

# Thermodynamic Overfitting and Generalization: Energetic Limits on Predictive Complexity

Alexander B. Boyd,<sup>1,2,3,\*</sup> James P. Crutchfield,<sup>4,†</sup> Mile Gu,<sup>5,6,7,‡</sup> and Felix C. Binder<sup>1,2,§</sup>

<sup>1</sup>*School of Physics, Trinity College Dublin, Dublin 2, Ireland*

<sup>2</sup>*Trinity Quantum Alliance, Unit 16, Trinity Technology and Enterprise Centre, Pearse Street, Dublin 2, Ireland*

<sup>3</sup>*Beyond Institute for Theoretical Science, San Francisco, California, USA*

<sup>4</sup>*Complexity Sciences Center and Physics Department,*

*University of California at Davis, One Shields Avenue, Davis, CA 95616*

<sup>5</sup>*Nanyang Quantum Hub, School of Physical and Mathematical Sciences,*

*Nanyang Technological University, 637371, Singapore*

<sup>6</sup>*Centre for Quantum Technologies, National University of Singapore, 3 Science Drive 2, 117543, Singapore*

<sup>7</sup>*MajuLab, CNRS-UNS-NUS-NTU International Joint Research Unit, UMI 3654, 117543, Singapore*

(Dated: February 28, 2024)

Efficiently harvesting thermodynamic resources requires a precise understanding of their structure. This becomes explicit through the lens of information engines—thermodynamic engines that use information as fuel. Maximizing the work harvested using available information is a form of physically-instantiated machine learning that drives information engines to develop complex predictive memory to store an environment’s temporal correlations. We show that an information engine’s complex predictive memory poses both energetic benefits and risks. While increasing memory facilitates detection of hidden patterns in an environment, it also opens the possibility of thermodynamic overfitting, where the engine dissipates additional energy in testing. To address overfitting, we introduce thermodynamic regularizers that incur a cost to engine complexity in training due to the physical constraints on the information engine. We demonstrate that regularized thermodynamic machine learning generalizes effectively. In particular, the physical constraints from which regularizers are derived improve the performance of learned predictive models. This suggests that the laws of physics jointly create the conditions for emergent complexity and predictive intelligence.

Keywords: nonequilibrium thermodynamics, Maxwell’s demon, Landauer’s Principle, machine learning, generalization

## I. INTRODUCTION

Modern machine learning has made remarkable advances mimicking our understanding of biological intelligence by incorporating biology’s fundamental building blocks, e.g., in the form of neural networks [1]. Similarly, thermodynamics and statistical mechanics have made considerable contributions to machine learning [2, 3]. Most recently, this has come to include proposals to directly use thermodynamic systems for machine learning tasks [4, 5]. Meanwhile, thermodynamic concepts in machine learning found their way back to elucidating biological intelligence [6, 7]. This cross-fertilisation falls in line with a larger concern about how biological intelligence is embodied [8]. That is, the physicality of an intelligent agent, as manifested through its environment interactions, determines how it learns about the world. In this way, it has been recognised that machine learning and pattern recognition “can be viewed as two facets of the same field” [9].

With this in mind, the following explores the thermodynamics of pattern prediction by analyzing a learning

agent who utilises information provided by its environment for the aim of maximal energy extraction. This *thermodynamic machine learning* (TML) implements the *principle of maximum work production* [10] that expresses how thermodynamic efficiency is tied to optimal prediction [11]. Specifically, selecting the model that harvests the most work from an information source is equivalent to performing *maximum likelihood estimation* (MLE) on predictive models of data. This highlights the thermodynamic roots of *computational mechanics* [12–14]—the information theory of time-series prediction and structural complexity. While maximum work production is by no means guaranteed in general nonequilibrium physical processes [15], the equivalence between machine learning and thermodynamic resource maximization does offer a mechanism that drives the emergence of structural complexity and intelligence.

Here, we explore the central role of model complexity in thermodynamic machine learning by showing that principles of prediction arise from simple physical principles—work maximization, autocorrection, and model initialization. Therein, a concern arises about overfitting which occurs in machine learning when a model performs well on training data, but fails to effectively predict further inputs (test data). In this case, the estimated model is overly-specific to the training data and does not generalize to further samples [16]. Overfitting is closely tied to high model complexity, because

\* alecboy@gmail.com

† chaos@ucdavis.edu

‡ mgu@quantumcomplexity.org

§ quantum@felix-binder.net

it often happens when the number of model parameters exceed what is justified by a limited dataset.

In the thermodynamic setting, there is the (conversely) related *principle of requisite complexity* which states that an information engine must (at least) match the structural complexity of the environment to operate efficiently [17]. In contrast, we have *thermodynamic overfitting* which reflects a thermodynamic cost to excessively complex models. Practically, this cost appears to prohibit learning exceedingly large models that may exhibit “double descent” [18] where a regime of improved learning occurs as model complexity increases beyond the regime of overfitting.

Paralleling the strategy in conventional machine learning, to mitigate overfitting we turn to regularization by incurring a training penalty that reflects model complexity [19, 20]. Specifically, we introduce a *thermodynamic regularizer*: a thermodynamic cost to model complexity that is proportional to the work that is dissipated. Algorithmically, adding this regularizer results in Bayesian updates of a predictive model’s edge-weights—the weights that control the engine’s operation. We also introduce a cost to autocorrect the engine’s predictive states, which arises from starting in a uniform distribution over the engine’s memory states. The net result is an effective regularization function—essentially, a new complexity measure that quantifies the degree to which different causal predictive states make different predictions.

Notably, the following derives a thermodynamically-based prediction algorithm that parallels many modern machine learning methods for prediction. These include reservoir computing [21], backpropagation through unrolling time in recurrent neural networks [22], and transformers [23] which underlie the marked effectiveness of modern large language models, such as ChatGPT.

To begin the development, the following section briefly introduces thermodynamic machine learning from basic principles. Section III describes memory-constrained work optimisation. We then characterize the performance of information engines learned through maximum work production by calculating the work production rate for both training and test data. This leads to the following main results:

1. We derive an exact expression for the asymptotic work rate of the information engine in Thm. 1, relating the engine’s estimated predictive model and the true predictive model of the input process.
2. We greatly simplify the search for the maximum-work engine by analytically deriving the engine parameters in Thm. 2.
3. We identify thermodynamic overfitting through divergent dissipated work in Fig. 6.
4. We introduce thermodynamic regularizers that add a cost to model complexity during training. This results in thermodynamic learning that generalizes and avoids overfitting; see Figs. 7 and 8.

Altogether, these results demonstrate that thermody-

amic principles spontaneously produce effective predictive learning.

## II. THERMODYNAMIC MACHINE LEARNING

Thermodynamic machine learning arises when a physical agent tries to extract energy from a complex, noisy environment. Facing such an environment the agent predicts the environment’s behavior and then converts that knowledge into useful work. This section describes the salient aspects of this process. Here, the “agent” may be understood as a version of Maxwell’s Demon [24–26] confronted with correlated patterns.

The work value of information has been thoroughly explored since Szilard’s proposal of his eponymous information engine [27] and the introduction of Landauer’s erasure principle [28]. In essence, given information-bearing degrees of freedom characterised by a random variable  $Y$  we may, on average, extract an amount of work upper-bounded by  $W_{ext}$  [29]:

$$\beta W_{ext} = \Delta H \equiv H[Y'] - H[Y] \quad (1)$$

by converting  $Y$  to an output random variable  $Y'$  under coupling with an external heat bath at inverse temperature  $\beta$ . Here,  $H[Y] = -\sum_{y \in \mathcal{Y}} \Pr(Y = y) \ln \Pr(Y = y)$  is the Shannon entropy of random variable  $Y$ . For clarity,  $\mathcal{Y}$  is the physical system,  $y \in \mathcal{Y}$  are realizations of states from that system, and  $Y$  is the random variable that determines the distribution over those system states via  $\{\Pr(Y = y)\}_{y \in \mathcal{Y}}$ .

We leave aside the ongoing investigation into the conditions for saturating this relationship [30–32] and consider a procedure that saturates the bound. The maximal value for  $\Delta H$  is attained by maximally randomizing  $Y'$  such that  $H[Y'] = \ln |\mathcal{Y}|$ , where  $|\mathcal{Y}|$  is the number of distinct values that  $Y$  may take. An agent that fully extracts this maximal amount of work from an information source  $Y$  is thermodynamically efficient, because it has harvested all available free energy.

We can further decompose the average work extracted by an efficient agent expressed in Eq. (1) into its single-shot elements:

$$\beta \langle W(y) \rangle = \ln \Pr(Y = y) + \ln |\mathcal{Y}|. \quad (2)$$

This quantifies the average extractable work from a probabilistically occurring realization  $y$  in the maximum-extraction limit of full randomization.

This is the necessary single-shot work extraction for an *efficient* agent, because it must have zero total entropy production  $\langle \Sigma \rangle \equiv \langle Q \rangle / T + k_B \Delta H = 0$  on average and therefore must also have zero fluctuations in entropy production  $\Sigma(y) \equiv Q(y) / T + k_B \ln \frac{p_y}{|\mathcal{Y}|} = 0$  [10, 33]. Here,  $Q/T$  is the entropy change in the environment due to heat  $Q$  and  $\Delta H$  is the change in entropy of the system  $\mathcal{Y}$ , which together account for the total entropy pro-

duced  $\Sigma$  [34, 35].

In the case where the process transforms an information reservoir, the system starts and ends in an energetically degenerate configuration such that the work production and heat are equal with opposite signs. To minimize the entropy production when harvesting energy from  $\mathcal{Y}$  through a quasistatic protocol, the probability of each outcome must be encoded in the initial energy landscape [10, 36]. The parameters of an agent's estimated input are explicitly encoded in its evolving energy landscape.

Here, we are here interested in the situation where  $Y$  represents a stochastic process taking values  $y_{0:L} \equiv y_0 y_1 \cdots y_{L-1}$  generated by a model  $\theta$  with probability:

$$\Pr(Y_{0:L}^\theta = y_{0:L}). \quad (3)$$

That is, we consider sequences of length  $L$ , indicated by subscript  $0:L$ . The work extracted by an efficient agent is then a specific case of Eq. (2) with  $Y = Y_{0:L}^\theta$  [10, 37]:

$$\beta \langle W^\theta(y_{0:L}) \rangle = \ln \Pr(Y_{0:L}^\theta = y_{0:L}) + L \ln |\mathcal{Y}|, \quad (4)$$

However, there is an important conceptual addition. To extract work from  $Y_{0:L}^\theta$ , the agent must interact with each element in sequence. Thus, the agent cannot establish a simultaneous energy landscape over all elements. Instead, the engine requires a memory system  $\mathcal{X}$  that tracks what has already been observed in the sequence such that it can make optimal estimates of the next input. A work extraction device that has access to such an internal memory is called an *information engine*. Such a device is a type of information ratchet [38, 39], a physically instantiated type of stochastic Turing machine, also known as a ‘‘Brownian computer’’ [40]. Figure 1 provides an illustration. Specifically, an information engine is a ratchet whose functionality is reduced to maximally randomizing the outputs while storing all relevant correlations in the machine memory.

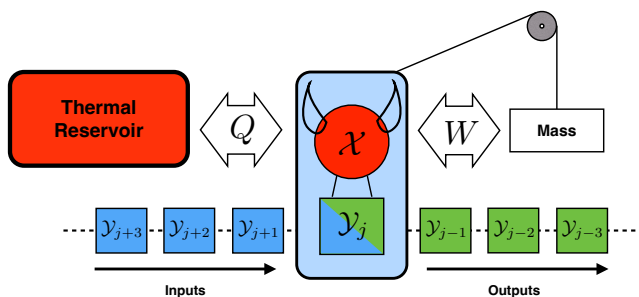


FIG. 1. An information engine: A physical device that can maximize work production using correlations on an information tape. It transforms an input sequence  $Y_{0:L}$  to an output sequence  $Y'_{0:L}$  while exchanging energy between a work reservoir (represented by the mass on a string) and thermal reservoir at inverse temperature  $\beta = 1/k_B T$ . The input  $Y_j$  and output  $Y'_j$  are stored in the physical system  $\mathcal{Y}_j$ .

Faced with the output of a stochastic process  $\mathcal{P}$ , a thermodynamic agent must find a good process model  $\theta$  among a candidate family of models  $\Theta$ . Denoting a given model  $\theta$ 's output distribution as  $Y_{0:L}^\theta$ , the likelihood of  $\theta$  being the generator of a sequence  $y_{0:L}$  is:

$$\ell(\theta|y_{0:L}) \equiv \Pr(Y_{0:L}^\theta = y_{0:L}). \quad (5)$$

To learn a model of the process  $\mathcal{P}$  the agent may then employ *maximum-likelihood estimation* (MLE) to select the model  $\theta$  that maximizes the likelihood:

$$\Theta^{\text{MLE}}(y_{0:L}) = \operatorname{argmax}_{\theta \in \Theta} \ell(y_{0:L}|\theta). \quad (6)$$

MLE is one of the most general techniques in machine learning. The resulting estimated distribution can be used for a variety of other learning tasks, including classification [41].

Next, let's consider in more detail how the agent models the process from which work is to be extracted. Note that generally any such process can be described by a *hidden Markov model* (HMM). Likewise, as the agent observes the process it must build its own model of the process. To minimize dissipation the information engine's memory must use the predictive states of the input process [17]. The agent's memory dynamics and energy landscape directly match a predictive model of  $Y_{0:L}^\theta$  [10]. The agent's task is thus to match its internal model to the process' true model.

A HMM is defined by hidden states  $s$  and transitions between them according to the probabilities:

$$T_{s \rightarrow s'}^{(y)} \equiv \Pr(Y_i = y, S_{i+1} = s' | S_i = s), \quad (7)$$

outputting a symbol  $y$  with the transition. Here,  $Y_i$  and  $S_i$  label the random variables corresponding to the output symbol and hidden state, respectively.  $y, s, s'$  denote their realizations.

There are many ways to predictively model a time series  $Y_{0:L}$ . Among different procedures we choose minimal, predictive HMMs called  $\epsilon$ -*machines*; see App. A along with examples there. For a given process there is no alternative predictive HMM that requires fewer hidden states and  $\epsilon$ -machines are sufficient for describing *any* stochastic process giving rise to it [13], even if nonstationary [10, 42]. In addition, they provide a prescription for designing an information engine that harvests all available free energy from that process [17]. Thus, we choose our set of candidate models  $\Theta \equiv \{\theta\}$  to be a subset of the class of  $\epsilon$ -machines.

Finding the engine that produces the most work from given data is then equivalent to finding the engine whose model produces that information with maximum likelihood. Thus, as illustrated in Fig. 2 thermodynamic machine learning is equivalent to maximum likelihood estimation over predictive models [10]. The input system encodes data, the information engine contains a model, and the performance of that engine (work production)

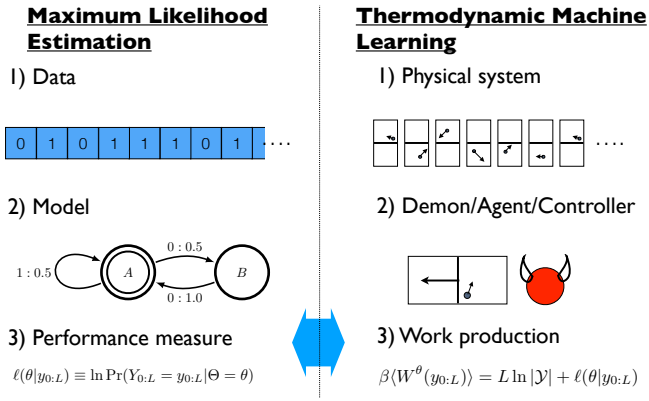


FIG. 2. Equivalence of thermodynamic machine learning and maximum likelihood estimation over predictive models. 1) MLE: Data (input information) is a realization of a random variable, such as a sequence of bits. TML: information is realized within a physical system, such as particles in partitioned boxes, each of which encodes a bit. 2) MLE: A model of the input information specifies the estimated probability of realizing data. TML: An efficient demon, agent, or controller of the physical system that contains the input information must have an internal model of its estimated input. 3) MLE: The performance of the model for a particular input data is given by the log-likelihood. TML: The performance of the efficient agent is measured by work production. The latter is linearly related to the log-likelihood. MLE and TML are equivalent estimation processes.

scales proportionally to the log-likelihood. As illustrated in Fig. 4, when performing thermodynamic learning on a collection of engines with models  $\Theta$ , we denote the maximum-work model for input  $y_{0:L}$ :

$$\Theta^{\max}(y_{0:L}) \equiv \operatorname{argmax}_{\theta \in \Theta} \langle W^\theta(y_{0:L}) \rangle. \quad (8)$$

The equivalence between TML and MLE means that the inferred models from both learning strategies are the same:

$$\Theta^{\max}(y_{0:L}) = \Theta^{\text{MLE}}(y_{0:L}). \quad (9)$$

Finally,  $\epsilon$ -machines can model any process and yield the same work production as any predictive model, meaning that it is sufficient to limit our model class  $\Theta$  to them. As described in App. A and Fig. 3, such machines are described by a causal update map on the hidden states (also called *predictive states*):

$$S_{i+1} = \epsilon(Y_i, S_i), \quad (10)$$

and edge-weights:

$$\theta(y|s) = \Pr(Y_i = y | S_i = s). \quad (11)$$

These are the explicit parameters we must explore through training [43]. The memory of the engine  $\mathcal{X}$  is a direct copy of its model's predictive state space

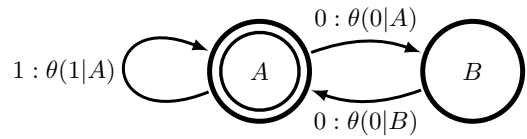


FIG. 3. Example  $\epsilon$ -machine: The Even Process with start state  $A$  is described by  $\theta$ , which is composed of the causal states  $\mathcal{S}$ , outputs  $\mathcal{Y}$ , start state  $s^*$ , topology  $S_{i+1} = \epsilon(Y_i, S_i)$ , and edge-weights  $\theta(Y_i|S_i)$ . The Even Process produces sequences of zeros in even numbers. The hidden states evolve according to  $\epsilon(0, A) = B$ ,  $\epsilon(1, A) = A$ ,  $\epsilon(0, B) = A$ , and transitions are taken with probabilities given by the edge weights  $\theta(1|A) = 0.5$ ,  $\theta(0|A) = 0.5$ , and  $\theta(0|B) = 1.0$ . Outputting a 0 from  $B$  has zero probability  $\theta(1|B) = 0.0$ , so we leave  $\epsilon(0, B)$  undefined. In this case, we chose the start state  $s^* = A$ . Altogether, we can describe this model with the set  $\theta = \{\mathcal{S}, \mathcal{Y}, s^*, \{\theta(y|s)\}_{y \in \mathcal{Y}, s \in \mathcal{S}}, \{\epsilon(y, s)\}_{y \in \mathcal{Y}, s \in \mathcal{S}}\}$ .

