attention is paid to results about "probably approximately correct" (PAC) procedures [113]: those having a high probability of finding members of a fixed "representation class" (e.g., neural nets, Boolean functions in disjunctive normal form, and deterministic finite automata). The key word here is "fixed"; as in contemporary time-series analysis, practitioners of this discipline acknowledge the importance of getting the representation class right. (Getting it wrong can make easy problems intractable.) In practice, however, they simply take the representation class as a given, even assuming that we can always count on it having at least one representation which exactly captures the target concept. Although this is in line with implicit assumptions in most of mathematical statistics, it seems dubious when analyzing learning in the real world [5,114,115].

In any case, the preceding development made no such assumption. One of the goals of computational mechanics is, exactly, discovering the best representation. This is not to say that the results of computational learning theory are not remarkably useful and elegant, nor that one should not take every possible advantage of them in implementing ϵ -machine reconstruction. In our view, though, these theories belong more to statistical inference, particularly to algorithmic parameter estimation, than to foundational questions about the nature of pattern and the dynamics of learning.

Finally, in a sense computational mechanics' focus on

causal states is a search for a particular kind of structural decomposition for a process. That decomposition is most directly reflected in the conditional independence of past and future that causal states induce. This decomposition reminds one of the important role that conditional independence plays in contemporary methods for artificial intelligence, both for developing systems that reason in fluctuating environments [116] and the more recently developed algorithmic methods of graphical models [117,118].

6. Description-Length Principles and Universal Coding Theory

Rissanen's minimum description length (MDL) principle, most fully described in Ref. [46], is a procedure for selecting the most concise generative model out of a family of models that are all statistically consistent with given data. The MDL approach starts from Shannon's results on the connection between probability distributions and codes. Rissanen's development follows the inductive framework introduced by Solomonoff [43].

Suppose we choose a representation that leads to a class \mathcal{M} of models and are given data set X. The MDL principle enjoins us to pick the model $M \in \mathcal{M}$ that minimizes the sum of the length of the description of X given M, plus the length of description of M given M. The description length of X is taken to be $-\log P(X|M)$; cf. Eq. (5). The description length of M may be regarded as either given by some coding scheme or, equivalently, by some distribution over the members of M. (Despite the similarities to model estimation in a Bayesian framework [119], Rissanen does not interpret this distribution as a Bayesian prior or regard description length as a measure of evidential support.)

The construction of causal states is somewhat similar to the states estimated in Rissanen's context algorithm [46,120] (and to the "vocabularies" built by universal coding schemes, such as the popular Lempel-Ziv algorithm [121,122]). Despite the similarities, there are significant differences. For a random source—for which there is a single causal state—the context algorithm estimates a number of states that diverges (at least logarithmically) with the length of the data stream, rather than inferring a single state, as ϵ -machine reconstruction would. Moreover, we avoid any reference to encodings of rival models or to prior distributions over them; $C_{\mu}(\mathcal{R})$ is not a description length.

7. Measure Complexity

Ref. [75] proposed that the appropriate measure of the complexity of a process was the "minimal average Shannon information needed" for optimal prediction. This true measure complexity was to be taken as the Shannon

entropy of the states used by some optimal predictor. The same paper suggested that it could be approximated (from below) by the excess entropy; there called the *effective measure complexity*, as noted in Sec. VI above. This is a position closely allied to that of computational mechanics, to Rissanen's MDL principle, and to the minimal embeddings introduced by the "geometry of a time series" methods [90] just described.

In contrast to computational mechanics, however, the key notion of "optimal prediction" was left undefined, as were the nature and construction of the states of the optimal predictor. In fact, the predictors used required knowing the process's underlying equations of motion. Moreover, the statistical complexity $C_{\mu}(\mathcal{S})$ differs from the measure complexities in that it is based on the well defined causal states, whose optimal predictive powers are in turn precisely defined. Thus, computational mechanics is an operational and constructive formalization of the insights expressed in Ref. [75].

