

Critical Behavior from Deep Dynamics: A Hidden Dimension in Natural Language

Henry W. Lin and Max Tegmark

*Dept. of Physics, Harvard University, Cambridge, MA 02138 and
Dept. of Physics & MIT Kavli Institute, Massachusetts Institute of Technology, Cambridge, MA 02139*

(Dated: July 12, 2016)

We show that in many data sequences — from texts in different languages to melodies and genomes — the mutual information between two symbols decays roughly like a power law with the number of symbols in between the two. In contrast, we prove that Markov/hidden Markov processes generically exhibit exponential decay in their mutual information, which explains why natural languages are poorly approximated by Markov processes. We present a broad class of models that naturally reproduce this critical behavior. They all involve deep dynamics of a recursive nature, as can be approximately implemented by tree-like or recurrent deep neural networks. This model class captures the essence of probabilistic context-free grammars as well as recursive self-reproduction in physical phenomena such as turbulence and cosmological inflation. We derive an analytic formula for the asymptotic power law and elucidate our results in a statistical physics context: 1-dimensional “shallow” models (such as Markov models or regular grammars) will fail to model natural language, because they cannot exhibit criticality, whereas “deep” models with one or more “hidden” dimensions representing levels of abstraction or scale can potentially succeed.

I. INTRODUCTION

Critical behavior, where long-range correlations decay as a power law with distance, has many important physics applications ranging from phase transitions in condensed matter experiments to turbulence and inflationary fluctuations in our early Universe. It has important applications beyond the traditional purview of physics as well [1–5] including new results which we report in Figure I: the number of bits of information provided by a symbol about another drops roughly as a power-law¹ with distance in sequences as diverse as the human genome, music by Bach, and text in English and French. Why is this, when so many other correlations in nature instead drop exponentially [9]?

Better understanding such statistical properties of natural languages (in the broad sense of information-transmitting sequences) is interesting not only for geneticists, musicologists and linguists, but also for the machine learning community. Consider how your phone can auto-correct your typing, how you can free up disk space with data compression software and how speech-to-text conversion enables you to talk to digital personal assistants such as Siri, Cortana and Google Now: these technolo-

gies all exploit statistical properties of language, and can all be further improved if we can better understand these properties. Such deepened understanding is the goal of the present paper, focusing on critical behavior.

Natural languages are difficult for machines to understand. This has been known at least as far back as Turing, whose eponymous test [14] relies upon this key fact. A tempting explanation is that natural language is something uniquely human. But this is far from a satisfactory one, especially given the recent successes of machines at performing tasks as complex and as “human” as playing *Jeopardy!* [15], chess [16], Atari games [17] and Go [18]. We will show that computer descriptions of language tend to suffer from a much simpler problem that has nothing to do with meaning, understanding or being non-human: they tend to get the basic statistical properties wrong. We will prove that Markov processes, the workhorse of modeling any sequential data with translational symmetry and one of the handful of models that are analytically tractable, fail epically by predicting exponentially decaying mutual information. On the other hand, impressive progress has been made by using deep neural networks for natural language processing (see, e.g., [19–22]); for recent reviews of deep neural networks, see [23, 24]. Unfortunately, unlike Markov and related n -gram models, these deep networks are often treated like inscrutable black boxes, given their enormous complexity. This has triggered many recent efforts to understand their advantages analytically [25–28], from a functional [29], topological [30], and geometric [31, 32] perspective; this paper explores the advantages from a statistical physics perspective [33].

We will see that a key reason that currently popular recurrent neural networks with long-short-term memory (LSTM) [34] do much better is that they can replicate critical behavior, but that even they can be further improved, since they can under-predict long-range mutual

¹ The power law discussed here should not be confused with another famous power law that occurs in natural languages: Zipf’s law [6]. Zipf’s law implies power law behavior in one-point statistics (in the histogram of word frequencies), whereas we are interested in two-point statistics. In the former case, the power law is in the frequency of words; in the latter case, the power law is in the separation between characters. One can easily cook up sequences which obey Zipf’s law but are not critical and do not exhibit a power law in the mutual information. However, there are models of certain physical systems where Zipf’s law follows from criticality [7, 8].

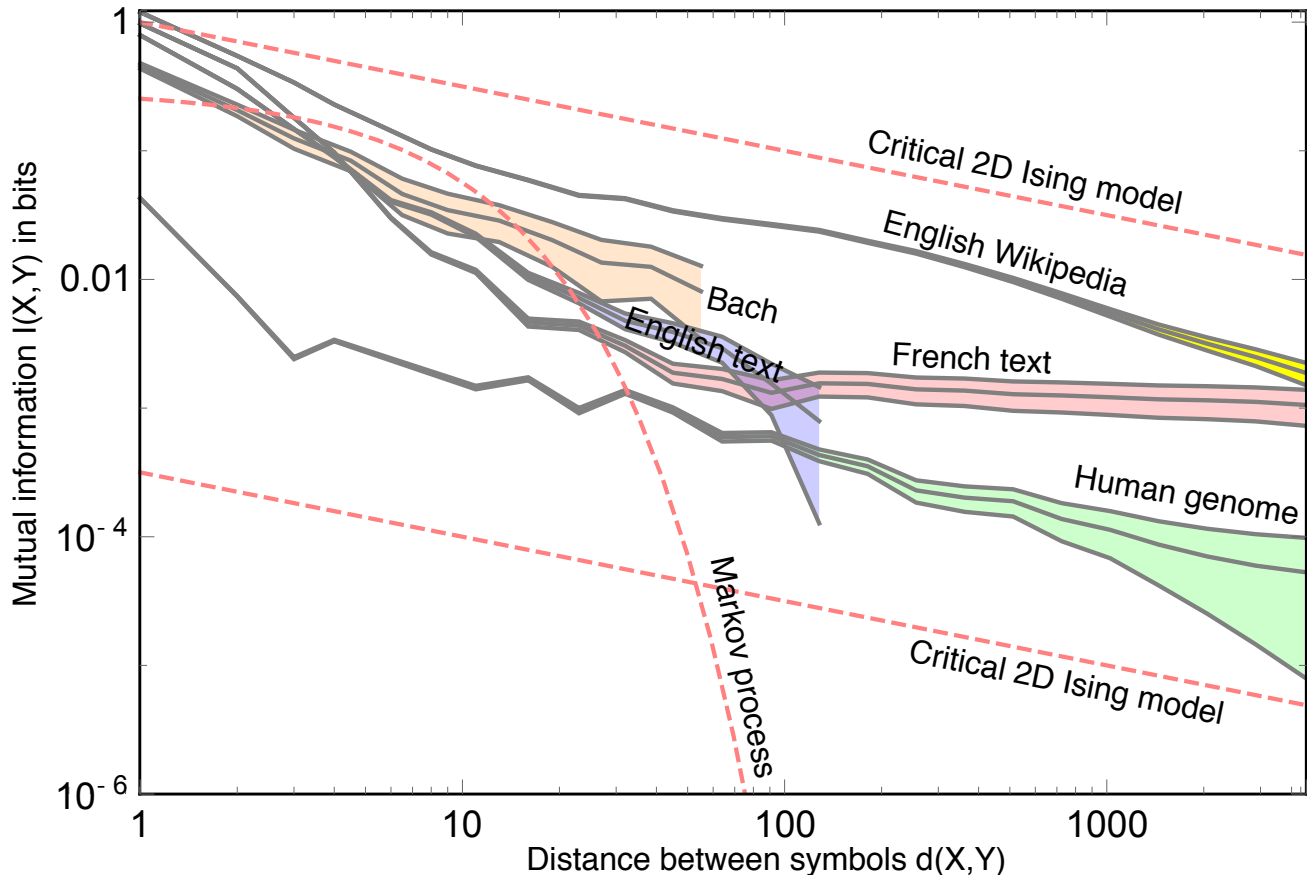


FIG. 1: Decay of mutual information with separation. Here the mutual information in bits per symbol is shown as a function of separation $d(X, Y) = |i - j|$, where the symbols X and Y are located at positions i and j in the sequence in question, and shaded bands correspond to $1 - \sigma$ error bars. All measured curves are seen to decay roughly as power laws, explaining why they cannot be accurately modeled as Markov processes — for which the mutual information instead plummets exponentially (the example shown has $I \propto e^{-d/6}$). The measured curves are seen to be qualitatively similar to that of a famous critical system in physics: a 1D slice through a critical 2D Ising model, where the slope is $-1/2$. The human genome data consists of 177,696,512 base pairs $\{A, C, T, G\}$ from chromosome 5 from the National Center for Biotechnology Information [10], with unknown base pairs omitted. The Bach data consists of 5727 notes from Partita No. 2 [11], with all notes mapped into a 12-symbol alphabet consisting of the 12 half-tones $\{C, C\#, D, D\#, E, F, F\#, G, G\#, A, A\#, B\}$, with all timing, volume and octave information discarded. The three text corpuses are 100 MB from Wikipedia [12] (206 symbols), the first 114 MB of a French corpus [13] (185 symbols) and 27 MB of English articles from `slate.com` (143 symbols). The large long range information appears to be dominated by poems in the French sample and by html-like syntax in the Wikipedia sample.

information.

A final goal of this paper is to ameliorate the following problem: machine learning typically involves using something we do not fully understand (neural nets, *etc.*) to study something we also do not fully understand (English, *etc.*). If we are ever to understand how some learning algorithm works, we must first understand what we are trying to learn. For this reason, we will construct a simple class of analytically tractable models which qualitatively reproduce (some of) the statistics of natural languages — specifically, critical behavior.

This paper is organized as follows. In Section II, we show how Markov processes exhibit exponential decay in mutual information with scale; we give a rigorous proof

of this and other results in a series of appendices. To enable such proofs, we introduce a convenient quantity that we term *rational mutual information*, which bounds the mutual information and converges to it in the near-independence limit. In Section III, we define a subclass of generative grammars and show that they exhibit critical behavior with power law decays. We then generalize our discussion using Bayesian nets and relate our findings to theorems in statistical physics. In Section IV, we discuss our results and explain how LSTM RNNs can reproduce critical behavior by emulating our generative grammar model.