$\mathcal{S}$ . Once the maximum work model is found, we evaluate its complexity as the size of the engine's memory:  $C = \ln |\mathcal{S}| = \ln |\mathcal{X}|$ .

### III. CONSTRAINED MEMORY WORK MAXIMIZATION

Starting from the equivalence between work maximization and MLE, the following details an algorithm for discovering predictive models of training data  $y_{0:L}$  via TML. The class of  $\epsilon$ -machines is a particularly appropriate model class for learning as they can produce any process  $\Pr(Y_{0:\infty})$  given sufficiently many memory states (causal states) through the causal equivalence relation [12, 14]. Since  $\epsilon$ -machines are the most general model class, potentially any pattern is discoverable via TML.

However, the extreme generality of  $\epsilon$ -machines comes with a downside in learning. For any sequence  $y_{0:L}$   $\epsilon$ -machines include the process that produces it with unit probability. This means that if we allow our engine arbitrarily large memory  $n = |\mathcal{X}|$ , we can trivially maximize work production  $W^\theta(y_{0:L})$ . This corresponds to simply storing the training data in the engine's memory, rather than trying to discover the underlying pattern. In this case, any other word besides the training word would be expected with zero likelihood, making this the extreme limit of overfitting. This limiting case demonstrates that there are learning algorithms for which double-descent [18] does not apply. Moreover, it is natural to constrain memory, because it is an informational resource.

We consider the maximum-work model from the set  $\Theta_n$  of  $\epsilon$ -machines with  $n$  predictive states. This means that the engine is limited to  $n$  memory states:

$$\Theta_n^{\max}(y_{0:L}) \equiv \operatorname{argmax}_{\theta \in \Theta_n} \langle W^\theta(y_{0:L}) \rangle. \quad (12)$$

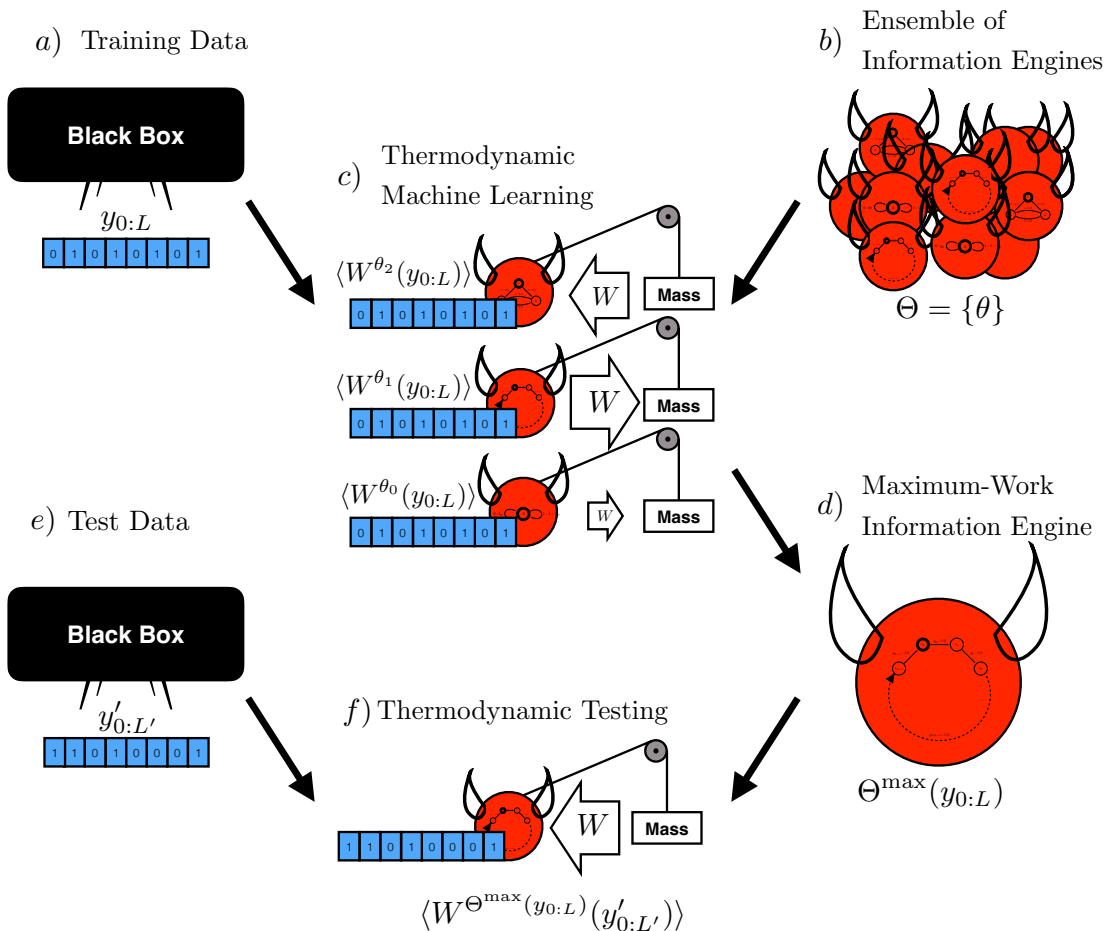


FIG. 4. Thermodynamic machine learning discovers the maximum-work information engine. This is followed by thermodynamic testing to validate its capacity to harvest energy using previously-unseen data: a) We start with training data produced from a black box. b) Each candidate information engine from our ensemble has an internal model, represented faintly by the  $\epsilon$ -machine within their red body. c) We search through the engines to find that with the best model by determining how much work each produces from the training data. d) Thermodynamic machine learning converges on the maximum-work engine. e) We take further test data from the black box of our environment. f) We feed the test data into our selected engine and track the work production to evaluate the effectiveness of the engine’s model at capturing the process generated by the black box.

This engine produces work during training equivalent to:

$$\langle W_n^{\max}(y_{0:L}) \rangle \equiv \max_{\theta \in \Theta_n} \langle W^\theta(y_{0:L}) \rangle. \quad (13)$$

This is the training we execute throughout our development here: Find the maximum-work  $n$ -state engine that corresponds to the maximum-likelihood model from the class of predictive  $n$ -state HMMs. Unlike other training algorithms—that rely on the convergence of numerical estimators—this training is essentially analytic. Given our chosen class of models, as long as we successfully enumerate all accessible topologies, we directly find the maximum-work engine among the candidates.

Recall that the outcomes of many machine learning algorithms change based on how the learning parameters are selected and whether local maxima can be escaped. In contrast, the TML algorithm always arrives at the same model for the same input data and definition of

work production. The following explains how this is implemented.

The challenge of finding the maximum-work model from such a general class of processes may at first seem like a daunting optimization process over the high dimensional space of symbol-labeled stochastic matrices that specify an  $\epsilon$ -machine’s Hidden Markov Model (HMM). For general classical and quantum HMMs, evaluating the likelihood of  $y_{0:L}$  requires a series of  $L$  linear operations [44]. However, the properties of unifilarity allow for a useful simplified expression for the work production in terms of the  $\epsilon$ -map and edge-weights, as shown in Appendix D:

$$\beta \langle W^\theta(y_{0:L}) \rangle = L \ln |\mathcal{Y}| + \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, s^*)). \quad (14)$$

For a particular topology  $\epsilon$ , start state  $s^*$ , and input word

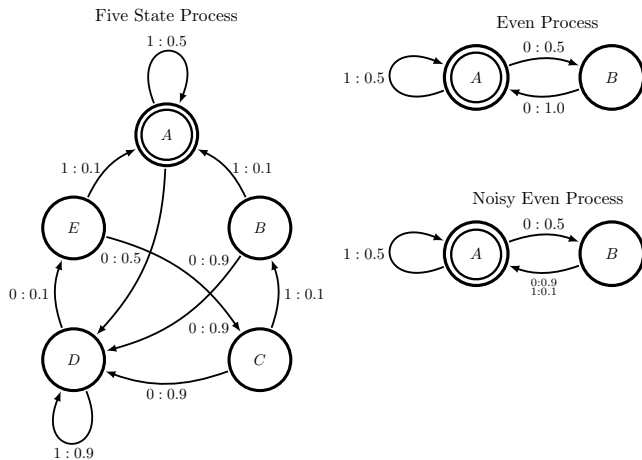


FIG. 5. Example  $\epsilon$ -machines: the “Five-State Process,” which is a randomly generated process with five causal states and full support, the “Even Process,” which produces sequences of 0s in even numbers, and the “Noisy Even Process” which adds some noise to the Even Process such that it has full support.

$y_{0:L}$ , we can analytically find the maximum-work edge-weights by counting the number of times that the input  $y$  is received by predictive state  $s$  due to the input-driven dynamics of the predictive state:

$$N(y, s|s^*, \epsilon, y_{0:L}) \equiv \sum_{i=0}^{L-1} \delta_{y, y_i} \delta_{s, \epsilon(y_{0:i}, s^*)}. \quad (15)$$

As shown in App. D, the work production can be rewritten:

$$\beta \langle W^\theta(y_{0:L}) \rangle = L \ln |\mathcal{Y}| + \sum_{s, y} N(y, s|s^*, \epsilon, y_{0:L}) \ln \theta(y|s). \quad (16)$$

The resulting maximum-work edge-weights are derived using the method of Lagrange multipliers. They are simply the fraction of times that predictive state  $s$  receives  $y$  when driven by  $y_{0:L}$ :

$$\Theta_{s^*, \epsilon, y_{0:L}}^{\max}(y|s) = \frac{N(y, s|s^*, \epsilon, y_{0:L})}{\sum_{y'} N(y', s|s^*, \epsilon, y_{0:L})}. \quad (17)$$

This is purely a frequentist estimate of the edge-weights based on how often each input visits each predictive causal state. This thermodynamically motivated engine design directly reflects the results of training a transformer, as described in Ref. [45].

Returning to analytically finding the maximum-likelihood model and maximum-work engine with  $n$  memory states, this reduces to checking every allowed topology [43]. In Fig. 6 we do this exactly for  $n \in \{1, 2, 3\}$  for a word  $y_{0:100}$ , for which we train on all intermediate length strings  $\{y_{0:L}\}_{L \in \{1, 2, \dots, 99, 100\}}$ . The word was

generated from the Five-State  $\epsilon$ -machine shown in Fig. 5. Since we are limited to models  $\theta$  with three or fewer causal states and the underlying model  $\theta'$  has five, we have technically misspecified our model class [46].

We consider the work production rate of the maximum work model  $\beta \langle W^{\max}(y_{0:L}) \rangle / L$  in Fig. 6, since it determines how much energy each symbol contributes on average during training. Taking the limit  $L \rightarrow \infty$ , this should approach  $\langle W^\theta \rangle_\infty$ —the average work extraction per bit for the engine, where  $\theta$  is the engine’s model.

While the next section derives an explicit expression for  $\langle W^\theta \rangle_\infty$ , for now we discuss the length- $L$  work production rate of the maximum-work engine. This is plotted with respect to three regions of asymptotic functionality:

1. Dud: An information engine that has nonpositive work production rate is a Dud, because it produces no useful energy.
2. Engine: An information engine that produces positive work is a functional engine.
3. Negative entropy production ( $\Sigma < 0$ ):  $\Sigma$  denotes the total entropy production in the thermal environment and information bearing degrees of freedom combined [35]. If the engine produces more work per bit than an engine that estimates a true model  $\langle W^{\theta'} \rangle_\infty$ , then it has extracted more energy than is available in the form of free energy per symbol. If it continues at this work production rate, the total entropy production will be negative on average and violate the Second Law of thermodynamics.

We also plot a dashed black line at  $\beta \langle W^\theta \rangle = \ln 2$ , which is the Landauer benefit of randomizing a bit [28].

Figure 6 shows that all the maximum-work information engines produce positive work and do not exceed Landauer’s bound. This is because it is always possible to find an engine that gains energy by interacting with a singular realization of the true model (the data). We also see that early in training (for short length strings  $L < \sim 20$ ) the training work rate often achieves the  $k_B T \ln 2$  Landauer limit on work production. This corresponds to discovering a predictive model that deterministically produces the training sequence with unit probability [47]. However, this is well into the regime of asymptotic Second Law violation if the true source is anything other than a deterministic repetition of the training string, meaning that entropy production for these instances is negative. This is possible since the data is a single realization: negative entropy fluctuations are allowed as long as the detailed fluctuation theorem is satisfied and entropy is nondecreasing *on average* [33].

As the training length increases in Fig. 6, a general trend appears: The training work-rate slowly ramps up, interspersed by sudden sharp dips. These dips happen at roughly the same point for all three curves ( $n = 1$ ,  $n = 2$ , and  $n = 3$ ). The slow ramps of increasing work rate correspond to the slow accumulation of low-suprival

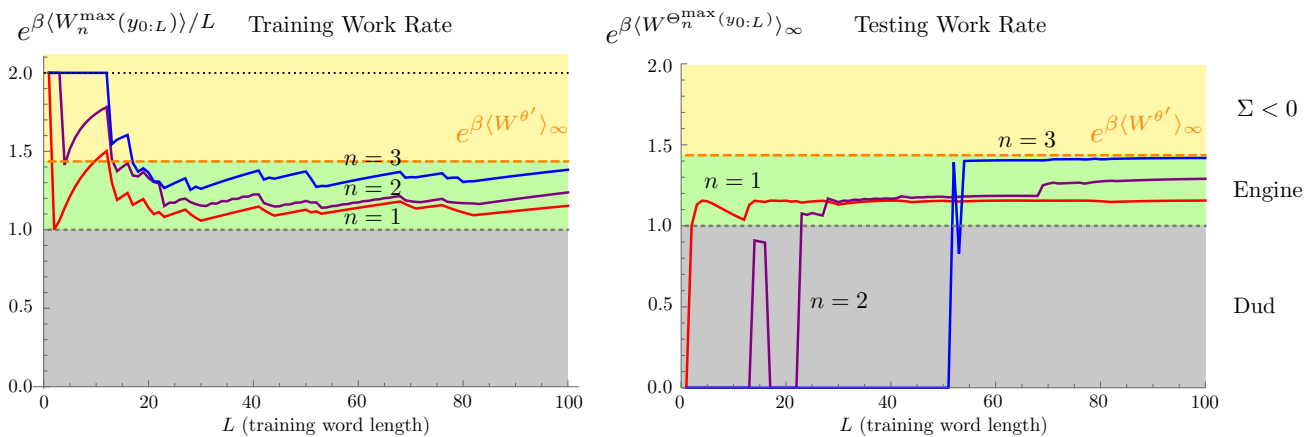


FIG. 6. Thermodynamic Overfitting: The work rate converges quickly on the training data (left), but this hides the divergent energy dissipation during testing (right) that corresponds to thermodynamic overfitting. Work production from training data increases as the number of the engine’s predictive states  $n$  increases. However, we see that asymptotic work production during testing is minimal ( $-\infty$ ) for small amounts of training data, especially for more complex engines, which corresponds to overfitting. For reference, we plot (i) the grey dashed line for zero work production  $\beta \langle W^{\theta} \rangle_{\infty} = 0$ , (ii) the orange dashed line for work production that corresponds to correctly estimating the true model  $\beta \langle W^{\theta'} \rangle_{\infty}$ , and (iii) the black dotted line for  $\beta \langle W^{\theta}(y_{0:L}) \rangle / L = \ln 2$  is the maximum work that can be harvested per bit from a single sequence. In the grey region below the grey dashed line, the engine is unable to harvest work on average, so it is a “Dud”. In the green region just above the grey dashed line, positive work is produced, effectively functioning as an “Engine”. Last, in the yellow region above the orange dashed line, we labeled the region with  $\Sigma < 0$  to indicate a negative entropy fluctuation.

symbols, where we can evaluate the surprisal of the next symbol  $-\ln \Pr(Y_L^{\theta} = y_L | Y_{0:L}^{\theta} = y_{0:L})$  via the probability conditioned on the past sequence. The ensuing dip corresponds to breaking the sequence with a rare high-surprisal element, significantly increasing the estimated input entropy rate and reducing the work rate.

Figure 6 also shows that memory size significantly affects work production. There is a consistent thermodynamic advantage to additional memory in training. Work production increases as memory increases  $\langle W_{n+1}^{\max}(y_{0:L}) \rangle \geq \langle W_n^{\max}(y_{0:L}) \rangle$ . This is as expected, because every  $n$ -state model can be described using  $n$  or more states.

The thermodynamic advantage of higher memory parallels the thermodynamic Principle of Requisite Complexity [17]—an engine’s memory must match the predictive complexity of its fuel to optimally leverage all correlations. The challenge of harvesting energy from an information source  $\theta'$  closely parallels the challenge of predicting that same source, and more memory yields better prediction of temporally correlated information [48]. The advantage of larger memory is especially pronounced for short training lengths, where we see that 3-state machines can produce work at the limit of the Second Law for far longer training lengths.

One might infer from the advantage of additional memory in training that more memory is always thermodynamically advantageous. However, we will now explore a more nuanced picture of the costs and benefits of engine complexity by analyzing the effectiveness of the engine harvesting energy from further inputs from the information source.

#### IV. ASYMPTOTIC WORK HARVESTING AND OVERFITTING

After training to identify the maximum-work information engine, we examine its thermodynamic performance on test data as shown in Fig. 4. The purpose of training is to find an engine whose internal model approximates the true model  $\theta'$  that produced the training data. TML has found a “good model” if the resulting engine is able to effectively produce work when it receives further inputs from the true source. Thus, we consider the average work that would be produced in the testing stage.