8. Hierarchical Scaling Complexity

Introduced in Ref. [123, ch. 9], this approach seeks, like computational mechanics, to extend certain traditional ideas of statistical physics. In brief, the method is to construct a hierarchy of n^{th} -order Markov models and examine the convergence of their predictions with the real distribution of observables as $n \to \infty$. The discrepancy between prediction and reality is, moreover, defined information theoretically, in terms of the relative entropy or Kullback-Leibler distance [62,71]. (We have not used this quantity.) The approach implements Weiss's discovery that for finite-state sources there is a structural distinction between block-Markovian sources (subshifts of finite type) and sofic systems. Weiss showed that, despite their finite memory, sofic systems are the limit of an infinite series of increasingly larger block-Markovian sources [124].

The hierarchical-scaling-complexity approach has several advantages, particularly its ability to handle issues of scaling in a natural way (see Ref. [123, sec. 9.5]). Nonetheless, it does not attain all the goals set in Sec. IIF. Its Markovian predictors are so many black boxes, saying little or nothing about the hidden states of the process, their causal connections, or the intrinsic computation carried on by the process. All of these properties, as we have shown, are manifest from the ϵ machine. We suggest that a productive line of future work would be to investigate the relationship between hierarchical scaling complexity and computational mechanics, and to see whether they can be synthesized. Along these lines, hierarchical scaling complexity reminds us somewhat of hierarchical ϵ -machine reconstruction described in Ref. [5].

9. Continuous Dynamical Computing

Using dynamical systems as computers has become increasingly attractive over the last ten years or so among physicists, computer scientists, and others exploring the physical basis of computation [125–128]. These proposals have ranged from highly abstract ideas about how to embed Turing machines in discrete-time nonlinear continuous maps [7,129] to, more recently, schemes for specialized numerical computation that could in principle be implemented in current hardware [130]. All of them, however, have been synthetic, in the sense that they concern designing dynamical systems that implement a given desired computation or family of computations. In contrast, one of the central questions of computational mechanics is exactly the converse: given a dynamical system, how can one detect what it is intrinsically computing?

We believe that having a mathematical basis and a set of tools for answering this question are important to the synthetic, engineering approach to dynamical computing. Using these tools we may be able to discover, for example, novel forms of computation embedded in natural processes that operate at higher speeds, with less energy, and with fewer physical degrees of freedom than currently possible.

- [1] Julia M. Yeomans. Statistical Mechanics of Phase Transitions. Clarendon Press, Oxford, 1992.
- [2] Paul Manneville. Dissipative Structures and Weak Turbulence. Academic Press, Boston, Massachusetts, 1990.
- [3] P. M. Chaikin and T. C. Lubensky. Principles of Condensed Matter Physics. Cambridge University Press, Cambridge, England, 1995.
- [4] Mark C. Cross and Pierre Hohenberg. Pattern Formation Out of Equilibrium. Reviews of Modern Physics, 65:851–1112, 1993.
- [5] James P. Crutchfield. The calculi of emergence: Computation, dynamics, and induction. *Physica D*, 75:11–54, 1994.
- [6] James P. Crutchfield and Karl Young. Inferring statistical complexity. *Physical Review Letters*, 63:105–108, 1989.
- [7] James P. Crutchfield and Karl Young. Computation at the onset of chaos. In Zurek [131], pages 223–269.
- [8] Nicolás Perry and P.-M. Binder. Finite statistical complexity for sofic systems. *Physical Review E*, 60:459–463, 1999.
- [9] James E. Hanson and James P. Crutchfield. Computational mechanics of cellular automata: An example. Physica D, 103:169–189, 1997.
- [10] Daniel R. Upper. Theory and Algorithms for Hidden Markov Models and Generalized Hidden Markov Models. PhD thesis, University of California, Berkeley, 1997.