II. MARKOV IMPLIES EXPONENTIAL DECAY

For two random variables X and Y , the following definitions of mutual information are all equivalent:

$$\begin{aligned} I(X, Y) &\equiv S(X) + S(Y) - S(X, Y) \\ &= D(p(XY) || p(X)p(Y)) \\ &= \left\langle \log_B \frac{P(a, b)}{P(a)P(b)} \right\rangle \\ &= \sum_{ab} P(a, b) \log_B \frac{P(a, b)}{P(a)P(b)}, \end{aligned} \quad (1)$$

where $S \equiv \langle -\log_B P \rangle$ is the Shannon entropy [35] and $D(p(XY) || p(X)p(Y))$ is the Kullback-Leibler divergence [36] between the joint probability distribution and the product of the individual marginals. If the base of the logarithm is taken to be $B = 2$, then $I(X, Y)$ is measured in bits. The mutual information can be interpreted as how much one variable knows about the other: $I(X, Y)$ is the reduction in the number of bits needed to specify for X once Y is specified. Equivalently, it is the number of encoding bits saved by using the true joint probability $P(X, Y)$ instead of approximating X and Y are independent. It is thus a measure of statistical dependencies between X and Y . Although it is more conventional to measure quantities such as the correlation coefficient ρ in statistics and statistical physics, the mutual information is more suitable for generic data, since it does not require that the variables X and Y are numbers or have any algebraic structure, whereas ρ requires that we are able to multiply $X \cdot Y$ and average. Whereas it makes sense to multiply numbers, is meaningless to multiply or average two characters such as “!” and “?”.

The rest of this paper is largely a study of the mutual information between two random variables that are realizations of a discrete stochastic process, with some separation in τ in time. More concretely, we can think of sequences $\{X_1, X_2, X_3, \dots\}$ of random variables, where each one might take values from some finite alphabet. For example, if we model English as a discrete stochastic process and take $\tau = 2$, X could represent the first character (“F”) in this sentence, whereas Y could represent the third character (“r”) in this sentence.

In particular, we start by studying the mutual information function of a Markov process, which is analytically tractable. Let us briefly recapitulate some basic facts about Markov processes (see, *e.g.*, [37] for a pedagogical review). A Markov process is defined by a matrix \mathbf{M} of conditional probabilities $M_{ab} = P(X_{t+1} = a | X_t = b)$. Such Markov matrices (also known as stochastic matrices) thus have the properties $M_{ab} \geq 0$ and $\sum_a M_{ab} = 1$. They fully specify the dynamics of the model:

$$\mathbf{p}_{t+1} = \mathbf{M} \mathbf{p}_t, \quad (2)$$

where \mathbf{p}_t is a vector with components $P(X_t = a)$ that specifies the probability distribution at time t . Let λ_i

denote the eigenvalues of \mathbf{M} , sorted by decreasing magnitude: $|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \dots$. All Markov matrices have $|\lambda_i| \leq 1$, which is why blowup is avoided when equation (2) is iterated, and $\lambda_1 = 1$, with the corresponding eigenvector giving a stationary probability distribution $\boldsymbol{\mu}$ satisfying $\mathbf{M}\boldsymbol{\mu} = \boldsymbol{\mu}$.

In addition, two mild conditions are usually imposed on Markov matrices: \mathbf{M} is *irreducible*, meaning that every state is accessible from every other state (otherwise, we could decompose the Markov process into separate Markov processes). Second, to avoid processes like $1 \rightarrow 2 \rightarrow 1 \rightarrow 2 \dots$ that will never converge, we take the Markov process to be *aperiodic*. It is easy to show using the Perron-Frobenius theorem that being irreducible and aperiodic implies $|\lambda_2| < 1$.

This section is devoted to the intuition behind the following theorem, whose full proof is given in Appendix A and B. The theorem states roughly that for a Markov process, the mutual information between two points in time t_1 and t_2 decays exponentially for large separation $|t_2 - t_1|$:

Theorem 1: Let \mathbf{M} be a Markov matrix that generates a Markov process. If \mathbf{M} is irreducible and aperiodic, then the asymptotic behavior of the mutual information $I(t_1, t_2)$ is exponential decay toward zero for $|t_2 - t_1| \gg 1$ with decay timescale $\log \frac{1}{|\lambda_2|}$, where λ_2 is the second largest eigenvalue of \mathbf{M} . If \mathbf{M} is reducible or periodic, I can instead decay to a constant; no Markov process whatsoever can produce power-law decay. Suppose \mathbf{M} is irreducible and aperiodic so that $\mathbf{p}_t \rightarrow \boldsymbol{\mu}$ as $t \rightarrow \infty$ as mentioned above. This convergence of one-point statistics, *e.g.*, \mathbf{p}_t , has been well-studied [37]. However, one can also study higher order statistics such as the joint probability distribution for two points in time. For succinctness, let us write $P(a, b) \equiv P(X = a, Y = b)$, where $X = X_{t_1}$ and $Y = X_{t_2}$ and $\tau \equiv |t_2 - t_1|$. We are interested in the asymptotic situation where the Markov process has converged to its steady state, so the marginal distribution $P(a) \equiv \sum_b P(a, b) = \mu_a$, independently of time.

If the joint probability distribution approximately factorizes as $P(a, b) \approx \mu_a \mu_b$ for sufficiently large and well-separated times t_1 and t_2 (as we will soon prove), the mutual information will be small. We can therefore Taylor expand the logarithm from equation (1) around the point $P(a, b) = P(a)P(b)$, giving

$$\begin{aligned} I(X, Y) &= \left\langle \log_B \left(\frac{P(a, b)}{P(a)P(b)} \right) \right\rangle \\ &= \left\langle \log_B \left[1 + \frac{P(a, b)}{P(a)P(b)} - 1 \right] \right\rangle \\ &\approx \left\langle \frac{P(a, b)}{P(a)P(b)} - 1 \right\rangle \frac{1}{\ln B} = \frac{I_R(X, Y)}{\ln B}, \end{aligned} \quad (3)$$

where we have defined the *rational mutual information*

$$I_R \equiv \left\langle \frac{P(a,b)}{P(a)P(b)} - 1 \right\rangle. \quad (4)$$

For comparing the rational mutual information with the usual mutual information, it will be convenient to take e as the base B of the logarithm. We derive useful properties of the rational mutual information in Appendix A. To mention just one, we note that the rational mutual information is not just asymptotically equal to the mutual information in the limit of near-independence, but it also provides a strict upper bound on it: $0 \leq I \leq I_R$.

Let us without loss of generality take $t_2 > t_1$. Then iterating equation (2) τ times gives $P(b|a) = (\mathbf{M}^\tau)_{ba}$. Since $P(a,b) = P(a)P(b|a)$, we obtain

$$\begin{aligned} I_R + 1 &= \left\langle \frac{P(a,b)}{P(a)P(b)} \right\rangle = \sum_{ab} P(a,b) \frac{P(a,b)}{P(a)P(b)} \\ &= \sum_{ab} \frac{P(b|a)^2 P(a)^2}{P(a)P(b)} = \sum_{ab} \frac{\mu_a}{\mu_b} [(\mathbf{M}^\tau)_{ba}]^2. \end{aligned}$$

We will continue the proof by considering the typical case where the eigenvalues of \mathbf{M} are all distinct (non-degenerate) and the Markov matrix is irreducible and aperiodic; we will generalize to the other cases (which form a set of measure zero) in Appendix B. Since the eigenvalues are distinct, we can diagonalize \mathbf{M} by writing

$$\mathbf{M} = \mathbf{BDB}^{-1} \quad (5)$$

for some invertible matrix \mathbf{B} and some a diagonal matrix \mathbf{D} whose diagonal elements are the eigenvalues: $D_{ii} = \lambda_i$. Raising equation (5) to the power τ gives $\mathbf{M}^\tau = \mathbf{BD}^\tau \mathbf{B}^{-1}$, *i.e.*,

$$(\mathbf{M}^\tau)_{ba} = \sum_c \lambda_c^\tau \mathbf{B}_{bc} (\mathbf{B}^{-1})_{ca}. \quad (6)$$

Since \mathbf{M} is non-degenerate, irreducible and aperiodic, $1 = \lambda_1 > |\lambda_2| > \dots > |\lambda_n|$, so all terms except the first in the sum of equation (6) decay exponentially with τ , at a decay rate that grows with c . Defining $r = \lambda_3/\lambda_2$, we have

$$\begin{aligned} (\mathbf{M}^\tau)_{ba} &= B_{b1} B_{1a}^{-1} + \lambda_2^\tau [B_{b2} B_{2a}^{-1} + \mathcal{O}(r^\tau)] \\ &= \mu_b + \lambda_2^\tau A_{ba}, \end{aligned} \quad (7)$$

where we have made use of the fact that an irreducible and aperiodic Markov process must converge to its stationary distribution for large τ , and we have defined \mathbf{A} as the expression in square brackets above, satisfying $\lim_{\tau \rightarrow \infty} A_{ba} = B_{b2} B_{2a}^{-1}$. Note that $\sum_b A_{ba} = 0$ in order for \mathbf{M} to be properly normalized.

Substituting equation (7) into equation (8) and using the facts that $\sum_a \mu_a = 1$ and $\sum_b A_{ba} = 0$, we obtain

$$\begin{aligned} I_R &= \sum_{ab} \frac{\mu_a}{\mu_b} [(\mathbf{M}^\tau)_{ba}]^2 - 1 \\ &= \sum_{ab} \frac{\mu_a}{\mu_b} (\mu_b^2 + 2\mu_b \lambda_2^\tau A_{ba} + \lambda_2^{2\tau} A_{ba}^2) - 1 \\ &= \sum_{ab} \lambda_2^{2\tau} (\mu_b^{-1} A_{ba}^2 \mu_a) = \mathcal{C} \lambda_2^{2\tau}, \end{aligned} \quad (8)$$

where the term in the last parentheses is of the form $\mathcal{C} = \mathcal{C}_0 + \mathcal{O}(r^\tau)$.

In summary, we have shown that an irreducible and aperiodic Markov process with non-degenerate eigenvalues cannot produce critical behavior, because the mutual information decays exponentially. In fact, *no* Markov processes can, as we show in Appendix B.

To hammer the final nail into the coffin of Markov processes as models of critical behavior, we need to close a final loophole. Their fundamental problem is lack of long-term memory, which can be superficially overcome by redefining the state space to include symbols from the past. For example, if the current state is one of n and we wish the process to depend on the the last τ symbols, we can define an expanded state space consisting of the n^τ possible sequences of length τ , and a corresponding $n^\tau \times n^\tau$ Markov matrix (or an $n^\tau \times n$ table of conditional probabilities for the next symbol given the last τ symbols). Although such a model could fit the curves in Figure I in theory, it cannot in practice, because \mathbf{M} requires way more parameters than there are atoms in our observable universe ($\sim 10^{78}$): even for as few as $n = 4$ symbols and $\tau = 1000$, the Markov process involves over $4^{1000} \sim 10^{602}$ parameters. Scale-invariance aside, we can also see how Markov processes fail simply by considering the structure of text. To model English well, \mathbf{M} would need to correctly close parentheses even if they were opened more than $\tau = 100$ characters ago, requiring an \mathbf{M} -matrix with than n^{100} parameters, where $n > 26$ is the number of characters used.