If the  $\epsilon$ -machine of the true input process is described by:

$$\theta' = \{\mathcal{S}', \mathcal{Y}, s^{*'}, \{\theta'(y|s')\}_{y \in \mathcal{Y}, s' \in \mathcal{S}'}, \{\epsilon'(y, s')\}_{y \in \mathcal{Y}, s' \in \mathcal{S}'}\},$$

then an efficient information engine with internal model:

$$\theta = \{\mathcal{S}, \mathcal{Y}, s^*, \{\theta(y|s)\}_{y \in \mathcal{Y}, s \in \mathcal{S}}, \{\epsilon(y, s)\}_{y \in \mathcal{Y}, s \in \mathcal{S}}\},$$

will on average produce work over  $L$  time steps:

$$\beta \langle W^{\theta} \rangle_{0:L} = \sum_{y_{0:L}} \Pr(Y^{\theta'} = y_{0:L}) \ln \Pr(Y^{\theta} = y_{0:L}) + L \ln |\mathcal{Y}|. \quad (18)$$

$Y_{0:L}^{\theta'}$  is the random variable of the actual input process over length  $L$ . The average work produced during the

$L$ th time step is:

$$\beta\langle W^\theta \rangle_L \equiv \beta\langle W^\theta \rangle_{0:L+1} - \beta\langle W^\theta \rangle_{0:L}. \quad (19)$$

We consider the asymptotic rate of work production:

$$\langle W^\theta \rangle_\infty \equiv \lim_{L \rightarrow \infty} \langle W^\theta \rangle_L, \quad (20)$$

as the measure of an engine's effectiveness.

**Theorem 1.** *The asymptotic work rate for an efficient engine with internal predictive model  $\theta$  when harvesting information from a source with predictive model  $\theta'$  is:*

$$\beta\langle W^\theta \rangle_\infty = \ln |\mathcal{Y}| + \sum_{s,s',y} \pi_{s,s'} \theta'(y|s') \ln \theta(y|s). \quad (21)$$

Here,  $\pi_{s,s'}$  is the steady-state of the joint hidden states  $\mathcal{S} \otimes \mathcal{S}'$  if the causal update of  $\theta$  is driven by  $\theta'$ :

$$\pi_{s_1,s'_1} = \sum_{s_0,s'_0,y} \delta_{s_1,\epsilon(s_0,y)} \delta_{s'_1,\epsilon'(s'_0,y)} \theta'(y|s'_0) \pi_{s_0,s'_0}. \quad (22)$$

*Proof.* See Appendix B.  $\square$

Using the expression in Thm. 1, we can calculate the rate of entropy production (dissipated work) [29] by comparing the work rate to the rate of nonequilibrium free energy change:

$$\langle \Sigma^\theta \rangle_\infty / k_B = \beta(-\langle W^\theta \rangle_\infty - \Delta F_\infty^{\text{NEQ}}). \quad (23)$$

This is the amount of free energy wasted per symbol, which quantifies the irreversibility of the information engine. As shown in App. E, the entropy production can be reduced to the average relative entropy between the next-input prediction of the true model's  $\theta'$  hidden state  $s'$  and the prediction of the estimated model's  $\theta$  hidden state  $s$ :

$$\langle \Sigma^\theta \rangle_\infty / k_B = \sum_{s,s'} \pi_{s,s'} D_{KL}(Y_i^{\theta'} | S'_i = s' || Y_i^\theta | S_i = s), \quad (24)$$

where:

$$D_{KL}(X'|Z' = z' || X|Z = z) \equiv \sum_x \Pr(X' = x | Z' = z') \ln \frac{\Pr(X' = x | Z' = z')}{\Pr(X = x | Z = z)}$$

denotes the relative entropy between the distribution on  $X'$  induced by the condition that  $Z'$  realizes  $z'$  and the distribution on  $X$  induced by the condition that  $Z$  realizes the element  $z$ . Such relative entropies appear as the additional dissipation incurred by misestimating the input distribution [49, 50]. If a learning process refines and improves the estimator  $\theta$ , its divergence from the actual process  $Y_{0:L}^{\theta'}$  should diminish, monitoring how much learning reduces entropy production [15, 51].

Figure 6 plots the asymptotic work rate for the maximum-work models  $\{\Theta_n^{\text{max}}(y_{0:L})\}_{L \in \{1,2,\dots,99,100\}}$  that result from training on the words  $\{y_{0:L}\}_{L \in \{1,2,\dots,99,100\}}$  appearing in  $y_{0:100}$  which was randomly generated by the Five-State model shown in Fig. 5. The result is compared to the engine's maximum possible work rate with the true model  $\langle W^{\theta'} \rangle_\infty$ . The difference between these work rates gives the entropy production rate.

Again, we decompose the regions of functionality into Dud (negative work production), Engine (positive work production), and  $\Sigma < 0$  (Second Law violation). We see that all learned engines respect the Second Law and extract less work than if they had estimated the true model  $\theta'$ . The upper bound for engines is the work rate that results from guessing the true model, as this is the difference in entropy rates between inputs and outputs. According to the Information Processing Second Law of thermodynamics, this entropy difference is the accessible free energy per symbol and bounds the work production [39].

Figure 6 also shows a thermodynamic advantage in the testing work rate for engines with larger memories, if trained on long words. While it is unclear how close to the optimal  $n$ -state engine these results are, we see that thermodynamic learning discovers enough of the hidden temporal structure to harvest much of the available free energy. In fact, the three-state information engine nearly achieves the thermodynamic limit of perfect efficiency for training length  $L \approx 100$ . Rate-distortion theory [52–54] provides a prescription for the most predictive model given a limitation on available memory.

This noted, training engines with additional memory produces less effective engines for short training words. In fact, the maximum-work three-state engine produces  $-\infty$  work in testing for training lengths up to  $L \approx 50$ . Divergent dissipated work in Fig. 6 corresponds to *overfitting*, where the number of available model parameters is much larger than can be reasonably deduced from frequentist estimates from the data. Each memory state must make a prediction of its input. For short training words and larger memories, it becomes more probable that one of the memory states will not receive any copies of one of the input symbols  $N(y, s|s^*, \epsilon, y_{0:L}) = 0$ . In this case, the engine estimates  $\theta(y|s) = 0$ , which comes at the cost of negative divergent work production if such an observation is actually possible for the input process  $\theta'$  [49, 50]. The size of the available parameter space within even 3-state models is large enough that it realizes the main feature of overfitting for this particular training word—good performance on training data, while failing to effectively predict and harvest energy from test data. There is an asymptotic benefit to memory, but there are also dire costs to using an excessively complex engine when training data does not justify it.

In this way, we identified thermodynamic overfitting, showing that it is a considerable hurdle for TML. We now turn to resolve this challenge.



## V. THERMODYNAMIC GENERALIZATION AND REGULARIZATION

Overfitting is commonly encountered in machine learning, as models with many degrees of freedom can encode a dataset explicitly in model parameters [55]. To circumvent it, it is commonplace to implement regularization techniques that allow models to generalize and better predict unseen data. These strategies include restrictions that limit the model complexity that can be discovered through training. See, for example, “dropout” in deep neural networks [56] and regularizers that add an explicit penalty to model complexity to modify the performance measure [19, 20].

The following considers two physically-motivated methods of *thermodynamic regularization*:

1. *Autocorrection*: An engine cannot start synchronized with the true predictive state of the input. Through the influence of the input symbols on the engine’s predictive state dynamics, the engine must autocorrect in order to synchronize its estimated predictive state. Not surprisingly, there is a transient energy cost as the engine autocorrects its predictive state and approaches its steady-state dynamics [57].
2. *Engine Initialization*: There is an energy cost to initializing the energy landscape of a predictive information engine.

We will now show that autocorrection incurs an additional cost to complex models with unnecessarily distinct predictions from each memory state. In addition, we find that the cost of engine initialization leads to Bayesian updates of the edge-weights according to Laplace’s rule of succession [58].

Let  $C(\theta)$  denote the energy penalty of initializing an engine with model  $\theta$ . It functions as a *regularizer*. If we also include the cost of autocorrecting from an initial distribution  $p(s_0) \equiv \Pr(S_0 = s_0)$ , as shown in App. D, the average work production can be expressed:

$$\beta \langle W_G^\theta(y_{0:L}) \rangle = L \ln |\mathcal{Y}| - C(\theta) + \sum_{s_0} p(s_0) \ln \Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = s_0),$$

where  $\Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = s_0)$  is the probability of model  $\theta$  producing  $y_{0:L}$  when starting from predictive state  $s_0$ . As with the case of unregularized TML, we can re-express the work production by tracking the state dynamics induced within the predictive states by the input word. Again, we need only count the number of times  $N(y, s | s_0, \epsilon, y_{0:L})$  that  $y$  is input to predictive state  $s$  for every the initial state  $s_0$ , then we obtain the work pro-

duction:

$$\beta \langle W_G^\theta(y_{0:L}) \rangle = L \ln |\mathcal{Y}| - C(\theta) + \sum_{s_0, s, y} p(s_0) N(y, s | s_0, \epsilon, y_{0:L}) \ln \theta(y | s). \quad (25)$$

The penalty  $C(\theta)$  for model  $\theta$  incurred through training is motivated by physical constraints, such as the energetic cost of initializing the model. If TML evaluates a model that is too complex and costly to initialize, then that cost should counteract the energetic benefit of predicting the training word. This echoes the *minimum description length principle*, in which the benefit gained through prediction is offset by the cost of describing the model [59]. In this spirit, we introduce an energy penalty of preparing the edge-weights  $\theta(y | s)$  of a particular model  $\theta$ .

The penalty originates from the fact that equilibrium probabilities  $\pi$  are directly related to energies via:

$$\beta E(z) = \beta F^{\text{eq}} - \ln \pi(z).$$

The equilibrium free energy  $-\beta F^{\text{eq}} \equiv \ln \sum_z e^{-\beta E(z)}$  is the upper limit on work that can be produced from an equilibrium distribution. The equilibrium distribution is the necessary starting point for an efficient quasistatic transformation of information [17, 36, 60]. This means that if we prepare a distribution  $q(z)$  that differs from the initial equilibrium distribution, we can calculate the dissipated work  $\langle W_{\text{diss}} \rangle$  as it partially relaxes to  $q'(z)$  via a difference in relative entropies [29, 49, 50]:

$$\beta \langle W_{\text{diss}} \rangle = D_{KL}(q | \pi) - D_{KL}(q' | \pi).$$

The dissipation associated with preparing every edge-weight of the model  $\theta$  should contribute to a Thermodynamic Regularizer.

As discussed in App. D, preparing the edge-weight  $\theta(y | s)$  of a particular combination of predictive state  $s$  and input  $y$  incurs the dissipated work of:

$$\beta \langle W_{\text{diss}}^{\text{prepare}}(s, y) \rangle = -\ln \theta(y | s).$$

We propose that the cost of implementing a model is proportional to the dissipation associated with preparing every edge-weight:

$$C(\theta) = \alpha \sum_{s, y} \beta \langle W_{\text{diss}}^{\text{prepare}} \rangle = -\alpha \sum_{s, y} \ln \theta(y | s),$$

where  $\alpha$  is a regularization parameter of our choosing. For the parameter  $\alpha = 1$ , this is the total energetic excess beyond the free energy for each combination engine

memory state and input:

$$C(\theta) = \sum_{s,y} \beta(E^\theta(s,y) - F^\theta(s)).$$

We associate this additional cost with examining the relaxation to check the estimated probability of every edge-weight. This is a *thermodynamic regularizer* in that the penalty originates from a thermodynamic implementation of the engine with model  $\theta$ .

The resulting regularized work production is

$$\begin{aligned} \beta \langle W_{p,\alpha,\epsilon}^\theta(y_{0:L}) \rangle &= L \ln |\mathcal{Y}| - \alpha \langle W_{\text{diss}}^{\text{prepare}} \rangle \\ &+ \sum_{s_0} p(s_0) N(y, s|s_0, \epsilon, y_{0:L}) \ln \theta(y|s). \end{aligned} \quad (26)$$

This work production can be maximized analytically.

**Theorem 2.** *The maximum-work edge-weights of an engine with input  $y_{0:L}$ , topology  $\epsilon$ , and regularization parameters  $\alpha$  and  $p$  can be analytically calculated by tracking the dynamics of the  $\epsilon$ -map from each start state:*

$$\Theta_{p,\alpha,\epsilon,y_{0:L}}^{\text{max}}(y|s) = \frac{\sum_{s_0} p(s_0)(\alpha + N(y, s|s_0, \epsilon, y_{0:L}))}{\sum_{y',s_0} p(s_0)(\alpha + N(y', s|s_0, \epsilon, y_{0:L}))}.$$

*Proof.* See App. D  $\square$

Theorem 2 gives a shortcut to regularized TML. Once given the topology of a candidate information engine, we need only track the memory dynamics and directly calculate the maximum-work edge-weights from  $\Theta_{p,\alpha,\epsilon,y_{0:L}}^{\text{max}}(y|s)$ . Thus, given a memory constraint, we scan through the available engine topologies and select that which produces the most work. While the set of topologies grows super-exponentially with memory, it is vastly simpler than discovering the edge-weights through numerical optimization of edge-weights.

## A. Autocorrection

When harvesting temporal correlations from a time series, the engine may start out of sync with the input process. As a result, it must *autocorrect* to the predictive states input. Autocorrection requires additional energy [57], as does synchronizing with the inputs [61]. To address this, consider the average work production if the engine starts in distribution  $p(s) \equiv \Pr(X_0^\theta = s)$  over its memory states. Appendix C shows that the resulting average work production is the weighted log-likelihood of

each predictive state:

$$\begin{aligned} &\beta \langle W_{\text{AC}}^\theta(y_{0:L}) \rangle \\ &= L \ln |\mathcal{Y}| + \sum_s p(s) \ln \Pr(Y_{0:L}^\theta = y_{0:L} | S_0^\theta = s) \\ &= L \ln |\mathcal{Y}| + \sum_s p(s) \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, s)). \end{aligned}$$

As an example, consider the case where the initial distribution over memory states is uniform  $p(s) = 1/|\mathcal{X}|$ . Rather than finding the  $n$ -state model that maximizes the work production  $\Theta_n^{\text{max}}(y_{0:N})$ , we find the model that produces maximum work when we take into account the cost of autocorrection:

$$\Theta_n^{\text{AC}}(y_{0:L}) \equiv \operatorname{argmax}_{\theta \in \Theta_n} \langle W_{\text{AC}}^\theta(y_{0:L}) \rangle. \quad (27)$$

This sets us up to introduce a thermodynamic complexity measure for autocorrection. Note that requiring the agent to autocorrect introduces unavoidable energy inefficiencies that are not present in the unregularized MLE/TML strategy. Given an initial density over the input's predictive states  $p(s_0) = \Pr(S'_0 = s_0)$ , we have the true input distribution:

$$\Pr(Y_{0:L}^{\theta'}) = \sum_{s_0} p(s_0) \Pr(Y_{0:L}^{\theta'} | S_0 = s_0).$$

The average entropy production is the difference between the average work production and the change in free energy:

$$\begin{aligned} \langle \Sigma_{\text{AC}}^\theta \rangle_{0:L} / k_B &= \beta(-\langle W_{\text{AC}}^\theta \rangle_{0:L} - \Delta F_{0:L}^{\text{NEQ}}) \\ &= \sum_{s_0} p(s_0) D_{KL}(Y_{0:L}^{\theta'} || Y_{0:L}^\theta | S_0 = s_0). \end{aligned}$$

The average dissipation is proportional to the average divergence between the actual input and the estimated input from each predictive state of the estimated model  $\theta$ .

Even if the agent correctly guesses the underlying  $\epsilon$ -machine model, such that the state transitions  $\epsilon' = \epsilon$  and the edge-weights  $\theta'(y|s) = \theta(y|s)$  are the same and the prediction from each causal state is the same:

$$\Pr(Y_{0:L}^{\theta'} | S'_0 = s_0) = \Pr(Y_{0:L}^\theta | S_0 = s_0), \quad (28)$$

the agent can still produce entropy due to the lack of model synchronisation.

The entropy produced through autocorrection is a new model complexity measure. Specifically, the complexity is proportional to the cost of autocorrecting to the correct predictive states of the process from an initially uniform

distribution over all predictive states:

$$\mathcal{C}_{\text{AC}}(\theta, L) \equiv \frac{\sum_{s_0} D_{KL}(Y_{0:L}^{\theta'} || Y_{0:L}^{\theta'} | S_0' = s_0)}{|\mathcal{S}|}.$$

Previous explorations of correlation-powered information engines show that they must autocorrect to synchronize with their inputs, incurring a thermodynamic cost along the way [57].

The autocorrection cost  $\mathcal{C}_{\text{AC}}(\theta, L)$  differs from a measure of stored information like statistical complexity  $C_\mu$  [13]. On the one hand, we can have many different possible configurations with the latter, with nearly identical predictions, but the difference between predictions of each causal state is immaterial to the complexity measure. On the other hand, the complexity measure  $\mathcal{C}_{\text{AC}}(\theta, L)$  reflects the thermodynamic cost of choosing a model that anticipates wildly divergent futures from different predictive states.

There are cases in which we know this will approach a finite quantity in the asymptotic limit  $L \rightarrow \infty$ . (For instance, when  $Y_{0:\infty}^{\theta'}$  has finite Markov order.) In such cases, we can refine the complexity measure to be the thermodynamic cost of exactly synchronizing to the input process:

$$\begin{aligned} \mathcal{C}_{\text{AC}}(\theta) &\equiv \lim_{L \rightarrow \infty} \mathcal{C}_{\text{AC}}(\theta, L) \\ &= \frac{\sum_{s_0} D_{KL}(Y_{0:\infty}^{\theta'} || Y_{0:\infty}^{\theta'} | S_0' = s_0)}{|\mathcal{S}|}. \end{aligned}$$

We define  $\mathcal{C}_{\text{AC}}(\theta)$  to be the *autocorrection complexity*.

The entropy production associated with autocorrection can only be minimized to zero if the engine is restricted to start in a single predictive state. Relative entropies are all non-negative and zero if and only if  $\Pr(Y_{0:\infty}^{\theta'}) = \Pr(Y_{0:\infty}^{\theta'} | S_0 = s_0)$ . This means that for all predictive states  $s_0$  with nonzero  $p(s_0)$ , to minimize entropy production it must be true that  $\Pr(Y_{0:\infty}^{\theta'}) = \Pr(Y_{0:\infty}^{\theta'} | S_0 = s_0)$ . However, by the definition of causal states, different memory states must give rise to different future predictions. Thus, zero dissipation is impossible unless the initial distribution is a unique start state  $p(s_0) = \delta_{s_0, s^*}$ .