- [11] James P. Crutchfield and Melanie Mitchell. The evolution of emergent computation. *Proceedings of the National Academy of Sciences*, 92:10742–10746, 1995.
- [12] A. Witt, A. Neiman, and J. Kurths. Characterizing the dynamics of stochastic bistable systems by measures of complexity. *Physical Review E*, 55:5050–5059, 1997.
- [13] Jordi Delgado and Ricard V. Solé. Collective-induced computation. *Physical Review E*, 55:2338–2344, 1997.
- [14] W. M. Gonçalves, R. D. Pinto, J. C. Sartorelli, and M. J. de Oliveira. Inferring statistical complexity in the dripping faucet experiment. *Physica A*, 257:385–389, 1998.
- [15] James P. Crutchfield and Cosma Rohilla Shalizi. Thermodynamic depth of causal states: Objective complexity via minimal representations. *Physical Review E*, 59:275–283, 1999.
- [16] Jorge Luis Borges. Other Inquisitions, 1937–1952. University of Texas Press, Austin, 1964. Trans. Ruth L. C. Simms.
- [17] James P. Crutchfield. Semantics and thermodynamics. In Martin Casdagli and Stephen Eubank, editors, Non-linear Modeling and Forecasting, volume 12 of Santa Fe Institute Studies in the Sciences of Complexity, pages 317–359, Reading, Massachusetts, 1992. Addison-Wesley.
- [18] Plato. Phaedrus.
- [19] A. R. Luria. The Working Brain: An Introduction to Neuropsychology. Basic Books, New York, 1973.
- [20] Norma Van Surdam Graham. Visual Pattern Analyzers, volume 16 of Oxford Psychology Series. Oxford University Press, Oxford, 1989.
- [21] Sara J. Shettleworth. Cognition, Evolution and Behavior. Oxford University Press, Oxford, 1998.
- [22] Julius T. Tou and Rafael C. Gonzalez. Pattern Recognition Principles. Addison-Wesley, Reading, Mass, 1974.
- [23] Stephen P. Banks. Signal Processing, Image Processing, and Pattern Recognition. Prentice Hall, New York, 1990.
- [24] Jae S. Lim. Two-Dimensional Signal and Image Processing. Prentice Hall, New York, 1990.
- [25] Plato. Meno. In Sec. 80D Meno says: "How will you look for it, Socrates, when you do not know at all what it is? How will you aim to search for something you do not know at all? If you should meet it, how will you know that this is the thing that you did not know?" The same difficulty is raised in Theaetetus, Sec. 197 et sea.
- [26] Alfred North Whitehead and Bertrand Russell. Principia Mathematica. Cambridge University Press, Cambridge, England, 2nd edition, 1925–27.
- [27] Bertrand Russell. Introduction to Mathematical Philosophy. The Muirhead Library of Philosophy. George Allen and Unwin, London, revised edition, 1920. First edition, 1919. Reprinted New York: Dover Books, 1993.
- [28] James P. Crutchfield. Information and its metric. In L. Lam and H. C. Morris, editors, Nonlinear Structures in Physical Systems—Pattern Formation, Chaos and Waves, page 119, New York, 1990. Springer-Verlag.
- [29] Bertrand Russell. Human Knowledge: Its Scope and Limits. Simon and Schuster, New York, 1948.
- [30] John Rhodes. Applications of Automata Theory and Algebra via the Mathematical Theory of Complexity to

- Biology, Physics, Psychology, Philosophy, Games, and Codes. University of California, Berkeley, California, 1971.
- [31] Chrystopher L. Nehaniv and John L. Rhodes. Krohn-Rhodes theory, hierarchies, and evolution. In Boris Mirkin, F. R. McMorris, Fred S. Roberts, and Andrey Rzhetsky, editors, *Mathematical Hierarchies and Biology: DIMACS workshop, November 13–15, 1996*, volume 37 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, Providence, Rhode Island, 1997. American Mathematical Society.
- [32] Ulf Grenander. Elements of Pattern Theory. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, Maryland, 1996.
- [33] Ulf Grenander, Y. Chow, and D. M. Keenan. Hands: A Pattern Theoretic Study of Biological Shapes, volume 2 of Research Notes in Neural Computing. Springer-Verlag, New York, 1991.
- [34] Ulf Grenander and K. Manbeck. A stochastic shape and color model for defect detection in potatoes. American Statistical Association, 2:131–151, 1993.
- [35] A. N. Kolmogorov. Three approaches to the quantitative definition of information. *Problems of Information Transmission*, 1:1–7, 1965.
- [36] Gregory Chaitin. On the length of programs for computing finite binary sequences. *Journal of the Association for Computing Machinery*, 13:547–569, 1966.
- [37] A. N. Kolmogorov. Combinatorial foundations of information theory and the calculus of probabilities. Russ. Math. Surveys, 38:29, 1983.
- [38] Ming Li and Paul M. B. Vitanyi. An Introduction to Kolmogorov Complexity and its Applications. Springer-Verlag, New York, 1993.
- [39] Marvin Minsky. Computation: Finite and Infinite Machines. Prentice-Hall, Englewood Cliffs, New Jersey, 1967.
- [40] P. Martin-Löf. The definition of random sequences. Information and Control, 9:602–619, 1966.
- [41] L. A. Levin. Laws of information conservation (non-growth) and aspects of the foundation of probability theory. Problemy Peredachi Informatsii, 10:30–35, 1974. Translation: Problems of Information Transmission 10 (1974) 206–210.
- [42] V. G. Gurzadyan. Kolmogorov complexity as a descriptor of cosmic microwave background maps. *Europhysics Letters*, 46:114–117, 1999. Also available as an electronic preprint, LANL archive, astro-phy/9902123.
- [43] Raymond J. Solomonoff. A formal theory of inductive inference. *Information and Control*, 7:1–22 and 224–254, 1964.
- [44] Paul Vitányi and Ming Li. Minimum description length induction, Bayesianism, and Kolmogorov complexity. Electronic pre-print, LANL Archive, cs.LG/9901014, 1999.
- [45] Gary William Flake. The Computational Beauty of Nature: Computer Explorations of Fractals, Chaos, Complex Systems and Adaptation. MIT Press, Cambridge, Massachusetts, 1998.
- [46] Jorma Rissanen. Stochastic Complexity in Statistical Inquiry. World Scientific, Singapore, 1989.