We can significantly generalize Theorem 1 into a theorem about hidden Markov models (HMM). In an HMM, the observed sequence X_1, \dots, X_n is only part of the picture: there are hidden variables Y_1, \dots, Y_n that themselves form a Markov chain. We can think of an HMM as follows: imagine a machine with an internal state space Y that updates itself according to some Markovian dynamics. The internal dynamics are never observed, but at each time-step, it also produces some output $Y_i \rightarrow X_i$ that form the sequence which we can observe. These models are quite general and are used to model a wealth of empirical data (see, e.g., [38]).

Theorem 2: Let \mathbf{M} be a Markov matrix that generates the transitions between hidden states Y_i in an HMM. If \mathbf{M} is irreducible and aperiodic, then the asymptotic behavior of the mutual information $I(t_1, t_2)$ is exponential

decay toward zero for $|t_2 - t_1| \gg 1$ with decay timescale $\log \frac{1}{|\lambda_2|}$, where λ_2 is the second largest eigenvalue of \mathbf{M} . This theorem is a strict generalization of Theorem 1, since given any Markov process \mathcal{M} with corresponding matrix \mathbf{M} , we can construct an HMM that reproduces the exact statistics of \mathcal{M} by using \mathcal{M} as the transition matrix between the Y 's and generating X_i from Y_i by simply setting $x_i = y_i$ with probability 1.

The proof is very similar in spirit to the proof of Theorem 1, so we will just present a sketch here, leaving a full proof to Appendix B. Let \mathbf{G} be the Markov matrix that governs $X_i \rightarrow Y_i$. To compute the joint probability between two random variables X_{t_1} and X_{t_2} , we simply compute the joint probability distribution between Y_{t_1} and Y_{t_2} , which again involves a factor of \mathbf{M}^τ and then use two factors of \mathbf{G} to convert the joint probability on Y_{t_1}, Y_{t_2} to a joint probability on X_{t_1}, X_{t_2} . These additional two factors of \mathbf{G} will not change the fact that there is an exponential decay given by \mathbf{M}^τ .

A simple, intuitive bound from information theory (namely the data processing inequality [37]) gives $I(Y_{t_1}, Y_{t_2}) \geq I(Y_{t_1}, X_{t_2}) \geq I(X_{t_1}, X_{t_2})$. However, Theorem 1 implies that $I(Y_{t_1}, Y_{t_2})$ decays exponentially. Hence $I(X_{t_1}, X_{t_2})$ must also decay at least as fast as exponentially.

There is a well-known correspondence between so-called *probabilistic regular grammars* [39] (sometimes referred to as stochastic regular grammars) and HMMs. Given a probabilistic regular grammar, one can generate an HMM that reproduces all statistics and vice versa. Hence, we can also state Theorem 2 as follows:

Corollary: No probabilistic regular grammar exhibits criticality.

In the next section, we will show that this statement is not true for context-free grammars.

III. POWER LAWS FROM GENERATIVE GRAMMAR

If computationally feasible Markov processes cannot produce genomes, melodies, texts or other sequences with roughly critical behavior, then how do such sequences arise? What sort of alternative processes can generate them? This question is not only interesting theoretically, but also important in practically: models which can approximate such critical sequences could explain how some machine learning algorithms perform better than others in tasks like such as language processing, and might suggest ways to improve existing algorithms. In the best case scenario, theoretical models may even shed light on how human brains can efficiently generate English sentences without storing googols of parameters.

One answer advanced by Chomsky [40] and others which will loosely inspire our work is the idea of deep structure.

Roughly speaking, the idea is that language is generated in a hierarchical rather than linear fashion. When we write an essay, we do so by thinking about some big idea, and then breaking it into sub-ideas, and sub-sub-ideas, etc. Similarly, we generate a sentence by first choosing its basic structure and then fleshing-out each part of the sentence with more modifiers, etc. For example, the two sentences “Bob loves Alice” and “Alice is loved by Bob” are “close” in meaning but there could be nothing similar about them if English were Markov. On the other hand, if English is generated hierarchically, these sentences might be close in the sense that they diverged close to the leaves of the generative tree.

A. A simple recursive grammar model

We can formalize the above considerations by giving production rules for a toy language L over an alphabet A . In the parlance of theoretical linguistics, our language is generated by a *stochastic* or *probabilistic context-free grammar* (PCFG) [41–44]. We will discuss the relationship between our model and a generic PCFG in Section C. The language is defined by how a native speaker of L produces sentences: first, she draws one of the $|A|$ characters from some probability distribution μ on A . She then takes this character x_0 and replaces it with q new symbols, drawn from a probability distribution $P(b|a)$, where $a \in A$ is the first symbol and $b \in A$ is any of the second symbols. This is repeated over and over. After u steps, she has a sentence of length q^u .²

One can ask for the character statistics of the sentence at production step u given the statistics of the sentence at production step $u - 1$. The character distribution is simply

$$P_u(b) = \sum_a P(b|a)P_u(a). \quad (9)$$

Of course this equation does *not* imply that the process is a Markov process when the sentences are read left to right. To characterize the statistics as read from left to right, we really want to compute the statistical dependencies *within* a given sequence, *e.g.*, at fixed u .

To see that the mutual information decays like a power law rather than exponentially with separation, consider two random variables X and Y separated by τ . One can ask how many generations took place between X and the nearest ancestor of X and Y . Typically, this will be about $\log_q \tau$ generations. Hence in the tree graph shown

² This exponential blow-up is reminiscent of de Sitter space in cosmic inflation. There is actually a much deeper mathematical analogy involving conformal symmetry and p -adic numbers that has been discussed by Harlow *et al.*[45].

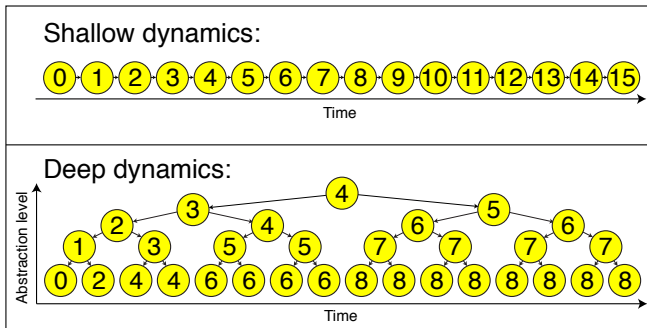


FIG. 2: Both a traditional Markov process (top) and our recursive generative grammar process (bottom) can be represented as Bayesian networks, where the random variable at each node depends only on the node pointing to it with an arrow. The numbers show the geodesic distance Δ to the leftmost node, defined as the smallest number of edges that must be traversed to get there. Our results show that the mutual information decays exponentially with Δ . Since this geodesic distance Δ grows only logarithmically with the separation in time in a hierarchical generative grammar (the hierarchy creates very efficient shortcuts), the exponential kills the logarithm and we are left with power-law decays of mutual information in such languages.

in Figure 2, which illustrates the special case $q = 2$, the number of edges Δ between X and Y is about $2 \log_q \tau$. Hence by the previous result for Markov processes, we expect an exponential decay of the mutual information in the variable $\Delta \sim 2 \log_q \tau$. This means that $I(X, Y)$ should be of the form

$$I(X, Y) \sim q^{-\gamma \Delta} = q^{2 \log_q(\tau)^{-\gamma}} = \tau^{-2\gamma}, \quad (10)$$

where γ is controlled by the second-largest eigenvalue of \mathbf{G} , the matrix of conditional probabilities $P(b|a)$. But this exponential decay in Δ is exactly a power-law decay in τ ! This intuitive argument is transformed into a rigorous proof in Appendix C.

B. Further Generalization: strongly correlated characters in words

In the model we have been describing so far, all nodes emanating from the same parent can be freely permuted since they are conditionally independent. In this sense, characters within a newly generated word are uncorrelated. We call models with this property *weakly correlated*. There are still arbitrarily large correlations between words, but not inside of words. If a weakly correlated grammar allows $a \rightarrow ab$, it must allow for $a \rightarrow ba$ with the same probability. We now wish to relax this property to allow for the *strongly-correlated* case where variables may not be conditionally independent given the parents. This allows us to take a big step towards mod-

eling realistic languages: in English, *god* significantly differs in meaning and usage from *dog*.

In the previous computation, the crucial ingredient was the joint probability $P(a, b) = P(X = a, Y = b)$. Let us start with a seemingly trivial remark. This joint probability can be re-interpreted as a conditional joint probability. Instead of X and Y being random variables at *specified* sites t_1 and t_2 , we can view them as random variables at randomly chosen locations, conditioned on their locations being t_1 and t_2 . Somewhat pedantically, we write $P(a, b) = P(a, b|t_1, t_2)$. This clarifies the important fact that the only way that $P(a, b|t_1, t_2)$ depends on t_1 and t_2 is via a dependence on $\Delta(t_1, t_2)$. Hence

$$P(a, b|t_1, t_2) = P(a, b|\Delta). \quad (11)$$

This equation is specific to weakly correlated models and does not hold for generic strongly correlated models.

In computing the mutual information as a function of separation, the relevant quantity is the right hand side of equation (7). The reason is that in practical scenarios, we estimate probabilities by sampling a sequence at fixed separation $t_1 - t_2$, corresponding to $\Delta \approx 2 \log_q |t_2 - t_1| + \mathcal{O}(1)$, but varying t_1 and t_2 . (The $\mathcal{O}(1)$ term is discussed in Appendix E).

Now whereas $P(a, b|t_1, t_2)$ will change when strong correlations are introduced, $P(a, b|\Delta)$ will retain a very similar form. This can be seen as follows: knowledge of the geodesic distance corresponds to knowledge of how high up the closest parent node is in the hierarchy (see Figure 1). Imagine flowing down from the parent node to the leaves. We start with the stationary distribution μ_i at the parent node. At the first layer below the parent node (corresponding to a causal distance $\Delta - 2$), we get $Q_{rr'} \equiv P(rr') = \sum_i P_S(rr'|i)P(i)$, where the symmetrized probability $P_S = \frac{1}{2} \sum_i [P(rr'|i) + P(r'r|i)]$ comes into play because knowledge of the fact that r, r' are separated by $\Delta - 2$ gives no information about their order. To continue this process to the second stage and beyond, we only need the matrix $G_{sr} = P(s|r) = \sum_{s'} P_S(ss'|r)$. The reason is that since we only wish to compute the two-point function at the bottom of the tree, the only place where a three-point function is ever needed is at the very top of the tree, where we need to take a single parent into two children nodes. After that, the computation only involves evolving a child node into a grand-child node, and so forth. Hence the overall two-point probability matrix $P(ab|\Delta)$ is given by the simple equation

$$\mathbf{P}(\Delta) = \left(\mathbf{G}^{\Delta/2-1} \right) \mathbf{Q} \left(\mathbf{G}^{\Delta/2-1} \right)^t. \quad (12)$$

As we can see from the above formula, changing to the strongly correlated case essentially reduces to the weakly correlated case where

$$\mathbf{P}(\Delta) = \left(\mathbf{G}^{\Delta/2} \right) \text{diag}(\boldsymbol{\mu}) \left(\mathbf{G}^{\Delta/2} \right)^t, \quad (13)$$

except for a perturbation near the top of the tree. We can think of the generalization as equivalent to the old model except for a different initial condition. We thus expect on intuitive grounds that the model will still exhibit power law decay. This intuition is correct, as we will prove rigorously in Appendix C.