## B. Bayesian Edge-Weights

We introduce another regularization technique by allowing the agent's memory to start in a single start state  $p(s) = \delta_{s, s^*}$ , but insisting on a complexity cost that is proportional to the work dissipated by initializing every edge-weight  $C(\theta) = \langle W_{\text{diss}}^{\text{prepare}} \rangle$  ( $\alpha = 1$ ). Applying Thm. 2, the resulting maximum-work edge-weights are given by:

$$\Theta_{\delta_{s, s^*}, 1, \epsilon, y_{0:L}}^{\max}(y|s) = \frac{1 + N(y, s | s^*, \epsilon, y_{0:L})}{|\mathcal{Y}| + \sum_{y'} N(y', s | s^*, \epsilon, y_{0:L})}.$$

This is simply Laplace's rule of succession when  $Y_t$  is a binary variable [58]. It follows from considering a uniform prior over the inputs for each causal state, then using Bayes' theorem to update the distribution using input counts  $N(y, s | s^*, \epsilon, y_{0:L})$  for each causal state. The generalization of this to larger alphabets  $\mathcal{Y}$  also comes from Bayesian inference applied to a prior given by a  $|\mathcal{Y}|$ -dimensional uniform distribution (Dirichlet distribution) over the possible edge-weights [58].

## C. Combined Regularization

If we choose to both incur a cost of initializing the edge-weights ( $\alpha = 1$ ) and a cost of autocorrection ( $p(s) = 1/|\mathcal{S}|$ ), then we obtain a "combined" regularization strategy for which the maximum-work edge-weights are:

$$\Theta_{1/|\mathcal{Y}|, 1, \epsilon, y_{0:L}}^{\max}(y|s) = \frac{\sum_{s_0} (1 + N(y, s | s_0, \epsilon, y_{0:L}))}{\sum_{y', s_0} (1 + N(y', s | s_0, \epsilon, y_{0:L}))}. \quad (29)$$

as shown in App. D.

The following numerically compares the advantages of each regularization strategy:

1. MLE ( $\alpha = 0$ ,  $p(s) = \delta_{s, s^*}$ ): Unregularized, where we simply maximize the work production, which yields the same inference as Maximum Likelihood Estimation.
2. BAYES ( $\alpha = 1$ ,  $p(s) = \delta_{s, s^*}$ ): We incur the cost of initializing edge-weights, resulting in Bayesian updates of the edge-weights.
3. AC ( $\alpha = 0$ ,  $p(s) = 1/|\mathcal{S}|$ ): We incur the cost of autocorrection, meaning that divergent predictions from different causal states must be strongly justified by data.
4. CMBD ( $\alpha = 1$ ,  $p(s) = 1/|\mathcal{S}|$ ): We combine the costs of autocorrection and initializing edge-weights.

## VI. TESTING GENERALIZATION STRATEGIES

As we attempt to determine the effectiveness of thermodynamic learning with an additional work penalty from complexity, it is worth noting that our original measure of testing performance for our model (the asymptotic work rate  $\langle W^\theta \rangle_\infty$ ) is unchanged. As long as the machines are ergodic, the initial distribution  $p(s)$  doesn't affect the steady-state distribution  $\pi_{s, s'}$  in Thm. 1. In addition, the complexity costs  $C(\theta)$  are transient, and don't affect the asymptotic dynamics. Thus, we can exactly calculate the asymptotic performance of regularized maximum-work models using Thm. 1, just as before. In the following, we evaluate the performance of the four different generalization strategies by calculating

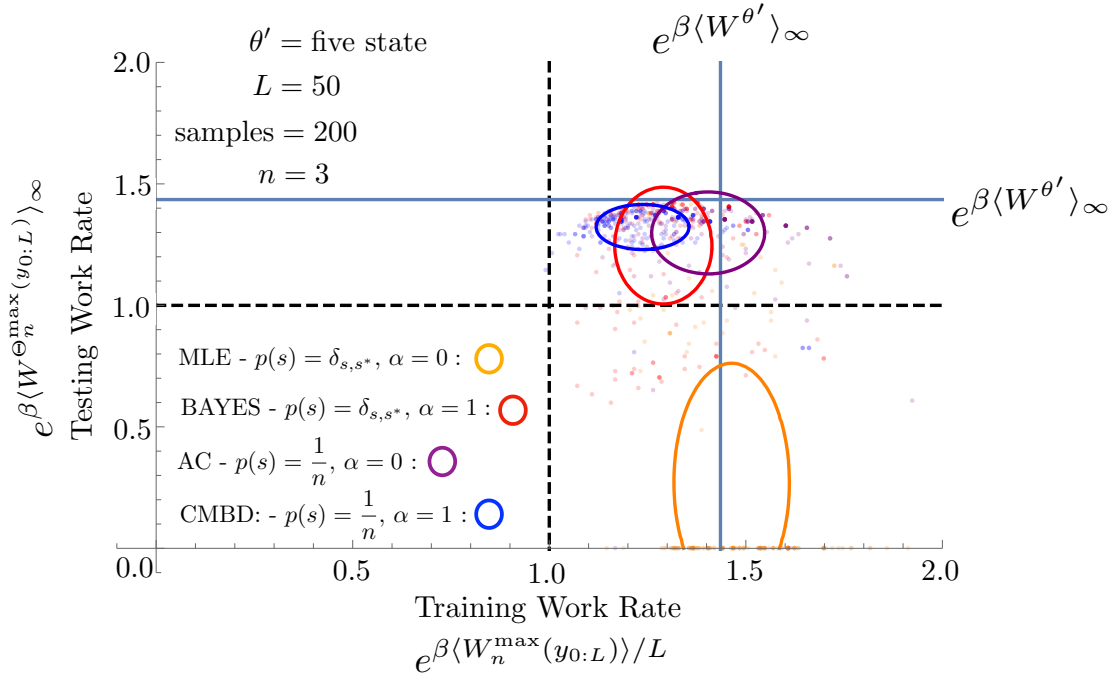


FIG. 7. Ensemble performance of thermodynamic regularization in training compared to testing: For the input process  $\theta'$  we randomly sample 200 length  $L = 50$  words and train 3-state information engines on each to find the exponential training work rate  $e^{\beta \langle W_{n=3}^{\max}(y_{0:L}) \rangle / L}$  and the exponential testing work rate  $e^{\beta \langle W_{n=3}^{\Theta_n^{\max}}(y_{0:L}) \rangle_{\infty}}$  for each regularization strategy: MLE (orange), BAYES (red), AC (purple), and CMBD (blue). The ovals are centered around the average work rates of these 200 samples, and their dimensions are given by the variance of the work rates. The dashed black lines represent work rates of zero along each dimension, and the blue lines represent the theoretical limit on the asymptotic work rate, given by  $e^{\beta \langle W^{\theta'} \rangle_{\infty}}$ . The strict MLE strategy performs best in training, but worst in testing. The CMBD strategy performs worst in training and best in testing.

the maximum-work model

$$\Theta_n^{\max}(y_{0:L}) \equiv \operatorname{argmax}_{\theta \in \Theta_n} \langle W_G^{\theta}(y_{0:L}) \rangle, \quad (30)$$

then evaluating the testing work rate  $\langle W_n^{\Theta_n^{\max}}(y_{0:L}) \rangle_{\infty}$ .

### A. Ensemble Performance of Regularized Thermodynamic Learning

The performance of thermodynamic learning from individual words as shown in Fig. 6 gives some insight into how patterns are discovered. However, we see enough variety in random word realizations in App. F that we must look at ensemble averages to determine whether Thermodynamic Regularization effectively generalizes.

Take the process generated by the Five-State model shown in Fig. 5 as our true source. This process requires five or more memory states for perfect prediction, while we have at most three memory states available. Examining learning from this process elucidates the case where there is a force towards more engine complexity, but practical constraints of overfitting limit complexity.

This illustrates how thermodynamic learning performs when the accessible model class is misspecified by not including sufficient memory states to fully capture the process.

Figure 7 plots the asymptotic testing work rate against the training work rate for engines with memory size  $n = 3$  from 200 different words of length 50 generated from the Five-State model. Plotted against the theoretical limit on the work rate  $\beta \langle W^{\theta'} \rangle_{\infty}$ , we see how close the learned engines are to optimal. Note that the work rates are exponentiated, so that the graph can accommodate infinitely divergent outcomes. We identify the behaviour of the ensemble of words by plotting an oval whose center is our numerical estimate of the average work rates for the learning process:

$$\begin{aligned} & \langle \langle x \rangle, \langle y \rangle \rangle \\ &= \left( \langle e^{\beta \langle W_{n=3}^{\max}(y_{0:L}) \rangle / L} \rangle_{Y_{0:L}^{\theta'}}, \langle e^{\beta \langle W_{n=3}^{\Theta_n^{\max}}(y_{0:L}) \rangle_{\infty}} \rangle_{Y_{0:L}^{\theta'}} \right), \end{aligned}$$

and whose radial dimensions are given by the variance of the exponential work rates:

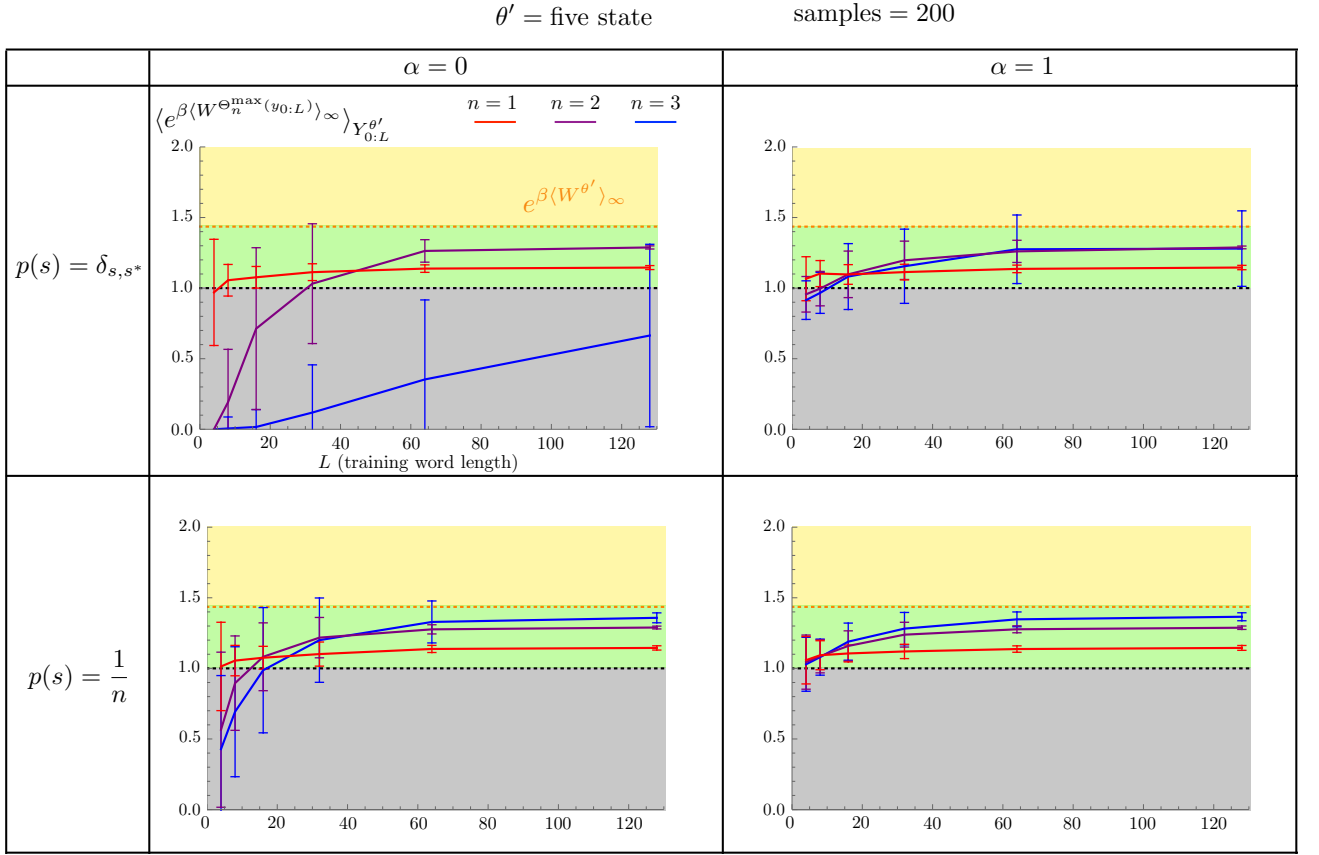


FIG. 8. CMBD regularization mitigates overfitting: The average and variance of the exponential asymptotic testing work rate that results from the four different learning strategies. We generate 200 words from the Five-State process for each length  $L \in \{4, 8, 16, 32, 64, 128\}$ , then train on each using MLE, BAYES, AC, and CMBD. The unregularized thermodynamic machine learning does an extremely poor job of discovering the underlying pattern in the word for large memories. The BAYES method does well for small memories, but has poor performance for large memories, because of the large variance. Using AC instead helps, but leads to divergent dissipation for small training sets. For MLE, BAYES, and AC, there is an advantage to using a small amount of memory, for small data. The CMBD strategy, by contrast, doesn't seem to have this suffer from using more memory, suggesting that it effectively mitigates overfitting for this case.

$$(\text{var}(x), \text{var}(y)) = \left( \text{var} \left( e^{\beta \langle W_{n=3}^{\max}(y_{0:L})/L \rangle} \right)_{Y_{0:L}^{\theta'}}, \text{var} \left( e^{\beta \langle W_{n=3}^{\max}(y_{0:L}) \rangle} \right)_{Y_{0:L}^{\theta'}} \right).$$

This uses the notation:

$$\langle f(x) \rangle_X \equiv \sum_x \Pr(X = x) f(x)$$

$$\text{var}(f(x))_X \equiv \sum_x \Pr(X = x) (f(x) - \langle f(x) \rangle)^2.$$

To interpret Fig. 7, note that it plots the average and variance of the exponential work  $e^W$  to accommodate cases of extreme overfitting, with infinite dissipation and  $-\infty$  work. If any elements of our training sample overfit in this way, then the resulting estimated average  $\langle W \rangle$  and variance  $\text{var}(W)$  of the work would diverge. For the

case where  $\alpha = 0$ , this is always a possibility, since a given training word may not realize a transition that is allowed by the input process.

We estimate a rough average from the plot via Jensen's inequality:

$$e^{\langle W \rangle} \leq \langle e^W \rangle.$$

Equality is only satisfied when the work always realizes the average. This means that as the variance of the distribution increases, so should the difference below the average exponential  $\langle e^W \rangle$ . Thus, when reading the figure, both higher average and lower variance correspond

to better engine performance.

We see that the training work rate is often above the theoretical limit set but  $\langle W^{\theta'} \rangle_{\infty}$ , but the testing work rate is always below it, as should be the case. We also see that a higher training work rate seems to correspond to worse testing performance. Beyond this, we focus on the testing work rate in the following analysis.

For the Five-State process, we see that the standard MLE technique without regularization badly overfits. The average exponential asymptotic work rate after learning is well below unity. This reflects the fact that many realizations dissipate infinite work. The BAYES strategy, by contrast, improves on this work rate. However, the variance is lower, and the average is higher for the AC and CMBD techniques, implying better test data performance for these two regularization strategies. The CMBD strategy does best of all.

Figure 8 shows the performance of different regularization strategies in greater detail by considering varying lengths. For each length  $L \in \{4, 8, 16, 32, 64, 128\}$ , we randomly generate 200 words from the Five-State model and train on that model using MLE, BAYES, AC, and CMBD regularization. We also compare different memory sizes for the information engine. We see that all strategies are relatively effective for engines with small memories. However, for engines with three memory states, unregularized MLE results in divergent dissipation even for long training words. We see better performance in BAYES and AC training, but for short words, we still see an advantage to training on models with less memory. However, the CMBD technique seems to derive a consistent advantage from additional memory. This suggests that the CMBD regularization technique could be used to reliably discover complex patterns from small amounts of data.

Appendix G goes into a similar analysis of the ‘‘Even Process’’ and ‘‘Noisy Even Process’’. These are interesting since the true source is not misspecified by the candidate model class. In both cases, since two memory states are sufficient to predict the process the  $n = 3$  curve lies at or below the  $n = 2$  for all samples, extracting less work on average. Thus, we still see a cost to having an engine that is more complex than necessary for this two-state input. However, CMBD regularization does the best in mitigating overfitting for larger engines.

## VII. SUMMARY

Thermodynamics mandates both a drive towards complexity and an impetus for simplicity. In this article, we see this through the lens of thermodynamic machine learning (TML), which discovers patterns in data by maximizing work production.

Our framework implies a fundamental energy dissipation cost when an engine’s internal model deviates from that of the true underlying distribution. To minimize this mismatch, an information engine must be at least as

complex as the input - illustrating the principle of Requisite Complexity [17, 62]. Our first contribution is an exact expression for the asymptotic work rate of the information engine (see Thm. 1), allowing us to directly evaluate the effectiveness of an engine in utilizing a particular information fuel.

Next, we greatly simplified the search for the maximum-work engine by analytically deriving the engine parameters in Thm. 2. This allowed us to implement TML over a wide class of engines, and demonstrate that solely maximizing work harvested from training data can lead to overfitting, with potentially dire energetic consequences (as seen in Fig. 6). Moreover, this risk increases as engines scale in complexity. While the work production improves uniformly with engine memory during training, the resulting information engine anticipates phantom causal relations that are not present during testing. In extreme cases, such agents become immensely dissipative on unexpected inputs, and work dissipation can diverge towards infinity — the key signature of thermodynamic overfitting.

Our final contribution involves introducing two thermodynamically means of regularization. The first adds a cost to initializing the information engine’s edge-weights that is proportional to the dissipated work arising from relaxing to equilibrium with the prediction of each memory state. The second includes the cost of auto-correction - the energetic cost of discovering the correct predictive state when the engine starts out of sync with its information fuel. Combining the two techniques, we introduced analytical methods to construct regularized information engines that perform significantly better in harnessing free energy outside the training phase (see in Figs. 7 and 8). In fact, our resulting engines perform similarly during testing and training phases, indicating that they have effectively mitigated overfitting. Meanwhile, the correspondence between work cost and maximum-likelihood estimation (MLE) together with the consistency of MLE [58], imply our agents converge to the correct model in the long-time limit. We thus illustrate how thermodynamic considerations can lead to effective means of pattern discovery without the risk of overfitting.