- [47] Charles H. Bennett. How to define complexity in physics, and why. In Zurek [131], pages 137–148.
- [48] Moshe Koppel. Complexity, depth, and sophistication. Complex Systems, 1:1087–1091, 1987.
- [49] Moshe Koppel and Henri Atlan. An almost machineindependent theory of program-length complexity, sophistication and induction. *Information Sciences*, 56:23–44, 1991.
- [50] Daniel C. Dennett. Real patterns. Journal of Philosophy, 88:27–51, 1991. Reprinted in Dennett (1997).
- [51] James P. Crutchfield. Is anything ever new? Considering emergence. In G. Cowan, D. Pines, and D. Melzner, editors, Complexity: Metaphors, Models, and Reality, volume 19 of Santa Fe Institute Studies in the Sciences of Complexity, pages 479–497, Reading, Massachusetts, 1994. Addison-Wesley.
- [52] John H. Holland. Emergence: From Chaos to Order. Addison-Wesley, Reading, Massachusetts, 1998.
- [53] Ludwig Boltzmann. Lectures on Gas Theory. University of California Press, Berkeley, 1964.
- [54] Harald Cramér. Mathematical Methods of Statistics. Almqvist and Wiksells, Uppsala, 1945. Republished by Princeton University Press, 1946, as vol. 9 in the Princeton Mathematics Series, and as a paperback, in the Princeton Landmarks in Mathematics and Physics series, 1999.
- [55] Claude E. Shannon. A mathematical theory of communication. Bell System Technical Journal, 27:379–423, 1948.
- [56] David Hume. A Treatise of Human Nature: Being an Attempt to Introduce the Experimental Method of Reasoning into Moral Subjects. John Noon, London, 1739. Reprint (Oxford: Clarendon Press, 1951) of original edition, with notes and analytical index.
- [57] Mario Bunge. Causality: The Place of the Causal Princple in Modern Science. Harvard University Press, Cambridge, Massachusetts, 1959. Reprinted as Causality and Modern Science, NY: Dover Books, 1979.
- [58] Wesley C. Salmon. Scientific Explanation and the Causal Structure of the World. Princeton University Press, Princeton, 1984.
- [59] Patrick Billingsley. Ergodic Theory and Information. Tracts on Probablity and Mathematical Statistics. Wiley, New York, 1965.
- [60] Patrick Billingsley. Probability and Measure. Wiley Series in Probability and Mathematical Statistics. Wiley, New York, 1979.
- [61] Bernard F. Schutz. Geometrical Methods of Mathematical Physics. Cambridge University Press, Cambridge, England, 1980.
- [62] Thomas M. Cover and Joy A. Thomas. Elements of Information Theory. Wiley, New York, 1991.
- [63] William of Ockham. Philosophical Writings: A Selection, Translated, with an Introduction, by Philotheus Boehner, O.F.M., Late Professor of Philosophy, The Franciscan Institute. Bobbs-Merrill, Indianapolis, 1964. first pub. various European cities, early 1300s.
- [64] Anonymous. Kuan Yin Tzu, T'ang Dynasty. Written in China during the T'ang dynasty. Partial translation in Joseph Needham, *Science and Civilisation in China*,