C. Further Generalization: Bayesian networks and context-free grammars

Just how generic is the scaling behavior of our model? What if the length of the words is not constant? What about more complex dependencies between layers? If we retrace the derivation in the above arguments, it becomes clear that the only key feature of all of our models considered so far is that the rational mutual information decays exponentially with the causal distance Δ :

$$I_R \propto e^{-\gamma\Delta}. \quad (14)$$

This is true for (hidden) Markov processes and the hierarchical grammar models that we have considered above. So far we have defined Δ in terms of quantities specific to these models; for a Markov process, Δ is simply the time separation. Can we define Δ more generically? In order to do so, let us make a brief aside about *Bayesian networks*. Formally, a Bayesian net is a directed acyclic graph (DAG), where the vertices are random variables and conditional dependencies are represented by the arrows. Now instead of thinking of X and Y as living at certain times (t_1, t_2) , we can think of them as living at vertices (i, j) of the graph.

We define $\Delta(i, j)$ as follows. Since the Bayesian net is a DAG, it is equipped with a partial order \leq on vertices. We write $k \leq l$ iff there is a path from k to l , in which case we say that k is an *ancestor* of l . We define the $L(k, l)$ to be the number of edges on the shortest directed path from k to l . Finally, we define the causal distance $\Delta(i, j)$ to be

$$\Delta(i, j) \equiv \min_{x \leq i, x \leq j} L(x, i) + L(x, j). \quad (15)$$

It is easy to see that this reduces to our previous definition of Δ for Markov processes and recursive generative trees (see Figure 2).

Is it true that our exponential decay result from equation (14) holds even for a generic Bayesian net? The answer is yes, under a suitable approximation. The approximation is to ignore long paths in the network when computing the mutual information. In other words, the mutual information tends to be dominated by the shortest paths via a common ancestor, whose length is Δ . This is a generally a reasonable approximation, because these longer paths will give exponentially weaker correlations,

so unless the number of paths increases exponentially (or faster) with length, the overall scaling will not change.

With this approximation, we can state a key finding of our theoretical work. Deep models are important because without the extra “dimension” of depth/abstraction, there is no way to construct “shortcuts” between random variables that are separated by large amounts of time with short-range interactions; 1D models will be doomed to exponential decay. *Hence the ubiquity of power laws explains the success of deep learning.* In fact, this can be seen as the Bayesian net version of the important result in statistical physics that there are no phase transitions in 1D [46, 47].

There are close analogies between our deep recursive grammar and more conventional physical systems. For example, according to the emerging standard model of cosmology, there was an early period of cosmological inflation when density fluctuations get getting added on a fixed scale as space itself underwent repeated doublings, combining to produce an excellent approximation to a power-law correlation function. This inflationary process is simply a special case of our deep recursive model (generalized from 1 to 3 dimensions). In this case, the hidden “depth” dimension in our model corresponds to cosmic time, and the time parameter which labels the place in the sequence of interest corresponds to space. A similar physical analogy is turbulence in a fluid, where energy in the form of vortices cascades from large scales to ever smaller scales through a recursive process where larger vortices create smaller ones, leading to a scale-invariant power spectrum. There is also a close analogy to quantum mechanics: in equation (13) expresses the exponential decay of the mutual information with geodesic distance through the Bayesian network; in quantum mechanics, the correlation function of a many body system decays exponentially with the geodesic distance defined by the tensor network which represents the wavefunction [48].

It is also worth examining our model using techniques from linguistics. A generic PCFG \mathcal{G} consists of three ingredients:

1. An alphabet $\mathcal{A} = A \cup T$ which consists of non-terminal symbols A and terminal symbols T .
2. A set of production rules of the form $a \rightarrow B$, where the left hand side $a \in A$ is always a single non-terminal character and B is a string consisting of symbols in \mathcal{A} .
3. Probabilities associated with each production rule $P(a \rightarrow B)$, such that for each $a \in A$, $\sum_B P(a \rightarrow B) = 1$.

It is a remarkable fact that any stochastic-context free grammars can be put in *Chomsky normal form* [43, 49]. This means that given \mathcal{G} , there exists some other grammar $\tilde{\mathcal{G}}$ such that all the production rules are either of

the form $a \rightarrow bc$ or $a \rightarrow \alpha$, where $a, b, c \in A$ and $\alpha \in T$ and the corresponding languages $L(\mathcal{G}) = L(\bar{\mathcal{G}})$. In other words, given some complicated grammar \mathcal{G} , we can always find a grammar $\bar{\mathcal{G}}$ such that the corresponding statistics of the languages are identical and all the production rules replace a symbol by at most two symbols (at the cost of increasing the number of production rules in $\bar{\mathcal{G}}$).

This formalism allows us to strengthen our claims. Our model with a branching factor $q = 2$ is precisely the class of all context-free grammars that are generated by the production rules of the form $a \rightarrow bc$. While this might naively seem like a very small subset of all possible context-free grammars, the fact that *any* context-free grammar can be converted into Chomsky normal form shows that our theory deals with a generic context-free grammar, except for the additional step of producing terminal symbols from non-terminal symbols. Starting from a single symbol, the deep dynamics of the PCFG in normal form are given by a strongly-correlated branching process with $q = 2$ which proceeds for a characteristic number of productions before terminal symbols are produced. Before most symbols have been converted to terminal symbols, our theory applies, and power-law correlations will exist amongst the non-terminal symbols. To the extent that the terminal symbols that are then produced from non-terminal symbols reflect the correlations of the non-terminal symbols, we expect context-free grammars to be able to produce power law correlations.

From our corollary to Theorem 2, we know that regular grammars cannot exhibit power-law decays in mutual information. Hence context-free grammars are the simplest grammars which support criticality, e.g., they are the lowest in the Chomsky hierarchy that supports criticality. Note that our corollary to Theorem 2 also implies that not all context-free grammars exhibit criticality since regular grammars are a strict subset of context-free grammars. Whether one can formulate an even sharper criterion should be the subject of future work.

IV. DISCUSSION

We have shown that many data sequences generated for the purposes of communications — from English and French text to Bach and the human genome — exhibit critical behavior, where the mutual information between symbols decays roughly like a power law with separation. By introducing a quantity we term rational mutual information, we have proved that (hidden) Markov processes generically exhibit exponential decay, whereas deep generative grammars exhibit power law decays. This explains why natural languages are poorly approximated by Markov processes, but better approximated by deep recurrent neural networks now widely used in machine learning for natural language processing, as we will discuss in detail below. Furthermore, we have identified

a crucial ingredient of any successful natural language model: it must have one or more “hidden” dimensions, which can be used to provide shortcuts between distant parts of the network; hence leading to longer-range correlations that are crucial for critical behavior.

Let us now explore some useful implications of these results, both for understanding the success of certain neural network architectures and for using the mutual information function as a tool for validating machine learning algorithms.

A. Connection to Recurrent Neural Networks

While the generative grammar model is appealing from a linguistic perspective, it may superficially appear to have little to do with machine learning algorithms that are implemented in practice. However, as we will now see, this model can in fact be viewed an idealized version of a long-short term memory (LSTM) recurrent neural network (RNN) that is generating (“hallucinating”) a sequence.

First of all, Figure 4 shows that an LSTM RNN can in fact reproduce critical behavior. In this example, we trained an RNN (consisting of three hidden LSTM layers of size 256 as described in [20]) to predict the next character in the 100MB Wikipedia sample known as enwik8 [12]. We then used the LSTM to hallucinate 1 MB of text and measured the mutual information as a function of distance. Figure 4 shows that not only is the resulting mutual information function a rough power law, but it also has a slope that is relatively similar to the original.

We can understand this success by considering a simplified model that is less powerful and complex than a full LSTM, but retains some of its core features — such an approach to studying deep neural nets has proved fruitful in the past (e.g., [27]).

The usual implementation of LSTMs consists of multiple cells stacked one on top of each other. Each cell of the LSTM (depicted as a yellow circle in Fig. 3) has a state that is characterized by a matrix of numbers \mathbf{C}_t and is updated according to the following rule

$$\mathbf{C}_t = \mathbf{f}_t \circ \mathbf{C}_{t-1} + \mathbf{i}_t \circ \mathbf{D}_t, \quad (16)$$

where \circ denotes element wise multiplication, and $\mathbf{D}_t = \mathbf{D}_t(\mathbf{C}_{t-1}, \mathbf{x}_t)$ is some function of the input \mathbf{x}_t from the cell from the layer above (denoted by downward arrows in Figure 3, the details of which do not concern us. Generically, a graph of this picture would look like a rectangular lattice, with each node having an arrow to its right (corresponding to the first term in the above equation), and an arrow from above (corresponding to the second term in the equation). However, if the forget weights \mathbf{f} weights decay rapidly with depth (e.g., as we go from the bottom

cell to the towards the top) so that the timescales for forgetting grow exponentially, we will show that a reasonable approximation to the dynamics is given by Figure 3.

If we neglect the dependency of \mathbf{D}_t on \mathbf{C}_{t-1} , the forget gate \mathbf{f}_t leads to exponential decay of \mathbf{C}_{t-1} e.g., $\mathbf{C}_t = \mathbf{f}^t \circ \mathbf{C}_0$; this is how LSTM’s forget their past. Note that all operations including exponentiation are performed element-wise in this section only.

In general, a cell will smoothly forget its past over a timescale of $\sim \log(1/f) \equiv \tau_f$. On timescales $\gtrsim \tau_f$, the cells are weakly correlated; on timescales $\lesssim \tau_f$, the cells are strongly correlated. Hence a discrete approximation to this above equation is the following:

$$\begin{aligned} \mathbf{C}_t &= \mathbf{C}_{t-1}, \text{ for } \tau_f \text{ timesteps} \\ &= \mathbf{D}_t(\mathbf{x}_t), \text{ on every } \tau_f + 1 \text{ timestep.} \end{aligned} \quad (17)$$

This simple approximation leads us right back to the hierarchical grammar! The first line of the above equation is labeled “remember” in Figure 2 and the second line is what we refer to as “Markov,” since the next state depends only on the previous. Since each cell perfectly remembers its pervious state for τ_f time-steps, the tree can be reorganized so that it is exactly of the form shown in Figure 3, by omitting nodes which simply copy the previous state. Now supposing that τ_f grows exponentially with depth $\tau_f(\text{layer } i) \propto q \tau_f(\text{layer } i + 1)$, we see that the successive layers become exponentially sparse, which is exactly what happens in our deep grammar model, identifying the parameter q , governing the growth of the forget timescale, with the branching parameter in the deep grammar model. (Compare Figure 2 and Figure 3.)