## VIII. OUTLOOK

Our results provide natural links between thermodynamics and machine learning and, in this way, establish a number of interesting connections worth further investigation. On a technical side, many of the results may be recast in the language of Fisher Information. Consider an ensemble of training words selected from the true distribution  $\Pr(Y_{0:L}^{\theta'} = y_{0:L})$ . The variance in the estimator  $\Theta_n^{\max}(y_{0:L})$  can be bounded by using the Cramér-Rao bound. Asymptotically, this quantity scales linearly with word-length  $L$ , with the constant multiplier being the Fisher information rate [63, 64]. Maximum-likelihood estimators are asymptotically efficient, and so achieve this

rate of learning with sufficient data [65]. Since the regularized TML methods presented here match MLE in the asymptotic limit, the variance of estimated parameters will follow the optimal  $1/L$  scaling determined by the Fisher information rate.

Meanwhile, calculating the maximum-work engine edge-weights shown in Thm. 2 strongly echoes reservoir computing in its computational simplicity [66]. Reservoir computing leverages the inherent information processing of a complex system—the reservoir—driving it with time-series data. The only training that happens through this process occurs via an output layer, which homes in on the subdynamics of the reservoir that carry the relevant temporal correlations from the input sequence. This process is computationally simple, corresponding to a matrix inverse, with the weights of the output layers paralleling the edge-weights  $\theta(y|s)$  of the predictive machine in TML. The question of what makes an effective reservoir remains open [67], with only heuristic design guides (e.g., good reservoirs are often thought to be on the “edge of chaos” [68]). It has been shown that the most effective reservoirs are deterministic (i.e., unifilar) [69]. An engine’s memory, viewed as a reservoir, satisfies this condition. Our framework may thus help provide thermodynamic guiding principles for finding effective reservoir computers.

More broadly, state-of-the-art time-series prediction and manipulation generally involve the use of recurrent neural networks (RNNs) [22] and transformers [23]. In RNNs, the process of training is generally computationally intensive [70], and recent works have suggested improved performance when such training makes use of causal discovery techniques in  $\epsilon$ -machines [71]. Meanwhile, transformer functionality is rooted in “next token prediction”: finding a function that maps past inputs of some context length to a probability for the next “token” (input) [72, 73]. Recent results showed that Transformers are indeed universal predictors, recovering a mapping from past inputs to hidden states and edge-weights [45]. The  $\epsilon$ -map along with the edge-weights does just this, with the advantage that it does not require infinitely large memory to exactly model infinite Markov-order processes. Such processes would require infinite context-length for a Transformer to exactly predict them. By contrast, the memory states of a prediction engine can capture information contained in inputs arbitrarily far in the past with relatively small memory for many non-Markovian processes [54]. As such, our methodologies may well provide new tools to tackle the unsustainable energetic cost of current AI models.

Finally, information is ultimately quantum-mechanical—leading to recurrent quantum models and quantum reservoir computers [8, 74, 75]. Our thermodynamic toolkit thus provides a physical means to compare quantum and classical models operationally. Indeed, relations between work dissipation and imperfect modelling extend to the quantum regime [76]. Meanwhile, there is mounting evidence that such quantum

machines can exhibit certain target behaviours in various contexts — stochastic modelling, string generation, adaptive strategies — while using less memory than any classical counterparts [47, 77–80]. In stochastic modelling, such memory advantages induce energetic advantages [81], while model memory can bound generalisation error in classification tasks [82]. It would thus be exciting to determine how thermodynamic overfitting applies to quantum models, and thus determine whether the pressure for energetically efficient learning naturally motivates quantum-enhanced artificial intelligence.

We see an encouraging, evolving picture of how thermodynamic resources govern the emergence of predictive agents. On the one hand, maximizing the energy extracted from information corresponds to discovering the maximum-likelihood predictive model of that information. However, reckless energy extraction leads to overly-precise probability estimates, resulting in elevated energetic cost downstream. Fortunately, physical constraints, such as the cost of instantiating the information engine and autocorrecting to the correct predictive state, prevent such glaring pitfalls. It appears that nature conspires to bring about predictive machines through thermodynamic resource optimization.

## ACKNOWLEDGMENTS

This work was supported by the Irish Research Council under grant number IRCLA/2022/3922, and by the Foundational Questions Institute and Fetzer Franklin Fund, a donor advised fund of the Silicon Valley Community Foundation, grant number FQXi-RFP-IPW-1910. ABB and JPC thank the Telluride Science Research Center for hospitality during visits and the participants of the Information Engines Workshops there. ABB acknowledges support from the Templeton World Charity Foundation Power of Information fellowships TWCF0337 and TWCF0560. This material is also based upon work supported by, or in part by, U.S. Army Research Laboratory, U.S. Army Research Office grant W911NF-21-1-0048, the National Research Foundation, Singapore, and Agency for Science, Technology and Research (A\*STAR) under its QEP2.0 programme (NRF2021-QEP2-02-P06), and the Singapore Ministry of Education Tier 1 Grants RG146/20 and RG77/22, and the John Templeton Foundation grant no. 62423.

## Appendix A: Predictive Models

The predictive models encoded within an efficient information engine are Hidden Markov Models (HMMs). Specifically, they are edge-emitting such that they generate sequences through symbol-labeled transition matrices over hidden states  $\mathcal{S}$ :

$$T_{s \rightarrow s'}^{(y)} \equiv \Pr(Y_i = y, S_{i+1} = s' | S_i = s).$$

Predictive models are a subclass of HMMs such that the hidden states contain no information about the future beyond that which can be determined from the past: [17]

$$I[S_t; \vec{Y}_t | \overleftarrow{Y}_t] = 0.$$

The hidden states are called predictive states, because they contain all information in the past relevant for predicting the future. If we additionally require that a process's model memory is minimal, then we obtain the  $\epsilon$ -machine: the minimal predictive generator of that process [14]. The predictive states of this minimal machine are often referred to as *causal states*, because they describe the past's causal influence on the future.

The predictive complexity is described by the memory resources associated with the stationary causal state distribution  $\Pr(S_t)$ . The  $\epsilon$ -machine is the predictive model whose hidden states minimize the  $\alpha$ -Rényi entropy  $H_\alpha[S_t]$  for all  $\alpha$  values [80].  $\alpha = 1$  characterizes the statistical complexity [13]:

$$\begin{aligned} C_\mu &\equiv H_1[S_t] \\ &= - \sum_s \Pr(S_t = s) \ln \Pr(S_t = s), \end{aligned}$$

which is the channel capacity (measured in Nats) necessary to communicate from past to future. By comparison,  $\alpha = 0$  yields is the topological predictive complexity:

$$H_0[S_t] = \ln |\mathcal{S}|,$$

which is the log of the number of hidden states necessary to predictively model the process.  $H_0[S_t]$  is the measure that is most directly relevant in designing information engines, because it determines the number of memory states that an engine must have in order to efficiently harvest information.

The  $\epsilon$ -machine's minimal causal states are determined by a *causal equivalence relation* [12]. As a result, they have the convenient property of *unifilarity* [83], which means that the next causal state  $S_{i+1}$  is uniquely specified by the current one  $S_i$  and its output  $Y_i$  by an  $\epsilon$ -map:

$$S_{i+1} = \epsilon(Y_i, S_i).$$

This is also known as the *topology* of the  $\epsilon$ -machine [43]. In addition, an  $\epsilon$ -machine has a unique start causal state  $s^*$ . In the case of a bi-infinite process  $Y_{-\infty:\infty}$  which is stationary, this reflects the belief state of having seen nothing so far, but in the non-stationary semi-infinite  $Y_{0:\infty}$  case,  $s^*$  is simply the sufficient statistic of the (non-existent) past about the future [10].

Unifilarity and a unique start state  $s_0$  guarantee that an output sequence  $y_{0:i}$  will lead to a unique causal state:

$$s_i = \epsilon(y_{0:i}, s_0),$$

where we've re-used the update map notation, defining:

$$\epsilon(y_{0:i}, s_0) \equiv \epsilon(y_{i-1}, \dots \epsilon(y_1, \epsilon(y_0, s_0)) \dots).$$

Because of the isomorphism between the information engine and its  $\epsilon$ -machine model, the engine's memory state  $x_i$  will be the same function of its input  $x_i = \epsilon(y_{0:i}, s^*)$ . Thus, we can see an efficient engine's memory dynamics as an input-conditioned deterministic dynamical system, paralleling efficient reservoir computers [69].

Unifilarity means that we can characterize any  $\epsilon$ -machine in terms of its topology, start state, and edge-weights, as shown in the example of the Even Process in Fig. 3. For a particular topology, the edge-weights  $\theta(y|s)$  are the



probability of outputting  $y$  from causal state  $s$ :

$$\begin{aligned}\theta(y|s) &\equiv T_{s \rightarrow \epsilon(y,s)}^{(y)} \\ &= \Pr(Y_i^\theta = y | S_i = s).\end{aligned}$$

This meaningfully simplifies the task of Thermodynamic Machine Learning.

Finally,  $\epsilon$ -machines can produce any process and yield the same likelihood as any predictive model, meaning that it is sufficient to limit our model class  $\Theta$  to these models. Such machines are described by a causal update map on the predictive states and edge-weights, meaning these are the explicit parameters we must explore through training.

## Appendix B: Asymptotic Work Rate: Derivation

The work produced by an efficient engine with model  $\theta$  at the  $L$ th time step can be expressed:

$$\begin{aligned}\beta \langle W^\theta \rangle_L &\equiv \beta \langle W^\theta \rangle_{0:L+1} - \beta \langle W^\theta \rangle_{0:L} \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \Pr(Y_{0:L+1}^\theta = y_{0:L+1}) - \sum_{y_{0:L}} \Pr(Y_{0:L}^{\theta'} = y_{0:L}) \ln \Pr(Y_{0:L}^\theta = y_{0:L}) \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \Pr(Y_{0:L+1}^\theta = y_{0:L+1}) - \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \Pr(Y_{0:L}^\theta = y_{0:L}) \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \frac{\Pr(Y_{0:L+1}^\theta = y_{0:L+1})}{\Pr(Y_{0:L}^\theta = y_{0:L})} \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \Pr(Y_L^\theta = y_L | Y_{0:L}^\theta = y_{0:L}).\end{aligned}$$

Using the fact that, for the model  $\theta$ , the edge-weights are related to output probabilities:

$$\theta(y_L | \epsilon(y_{0:L})) = \Pr(Y_L^\theta = y_L | Y_{0:L}^\theta = y_{0:L}),$$

we can express the asymptotic work production:

$$\beta \langle W^\theta \rangle_L = \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \theta(y_L | \epsilon(y_{0:L})).$$

Another way of expressing this is to say the probability of the agent's causal state  $s$  at time  $L$  given the past inputs is  $\Pr(S_L = s | Y_{0:L}^{\theta'} = y_{0:L}) = \delta_{s, \epsilon(y_{0:L})}$ , so that the work production can be expressed:

$$\begin{aligned}\beta \langle W^\theta \rangle_L &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \theta(y_L | \epsilon(y_{0:L})) \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}, s} \Pr(Y_{0:L+1}^{\theta'} = y_{0:L+1}) \delta_{s, \epsilon(y_{0:L})} \ln \theta(y_L | \epsilon(y_{0:L})) \\ &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}, s} \Pr(S_L = s, Y_{0:L+1}^{\theta'} = y_{0:L+1}) \ln \theta(y_L | s).\end{aligned}$$

The causal states of the input's  $\epsilon$ -machine are given by the random variables  $S'_L$ , so we express the asymptotic work:

$$\begin{aligned}\beta \langle W^\theta \rangle_L &= \ln |\mathcal{Y}| + \sum_{y_{0:L+1}, s, s'} \Pr(S_L = s, S'_L = s', Y_L^{\theta'} = y_L, Y_{0:L}^{\theta'} = y_{0:L}) \ln \theta(y_L | s) \\ &= \ln |\mathcal{Y}| + \sum_{y_L, s, s'} \Pr(S_L = s, S'_L = s', Y_L^{\theta'} = y_L) \ln \theta(y_L | s) \\ &= \ln |\mathcal{Y}| + \sum_{y, s, s'} \Pr(S_L = s, S'_L = s') \theta'(y_L | s') \ln \theta(y_L | s).\end{aligned}$$

If we take the asymptotic limit, we obtain the steady-state distribution over the causal states of the input  $\epsilon$ -machine driving the estimated  $\epsilon$ -machine:

$$\pi_{s,s'} \equiv \lim_{L \rightarrow \infty} \Pr(S_L = s, S'_L = s'),$$

which can be found by solving the steady-state equation:

$$\pi_{s_1, s'_1} = \sum_{s_0, s'_0, y} \delta_{s_1, \epsilon(s_0, y)} \delta_{s'_1, \epsilon'(s'_0, y)} \theta'(y | s'_0) \pi_{s_0, s'_0}.$$

Thus, we obtain the expression for the asymptotic work rate:

$$\beta \langle W^\theta \rangle_\infty = \ln |\mathcal{Y}| + \sum_{s, s', y} \pi_{s, s'} \theta'(y | s') \ln \theta(y | s).$$

### Appendix C: Work Production From Distributed Start State

Reference [10] primarily considers the work production from a particular start-state. However, it also includes an equation for the work production when the hidden state starts in a distribution  $\Pr(X_0 = x_0)$ , given in Eq. (G1):

$$\beta \langle W^\theta(y_{0:L}) \rangle = \sum_{x_0} \Pr(X_0 = x_0) \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, x_0)) + L \ln |\mathcal{Y}| + \sum_{x_0} \Pr(X_0 = x_0) \ln \frac{\Pr(X_0^\theta = x_0)}{\Pr(X_L^\theta = \epsilon(y_{0:L}, x_0))},$$

where  $\Pr(X_i^\theta)$  is the estimated distribution of the agent memory  $\mathcal{X}$  at time  $i$ , and  $\Pr(X_i)$  is the actual distribution at the same time. This is the average work that is produced from the initial memory distribution  $\Pr(X_0)$  if one applies an efficient information engine based on the model  $\theta$ .

However, we should note that operating on the input string will transform the memory distribution from  $\Pr(X_0 = x_0)$  to  $\Pr(X_L = x_L) = \sum_{y_{0:L}, x_0} \Pr(X_0 = x_0, Y_{0:L} = y_{0:L}) \delta_{x_L, \epsilon(y_{0:L}, x_0)}$ , which the agent will estimate as transforming from  $\Pr(X_0^\theta = x_0)$  to  $\Pr(X_L^\theta = x_L) = \sum_{y_{0:L}, x_0} \Pr(X_0^\theta = x_0, Y_{0:L} = y_{0:L}) \delta_{x_L, \epsilon(y_{0:L}, x_0)}$ . To reset the protocol, we will reset to our estimated initial memory distribution  $\Pr(X_0^\theta)$ , which will involve an efficient transformation for which (according to Thm. 1 of Ref. [10]) the work production for input and output will be:

$$\langle W(x_L \rightarrow x') \rangle = \ln \frac{\Pr(X_L^\theta = x_L)}{\Pr(X_0^\theta = x')}.$$

Note that at the beginning of the reset, the distribution over  $\mathcal{X}$  is given by:

$$\Pr(X_L^\theta = x_L | Y_{0:L}^\theta = y_{0:L}) \equiv \sum_{x_0} \Pr(X_0 = x_0) \delta_{x_L, \epsilon(y_{0:L}, x_0)}.$$

After the reset, the distribution is  $\Pr(X_0^\theta = x')$ , and because we quasistatically evolve to the post-reset distribution, the final distribution is independent of the distribution before reset, meaning that the average work production is:

$$\begin{aligned} & \sum_{x_L, x'} \Pr(X_L^\theta = x_L | Y_{0:L}^\theta = y_{0:L}) \Pr(X_0^\theta = x') \langle W(x_L \rightarrow x') \rangle \\ &= \sum_{x_L, x'} \Pr(X_L^\theta = x_L | Y_{0:L}^\theta = y_{0:L}) \Pr(X_0^\theta = x') \ln \Pr(X_L^\theta = x_L) - \sum_{x_L, x'} \Pr(X_L^\theta = x_L | Y_{0:L}^\theta = y_{0:L}) \Pr(X_0^\theta = x') \ln \Pr(X_0^\theta = x') \\ &= \sum_{x_L} \Pr(X_L^\theta = x_L | Y_{0:L}^\theta = y_{0:L}) \ln \Pr(X_L^\theta = x_L) - \sum_{x'} \Pr(X_0^\theta = x') \ln \Pr(X_0^\theta = x') \\ &= \sum_{x_L} \sum_{x_0} \Pr(X_0 = x_0) \delta_{x_L, \epsilon(y_{0:L}, x_0)} \ln \Pr(X_L^\theta = x_L) - \sum_{x'} \Pr(X_0^\theta = x') \ln \Pr(X_0^\theta = x') \\ &= \sum_{x_0} \Pr(X_0 = x_0) \ln \Pr(X_L^\theta = \epsilon(y_{0:L}, x_0)) - \sum_{x_0} \Pr(X_0^\theta = x_0) \ln \Pr(X_0^\theta = x_0). \end{aligned}$$

When we add this reset cost to the work benefit of harvesting energy, we find the total work production:

$$\beta \langle W^\theta(y_{0:L}) \rangle = \sum_{x_0} \Pr(X_0 = x_0) \Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = x_0) + L \ln |\mathcal{Y}| + \sum_{x_0} (\Pr(X_0 = x_0) - \Pr(X_0^\theta = x_0)) \ln \Pr(X_0^\theta = x_0),$$

where the probability of input  $y_{0:L}$  given the initial state causal state is:

$$\Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = x_0) = \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, x_0)).$$

Let us posit that the agent correctly estimates the initial distribution to be  $p(s_0)$ , such that:

$$p(x_0) = \Pr(X_0 = x_0) = \Pr(X_0^\theta = x_0),$$

either because it prepares the initial memory distribution, or has reliably measured it in the past. Finally, since  $\mathcal{X} = \mathcal{S}$ , we can express the work production on this distributed state:

$$\beta \langle W^\theta(y_{0:L}) \rangle = \sum_{x_0} p(s_0) \Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = s_0) + L \ln |\mathcal{Y}|,$$

which is purely a function of underlying  $\epsilon$ -machine model  $\theta$ .