- vol. II (Cambridge University Press, 1956), p. 73.
- [65] David P. Feldman and James P. Crutchfield. Discovering non-critical organization: Statistical mechanical, information theoretic, and computational views of patterns in simple one-dimensional spin systems. *Journal of Statistical Physics*, submitted, 1998. Santa Fe Institute Working Paper 98-04-026, http://www.santafe.edu/projects/CompMech/papers/DNCO.html.
- [66] John E. Hopcroft and Jeffrey D. Ullman. Introduction to Automata Theory, Languages, and Computation. Addison-Wesley, Reading, 1979. 2nd edition of Formal Languages and Their Relation to Automata, 1969.
- [67] John G. Kemeny and J. Laurie Snell. Finite Markov Chains. Springer-Verlag, New York, 1976.
- [68] John G. Kemeny, J. Laurie Snell, and Anthony W. Knapp. Denumerable Markov Chains. Springer-Verlag, New York, 2nd edition, 1976.
- [69] James E. Hanson. Computational Mechanics of Cellular Automata. PhD thesis, University of California, Berkeley, 1993.
- [70] Gregory Bateson. Mind and Nature: A Necessary Unity.E. P. Dutton, New York, 1979.
- [71] Solomon Kullback. Information Theory and Statistics. Dover Books, New York, 2nd edition, 1968. First edition New York: Wiley, 1959.
- [72] Claude Bernard. Introduction a l'etude de la medecine experimentale. J. B. Bailliere, Paris, 1865. Trans. by Henry Copley Green as Introduction to the Study of Experimental Medicine, New York: Macmillian, 1927; reprinted New York: Dover, 1957.
- [73] James P. Crutchfield and Norman H. Packard. Symbolic dynamics of noisy chaos. *Physica D*, 7:201–223, 1983.
- [74] Robert Shaw. The Dripping Faucet as a Model Chaotic System. Aerial Press, Santa Cruz, California, 1984.
- [75] Peter Grassberger. Toward a quantitative theory of selfgenerated complexity. *International Journal of Theoret*ical Physics, 25:907–938, 1986.
- [76] Kristian Lindgren and Mats G. Nordahl. Complexity measures and cellular automata. *Complex Systems*, 2:409–440, 1988.
- [77] W. Li. On the relationship between complexity and entropy for Markov chains and regular languages. *Complex Systems*, 5:381–399, 1991.
- [78] Dirk Arnold. Information-theoretic analysis of phase transitions. Complex Systems, 10:143–155, 1996.
- [79] William Bialek and Naftali Tishby. Predictive information. Electronic pre-print, LANL archive, cond-mat/9902341, 1999.
- [80] James P. Crutchfield and David P. Feldman. Statistical complexity of simple one-dimensional spin systems. *Physical Review E*, 55:1239R–1243R, 1997.
- [81] W. Ross Ashby. An Introduction to Cybernetics. Chapman and Hall, London, 1956.
- [82] Hugo Touchette and Seth Lloyd. Information-theoretic limits of control. *Physical Review Letters*, 84:1156–1159, 1999.
- [83] Abraham Lempel and Jacob Ziv. Compression of twodimensional data. *IEEE Transactions in Information* Theory, IT-32:2–8, 1986.