B. A new diagnostic for machine learning

How can one tell whether an neural network can be further improved? For example, an LSTM RNN similar to the one we used in Figure 3 can predict Wikipedia text with a residual entropy ~ 1.4 bits/character [20], which is very close to the performance of current state of the art custom compression software — which achieves ~ 1.3 bits/character [50]. Is that essentially the best compression possible, or can significant improvements be made?

Our results provide a powerful diagnostic for shedding further light on this question: measuring the mutual information as a function of separation between symbols is a computationally efficient way of extracting much more meaningful information about the performance of a model than simply evaluating the loss function, usually given by the conditional entropy $H(X_t|X_{t-1}, X_{t-2}, \dots)$.

Figure 3 shows that even with just three layers, the LSTM-RNN is able to learn long-range correlations; the slope of the mutual information of hallucinated text is comparable to that of the training set. However, the figure also shows that the predictions of our LSTM-RNN

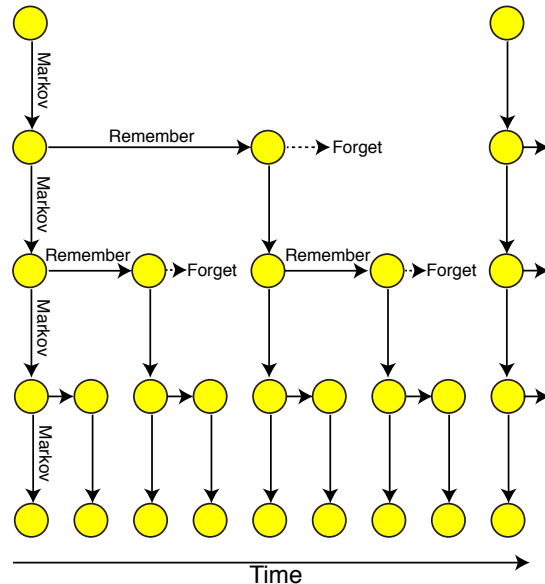


FIG. 3: Our deep generative grammar model can be viewed as an idealization of a long-short term memory (LSTM) recurrent neural net, where the “forget weights” drop with depth so that the forget timescales grow exponentially with depth. The graph drawn here is clearly isomorphic to the graph drawn in Figure 1. For each cell, we approximate the usual incremental updating rule by either perfectly remembering the previous state (horizontal arrows) or by ignoring the previous state and determining the cell state by a random rule depending on the node above (vertical arrows).

are far from optimal. Interestingly, the hallucinated text shows about the same mutual information for distances $\sim \mathcal{O}(1)$, but significantly less mutual information at large separation. Without requiring any knowledge about the true entropy of the input text (which is famously NP-hard to compute), this figure immediately shows that the LSTM-RNN we trained is performing sub-optimally; it is not able to capture all the long-term dependencies found in the training data.

As a comparison, we also calculated the bigram transition matrix $P(X_3 X_4 | X_1 X_2)$ from the data and used it to hallucinate 1 MB of text. Despite the fact that this higher order Markov model needs $\sim 10^3$ more parameters than our LSTM-RNN, it captures less than a fifth of the mutual information captured by the LSTM-RNN even at modest separations $\gtrsim 5$.

In summary, Figure 3 shows both the successes and shortcomings of machine learning. On the one hand, LSTM-RNN’s can capture long-range correlations much more efficiently than Markovian models; on the other hand, they cannot match the two point functions of training data, never mind higher order statistics!

One might wonder how the lack of mutual information at large scales for the bigram Markov model is manifested in the hallucinated text. Below we give a line from the

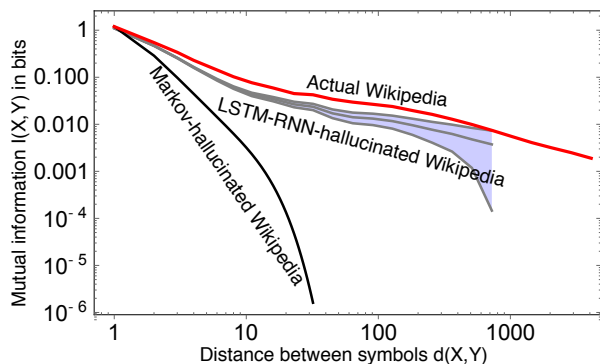


FIG. 4: Diagnosing different models with by hallucinating text and then measuring the mutual information as a function of separation. The red line is the mutual information of enwik8, a 100 MB sample of English Wikipedia. In shaded blue is the mutual information of hallucinated Wikipedia from a trained LSTM with 3 layers of size 256. We plot in solid black the mutual information of a Markov process on single characters, which we compute exactly. (This would correspond to the mutual information of hallucinations in the limit where the length of the hallucinations goes to infinity). This curve shows a sharp exponential decay after a distance of ~ 10 , in agreement with our theoretical predictions. We also measured the mutual information for hallucinated text on a Markov process for bigrams, which still underperforms the LSTMs in long-ranged correlations, despite having $\sim 10^3$ more parameters than

Markov hallucinations:

```
[[computhourgist, Flagesernmenserved whirequotes
or thand dy excommentaligmaktophy as
its:Fran at ||&lt;If ISBN 088;&ampategor
and on of to [[Prefung]]' and at them rector>
```

This can be compared with an example from the LSTM RNN:

```
Proudknow pop groups at Oxford
- [http://ccw.com/faqsisdaler/cardiffstwander
--helgar.jpg] and Cape Normans's first
attacks Cup rigid (AM).
```

Despite using many fewer parameters, the LSTM manages to produce a realistic looking URL and is able to close brackets correctly [51], something that the Markov model struggles with.

C. Outlook

Although great challenges remain to accurately model natural languages, our results at least allow us to improve on some earlier answers to key questions we sought to address :

1. *Why is natural language so hard?* The old answer was that language is uniquely human. The new answer is that at least part of the difficulty is that natural language is a critical system, with long-ranged correlations that are difficult for machines to learn.
2. *Why are machines bad at natural languages, and why are they good?* The old answer is that Markov models are simply not brain/human-like, whereas neural nets are more brain-like and hence better. The new answer is that Markov models or other 1-dimensional models cannot exhibit critical behavior, whereas neural nets and other deep models (where an extra hidden dimension is formed by the layers of the network) are able to exhibit critical behavior.
3. *How can we know when machines are bad or good?* The old answer is to compute the loss function. The new answer is to also compute the mutual information as a function of separation, which can immediately show how well the model is doing at capturing correlations on different scales.

Future studies could include more comprehensive measurements from multiple language/music corpuses and other one dimensional data. In addition, more theoretical work is required to explain why natural languages have slopes that are all comparable $\sim 1/2$. In statistical physics, “coincidences” of these sort are usually signs of universality: many seemingly unrelated systems have the same long-wavelength effective field theory at the critical point and hence display the same power law slopes. Perhaps something similar is happening here, though this is unexpected from the deep generative grammar model where any power law slope is possible.

Acknowledgments: This work was supported by the Foundational Questions Institute <http://fqxi.org>. The authors wish to thank Noam Chomsky and Greg Lessard for valuable comments on the linguistic aspects of this work, Taiga Abe, Meia Chita-Tegmark, Hanna Field, Esther Goldberg, Emily Mu, John Peurifoi, Tomaso Poggio, Luis Seoane, Leon Shen and David Theurel for helpful discussions and encouragement, Michelle Xu for help acquiring genome data and the Center for Brains Minds and Machines (CMBB) for hospitality.

Author contributions: HL proposed the project idea in collaboration with MT. HL and MT collaboratively formulated the proofs, performed the numerical experiments, analyzed the data, and wrote the manuscript.

Appendix A: Properties of rational mutual information

In this appendix, we prove the following elementary properties of rational mutual information:

1. **Symmetry:** for any two random variables X and Y , $I_R(X, Y) = I_R(Y, X)$. The proof is straightforward:

$$\begin{aligned} I_R(X, Y) &= \sum_{ab} \frac{P(X=a, Y=b)^2}{P(X=a)P(Y=b)} - 1 \\ &= \sum_{ba} \frac{P(Y=b, X=a)^2}{P(Y=b)P(X=a)} - 1 = I_R(Y, X). \end{aligned} \quad (\text{A1})$$

2. **Upper bound to mutual information:** The logarithm function satisfies $\ln(1+x) \leq x$ with equality if and only if (iff) $x = 0$. Therefore setting $x = \frac{P(a,b)}{P(a)P(b)} - 1$ gives

$$\begin{aligned} I(X, Y) &= \left\langle \log_B \frac{P(a,b)}{P(a)P(b)} \right\rangle \\ &= \frac{1}{\ln B} \left\langle \ln \left[1 + \left(\frac{P(a,b)}{P(a)P(b)} - 1 \right) \right] \right\rangle \\ &\leq \frac{1}{\ln B} \left\langle \frac{P(a,b)}{P(a)P(b)} - 1 \right\rangle = \frac{I_R(X, Y)}{\ln B}. \end{aligned} \quad (\text{A2})$$

Hence the rational mutual information $I_R \geq I \ln B$ with equality iff $I = 0$ (or simply $I_R \geq I$ if we use the natural logarithm base $B = e$).

3. **Non-negativity.** It follows from the above inequality that $I_R(X, Y) \geq 0$ with equality iff $P(a, b) = P(a)P(b)$, since $I_R = I = 0$ iff $P(a, b) = P(a)P(b)$. Note that this short proof is only possible because of the information inequality $I \geq 0$. From the definition of I_R , it is only obvious that $I_R \geq -1$; information theory gives a much tighter bound. Our findings 1-3 can be summarized as follows:

$$I_R(X, Y) = I_R(Y, X) \geq I(X, Y) \geq 0, \quad (\text{A3})$$

where both equalities occur iff $p(X, Y) = p(X)p(Y)$. It is impossible for one of the last two relations to be an equality while the other is an inequality.

4. **Generalization.** Note that if we view the mutual information as the divergence between two joint probability distributions, we can generalize the notion of rational mutual information to that of *rational divergence*:

$$D_R(p||q) = \left\langle \frac{p}{q} \right\rangle - 1, \quad (\text{A4})$$

where the expectation value is taken with respect to the “true” probability distribution p . Note that as it is written, p could be any probability measure on either a discrete or continuous space. The above results can be trivially modified to show that $D_R(p||q) \geq D_{KL}(p||q)$ and hence $D_R(p||q) \geq 0$, with equality iff $p = q$.