#### Appendix D: Maximum Work Edge-Weights

For a given engine topology  $\epsilon$ , a particular initial distribution  $p(s)$ , and an additional energy cost of initializing the machine  $C(\theta)$ , the work production from a particular input string  $y_{0:L}$  is:

$$\begin{aligned} \beta \langle W_G^\theta(y_{0:L}) \rangle &= L \ln |\mathcal{Y}| - C(\theta) + \sum_{s_0} p(s_0) \ln \Pr(Y_{0:L}^\theta = y_{0:L} | S_0 = s_0) \\ &= L \ln |\mathcal{Y}| - C(\theta) + \sum_s p(s_0) \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, s_0)). \end{aligned}$$

Here, the subscript  $G$  indicates that this work production is modified in an attempt to make maximum-work training generalize. We find the edge-weights for this topology by counting the input-memory state combinations of the engine when driven by  $y_{0:L}$ .  $N(y, s | s_0, \epsilon, y_{0:L})$  is the number of times predictive memory state  $s$  receives input  $y$  given that an engine with topology  $\epsilon$  started in  $s_0$  and received input word  $y_{0:L}$ . This allows us to rewrite the work production:

$$\begin{aligned} \beta \langle W_G^\theta(y_{0:L}) \rangle &= L \ln |\mathcal{Y}| - C(\theta) + \sum_{s_0} p(s_0) \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, s_0)) \\ &= L \ln |\mathcal{Y}| - C(\theta) + \sum_{s_0} p(s_0) \sum_{s', y'} N(y, s | s_0, \epsilon, y_{0:L}) \ln \theta(y' | s'). \end{aligned}$$

We can find the resulting maximum-work edge-weights by taking the set of constraints:

$$g_{s'}(\theta) \equiv \sum_{y'} \theta(y' | s') = 1,$$

and solving:

$$\begin{aligned}
\partial_{\theta(y|s)}\beta\langle W_G^\theta(y_{0:L})\rangle &= \sum_{s'} \lambda_{s'} \partial_{\theta(y|s)} g_{s'}(\theta) \\
-\partial_{\theta(y|s)}C(\theta) + \sum_{s_0, s', y'} p(s_0) \frac{N(y, s|s_0, \epsilon, y_{0:L})}{\theta(y'|s')} \delta_{sy, s'y'} &= \sum_{s'} \lambda_{s'} \sum_{y'} \delta_{sy, s'y'} \\
-\partial_{\theta(y|s)}C(\theta) + \frac{\sum_{s_0} p(s_0) N(y, s|s_0, \epsilon, y_{0:L})}{\theta(y|s)} &= \lambda_s \\
\frac{\sum_{s_0} p(s_0) N(y, s|s_0, \epsilon, y_{0:L}) - \theta(y|s) \partial_{\theta(y|s)}C(\theta)}{\lambda_s} &= \theta(y|s) \\
\frac{\sum_{s_0} p(s_0) (N(y, s|s_0, \epsilon, y_{0:L}) - \theta(y|s) \partial_{\theta(y|s)}C(\theta))}{\lambda_s} &= \theta(y|s).
\end{aligned}$$

The constraint of normalized edge-weights  $\sum_y \theta(y|s)$  allows us to solve for  $\lambda_s$  as the normalization constant:

$$\sum_{y', s_0} p(s_0) (N(y, s|s_0, \epsilon, y_{0:L}) - \theta(y'|s) \partial_{\theta(y'|s)}C(\theta)) = \lambda_s.$$

The resulting maximum-work edge-weights  $\Theta_{p, \epsilon, y_{0:L}}^{\max}(y|s)$  are the solution to the recursive relation:

$$\theta(y|s) = \frac{\sum_{s_0} p(s_0) (N(y, s|s_0, \epsilon, y_{0:L}) - \theta(y|s) \partial_{\theta(y|s)}C(\theta))}{\sum_{y', s_0} p(s_0) (N(y', s|s_0, \epsilon, y_{0:L}) - \theta(y'|s) \partial_{\theta(y'|s)}C(\theta))}.$$

Here, we consider the cost:

$$C(\theta) = -\alpha \sum_{y', s'} \ln \theta(y'|s') + \text{const.},$$

such that:

$$-\partial_{\theta(y|s)}C(\theta) = \frac{\alpha}{\theta(y|s)},$$

and we can remove the dependence on the right-hand side of the equation. Our resulting solution for edge-weights is:

$$\Theta_{p, \alpha, \epsilon, y_{0:L}}^{\max}(y|s) = \frac{\sum_{s_0} p(s_0) (\alpha + N(y, s|s_0, \epsilon, y_{0:L}))}{\sum_{y', s_0} p(s_0) (\alpha + N(y', s|s_0, \epsilon, y_{0:L}))}.$$

Note that if we are implementing the transformation quasistatically, the distribution  $\theta(y|s)$  must be the equilibrium distribution over  $\mathcal{Y}$  when conditioned on the engine memory state  $s$ . If we have the initial energy landscape  $E(y, s)$ , then the equilibrium distribution is:

$$\pi(y, s) = e^{\beta(F^{\text{eq}} - E(y, s))},$$

meaning:

$$\begin{aligned}
\theta(y|s) &= \frac{\pi^\theta(y, s)}{\pi^\theta(s)} \\
&= \frac{e^{-\beta E^\theta(y, s)}}{\sum_y e^{-\beta E^\theta(y, s)}},
\end{aligned}$$

where  $\pi^\theta(s) \equiv \sum_y \pi^\theta(y, s)$  is the marginal equilibrium distribution over the agent memory.

Note that we can also write the metastable free energy [29, 84, 85], of the memory state  $s$ :

$$\beta F^\theta(s) = -\ln \sum_y e^{-\beta E^\theta(y, s)}.$$

We can express the edge-weights in terms of the free energy of memory state  $s$  and the initial energy of:

$$\theta(y|s) = e^{\beta(F^\theta(s) - E^\theta(s,y))}.$$

Thus, the cost can be expressed as the thermodynamic quantity:

$$\begin{aligned} C(\theta) &= -\alpha \sum_{y',s'} \ln \theta(y'|s') \\ &= \alpha \sum_{y',s'} \beta(E^\theta(s,y) - F^\theta(s)). \end{aligned}$$

If we incur a cost that is proportional to the total energetic excess beyond the free energy for each memory state, then we will find the solution for the edge-weights given by  $\Theta_{p,\alpha,\epsilon,y_{0:L}}^{\max}(y|s)$ .

This can be recovered by considering the cost of initializing every edge-weight and relaxing into equilibrium with every predictive state  $s$ . Mechanically, this relates to the fact that the energy of a state is related to its equilibrium distribution:

$$E(z) = F^{\text{eq}} - \ln \Pr(Z^{\text{eq}} = z).$$

At the beginning of a quasistatic protocol, the energy landscape must be in equilibrium with the estimated distribution, meaning that:

$$E(y, s) = F^{\text{eq}} - \ln \theta(y|s)\pi(s),$$

where  $\pi(s) = \Pr(S^{\text{eq}} = s)$ . For each state, we must initialize it. This corresponds to confining the state to  $y$  and  $s$ , then letting it relax to equilibrium with output distribution  $\theta(y|s)$ . The amount of work dissipated in this relaxation is the change in nonequilibrium addition to free energy, which is the relative entropy:

$$\begin{aligned} \langle W_{\text{diss}}(y, s) \rangle &= D_{KL}(\delta_{y,y'}\delta_{s,s'} || \theta(y'|s')\pi(s)) - D_{KL}(\theta(y'|s')\delta_{s,s'} || \theta(y'|s')\pi(s')) \\ &= \sum_{s'y'} \delta_{y,y'}\delta_{s,s'} \ln \frac{\delta_{y,y'}\delta_{s,s'}}{\theta(y'|s')\pi(s')} - \sum_{s'y'} \theta(y'|s')\delta_{s,s'} \ln \frac{\theta(y'|s')\delta_{s,s'}}{\theta(y'|s')\pi(s')} \\ &= -\ln \theta(y|s)\pi(s) - \sum_{y'} \theta(y'|s) \ln \frac{1}{\pi(s)} \\ &= -\ln \theta(y|s)\pi(s) + \ln \pi(s) \\ &= -\ln \theta(y|s). \end{aligned}$$

If we incur this dissipation for every combination of predictive state and input, then the total dissipated work is:

$$\langle W_{\text{diss}}^{\text{prepare}} \rangle = - \sum_{s,y} \ln \theta(y|s).$$

Thus, our cost is proportional to our original cost function  $C(\theta)$ , meaning that we can set:

$$C(\theta) = \alpha \langle W_{\text{diss}}^{\text{prepare}} \rangle,$$

to solve for the the same maximum-work edge-weights  $\Theta_{p,\alpha,\epsilon,y_{0:L}}^{\max}(y|s)$  described in Eq. (D1).

We investigate four cases:

1. Without any attempt at generalization, we choose  $\alpha = 0$  and a peaked initial distribution  $p(s_0) = \delta_{s_0, s^*}$ , then we obtain the standard likelihood expression for work:

$$\begin{aligned} \beta \langle W_G^\theta(y_{0:L}) \rangle &= L \ln |\mathcal{Y}| + \ln \Pr(Y_{0:L} = y_{0:L} | S_0 = s^*) \\ &= L \ln |\mathcal{Y}| + \ln \prod_{i=0}^{L-1} \theta(y_i | \epsilon(y_{0:i}, s^*)). \end{aligned} \tag{D1}$$

which produces the familiar MLE estimate for the edge-weights:

$$\Theta_{\delta_{s,s^*},0,\epsilon,y_{0:L}}^{\max}(y|s) = \frac{N(y, s|s^*, \epsilon, y_{0:L})}{\sum_{y'} N(y', s|s^*, \epsilon, y_{0:L})}.$$

2. If we attempt to generalize through autocorrection, we choose a uniform initial state  $p(s) = 1/|\mathcal{S}|$  and zero contribution from the complexity cost  $\alpha = 0$ , yielding the work production:

$$\beta\langle W_G^\theta(y_{0:L}) \rangle = L \ln |\mathcal{Y}| + \frac{1}{|\mathcal{S}|} \sum_{s_0} \ln \Pr(Y_{0:L} = y_{0:L} | S_0 = s_0).$$

The modified edge-weight estimator includes the contributions from every start state:

$$\Theta_{1/|\mathcal{S}|,0,\epsilon,y_{0:L}}^{\max}(y|s) = \frac{\sum_{s_0} N(y, s|s_0, \epsilon, y_{0:L})}{\sum_{y',s_0} N(y', s|s_0, \epsilon, y_{0:L})}. \quad (\text{D2})$$

3. We might also try to generalize by only adding a cost from the dissipation of initializing the system at  $\alpha = 1$ , while allowing a unique start state  $p(s) = \delta_{s,s^*}$ , such that the work production is:

$$\beta\langle W_G^\theta(y_{0:L}) \rangle = \ln \Pr(Y_{0:L} = y_{0:L} | S_0 = s^*) + L \ln |\mathcal{Y}| - |\mathcal{Y}| \sum_{s'} \beta\langle W_{\text{diss}}^\theta(s') \rangle.$$

The resulting maximum-work estimator is:

$$\Theta_{\delta_{s,s^*},1,\epsilon,y_{0:L}}^{\max}(y|s) = \frac{1 + N(y, s|s^*, \epsilon, y_{0:L})}{|\mathcal{Y}| + \sum_{y'} N(y', s|s^*, \epsilon, y_{0:L})}.$$

This is precisely Laplace's rule of succession applied to each causal state, which is the result of Bayesian updating from a uniform distribution over output probabilities.

4. Last, we combine the complexity cost  $\alpha = 1$  with the cost of autocorrection  $p(s) = 1/|\mathcal{S}|$  such that the work production is:

$$\beta\langle W_G^\theta(y_{0:L}) \rangle = \frac{1}{|\mathcal{Y}|} \sum_s \ln \Pr(Y_{0:L} = y_{0:L} | S_0 = s) + L \ln |\mathcal{Y}| - |\mathcal{Y}| \sum_{s'} \beta\langle W_{\text{diss}}^\theta(s') \rangle.$$

The maximum-work edge-weights combine the benefits of both generalization strategies:

$$\Theta_{1/|\mathcal{Y}|,1,\epsilon,y_{0:L}}^{\max}(y|s) = \frac{\sum_{s_0} (1 + N(y, s|s_0, \epsilon, y_{0:L}))}{\sum_{y',s_0} (1 + N(y', s|s_0, \epsilon, y_{0:L}))}.$$

## Appendix E: Entropy Production as Divergence

We monitor engine inefficiency via entropy production (dissipated work) [29]. This is the net work invested minus the change in nonequilibrium free energy:

$$\langle \Sigma^\theta \rangle_{0:L} / k_B = \beta(-\langle W^\theta \rangle_{0:L} - \Delta F_{0:L}^{\text{NEQ}}).$$

Because an information reservoir has equal energy for all configurations, the contributions to the nonequilibrium free energy are only informational:

$$\begin{aligned} \beta \Delta F_{0:L}^{\text{NEQ}} &= -\Delta H_{0:L} \\ &= H[Y_{0:L}^\theta] - L \ln |\mathcal{Y}| \end{aligned}$$

where  $H[Z] \equiv -\sum_z \Pr(Z = z) \ln \Pr(Z = z)$  is the Shannon entropy of random variable  $Z$  measured in Nats. The average work, by contrast, is:

$$\langle W^\theta \rangle_{0:L} = L \ln |\mathcal{Y}| + \sum_{y_{0:L}} \Pr(Y_{0:L}^{\theta'} = y_{0:L}) \ln \Pr(Y_{0:L}^\theta = y_{0:L}).$$

Adding these terms together, we get the relative entropy between the true input process  $Y_{0:L}^{\theta'}$  and the estimated process  $Y_{0:L}^\theta$ :

$$\begin{aligned} \langle \Sigma^\theta \rangle_{0:L} / k_B &= \sum_{y_{0:L}} \Pr(Y_{0:L}^{\theta'} = y_{0:L}) \ln \frac{\Pr(Y_{0:L}^{\theta'} = y_{0:L})}{\Pr(Y_{0:L}^\theta = y_{0:L})} \\ &\equiv D_{KL}(Y_{0:L}^{\theta'} || Y_{0:L}^\theta), \end{aligned}$$

which is the additional dissipation that is incurred by misestimating the input distribution [49, 50].

If a learning process refines and improves the estimator  $\theta$ , its divergence from the actual process  $Y_{0:L}^{\theta'}$  should diminish, reflecting the idea that learning reduces entropy production [15, 51]. The consistency of MLE guarantees that, as long as the true process can be described by one of the available models, the learning process will discover the true distribution [58, 86]. Thus, given a sufficiently large class of  $\epsilon$ -machines to select from, thermodynamic machine learning will discover the hidden process and minimize the average entropy production to zero for future inputs.

For an information engine harvesting energy from an information reservoir, the asymptotic rate of change in free energy is the difference between the entropy rates of the input process and output process [39]:

$$\begin{aligned} \beta \Delta F_\infty^{\text{NEQ}} &\equiv \lim_{L \rightarrow \infty} \beta (\Delta F_{0:L+1}^{\text{NEQ}} - \Delta F_{0:L}^{\text{NEQ}}) \\ &= h_\mu^{\theta'} - \ln |\mathcal{Y}|. \end{aligned}$$

Here, the entropy rate  $h_\mu^{\theta'}$  of the inputs can be directly calculated from its  $\epsilon$ -machine [13]:

$$h_\mu^{\theta'} = - \sum_{s', y} \pi_{s'}' \theta'(y|s') \ln \theta'(y|s'),$$

where  $\pi_{s'}' = \sum_s \pi_{s, s'}$  is the steady-state distribution of the true  $\epsilon$ -machine's causal states. As a result, the asymptotic entropy production rate can be expressed as the average divergence between the edge-weights:

$$\begin{aligned} \langle \Sigma^\theta \rangle_\infty / k_B &= \beta (-\langle W^\theta \rangle_\infty - \Delta F_\infty^{\text{NEQ}}) \\ &= \sum_{s, s', y} \pi_{s, s'} \theta'(y|s') \ln \frac{\theta'(y|s')}{\theta(y|s)} \\ &= \sum_{s, s'} \pi_{s, s'} D_{KL}(Y_i^{\theta'} | S_i' = s' || Y_i^\theta | S_i = s). \end{aligned}$$

Here:

$$D_{KL}(Y_i^{\theta'} | S_i' = s' || Y_i^\theta | S_i = s) \equiv \sum_y \Pr(Y_i^{\theta'} = y | S_i' = s') \ln \frac{\Pr(Y_i^{\theta'} = y | S_i' = s')}{\Pr(Y_i^\theta = y | S_i = s)}$$

is the divergence between the prediction of the next input from state memory state  $s$  in model  $\theta$  and the prediction from memory state  $s'$  in model  $\theta'$ .

## Appendix F: Training and Testing Individual Words

As an illustrative example, consider two training words of length 100 generated from the Five-State machine shown in Fig. 5. Figure 9 shows the resulting (exponentiated) maximum training work production rate  $e^{\beta \langle W_n^{\text{max}}(y_{0:L}) \rangle} / L$  (solid lines) for memory sizes  $n \in \{1, 2, 3\}$  when we train on length-1 to length-100 for

two different training words. The dashed lines show the (exponentiated) asymptotic work production rate  $e^{\beta(W_n^{\Theta_{\text{max}}(y_0:L)})_{\infty}}$ , representing our performance in testing.

In the upper-left corner of each diagram, we see the result of un-regularized work-maximization ( $\alpha = 0$  and  $p(s) = \delta_{s,s^*}$ ), which yields improved work production in training with greater memory. But, it dangerously overfits, dissipating divergent work for large model memories.

By contrast, we see in the upper right corners that when  $\alpha = 1$  and  $p(s) = \delta_{s,s^*}$ , leading to Bayesian estimates of edge-weights, the model doesn't divergently overfit. For some training words (word 1 of Fig. 9), we see that the model learns how to extract almost all available work for the large-memory case ( $n = 3$ ). However, other training words (word 2 of Fig. 9) show that, even though the model doesn't divergently overfit, it dissipates considerable energy, producing negative work on average for the large-memory case.

The case of autocorrection in the bottom left, when  $\alpha = 0$  and  $p(s) = 1/|\mathcal{S}|$ , leads to error mitigation as well. There is still divergent dissipation for very short training words shown in Fig. 9, but the thermodynamic learning appears to discover useful patterns in the input data, with a growing advantage for larger memories. This suggests that autocorrection mitigates overfitting for complex processes.

Finally, in the bottom right ( $\alpha = 1$  and  $p(s) = 1/|\mathcal{S}|$ ), where we have combined regularization strategies to both require autocorrection and a complexity cost in initializing the model, we see that overfitting appears to be largely mitigated. Generalizing in this way appears to allow us to add memory without adding considerable risk of overfitting. This is promising, as we would like to be able to perform thermodynamic learning with much more complex  $\epsilon$ -machine models to discover complex patterns in data without the risk of projecting structure onto the data that isn't there.