- [84] David P. Feldman. Computational Mechanics of Classical Spin Systems. PhD thesis, University of California, Davis, 1998. Available on-line at http://hornacek.coa.edu/dave/Thesis/thesis.html.
- [85] Deborah Mayo. Error and the Growth of Experimental Knowledge. Science and Its Conceptual Foundations. University of Chicago Press, Chicago, 1996.
- [86] James P. Crutchfield and Cristopher Douglas. Imagined complexity: Learning a random process. in preparation, 1999.
- [87] Rudolf Lidl and Gunter Pilz. Applied Abstract Algebra. Springer, New York, 1984.
- [88] E. S. Ljapin. Semigroups, volume 3 of Translations of Mathematical Monographs. American Mathematical Society, Providence, Rhode Island, 1963.
- [89] Karl Young. The Grammar and Statistical Mechanics of Complex Physical Systems. PhD thesis, University of California, Santa Cruz, 1991.
- [90] Norman H. Packard, James P. Crutchfield, J. Doyne Farmer, and Robert S. Shaw. Geometry from a time series. *Physical Review Letters*, 45:712–716, 1980.
- [91] Floris Takens. Detecting strange attractors in fluid turbulence. In D. A. Rand and L. S. Young, editors, Symposium on Dynamical Systems and Turbulence, volume 898 of Lecture Notes in Mathematics, page 366, Berlin, 1981. Springer-Verlag.
- [92] James P. Crutchfield and Bruce S. McNamara. Equations of motion from a data series. *Complex Systems*, 1:417–452, 1987.
- [93] Jerzy Neyman. First Course in Probability and Statistics. Henry Holt, New York, 1950.
- [94] David Blackwell and M. A. Girshick. Theory of Games and Statistical Decisions. Wiley, New York, 1954. Reprinted New York: Dover Books, 1979.
- [95] R. Duncan Luce and Howard Raiffa. Games and Decisions: Introduction and Critical Survey. Wiley, New York, 1957.
- [96] I. M. Gel'fand and A. M. Yaglom. Calculation of the amount of information about a random function contained in another such function. *Uspekhi Matematich*eski Nauk, 12:3–52, 1956. Trans. in *American Math*ematical Society Translations, 2nd series, 12 (1959): 199–246.
- [97] Peter E. Caines. *Linear Stochastic Systems*. Wiley, New York, 1988.
- [98] Robert M. Gray. Entropy and Information Theory. Springer-Verlag, New York, 1990.
- [99] David Blackwell and Lambert Koopmans. On the identifiability problem for functions of finite Markov chains. Annals of Mathematical Statistics, 28:1011–1015, 1957.
- [100] H. Ito, S.-I. Amari, and K. Kobayashi. Identifiability of hidden Markov information sources and their minimum degrees of freedom. *IEEE Transactions on Information Theory*, 38:324–333, 1992.
- [101] H. Jaeger. Observable operator models for discrete stochastic time series. Neural Computation, forthcoming, 1999. ftp://ftp.gmd.de/GMD/ais/ publications/1999/.
- [102] Paul Algoet. Universal schemes for prediction, gambling and portfolio selection. The Annals of Probability,

- 20:901–941, 1992. See also an important Correction, The Annals of Probability, **23** (1995): 474–478.
- [103] A. N. Kolmogorov. Interpolation und extrapolation von stationären zufälligen folgen. Bull. Acad. Sci. U.S.S.R., Math., 3:3–14, 1941. In German.
- [104] Norbert Wiener. Extrapolation, Interpolation, and Smoothing of Stationary Time Series: With Engineering Applications. The Technology Press of the Massachusetts Institute of Technology, Cambridge, Massachusetts, 1949. "First published during the war as a classifed report to Section D₂, National Defense Research Council".
- [105] Norbert Wiener. Nonlinear Problems in Random Theory. The Technology Press of the Massachusetts Institute of Technology, Cambridge, Massachusetts, 1958.
- [106] Norbert Wiener. Cybernetics: Or, Control and Communication in the Animal and the Machine. MIT Press, Cambridge, Massachusetts, 2nd edition, 1961. First edition New York: Wiley, 1948.
- [107] Noam Chomsky. Three models for the description of language. IRE Transactions on Information Theory, 2:113, 1956
- [108] Noam Chomsky. Syntactic Structures, volume 4 of Janua linguarum, series minor. Mouton, The Hauge, 1957.
- [109] B. A. Trakhtenbrot and Ya. M. Barzdin. Finite Automata. North-Holland, Amsterdam, 1973.
- [110] Eugene Charniak. Statistical Language Learning. Language, Speech and Communication. MIT Press, Cambridge, Massachusetts, 1993.
- [111] Michael J. Kearns and Umesh V. Vazirani. An Introduction to Computational Learning Theory. MIT Press, Cambridge, Massachusetts, 1994.
- [112] V. N. Vapnik. The Nature of Statistical Learning Theory. Springer-Verlag, Berlin, 2nd edition, 2000.
- [113] Leslie G. Valiant. A theory of the learnable. Communications of the Association for Computing Machinery, 27:1134–1142, 1984.
- [114] Margaret A. Boden. Precis of The Creative Mind: Myths and Mechanisms. Behaviorial and Brain Sciences, 17:519–531, 1994.
- [115] Chris Thornton. Truth from Trash: How Learning Makes Sense. Complex Adaptive Systems. MIT Press, Cambridge, Massachusetts, 2000.
- [116] Judea Pearl. Causality: Models, Reasoning, and Inference. Cambridge University Press, Cambridge, England, 2000.
- [117] M. I. Jordan, editor. Learning in Graphical Models, volume 89 of NATO Science Series D: Behavioral and Social Sciences, Dordrecht, 1998.
- [118] Peter Spirtes, Clark Glymour, and Richard Scheines. Causation, Prediction, and Search. Adaptive Computation and Machine Learning. MIT Press, Cambridge, Massachusetts, 2000.
- [119] David V. Lindley. Bayesian Statistics, a Review. Society for Industrial and Applied Mathematics, Philadelphia, 1972.
- [120] Jorma Rissanen. Universal coding, information, prediction, and estimation. IEEE Transactions in Information Theory, IT-30:629–636, 1984.

