
STOCHASTIC THERMODYNAMICS OF LEARNING PARAMETRIC PROBABILISTIC MODELS

Shervin Sadat Parsi

Physics program at The Graduate Center
City University of New York
{shsparsi}@gmail.com

ABSTRACT

We have formulated a family of machine learning problems as the time evolution of Parametric Probabilistic Models (PPMs), inherently rendering a thermodynamic process. Our primary motivation is to leverage the rich toolbox of thermodynamics of information to assess the information-theoretic content of learning a probabilistic model. We first introduce two information-theoretic metrics: Memorized-information (M-info) and Learned-information (L-info), which trace the flow of information during the learning process of PPMs. Then, we demonstrate that the accumulation of L-info during the learning process is associated with entropy production, and parameters serve as a heat reservoir in this process, capturing learned information in the form of M-info.

Keywords Generative models, Machine Learning, Thermodynamics of Information, Entropy Production, Information Theory

Contents

1	Introduction	3
2	Information content of PPMs	4
3	The learning trajectory of a PPM	6
3.1	The model subsystem	6
3.2	The parameters subsystem	6
3.2.1	Lagged bipartite dynamics	7
3.3	Trajectory probabilities	8
3.4	Local Detailed Balance (LDB) for learning PPMs	9
3.5	L-info from fluctuation theorem	10
3.6	M-info and the role parameters subsystem	11
3.7	The ideal learning process	12

4	The parameters' reservoir	13
4.0.1	Naive parametric reservoir	15
4.0.2	Realistic parametric reservoir	15
5	Discussion	17

1 Introduction

Starting from nearly half a century ago, physicists began to learn that information is a physical entity [1, 2, 3]. Today, the information-theoretic perspective has significantly impacted various fields of physics, including quantum computing [4], cosmology [5], and thermodynamics [6]. Simultaneously, recent years have witnessed the remarkable success of an algorithmic approach known as machine learning, which is adept at learning information from data. This paper is propelled by a straightforward proposition: if "information is physical", then the process of learning information must inherently be a physical process.

The concepts of memory, prediction, and information exchange between subsystems have undergone extensive exploration within the realms of Thermodynamics of Information [6] and Stochastic Thermodynamics [7]. For instance, Still et al. [8] delved into the thermodynamics of prediction. And, the role of information exchange between thermodynamic subsystems has been studied by Sagawa and Ueda [9], and Esposito et al. [10]. This rich toolbox of thermodynamic of information is our main venue to study physics of machine learning process, with motivation to assess the information content of the learning process.

The type of machine learning problems we consider in this study encompasses any algorithmic approach that *evolves* a Parametric Probabilistic Model (PPM), or simply the model, towards a desirable target distribution through gradient-based loss function minimization. To establish our notation, consider a set of observations denoted by the training dataset B , drawn from an unknown target distribution p^* . The PPM, without lose of generality, can be written as the follows:

$$p_{\theta}(X = x) = e^{-\phi_{\theta}(x)} \quad (1)$$

This distribution is parameterized by a set of parameters $\theta \in \mathbb{R}^M$. The objective of learning is to find a set of parameters such that samples drawn from the model, $x \sim p(X|\theta)$, exhibit desirable statistical characteristics. In machine learning practice, one constructs the function $\phi_{\theta}(x)$ with a (deep) neural network and leave the parameter selection task to an optimizer that minimizes a loss function. Examples encompass Energy-based models [11], Large Language Models, Softmax classifiers, Variational Autoencoders (VAEs) [12], among others.

While the information-theoretic approach to this problem is prevalent in the field [13, 14, 15, 16], it has also faced criticisms [17]. Our primary motivation for framing learning in a PPM as a thermodynamic process is to facilitate the assessment of the information content inherent in the learning process. The structure of this paper is outlined as follows: Section 2 briefly discusses prior information-theoretic approaches to the learning problem with PPM, and the challenges they encounter. Subsequently, we introduce our own information-theoretic metrics. Finally, sections 3 and 4 employ the thermodynamic framework to address these information-theoretic inquiries.