On the whole, it appears that the result of training varies considerably depending on the training word. This is as expected for short word lengths, as the same word can come from many different processes. For this reason, we now explore in more detail as we vary across ensembles of training words.

### Appendix G: Learning the Even Process and Noisy Even Process

In contrast to the ‘‘Five-State Process,’’ the ‘‘Even Process’’ and ‘‘Noisy Even Process’’ require only two states for perfect prediction, so they reveal the case when we have more memory accessible ( $n = 3$ ) than is strictly necessary. In this case, the true model is contained within our class of model candidates.

The ‘‘Noisy Even Process’’ is an interesting case for the same reason, but it has full support in the possible words that it can produce. A good learning algorithm should robustly learn patterns that may contain rare fluctuations and noise, and the ‘‘Noisy Even Process’’ behaves much like the ‘‘Even process,’’ but its rare fluctuations are extremely important to thermodynamic behavior [87].

In Fig 10, we compare the training work rate to the testing work rate for 200 samples for all three processes. We see in all cases that the MLE strategy works best for training data and words for test data, while the CMBD method works worst for training and best for testing. However, the relative effectiveness of the BAYES and AC method changes based on which true process is being sampled. This suggests that autocorrection serves to address some features of processes, while Bayesian edge-weight updates serve to address others.

The even process appears to be easy to learn for all cases, as seen in Fig. 11. The unregularized MLE strategy appears to do worst, because its variance is the highest while its average is comparable to the regularized learning techniques. BAYES, AC, and CMBD all appear to have comparable results for the testing work rate, with BAYES perhaps performing the best. In this case, unlike for the Five-State process, we get to see what happens when our model has enough memory to fully capture the pattern in the data. Interestingly, additional memory then comes at a cost for all regularization techniques, which can be identified by noting that the blue curve ( $n = 3$ ) lies at or below the purple curve ( $n = 2$ ). This may be because of the fact that allowing three memory states simply creates more opportunity for overfitting when the process is fully described by two predictive states. However, the CMBD strategy seems to mostly mitigate the overfitting, because the blue curve is very close to the purple one.

Figure 12 shows that the noisy even process appears to be much more difficult to learn than the even process. We hypothesize that this is because we need to include words that are rare fluctuations to fully see the support of the process. When a model doesn't anticipate those words, it results in overfitting and negative work production. We see that the unregularized MLE and AC techniques encounter this problem frequently. By contrast, the BAYES and CMBD strategies seem to discover the underlying pattern after training on length-128. However,



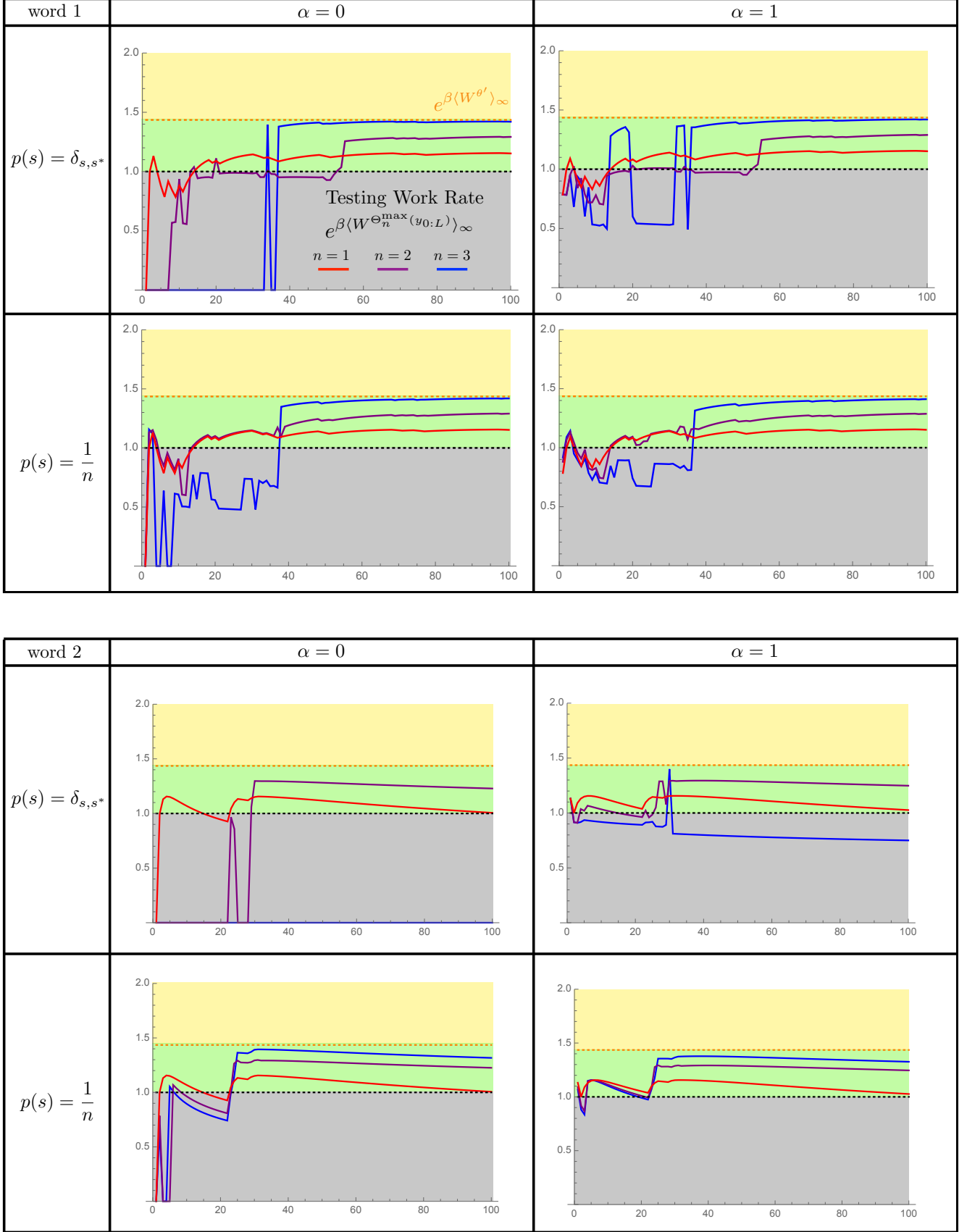
$\theta' = \text{five state}$ 

FIG. 9. The top and bottom correspond to two different training words (“word 1” and “word 2”) from the “Five-State”  $\epsilon$ -machine. We plot testing work rate with four different regularization strategies: 1) Unregularized TML ( $\alpha = 0$  and  $p(s) = \delta_{s,s^*}$ ) 2) Autocorrection ( $\alpha = 0$  and  $p(s) = 1/|\mathcal{S}|$ ) 3) Bayesian complexity cost ( $\alpha = 1$  and  $p(s) = \delta_{s,s^*}$ ) 4) Combined ( $\alpha = 1$  and  $p(s) = 1/|\mathcal{S}|$ ).

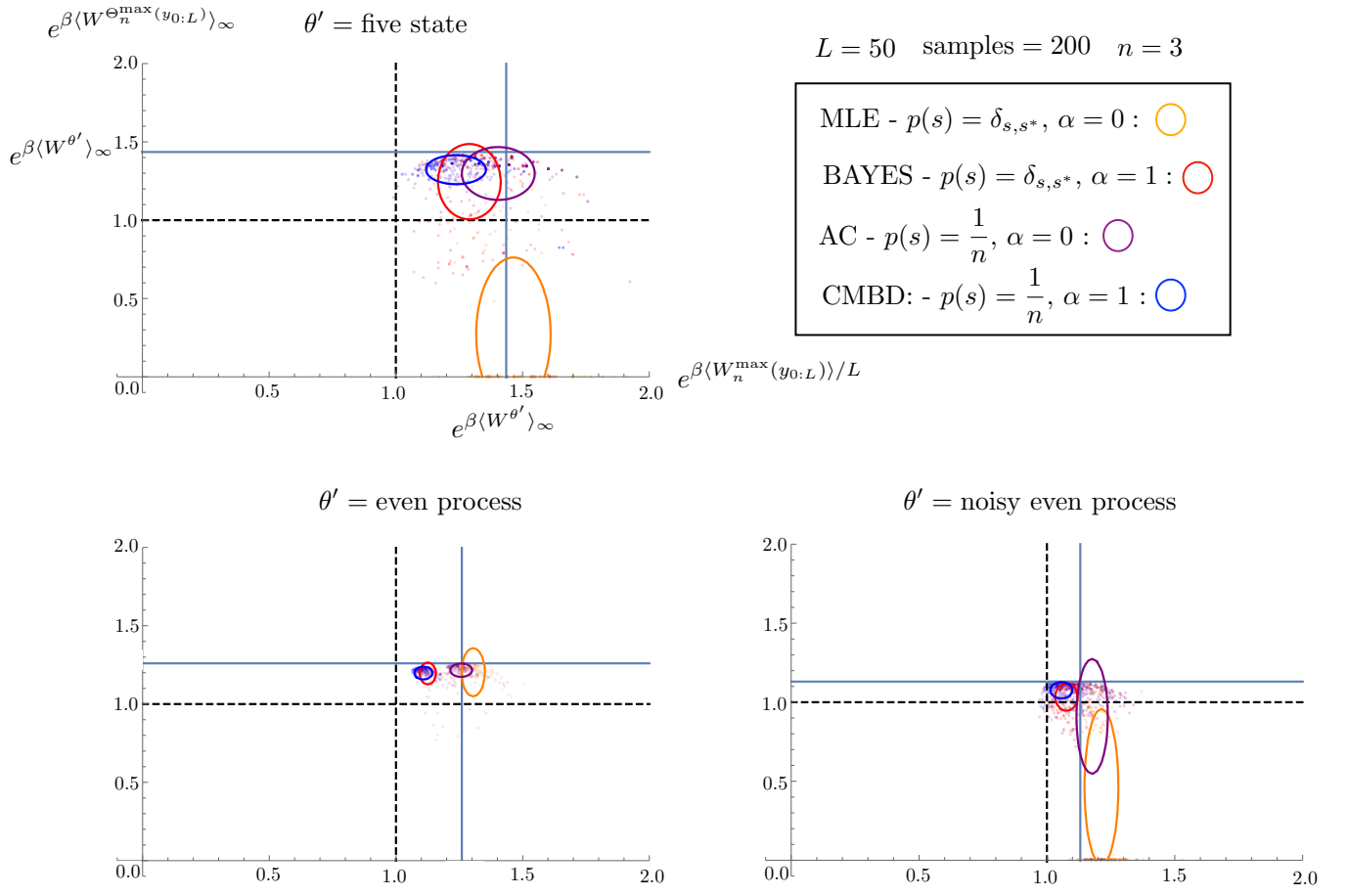


FIG. 10. Thermodynamic Machine Learning Performance for Different Processes: For each input process  $\theta'$  (Five-State, Even, and Noisy Even), we randomly sample 200 length  $L = 50$  words and train a 3-state  $\epsilon$ -machines on each to find the exponential training work rate  $e^{\beta\langle W_{n=3}^{\max}(y_{0:L})\rangle/L}$  and the exponential testing work rate  $e^{\beta\langle W_{n=3}^{\Theta_{\max}(y_{0:L})}\rangle_{\infty}}$  for each regularization strategy: MLE (orange), BAYES (red), AC (purple), and CMBD (blue). The ovals are centered around the average work rates of these 200 samples, and their dimensions are given by the variance of the work rates. The dashed black lines represent work rates of zero along each dimension, and the blue lines represent the theoretical limit on the asymptotic work rate, given by  $e^{\beta\langle W^{\theta'}\rangle}$ .

the BAYES technique still yields high variance, meaning that some outcomes are dissipating a lot of work. By contrast, the CMBD algorithm has low variance and its average is nearly the theoretical limit on work harvesting, indicating that it is reliably discovering the pattern.

$\theta' = \text{even process}$ 

samples = 200

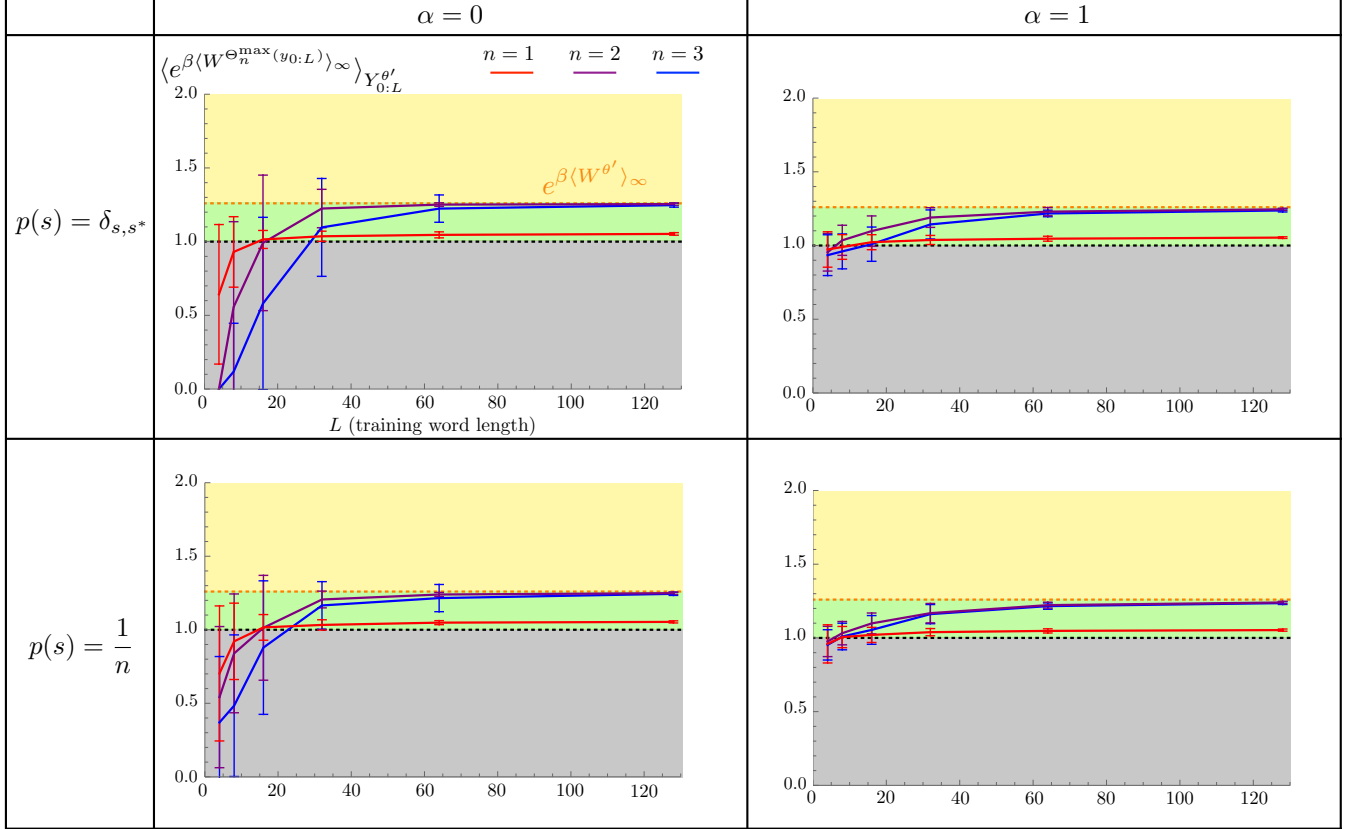


FIG. 11. The average and variance of the exponential asymptotic testing work rate that results from the four different learning strategies. We generate 200 words from the even process for each length  $L \in \{4, 8, 16, 32, 64, 128\}$ , then train on each using MLE, BAYES, AC, and CMBD. The un-regularized thermodynamic machine learning does a fairly good job of fitting, in comparison to the Five-State process. The AC method is a slight improvement over MLE, but not considerable. The BAYES and CMBD regularization techniques have less divergent work production for small training words, and seem to perform slightly better in general. One notable distinction between the BAYES and CMBD methods is that there is a clear disadvantage to using larger memory than necessary ( $n = 3$  instead of  $n = 2$ ) for the BAYES technique, but the CMBD technique does not seem to have this difference.