Appendix B: General proof for Markov processes

In this appendix, we drop the assumptions of non-degeneracy, irreducibility and non-periodicity made in the main body of the paper where we proved that Markov processes lead to exponential decay.

1. The degenerate case

First, we consider the case where the Markov matrix \mathbf{M} has degenerate eigenvalues. In this case, we cannot guarantee that \mathbf{M} can be diagonalized. However, any complex matrix can be put into Jordan normal form. In Jordan normal form, a matrix is block diagonal, with each $d \times d$ block corresponding to an eigenvalue with degeneracy d . These blocks have a particularly simple form, with block i having λ_i on the diagonal and ones right above the diagonal. For example, if there are only three distinct eigenvalues and λ_2 is threefold degenerate, the the Jordan form of \mathbf{M} would be

$$\mathbf{B}^{-1}\mathbf{M}\mathbf{B} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2 & 1 & 0 & 0 \\ 0 & 0 & \lambda_2 & 1 & 0 \\ 0 & 0 & 0 & \lambda_2 & 0 \\ 0 & 0 & 0 & 0 & \lambda_3 \end{bmatrix}. \quad (\text{B1})$$

Note that the largest eigenvalue is unique and equal to 1 for all irreducible and aperiodic \mathbf{M} . In this example, the matrix power \mathbf{M}^τ is

$$\mathbf{B}^{-1}\mathbf{M}^\tau\mathbf{B} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \lambda_2^\tau & \binom{\tau}{1}\lambda_2^{\tau-1} & \binom{\tau}{2}\lambda_2^{\tau-2} & 0 \\ 0 & 0 & \lambda_2^\tau & \binom{\tau}{1}\lambda_2^{\tau-1} & 0 \\ 0 & 0 & 0 & \lambda_2^\tau & 0 \\ 0 & 0 & 0 & 0 & \lambda_3^\tau \end{bmatrix}. \quad (\text{B2})$$

In the general case, raising a matrix to an arbitrary power will yield a matrix which is still block diagonal, with each block being an upper triangular matrix. The important point is that in block i , every entry scales $\propto \lambda_i^\tau$, up to a combinatorial factor. Each combinatorial factor grows only polynomially with τ , with the degree of the polynomials in the i th block bounded by the multiplicity of λ_i , minus one.

Using this Jordan decomposition, we can replicate equation (7) and write

$$M_{ij}^\tau = \mu_i + \lambda_2^\tau A_{ij}. \quad (\text{B3})$$

There are two cases, depending on whether the second eigenvalue λ_2 is degenerate or not. If not, then the equation

$$\lim_{\tau \rightarrow \infty} A_{ij} = B_{i2} B_{2j}^{-1} \quad (\text{B4})$$

still holds, since for $i \geq 3$, $(\lambda_i/\lambda_2)^\tau$ decays faster than any polynomial of finite degree. On the other hand, if the second eigenvalue is degenerate with multiplicity m_2 , we instead define \mathbf{A} with the combinatorial factor removed:

$$M_{ij}^\tau = \mu_i + \binom{\tau}{m_2} \lambda_2^\tau A_{ij}. \quad (\text{B5})$$

If $m_2 = 1$, this definition simply reduces to the previous definition of \mathbf{A} . With this definition,

$$\lim_{\tau \rightarrow \infty} A_{ij} = \lambda_2^{-m_2} B_{i2} B_{(2+m_2)j}^{-1}, \quad (\text{B6})$$

Hence in the most general case, the mutual information decays like a polynomial $\mathcal{P}(\tau)e^{-\gamma\tau}$, where $\gamma = 2 \ln \frac{1}{\lambda_2}$. The polynomial is non-constant if and only if the second largest eigenvalue is degenerate. Note that even in this case, the mutual information decays exponentially in the sense that it is possible to bound the mutual information by an exponential.

2. The reducible case

Now let us generalize to the case where the Markov process is reducible, *i.e.*, decomposable into a set of m non-interacting Markov processes for some integer $m > 1$. This means that the state space can be decomposed as

$$S = \bigcup_{i=1}^m S_i, \quad (\text{B7})$$

where restricting the Markov process to each S_i results in a well-defined Markov process on S_i that does not interact with any other S_i . If the system starts off in S_1 , it can never transition to S_2 , and so forth. Now if we restrict our model to the subset of interest, the mutual information will exponentially decay; however, for a generic initial state, the total probability within each set S_i remains constant, so the mutual information will be asymptotically approach the entropy of the probability distribution across sets, which can be at most $I = \log_2 m$. In the language of statistical physics, this is an example of topological order which leads to constant terms in the correlation functions; here, the Markov graph of \mathbf{M} is disconnected, so there are m degenerate equilibrium states.

3. The periodic case

If a Markov process is periodic, with a sequence that repeats forever, then the mutual information is a constant and never decays to zero, so the only power law that can be attained is the case of slope zero which does not correspond to critical behavior.

4. The $n > 1$ case

The following proof holds only for order $n = 1$ Markov processes, but we can easily extend the results for arbitrary n . Any $n = 2$ Markov process can be converted into an $n = 1$ Markov process on pairs of letters $X_1 X_2$. Hence our proof shows that $I(X_1 X_2, Y_1 Y_2)$ decays exponentially. But for any random variables X, Y , the data processing inequality [37] states that $I(X, g(Y)) \leq I(X, Y)$, where g is an arbitrary function of Y . Letting $g(Y_1 Y_2) = Y_1$, and then permuting and applying $g(X_1, X_2) = X_1$ gives

$$I(X_1 X_2, Y_1 Y_2) \geq I(X_1 X_2, Y_1) \geq I(X_1, Y_1). \quad (\text{B8})$$

Hence, we see that $I(X_1, Y_1)$ must exponentially decay. The preceding remarks can be easily formalized into a proof for an arbitrary Markov process by induction on n .

5. The detailed balance case

This asymptotic relation can be strengthened for a subclass of Markov processes which obey a condition known as detailed balance. This subclass arises naturally in the study of statistical physics [52]. For our purposes, this simply means that there exist some real numbers K_m and a symmetric matrix $S_{ab} = S_{ba}$ such that

$$M_{ab} = e^{K_a/2} S_{ab} e^{-K_b/2}. \quad (\text{B9})$$

Let us note the following facts. (1) The matrix power is simply $(M^\tau)_{ab} = e^{K_a/2} (S^\tau)_{ab} e^{-K_b/2}$. (2) By the spectral theorem, we can diagonalize S into an orthonormal basis of eigenvectors, which we label as v (or sometimes w), *e.g.*, $Sv = \lambda_i v$ and $v \cdot w = \delta_{vw}$. Notice that

$$\sum_n M_{ab} e^{K_n/2} v_n = \sum_n e^{K_m/2} S_{mn} v_n = \lambda_i e^{K_m/2} v_m.$$

Hence we have found an eigenvector of M for every eigenvector of S . Conversely, the set of eigenvectors of S forms a basis, so there cannot be any more eigenvectors of M . This implies that all the eigenvalues of M are given by $P_m^v = e^{K_m/2} v_m$, and the eigenvalues of P^v are λ_i . In other words, M and S share the same eigenvalues.

(3) $\mu_a = \frac{1}{Z} e^{K_a}$ is an eigenvector with eigenvalue 1, and hence is the stationary state:

$$\begin{aligned} \sum_b M_{ab} \mu_b &= \frac{1}{Z} \sum_b e^{(K_a + K_b)/2} S_{ab} \\ &= \frac{1}{Z} e^{K_a} \sum_b e^{K_b/2} S_{ba} e^{-K_a/2} = \mu_a \sum_b M_{ba} = \mu_a. \end{aligned} \quad (\text{B10})$$

The previous facts then let us finish the calculation:

$$\begin{aligned} \left\langle \frac{P(a,b)}{P(a)P(b)} \right\rangle &= \sum_{ab} \left(e^{K_a} (S^\tau)_{ab}^2 e^{-K_b} \right) (e^{K_b - K_a}) \\ &= \sum_{ab} \left(e^{K_a} (S^\tau)_{ab}^2 e^{-K_b} \right) (e^{K_b - K_a}) \\ &= \sum_{ab} (S^\tau)_{ab}^2 = \|S^\tau\|^2. \end{aligned} \quad (\text{B11})$$

Now using the fact that $\|A\|^2 = \text{tr}(A^T A)$ and is therefore invariant under an orthogonal change of basis, we find that

$$\left\langle \frac{P(a,b)}{P(a)P(b)} \right\rangle = \sum_i |\lambda_i|^{2\tau}. \quad (\text{B12})$$

Since the λ_i 's are both the eigenvalues of M and S , and since M is irreducible and aperiodic, there is exactly one eigenvalue $\lambda_1 = 1$, and all other eigenvalues are less than one. Altogether,

$$I_R(t_1, t_2) = \left\langle \frac{P(a,b)}{P(a)P(b)} \right\rangle - 1 = \sum_{i=2} |\lambda_i|^{2\tau}. \quad (\text{B13})$$

Hence one can easily estimate the asymptotic behavior of the mutual information if one has knowledge of the spectrum of M . We see that the mutual information exponentially decays, with a decay scale time-scale given by the second largest eigenvalue λ_2 :

$$\tau_{\text{decay}}^{-1} = 2 \log \frac{1}{\lambda_2}. \quad (\text{B14})$$

6. Hidden Markov Model

In this subsection, we generalize our findings to hidden Markov models and present a proof of Theorem 2. Based on the considerations in the main body of the text, the joint probability distribution between two visible states X_{t_1}, X_{t_2} is given by

$$P(a,b) = \sum_{cd} G_{bd} [(M^\tau)_{dc} \mu_c] G_{ac}, \quad (\text{B15})$$

where the term in brackets would have been there in an ordinary Markov model and the two new factors of G

are the result of the generalization. Note that as before, μ is the stationary state corresponding to \mathbf{M} . We will only consider the typical case where \mathbf{M} is aperiodic, irreducible, and non-degenerate; once we have this case, the other cases can be easily treated by mimicking our above proof for ordinary Markov processes. Using equation (7) and defining $\mathbf{g} = \mathbf{M}\mu$ gives

$$\begin{aligned} P(a,b) &= \sum_{cd} G_{bd} [(M^\tau)_{dc} \mu_c] G_{ac} \\ &= g_a g_b + \lambda_2^\tau \sum_{cd} (G_{bd} A_{dc} \mu_c G_{ac}). \end{aligned} \quad (\text{B16})$$

Plugging this in to our definition of rational mutual information gives

$$\begin{aligned} I_R + 1 &= \sum_{ab} \frac{P(a,b)^2}{g_a g_b} \\ &= \sum_{ab} \left(g_a g_b + \lambda_2^\tau \sum_{cd} G_{bd} A_{dc} \mu_c G_{ac} \right) \\ &\quad + \lambda_2^{2\tau} \mathcal{C} \\ &= 1 + \lambda_2^\tau \sum_{cd} A_{dc} \mu_c + \lambda_2^{2\tau} \mathcal{C} \\ &= 1 + \lambda_2^{2\tau} \mathcal{C}, \end{aligned} \quad (\text{B17})$$

where we have used the facts that $\sum_i G_{ij} = 1$, $\sum_i A_{ij} = 0$, and as before \mathcal{C} is asymptotically constant. This shows that $I_R \propto \lambda_2^{2\tau}$ exponentially decays.