- [121] Abraham Lempel and Jacob Ziv. On the complexity of finite sequences. *IEEE Transactions in Information Theory*, IT-22:75–81, 1976.
- [122] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *IEEE Transactions in Information Theory*, IT-23:337–343, 1977.
- [123] Remo Badii and Antonio Politi. Complexity: Hierarchical Structures and Scaling in Physics, volume 6 of Cambridge Nonlinear Science Series. Cambridge University Press, Cambridge, 1997.
- [124] Benjamin Weiss. Subshifts of finite type and sofic systems. Monatshefte für Mathematik, 77:462–474, 1973.
- [125] Cristopher Moore. Recursion theory on the reals and continuous-time computation. Theoretical Computer Science, 162:23–44, 1996.
- [126] Cristopher Moore. Dynamical recognizers: Real-time language recognition by analog computers. *Theoretical Computer Science*, 201:99–136, 1998.
- [127] Pekka Orponen. A survey of continuous-time computation theory. In D.-Z. Du and K.-I Ko, editors, Advances in Algorithms, Languages, and Complexity, pages 209– 224. Kluwer Academic, Dordrecht, 1997.
- [128] Lenore Blum, Michael Shub, and Steven Smale. On a theory of computation and complexity over the real numbers: NP-completeness, recursive functions and universal machines. *Bulletin of the American Mathe*matical Society, 21:1–46, 1989.
- [129] Cristopher Moore. Unpredictability and undecidability in dynamical systems. *Physical Review Letters*, 64:2354– 2357, 1990.
- [130] Sudeshna Sinha and William L. Ditto. Dynamics based computation. *Physical Review Letters*, 81:2156–2159, 1998.
- [131] Wojciech H. Zurek, editor. Complexity, Entropy, and the Physics of Information, volume 8 of Santa Fe Institute Studies in the Sciences of Complexity, Reading, Massachusetts, 1990. Addison-Wesley.

APPENDIX: GLOSSARY OF NOTATION

In the order of their introduction.