$\theta' = \text{noisy even process}$ 

samples = 200

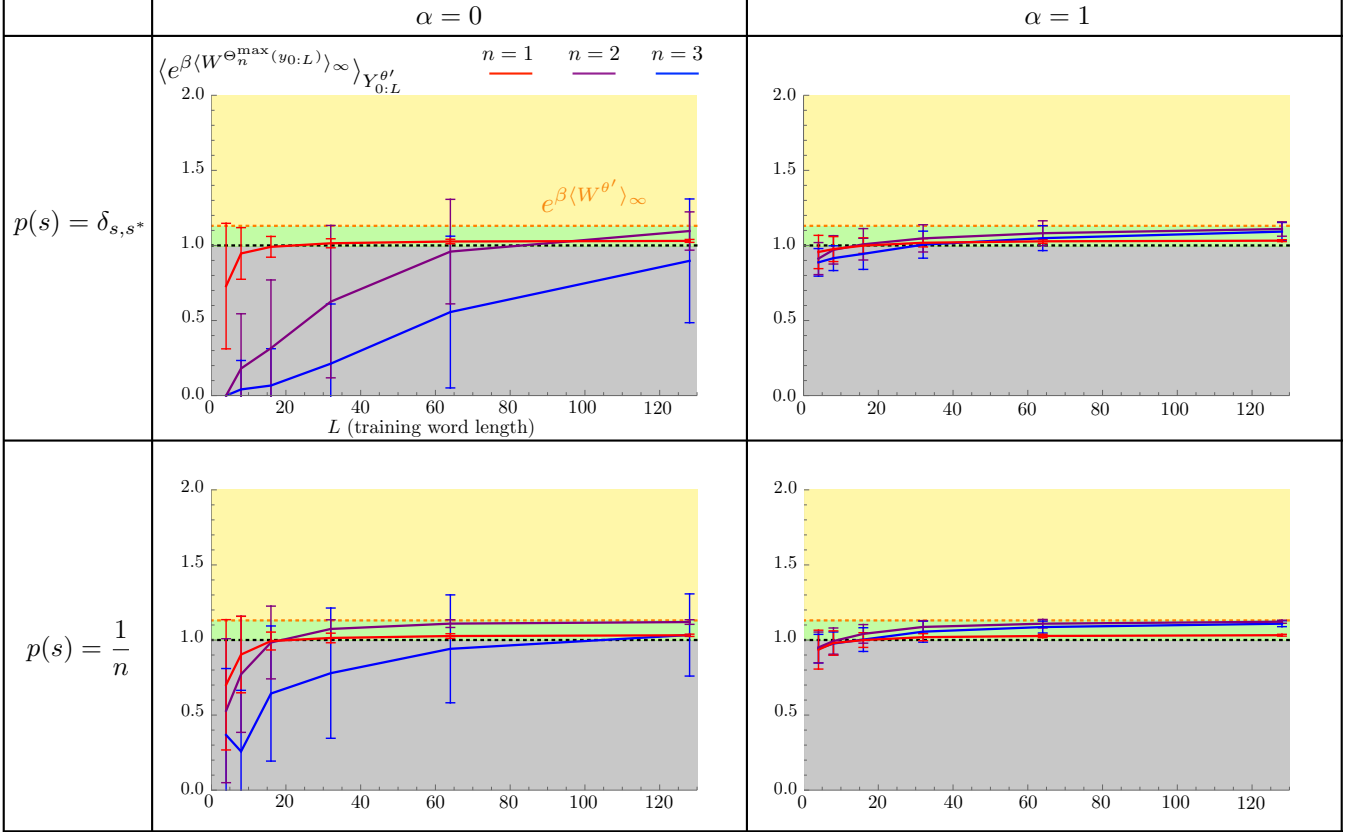


FIG. 12. The average and variance of the exponential asymptotic testing work rate that results from the four different learning strategies. We generate 200 words from the noisy even process for each length  $L \in \{4, 8, 16, 32, 64, 128\}$ , then train on each using MLE, BAYES, AC, and CMBD. This process appears much harder to learn than the even process for the unregularized MLE and AC techniques. The BAYES technique does better, even after 128 training words, the variance in the outcome is still relatively large, indicating that many samples are failing to effectively learn the process. By contrast, the CMBD technique appears to approach the upper bound on work production, with variance that's relatively small for memory sizes  $n = 2$  and  $n = 3$ .

- 
- [1] AD Dongare, RR Kharde, Amit D Kachare, et al. “Introduction to artificial neural network”. *International Journal of Engineering and Innovative Technology (IJEIT)* **2**, 189–194 (2012). url: <https://api.semanticscholar.org/CorpusID:212457035>.
- [2] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. “Deep unsupervised learning using nonequilibrium thermodynamics”. In *Proceedings of the 32nd International Conference on Machine Learning*. Volume 37, pages 2256–2265. PMLR (2015). url: <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- [3] Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein, and Surya Ganguli. “Statistical mechanics of deep learning”. *Annual Review of Condensed Matter Physics* **11**, 501–528 (2020).
- [4] Denis Melanson, Mohammad Abu Khater, Maxwell Aifer, Kaelan Donatella, Max Hunter Gordon, Thomas Ahle, Gavin Crooks, Antonio J Martinez, Faris Sbahi, and Patrick J Coles. “Thermodynamic computing system for ai applications” (2023). arXiv:2312.04836.
- [5] Patrick J Coles, Collin Szczepanski, Denis Melanson, Kaelan Donatella, Antonio J Martinez, and Faris Sbahi. “Thermodynamic ai and the fluctuation frontier”. In *2023 IEEE International Conference on Rebooting Computing (ICRC)*. Pages 1–10. IEEE (2023).
- [6] Karl Friston. “The free-energy principle : a unified brain theory?”. *Nature Reviews Neuroscience* **11**, 127–138 (2010).
- [7] Abdullahi Ali, Nasir Ahmad, Elgar de Groot, Marcel Antonius Johannes van Gerven, and Tim Christian Kietzmann. “Predictive coding is a consequence of energy efficiency in recurrent neural networks”. *Patterns* **3**, 100661 (2022).
- [8] Agrim Gupta, Silvio Savarese, Surya Ganguli, and Li Fei-Fei. “Embodied intelligence via learning and evolution”. *Nature communications* **12**, 5721 (2021).
- [9] Christopher Bishop. “Pattern Recognition and Machine Learning”. Springer. (2006). url: <https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/>.
- [10] Alexander B Boyd, James P Crutchfield, and Mile Gu. “Thermodynamic machine learning through maximum work production”. *New Journal of Physics* **24**, 083040 (2021).
- [11] Susanne Still, David A. Sivak, Anthony J. Bell, and Gavin E. Crooks. “Thermodynamics of Prediction”. *Physical Review Letters* **109**, 120604 (2012).
- [12] J. P. Crutchfield and K. Young. “Inferring statistical complexity”. *Physical Review Letters* **63**, 105–108 (1989).
- [13] J. P. Crutchfield and D. P. Feldman. “Regularities unseen, randomness observed: Levels of entropy convergence”. *CHAOS* **13**, 25–54 (2003).
- [14] J. P. Crutchfield. “Between order and chaos”. *Nature Physics* **8**, 17–24 (2012).
- [15] Jacob M Gold and Jeremy L England. “Self-organized novelty detection in driven spin glasses” (2019).
- [16] Xue Ying. “An overview of overfitting and its solutions”. *Journal of physics: Conference series* **1168**, 022022 (2019).
- [17] A. B. Boyd, D. Mandal, and J. P. Crutchfield. “Thermodynamics of modularity: Structural costs beyond the landauer bound”. *Physical Review X* **8**, 031036 (2018).
- [18] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. “Deep double descent: Where bigger models and more data hurt”. *Journal of Statistical Mechanics: Theory and Experiment* **2021**, 124003 (2021).
- [19] Robert Tibshirani. “Regression shrinkage and selection via the lasso”. *Journal of the Royal Statistical Society Series B: Statistical Methodology* **58**, 267–288 (1996).
- [20] Yiyun Zhang, Runze Li, and Chih-Ling Tsai. “Regularization parameter selections via generalized information criterion”. *Journal of the American statistical Association* **105**, 312–323 (2010).
- [21] Gouhei Tanaka, Toshiyuki Yamane, Jean Benoit Héroux, Ryosho Nakane, Naoki Kanazawa, Seiji Takeda, Hidetoshi Numata, Daiju Nakano, and Akira Hirose. “Recent advances in physical reservoir computing: A review”. *Neural Networks* **115**, 100–123 (2019).
- [22] Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urtasun, and Richard Zemel. “Reviving and improving recurrent back-propagation”. In *International Conference on Machine Learning*. Volume 80, pages 3082–3091. PMLR (2018). url: <https://proceedings.mlr.press/v80/liao18c.html>.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. “Attention is all you need”. *Advances in neural information processing systems* **30**, 5998–6008 (2017). url: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- [24] J. C. Maxwell. “Theory of heat”. Longmans, Green and Co. London, United Kingdom (1871).
- [25] W. Thomson. “Kinetic theory of the dissipation of energy”. *Nature* **9**, 441–444 (1874).
- [26] Harvey Leff and Andrew F. Rex, editors. “Maxwell’s Demon 2: Entropy, Classical and Quantum Information, Computing”. CRC Press. (2002).
- [27] Leo Szilard. “Über die entropieverminderung in einem thermodynamischen system bei eingriffen intelligenter wesen”. *Zeitschrift für Physik* **53**, 840–856 (1929).
- [28] R. Landauer. “Irreversibility and heat generation in the computing process”. *IBM J. Res. Develop.* **5**, 183–191 (1961).
- [29] J. M. R. Parrondo, J. M. Horowitz, and T. Sagawa. “Thermodynamics of information”. *Nature Physics* **11**, 131–139 (2015).
- [30] Massimiliano Esposito, Katja Lindenberg, and Christian Van den Broeck. “Entropy production as correlation between system and reservoir”. *New Journal of Physics* **12**, 013013 (2010).
- [31] David Reeb and Michael M Wolf. “An improved Landauer principle with finite-size corrections”. *New Journal of Physics* **16**, 103011 (2014).
- [32] Philip Taranto, Faraj Bakhshinezhad, Andreas Bluhm, Ralph Silva, Nicolai Friis, Maximilian P.E. Lock, Giuseppe Vitagliano, Felix C. Binder, Tiago Debarba,

- Emanuel Schwarzthans, Fabien Clivaz, and Marcus Huber. “Landauer Versus Nernst: What is the True Cost of Cooling a Quantum System”. *PRX Quantum* **4**, 010332 (2023).
- [33] G. E. Crooks. “Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences”. *Physical Review E* **60**, 2721 (1999).
- [34] Massimiliano Esposito and Christian Van den Broeck. “Second law and landauer principle far from equilibrium”. *Europhysics Letters* **95**, 40004 (2011).
- [35] Udo Seifert. “Entropy production along a stochastic trajectory and an integral fluctuation theorem”. *Physical Review Letters* **95**, 040602 (2005).
- [36] Andrew JP Garner, Jayne Thompson, Vlatko Vedral, and Mile Gu. “Thermodynamics of complexity and pattern manipulation”. *Physical Review E* **95**, 042140 (2017).
- [37] Léo Touzo, Matteo Marsili, Neri Merhav, and Édgar Roldán. “Optimal work extraction and the minimum description length principle”. *Journal of Statistical Mechanics: Theory and Experiment* **2020**, 093403 (2020).
- [38] D. Mandal and C. Jarzynski. “Work and information processing in a solvable model of Maxwell’s demon”. *Proc. Natl. Acad. Sci. USA* **109**, 11641–11645 (2012).
- [39] A. B. Boyd, D. Mandal, and J. P. Crutchfield. “Identifying functional thermodynamics in autonomous Maxwellian ratchets”. *New Journal of Physics* **18**, 023049 (2016).
- [40] Philipp Strasberg, Javier Cerrillo, Gernot Schaller, and Tobias Brandes. “Thermodynamics of stochastic turing machines”. *Physical Review E* **92**, 042104 (2015).
- [41] Henry W Lin, Max Tegmark, and David Rolnick. “Why does deep and cheap learning work so well?”. *Journal of Statistical Physics* **168**, 1223–1247 (2017).
- [42] Daniel Ray Upper. “Theory and algorithms for hidden markov models and generalized hidden markov models”. University of California, Berkeley. (1997).
- [43] Christopher C Strelhoff and James P Crutchfield. “Bayesian structural inference for hidden processes”. *Physical Review E* **89**, 042119 (2014).
- [44] Sandesh Adhikary, Siddharth Srinivasan, Geoff Gordon, and Byron Boots. “Expressiveness and learning of hidden quantum markov models”. In *International Conference on Artificial Intelligence and Statistics*. Pages 4151–4161. PMLR (2020). url: <https://proceedings.mlr.press/v108/adhikary20a.html>.
- [45] Sourya Basu, Moulik Choraria, and Lav R Varshney. “Transformers are universal predictors” (2023). arXiv:2307.07843.
- [46] Gerda Claeskens, Nils Lid Hjort, et al. “Model selection and model averaging”. Volume 330. Cambridge University Press Cambridge. (2008).
- [47] Lucas B Vieira and Costantino Budroni. “Temporal correlations in the simplest measurement sequences”. *Quantum* **6**, 623 (2022).
- [48] Alexandra M Jurgens and James P Crutchfield. “Shannon entropy rate of hidden markov processes”. *Journal of Statistical Physics* **183**, 32 (2021).
- [49] A. Kolchinsky and D. H. Wolpert. “Dependence of dissipation on the initial distribution over states”. *J. Stat. Mech.: Th. Expt.* **2017**, 083202 (2017).
- [50] P. M. Riechers and M. Gu. “Initial-state dependence of thermodynamic dissipation for any quantum process”. *Physical Review E* **Page 042145** (2020). arXiv:2002.11425.
- [51] G. J. Milburn. “Quantum learning machines” (2023). arXiv:2305.07801.
- [52] S. Still, J. P. Crutchfield, and C. J. Ellison. “Optimal causal inference: Estimating stored information and approximating causal architecture”. *CHAOS* **20**, 037111 (2010).
- [53] Felix Creutzig, Amir Globerson, and Naftali Tishby. “Past-future information bottleneck in dynamical systems”. *Physical Review E* **79**, 041925 (2009).
- [54] Sarah E Marzen and James P Crutchfield. “Predictive rate-distortion for infinite-order markov processes”. *Journal of Statistical Physics* **163**, 1312–1338 (2016).
- [55] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. “The elements of statistical learning: data mining, inference, and prediction”. Volume 2. Springer. (2009).
- [56] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. “Dropout: a simple way to prevent neural networks from overfitting”. *The journal of machine learning research* **15**, 1929–1958 (2014). url: <http://jmlr.org/papers/v15/srivastava14a.html>.
- [57] A. B. Boyd, D. Mandal, and J. P. Crutchfield. “Correlation-powered information engines and the thermodynamics of self-correction”. *Physical Review E* **95**, 012152 (2017).
- [58] Edwin T Jaynes. “Probability theory: The logic of science”. Cambridge university press. (2003).
- [59] Peter D Grünwald. “The minimum description length principle”. MIT press. (2007).
- [60] Y. Jun, M. Gavrilov, and J. Bechhoefer. “High-precision test of Landauer’s principle in a feedback trap”. *Physical Review Letters* **113**, 190601 (2014).
- [61] A. B. Boyd, D. Mandal, P. M. Riechers, and J. P. Crutchfield. “Transient dissipation and structural costs of physical information transduction”. *Physical Review Letters* **118**, 220602 (2017).
- [62] A. B. Boyd, D. Mandal, and J. P. Crutchfield. “Leveraging environmental correlations: The thermodynamics of requisite variety”. *Journal of Statistical Physics* **167**, 1555–1585 (2016).
- [63] Marco Radaelli, Gabriel T Landi, Kavan Modi, and Felix C Binder. “Fisher information of correlated stochastic processes”. *New Journal of Physics* **25**, 053037 (2023).
- [64] Paul M Riechers. “Ultimate limit on learning non-markovian behavior: Fisher information rate and excess information” (2023). arXiv:2310.03968.
- [65] Erich L Lehmann and George Casella. “Theory of point estimation”. Springer Science & Business Media. (2006).
- [66] Alexander Hsu and Sarah E Marzen. “Strange properties of linear reservoirs in the infinitely large limit for prediction of continuous-time signals”. *Journal of Statistical Physics* **190**, 32 (2023).
- [67] Thomas L Carroll. “Do reservoir computers work best at the edge of chaos?”. *Chaos: An Interdisciplinary Journal of Nonlinear Science* **30**, 121109 (2020).
- [68] Felix Schürmann, Karlheinz Meier, and Johannes Schemmel. “Edge of chaos computation in mixed-mode vlsi-a hard liquid”. *Advances in neural information processing systems* **17** (2004). url: [proceedings.neurips.cc/paper\\_files/paper/2004/file/dbab2adc8f9d078009ee3fa810bea142-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2004/file/dbab2adc8f9d078009ee3fa810bea142-Paper.pdf).
- [69] Sarah Marzen. “Infinitely large, randomly wired sensors

- cannot predict their input unless they are close to deterministic". *Plos one* **13**, e0202333 (2018).
- [70] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. "On the difficulty of training recurrent neural networks". In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*. Volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318. PMLR (2013).
- [71] Amy Zhang, Zachary C Lipton, Luis Pineda, Kamyar Azizzadenesheli, Anima Anandkumar, Laurent Itti, Joelle Pineau, and Tommaso Furlanello. "Learning causal state representations of partially observable environments" (2019). arXiv:1906.10437.
- [72] Seohyun Kim, Jinman Zhao, Yuchi Tian, and Satish Chandra. "Code prediction by feeding trees to transformers". In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. Pages 150–162. IEEE (2021).
- [73] Joe O'Connor and Jacob Andreas. "What context features can transformer language models use?" (2021). arXiv:2106.08367.
- [74] Mile Gu, Karoline Wiesner, Elisabeth Rieper, and Vlatko Vedral. "Quantum mechanics can reduce the complexity of classical models". *Nature communications* **3**, 762 (2012).
- [75] Yuto Takaki, Kosuke Mitarai, Makoto Negoro, Keisuke Fujii, and Masahiro Kitagawa. "Learning temporal data with a variational quantum recurrent neural network". *Physical Review A* **103**, 052414 (2021).
- [76] Paul M. Riechers and Mile Gu. "Initial-state dependence of thermodynamic dissipation for any quantum process". *Phys. Rev. E* **103**, 042145 (2021).
- [77] Felix C Binder, Jayne Thompson, and Mile Gu. "Practical unitary simulator for non-markovian complex processes". *Physical review letters* **120**, 240502 (2018).
- [78] Thomas J Elliott, Mile Gu, Andrew JP Garner, and Jayne Thompson. "Quantum adaptive agents with efficient long-term memories". *Physical Review X* **12**, 011007 (2022).
- [79] Hiroshi Yano, Yudai Suzuki, Kohei M Itoh, Rudy Raymond, and Naoki Yamamoto. "Efficient discrete feature encoding for variational quantum classifier". *IEEE Transactions on Quantum Engineering* **2**, 1–14 (2021).
- [80] Samuel P Loomis and James P Crutchfield. "Strong and weak optimizations in classical and quantum models of stochastic processes". *Journal of Statistical Physics* **176**, 1317–1342 (2019).
- [81] Samuel P Loomis and James P Crutchfield. "Thermal efficiency of quantum memory compression". *Physical review letters* **125**, 020601 (2020).
- [82] Leonardo Banchi, Jason Pereira, and Stefano Pirandola. "Generalization in quantum machine learning: A quantum information standpoint". *PRX Quantum* **2**, 040321 (2021).
- [83] C. J. Ellison, J. R. Mahoney, and J. P. Crutchfield. "Prediction, retrodiction, and the amount of information stored in the present". *Journal of Statistical Physics* **136**, 1005–1034 (2009).
- [84] Paul M Riechers, Alexander B Boyd, Gregory W Wimsatt, and James P Crutchfield. "Balancing error and dissipation in computing". *Physical Review Research* **2**, 033524 (2020).
- [85] Gregory Wimsatt, Alexander B Boyd, and James P Crutchfield. "Trajectory class fluctuation theorem" (2022). arXiv:2207.03612.
- [86] T. M. Cover and J. A. Thomas. "Elements of information theory". Wiley-Interscience. New York (2006). Second edition.
- [87] James P Crutchfield and Cina Aghamohammadi. "Not all fluctuations are created equal: Spontaneous variations in thermodynamic function" (2016). arXiv:1609.02519.