Appendix C: Power laws for generative grammars

In this appendix, we prove that the rational mutual information decays like a power law for a sub-class of generative grammars. We proceed by mimicking the strategy employed in the above appendix. Let \mathbf{G} be the linear operator associated with the matrix $P_{b|a}$, the probability that a node takes the value b given that the parent node has value a . We will assume that \mathbf{G} is irreducible and aperiodic, with no degeneracies. From the above discussion, we see that removing the degeneracy assumption does not qualitatively change things; one simply replaces the procedure of diagonalizing \mathbf{G} with putting \mathbf{G} in Jordan normal form.

Let us start with the weakly correlated case. In this case,

$$P(a,b) = \sum_r \mu_r \left(G^{\Delta/2} \right)_{ar} \left(G^{\Delta/2} \right)_{br}, \quad (\text{C1})$$

since as we have discussed in the main text, the parent node has the stationary distribution μ and $\mathbf{G}^{\Delta/2}$ give the conditional probabilities from transitioning from the parent node to the nodes at the bottom of the tree that we are interested in. We now employ our favorite trick of diagonalizing \mathbf{G} and then writing

$$(G^{\Delta/2})_{ij} = \mu_i + \lambda_2^{\Delta/2} A_{ij}, \quad (\text{C2})$$

which gives

$$\begin{aligned} P(a,b) &= \sum_r \mu_r \left(\mu_a + \lambda_2^{\Delta/2} A_{ar} \right) \left(\mu_b + \lambda_2^{\Delta/2} A_{br} \right), \\ &= \sum_r \mu_r \left(\mu_a \mu_b + \mu_a \epsilon A_{br} + \mu_b \epsilon A_{ar} + \epsilon^2 A_{ar} A_{br} \right) \end{aligned} \quad (\text{C3})$$

where we have defined $\epsilon = \lambda_2^{\Delta/2}$. Now note that $\sum_r A_{ar} \mu_r = 0$, since $\boldsymbol{\mu}$ is an eigenvector with eigenvalue 1 of $\mathbf{G}^{\Delta/2}$. Hence this simplifies the above to just

$$P(a,b) = \mu_a \mu_b + \epsilon^2 \sum_r \mu_r A_{ar} A_{br}. \quad (\text{C4})$$

From the definition of rational mutual information, and employing the fact that $\sum_i A_{ij} = 0$ gives

$$\begin{aligned} I_R + 1 &\approx \sum_{ab} \frac{(\mu_a \mu_b + \epsilon^2 \sum_r \mu_r A_{ar} A_{br})^2}{\mu_a \mu_b} \\ &= \sum_{ab} [\mu_a \mu_b + \epsilon^4 N_{ab}^2], \\ &= 1 + \epsilon^4 \|\mathbf{N}\|^2, \end{aligned} \quad (\text{C5})$$

where $N_{ab} \equiv (\mu_a \mu_b)^{-1/2} \sum_r \mu_r A_{ar} A_{br}$ is a symmetric matrix and $\|\cdot\|$ denotes the Frobenius norm. Hence

$$I_R = \lambda_2^{2\Delta} \|S\|^2. \quad (\text{C6})$$

Let us now generalize to the strongly correlated case. As discussed in the text, the joint probability is modified to

$$P(a,b) = \sum_{rs} Q_{rs} \left(G^{\Delta/2-1} \right)_{ar} \left(G^{\Delta/2-1} \right)_{bs}, \quad (\text{C7})$$

where Q is some symmetric matrix which satisfies $\sum_r Q_{rs} = \mu_s$. We now employ our favorite trick of diagonalizing \mathbf{G} and then writing

$$(G^{\Delta/2})_{ij} = \mu_i + \epsilon A_{ij}, \quad (\text{C8})$$

where $\epsilon \equiv \lambda_2^{\Delta/2-1}$. This gives

$$\begin{aligned} P(a,b) &= \sum_{rs} Q_{rs} (\mu_a + \epsilon A_{ar}) (\mu_b + \epsilon A_{bs}), \\ &= \mu_a \mu_b + \sum_{rs} Q_{rs} (\mu_a \epsilon A_{bs} + \mu_b \epsilon A_{ar} + \epsilon^2 A_{ar} A_{bs}) \\ &= \mu_a \mu_b + \sum_s \mu_a \epsilon A_{bs} \mu_s + \sum_r \mu_b \epsilon A_{ar} \mu_r \\ &\quad + \epsilon^2 \sum_{rs} Q_{rs} A_{ar} A_{bs} \\ &= \mu_a \mu_b + \epsilon^2 \sum_{rs} Q_{rs} A_{ar} A_{bs}. \end{aligned} \quad (\text{C9})$$

Now defining the symmetric matrices $R_{ab} \equiv \sum_{rs} Q_{rs} A_{ar} A_{bs} \equiv (\mu_a \mu_b)^{1/2} N_{ab}$, and noting that $\sum_a R_{ab} = 0$, we have

$$\begin{aligned} I_R + 1 &= \sum_{ab} \frac{(\mu_a \mu_b + \epsilon^2 R_{ab})^2}{\mu_a \mu_b} \\ &= \sum_{ab} [\mu_a \mu_b + \epsilon^4 N_{ab}^2], \\ &= 1 + \epsilon^4 \|\mathbf{N}\|^2, \end{aligned} \quad (\text{C10})$$

which gives

$$I_R = \lambda_2^{2\Delta-4} \|\mathbf{N}\|^2. \quad (\text{C11})$$

In either the strongly or the weakly correlated case, note that \mathbf{N} is asymptotically constant. We can write the second largest eigenvalue $|\lambda_2|^2 = q^{-k_2/2}$, where q is the branching factor,

$$I_R \propto q^{-\Delta k_2/2} \propto q^{-k_2 \log_q |i-j|} = \mathcal{C} |i-j|^{-k_2}. \quad (\text{C12})$$

Behold the glorious power law! We note that the normalization \mathcal{C} must be a function of the form $\mathcal{C} = m_2 f(\lambda_2, q)$, where m_2 is the multiplicity of the eigenvalue λ_2 . We evaluate this normalization in the next section.

As before, this result can be sharpened if we assume that \mathbf{G} satisfies detailed balance $G_{mn} = e^{K_m/2} S_{mn} e^{-K_n/2}$ where \mathbf{S} is a symmetric matrix and K_n are just numbers. Let us only consider the weakly correlated case. By the spectral theorem, we diagonalize S into an orthonormal basis of eigenvectors v . As before, G and S share the same eigenvalues. Proceeding,

$$P(a,b) = \frac{1}{Z} \sum_v \lambda_v^\Delta v_a v_b e^{(K_a + K_b)/2}, \quad (\text{C13})$$

where Z is a constant that ensures that P is properly normalized. Let us move full steam ahead to compute the rational mutual information:

$$\begin{aligned} \sum_{ab} \frac{P(a,b)^2}{P(a)P(b)} &= \sum_{ab} e^{-(K_a + K_b)} \left(\sum_v \lambda_v^\Delta v_a v_b e^{(K_a + K_b)/2} \right)^2 \\ &= \sum_{ab} \left(\sum_v \lambda_v^\Delta v_a v_b \right)^2. \end{aligned} \quad (\text{C14})$$

This is just the Frobenius norm of the symmetric matrix $H \equiv \sum_v \lambda_v^\Delta v_a v_b!$ The eigenvalues of the matrix can be read off, so we have

$$I_R(a,b) = \sum_{i=2} |\lambda_i|^{2\Delta}. \quad (\text{C15})$$

Hence we have computed the rational mutual information exactly as a function of Δ . In the next section, we use this result to compute the mutual information as a function of separation $|i - j|$, which will lead to a precise evaluation of the normalization constant \mathcal{C} in the equation

$$I(a, b) \approx \mathcal{C}|i - j|^{-k_2}. \quad (\text{C16})$$

1. Detailed evaluation of the normalization

For simplicity, we specialize to the case $q = 2$ although our results can surely be extended to $q > 2$. Define $\delta = \Delta/2$ and $d = |i - j|$. We wish to compute the expected value of I_R conditioned on knowledge of d . By Bayes rule, $p(\delta|d) \propto p(d|\delta)p(\delta)$. Now $p(d|\delta)$ is given by a triangle distribution with mean $2^{\delta-1}$ and compact support $(0, 2^\delta)$. On the other hand, $p(\delta) \propto 2^\delta$ for $\delta \leq \delta_{\max}$ and $p(\delta) = 0$ for $\delta \leq 0$ or $\delta > \delta_{\max}$. This new constant δ_{\max} serves two purposes. First, it can be thought of as a way to regulate the probability distribution $p(\delta)$ so that it is normalizable; at the end of the calculation we formally take $\delta_{\max} \rightarrow \infty$ without obstruction. Second, if we are interested in empirically sampling the mutual information, we cannot generate an infinite string, so setting δ_{\max} to a finite value accounts for the fact that our generated string may be finite.

We now assume $d \gg 1$ so that we can swap discrete sums with integrals. We can then compute the conditional expectation value of $2^{-k_2\delta}$. This yields

$$I_R \approx \int_0^\infty 2^{-k_2\delta} P(d|\delta) d\delta = \frac{(1 - 2^{-k_2}) d^{-k_2}}{k_2(k_2 + 1) \log(2)}, \quad (\text{C17})$$

or equivalently,

$$\mathcal{C}_{q=2} = \frac{1 - |\lambda_2|^4}{k_2(k_2 + 1) \log 2}. \quad (\text{C18})$$

It turns out it is also possible to compute the answer without making any approximations with integrals:

$$I_R = \frac{2^{-(k_2+1)\lceil \log_2(d) \rceil} ((2^{k_2+1} - 1) 2^{\lceil \log_2(d) \rceil} - 2d(2^{k_2} - 1))}{2^{k_2+1} - 1}. \quad (\text{C19})$$

The resulting predictions are compared in figure Figure 5.

Appendix D: Estimating (rational) mutual information from empirical data

Estimating mutual information or rational mutual information from empirical data is fraught with subtleties.