2 Information content of PPMs

Locating information within the parametric model, a.k.a. the neural network, remains a fundamental question in machine learning machine [18]. This challenge is central to any information-theoretic perspective on machine learning problems. In a pioneering study, Shwartz-Ziv et al. [19] quantified the internal information within neural networks by estimating the mutual information between inputs and the activities of hidden neurons. Moreover, they employed the information bottleneck theory to interpret the decrease in this mutual information as evidence of data compression during learning. This perspective garnered significant attention in the field [16, 15], reinforcing the view of neural networks as an information channel. However, the study encountered critiques [20, 17]. A primary problem was that the hidden neurons' activity in a neural network constitutes a deterministic function of the input. Such determinism inherently possesses a trivial mutual information value, even prior to any learning. The challenge of defining a well-defined and interpretable (Shannon) information metric in deterministic neural networks has prompted the proposition that neural network information processing is geometric in nature [21] (given that inputs are mapped to a latent space of differing dimensions), rather than information-theoretical.

In a distinct research direction, Ref. [22] addresses the significance of assessing the information content of the model's parameters. In our study, we echo this view, emphasizing that parameters are the primary carriers of learned information within neural networks. Consequently, any information-theoretic measure of learned information by the model should be grounded in parameters rather than the deterministic activity of hidden neurons. However, quantifying the information within parameters poses challenges, primarily due to the elusive nature of their distribution [23]. In this section, we introduce two information-theoretic metrics crafted to assess the information content within the learning process of a PPM. This paves the way for computing these quantities within the thermodynamic framework.

To avoid introducing new notation, we also denote B as the ground truth random variable associated with the target distribution p^* from which the training dataset is sampled. Subsequently, we represent the action of the optimizer as a map between this ground truth random variable and the desired set of parameters after n optimization steps:

$$\theta_{t_n} = \Lambda_n(B) \quad (2)$$

The map Λ_n incorporates the structure of the loss function, the optimization algorithm, and any hyperparameters related to the optimizer's action. We exclude the initial parameters' value from this map's argument, under the assumption that as n increases, the final set of parameters becomes independent of its initial condition. In Information Theory terminology, this map corresponds to a *statistic* of the ground truth random variable [24]. Moreover, the outcome of this map defines a model, from which the final model-generated sample is sampled: $x_{t_n} \sim p(x|\theta_{t_n})$. Considering that the model-generated sample becomes independent of the ground truth random variable given the parameters, we can express the following Markov chain governing the learning process:

$$B \rightarrow \theta_{t_n} \rightarrow x_{t_n}. \quad (3)$$

The Data Processing Inequality (DPI) associated with this Markov chain serves as our framework to define two information-theoretic metrics that gauge the information content of the model:

$$\underbrace{I_{\Theta;B}(t_n)}_{\text{M-info}} \geq \underbrace{I_{X;B}(t_n)}_{\text{L-info}}. \quad (4)$$

We have used notations presented in table 2.1. The left-hand side of this inequality quantifies the accumulation of mutual information between the parameters and the training dataset, while the right-hand side characterizes the performance of the generative model, it gauges the accumulation of mutual information between the model's generated samples and the training dataset. We refer to the former as Memorized Information (M-info) and the latter as Learned-information (L-info). We also note that both of these quantities start at zero before the training begins. Thus, their measurements at t_n , reveal accumulation of information during the learning process.

$\Delta_{t_n} f(t) := f(t_n) - f(0)$	Change over the interval $[0, t_n]$
$\langle f(x) \rangle_{p(x)} := \int dx p(x) f(x)$	Average over $p(x)$
$s_X(t) := s[p_t(x)] := -\ln p_t(x)$	Surprisal of $p_t(x)$
$S_X(t) := S[p_t(x)] := \langle -\ln p_t(x) \rangle_{p_t(x)}$	Shannon entropy of $p_t(x)$
$I_{X;\Theta}(t) := I[X_t; \Theta_t] := S_X(t) - S_{X \Theta}(t)$	Mutual information between X and Θ at time t

Table 2.1: A list of notations used in this paper

In the context of the learning problem, the DPI as referenced in 4 suggests that what is *Memorized* is always greater than or equal to what is *Learned*. The L-info metric is task-oriented. For example, in the realm of image generation, it quantifies the statistical resemblance between the model’s outputs and the genuine images. In the case of classification task, L-info would encapsulate only the pertinent information for label prediction. In contrast, M-info can encompass information not directly pertinent to the current task. For instance, it might capture intricate pixel configurations in an image dataset, which aren’t crucial for identifying distinct patterns like human faces. The DPI 4 neatly illustrate the risk of overfitting, when a model starts to incorporate extraneous information that doesn’t align with the learning objective. The necessity of constraining the information in a model’s parameters is highlighted in Ref. [22], echoing the Minimum Description Length Principle [25]. Additionally, studies suggest that the SGD optimizer tends to favor models with minimal information in their parameters [23]. Recent work by Ref. [26] has even proposed an upper limit for minimizing parameter information to bolster generalization capabilities. These findings suggest that the learning process seeks to minimize the left-hand side of the DPI inequality while simultaneously maximizing the right-hand side, that measures the model performance. This leads us to an ideal scenario where $I_{\Theta;B}(t_n) = I_{X;B}(t_n)$, signifying that all memorized information is relevant to the learning task.