Symbol	Description	Where Introduced
\mathcal{O}	Object in which we wish to find a pattern	Sec. II, p. 3
\mathcal{P}	Pattern in \mathcal{O}	Sec. II, p. 3
${\mathcal A}_{\leftrightarrow}$	Countable alphabet	Sec. III A, p. 6
$S \longrightarrow$	Bi-infinite, stationary, discrete stochastic process on \mathcal{A}	Def. 1, p. 6
$\stackrel{\searrow}{s}$	Particular realization of S	Def. 1, p. 6
$\begin{array}{c} \mathcal{A} \\ \stackrel{\hookrightarrow}{S} \\ \stackrel{\hookrightarrow}{S} \\ \stackrel{\rightarrow}{S}^{L} \\ \stackrel{\rightarrow}{S} \\ \stackrel{\downarrow}{S} \\ \stackrel{\leftarrow}{S} \\ \stackrel{\hookrightarrow}{S} \\ \stackrel{\hookrightarrow}{S} \\ \stackrel{\leftarrow}{S} \\ \stackrel{\leftarrow}{S} \\ \stackrel{\leftarrow}{\lambda} \\ \end{array}$	Random variable for the next L values of $\overset{\leftrightarrow}{S}$	Sec. III A, p. 6
$\overset{ ightarrow}{S}^{L}$	Particular value of \overrightarrow{S}	Sec. III A, p. 6
$\overrightarrow{S}_{\tau}^{1}$	Next observable generated by $\overset{\leftrightarrow}{S}$	Sec. III A, p. 6
$\overset{\leftarrow}{S}^L$	As \overrightarrow{S}^L , but for the last L values, up to the present	Sec. III A, p. 6
$\overset{\leftarrow}{s}^L$	Particular value of $\overset{\leftarrow}{S}^L$	Sec. III A, p. 6
\overrightarrow{S}	Semi-infinite future half of $\overset{\leftrightarrow}{S}$	Sec. III A, p. 6
$\stackrel{ ightarrow}{s}$	Particular value of \overrightarrow{S}	Sec. III A, p. 6
$\stackrel{\leftarrow}{S}$	Semi-infinite past half of $\overset{\leftrightarrow}{S}$	Sec. III A, p. 6
$\stackrel{\leftarrow}{s}$	Particular value of \overleftarrow{S}	Sec. III A, p. 6
λ	Null string or null symbol	Sec. III A, p. 6
$\overset{\leftarrow}{\mathbf{S}}$	Set of all pasts realized by the process S	Sec. IIIB, p. 6
${\cal R}$	Partition of S into effective states	Sec. IIIB, p. 6
ho	Member-class of \mathcal{R} ; a particular effective state	Sec. III B, p. 6
η	Function from S to R	Sec. III B, Eq. (4), p. 6
$\mathcal{R} \ \mathcal{R}'$	Current effective (η) state, as a random variable Next effective state, as a random variable	Sec. III B, p. 6 Sec. III B, p. 6
H[X]	Entropy of the random variable X	Sec. III C 1, p. 7
H[X,Y]	Joint entropy of the random variables X and Y	Sec. III C 2, p. 7
H[X Y]	Entropy of X conditioned on Y	Sec. III C 2, p. 7
I[X;Y]	Mutual information of X and Y	Sec. III C 3, p. 7
$h_{\mu}[\overrightarrow{S}]$	Entropy rate of \overrightarrow{S}	Sec. III D, Eq. (9), p. 8
$h_{\mu}[\overrightarrow{S} X]$	Entropy rate of \overrightarrow{S} conditioned on X	Sec. III D, Eq. (10), p. 8
$C_{\mu}(\mathcal{R})$	Statistical complexity of \mathcal{R}	Def. 4, p. 8
S	Set of the causal states of S Particular causal state	Def. 5, p. 9
$rac{\sigma}{\epsilon}$	Function from histories to causal states	Def. 5, p. 9 Def. 5, p. 9
$\overset{\circ}{\mathcal{S}}$	Current causal state, as a random variable	Def. 5, p. 9
\mathcal{S}'	Next causal state, as a random variable	Def. 5, p. 9
\sim_{ϵ}	Relation of causal equivalence between two histories	Sec. IV A, p. 9
$T_{ij}^{(s)}$	Probability of going from causal state i to j , emitting s	
$\mathcal{R}_{\hat{A}}$	Set of prescient rival states	Def. 11, p. 14
$\hat{ ho}$	Particular prescient rival state	Def. 11, p. 14
$T_{ij}^{(s)} \ \hat{oldsymbol{\mathcal{R}}} \ \hat{oldsymbol{\mathcal{R}}} \ \hat{oldsymbol{\mathcal{R}}} \ \hat{oldsymbol{\mathcal{R}}}'$	Current prescient rival state, as a random variable	Def. 11, p. 14
$C_{\mu}(\mathcal{O})$	Next prescient rival state, as a random variable Statistical complexity of the process \mathcal{O}	Def. 11, p. 14 Def. 12, p. 15
C_{μ}	Without an argument, short for $C_{\mu}(\mathcal{O})$	Def. 12, p. 15
\mathbf{E}^{μ}	Excess entropy	Def. 13, p. 16