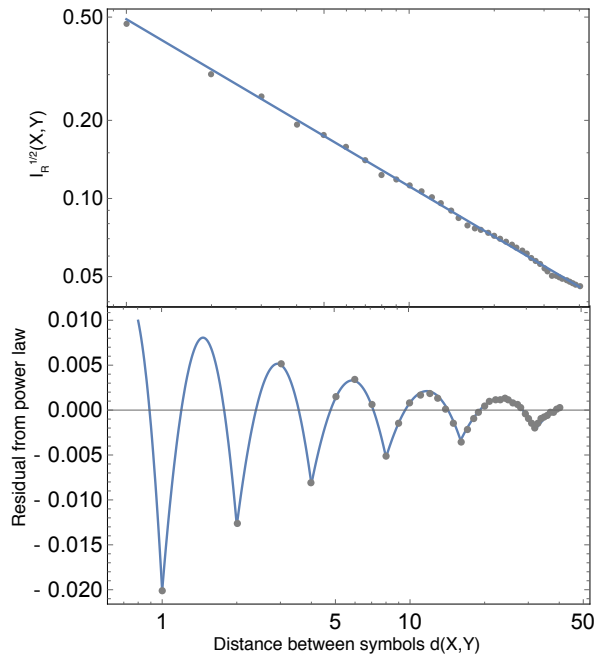


FIG. 5: Decay of rational mutual information with separation for a binary sequence from a numerical simulation with probabilities $p(0|0) = p(1|1) = 0.9$ and a branching factor $q = 2$. The blue curve is *not* a fit to the simulated data but rather an analytic calculation. The smooth power law displayed on the left is what is predicted by our “continuum” approximation. The very small discrepancies (right) are not random but are fully accounted for by more involved exact calculations with discrete sums.

It is well known that a naive estimate of the Shannon entropy obtained $\hat{S} = -\sum_{i=1}^K \frac{N_i}{N} \log \frac{N_i}{N}$ is biased, generally underestimating the true entropy from finite samples. For example, We use the estimator advocated by Grassberger [53]:

$$\hat{S} = \log N - \frac{1}{N} \sum_{i=1}^K N_i \psi(N_i), \quad (\text{D1})$$

where $\psi(x)$ is the digamma function, $N = \sum N_i$, and K is the number of characters in the alphabet. The mutual information estimator can then be estimated by $\hat{I}(X, Y) = \hat{S}(X) + \hat{S}(Y) - \hat{S}(X, Y)$. The variance of this estimator is then the sum of the variances

$$\text{var}(\hat{I}) = \text{varEnt}(X) + \text{varEnt}(Y) + \text{varEnt}(X, Y), \quad (\text{D2})$$

where the varEntropy is defined as

$$\text{varEnt}(X) = \text{var}(-\log p(X),) \quad (\text{D3})$$

where we can again replace logarithms with the digamma function ψ . The uncertainty after N measurements is then $\approx \sqrt{\text{var}(\hat{I})/N}$.

To compare our theoretical results with experiment in Fig. 4, we must measure the rational mutual information

for a binary sequence from (simulated) data. For a binary sequence with covariance coefficient $\rho(X, Y) = P(1, 1) - P(1)^2$, the rational mutual information is

$$I_R(X, Y) = \left(\frac{\rho(X, Y)}{P(0)P(1)} \right)^2. \quad (\text{D4})$$

This was essentially calculated in [54] by considering the limit where the covariance coefficient is small $\rho \ll 1$. In their paper, there is an erroneous factor of 2. To estimate covariance $\rho(d)$ as a function of d (sometimes confusingly referred to as the correlation function), we use the unbiased estimator for a data sequence $\{x_1, x_2, \dots, x_n\}$:

$$\hat{\rho}(d) = \frac{1}{n-d-1} \sum_{i=1}^{n-d} (x_i - \bar{x})(x_{i+d} - \bar{x}). \quad (\text{D5})$$

However, it is important to note that estimating the covariance function ρ by averaging and then squaring will generically yield a biased estimate; we circumvent this by simply estimating $I_R(X, Y)^{1/2} \propto \rho(X, Y)$.

-
- [1] P. Bak, Physical Review Letters **59**, 381 (1987).
- [2] P. Bak, C. Tang, and K. Wiesenfeld, Physical review A **38**, 364 (1988).
- [3] K. Linkenkaer-Hansen, V. V. Nikouline, J. M. Palva, and R. J. Ilmoniemi, The Journal of Neuroscience **21**, 1370 (2001), <http://www.jneurosci.org/content/21/4/1370.full.pdf+html>, URL <http://www.jneurosci.org/content/21/4/1370.abstract>.
- [4] D. J. Levitin, P. Chordia, and V. Menon, Proceedings of the National Academy of Sciences **109**, 3716 (2012).
- [5] M. Tegmark, ArXiv e-prints (2014), 1401.1219.
- [6] G. K. Zipf, *Human behavior and the principle of least effort* (Addison-Wesley Press, 1949).
- [7] H. W. Lin and A. Loeb, Physical Review E **93**, 032306 (2016).
- [8] L. Pietronero, E. Tosatti, V. Tosatti, and A. Vespignani, Physica A: Statistical Mechanics and its Applications **293**, 297 (2001), ISSN 0378-4371, URL <http://www.sciencedirect.com/science/article/pii/S0378437100006336>.
- [9] M. Kardar, *Statistical physics of fields* (Cambridge University Press, 2007).
- [10] URL ftp://ftp.ncbi.nih.gov/genomes/Homo_sapiens/.
- [11] URL http://www.jsbach.net/midi/midi_solo_violin.html.
- [12] URL <http://prize.hutter1.net/>.
- [13] URL <http://www.lexique.org/public/lisezmoi.corpatext.htm>.
- [14] A. M. Turing, Mind **59**, 433 (1950).
- [15] D. Ferrucci, E. Brown, J. Chu-Carroll, J. Fan, D. Gondek, A. A. Kalyanpur, A. Lally, J. W. Murdock, E. Nyberg, J. Prager, et al., AI magazine **31**, 59 (2010).
- [16] M. Campbell, A. J. Hoane, and F.-h. Hsu, Artificial intelligence **134**, 57 (2002).
- [17] V. Mnih, Nature **518**, 529 (2015).
- [18] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al., Nature **529**, 484 (2016), URL <http://dx.doi.org/10.1038/nature16961>.
- [19] Y. Kim, Y. Jernite, D. Sontag, and A. M. Rush (2015), 1508.06615, URL <https://arxiv.org/abs/1508.06615>.
- [20] A. Graves, ArXiv e-prints (2013), 1308.0850.
- [21] A. Graves, A.-r. Mohamed, and G. Hinton, in *2013 IEEE international conference on acoustics, speech and signal processing* (IEEE, 2013), pp. 6645–6649.
- [22] R. Collobert and J. Weston, in *Proceedings of the 25th international conference on Machine learning* (ACM, 2008), pp. 160–167.
- [23] J. Schmidhuber, Neural Networks **61**, 85 (2015).
- [24] Y. LeCun, Y. Bengio, and G. Hinton, Nature **521**, 436 (2015).
- [25] R. Meir and E. Domany, Phys. Rev. Lett. **59**, 359 (1987), URL <http://link.aps.org/doi/10.1103/PhysRevLett.59.359>.
- [26] K. Hornik, M. Stinchcombe, and H. White, Neural networks **2**, 359 (1989).
- [27] A. M. Saxe, J. L. McClelland, and S. Ganguli, arXiv preprint arXiv:1312.6120 (2013).
- [28] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio, in *Advances in neural information processing systems* (2014), pp. 2924–2932.
- [29] H. Mhaskar, Q. Liao, and T. Poggio, ArXiv e-prints (2016), 1603.00988.
- [30] M. Bianchini and F. Scarselli, IEEE transactions on neural networks and learning systems **25**, 1553 (2014).
- [31] B. Poole, S. Lahiri, M. Raghu, J. Sohl-Dickstein, and S. Ganguli, ArXiv e-prints (2016), 1606.05340.
- [32] M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein, ArXiv e-prints (2016), 1606.05336.
- [33] P. Mehta and D. J. Schwab, ArXiv e-prints (2014), 1410.3831.
- [34] S. Hochreiter and J. Schmidhuber, Neural computation **9**, 1735 (1997).
- [35] C. E. Shannon, ACM SIGMOBILE Mobile Computing and Communications Review **5**, 3 (1948).
- [36] S. Kullback and R. A. Leibler, Ann. Math. Statist. **22**, 79 (1951), URL <http://dx.doi.org/10.1214/aoms/1177729694>.
- [37] T. M. Cover and J. A. Thomas, *Elements of information theory* (John Wiley & Sons, 2012).
- [38] L. R. Rabiner, Proceedings of the IEEE **77**, 257 (1989).
- [39] R. C. Carrasco and J. Oncina, in *International Colloquium on Grammatical Inference* (Springer, 1994), pp. 139–152.
- [40] N. Chomsky, *Aspects of the Theory of Syntax*, vol. 11 (MIT press, 1965).
- [41] S. Ginsburg, *The Mathematical Theory of Context Free Languages.*[Mit Fig.] (McGraw-Hill Book Company, 1966).

- [42] T. L. Booth, in *Switching and Automata Theory, 1969., IEEE Conference Record of 10th Annual Symposium on* (IEEE, 1969), pp. 74–81.
- [43] T. Huang and K. Fu, *Information Sciences* **3**, 201 (1971), ISSN 0020-0255, URL <http://www.sciencedirect.com/science/article/pii/S0020025571800075>.
- [44] K. Lari and S. J. Young, *Computer speech & language* **4**, 35 (1990).
- [45] D. Harlow, S. H. Shenker, D. Stanford, and L. Susskind, *Physical Review D* **85**, 063516 (2012).
- [46] L. Van Hove, *Physica* **16**, 137 (1950).
- [47] J. A. Cuesta and A. Sánchez, *Journal of Statistical Physics* **115**, 869 (2004), cond-mat/0306354.
- [48] G. Evenbly and G. Vidal, *Journal of Statistical Physics* **145**, 891 (2011).
- [49] N. Chomsky, *Information and control* **2**, 137 (1959).
- [50] M. Mahoney, *Large text compression benchmark*.
- [51] A. Karpathy, J. Johnson, and L. Fei-Fei, *ArXiv e-prints* (2015), 1506.02078.
- [52] C. W. Gardiner et al., *Handbook of stochastic methods*, vol. 3 (Springer Berlin, 1985).
- [53] P. Grassberger, *ArXiv Physics e-prints* (2003), physics/0307138.
- [54] W. Li, *Journal of Statistical Physics* **60**, 823 (1990), ISSN 1572-9613, URL <http://dx.doi.org/10.1007/BF01025996>.