We now take one step further in our definition of M-info and L-info. First, the presence of the optimizer map, as referenced in 2, connecting the ground truth source of the training dataset to the parameters, allows us to simplify M-info as follows:

$$\text{M-info} := I_{\Theta;B}(t_n) = S(\Theta_{t_n}) \quad (5)$$

Thus, the parameters naturally emerge as the model’s *memory*, where its Shannon entropy measures the stored information during the learning process.

Second, we swap B for Θ_{t_n} , in the definition of L-info in cost of losing some information:

$$\begin{aligned} \text{L-info} &:= I_{X;B}(t_n) = I_{X;\Lambda_{t_n}(B)} + \epsilon \\ &= I_{X;\Theta}(t_n) + \epsilon \end{aligned} \quad (6)$$

where ϵ is a non-negative number that equals zero only when the map Λ_n outcome is a *sufficient statistic* for B . For the above expression, the condition of sufficient statistic can be eased as Θ to be sufficient with respect to X . This means the map Λ_n preserve all information in B that is also mutual in X . Indeed, in the problem, we are interested in this type of preservative maps that their action on training dataset preserve task-related information. Therefore, we consider $I_{X;\Theta}$ as a reasonable proxy to L-info, and we use the two interchangeably:

$$\text{L-info} := I_{X;\Theta}(t_n) \quad (7)$$

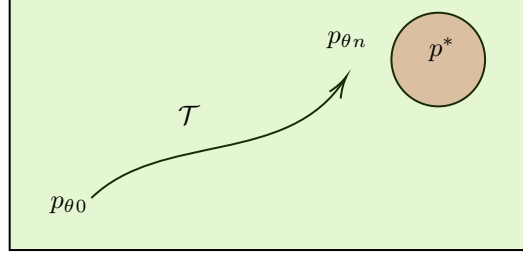


Figure 3.1: The learning trajectory \mathcal{T} depicts the thermodynamic process that take the initial model state to final state. The green area shows the space of family of distribution accessible to the PPM. The red area considers the possibility that the target distribution, p^* , is not in this family.

3 The learning trajectory of a PPM

The time evolution of the PPM is the first clue to frame the learning process as a thermodynamic process. To illustrate this, consider a discretized time interval $[0, t_n]$, which represents the time needed for n optimization steps of the parameters. During this time, the optimizer algorithm draws a sequence of i.i.d samples from the training dataset. We denote this sequence by $\mathbf{b}_n := \{b_{t_1}, b_{t_2}, \dots, b_{t_n}\}$, and refer to it as the "input trajectory". Then, the outcome of the optimization defines a sequence of parameters, call it the "parameters' trajectory": $\boldsymbol{\theta}_n := \{\theta_0, \theta_{t_1}, \theta_{t_2}, \dots, \theta_{t_n}\}$. Each realization of parameters defines a specific PPM. Consequently, the parameters' trajectory produces a sequence of PPMs:

$$\mathcal{T} := \{p(X|\theta_0), p(X|\theta_{t_1}), p(X|\theta_{t_2}), \dots, p(X|\theta_{t_n})\} \quad (8)$$

We refer to this sequence as the *learning trajectory*, depicted in figure 3.1. On the other hand, a thermodynamic process can be constructed solely from the time evolution of a distribution [27]. Therefore, we see \mathcal{T} as a thermodynamic process. The physics of this process is encoded in the transition rates governing the master equation of this time evolution. Finding the transition rate associated to learning a PPM, is our main task in this section.

3.1 The model subsystem

We refer to the subsystem that goes under the thermodynamic process \mathcal{T} as *the model subsystem*. This subsystem has X degrees of freedom, and its microscopic states' realization along the learning trajectory represent model-generated samples at each time step: $x_{t_i} \sim p(X|\theta_{t_i})$. Furthermore, we denote the stochastic trajectory of model-generated samples by $\mathbf{x}_n := \{x_{t_1}, x_{t_2}, \dots, x_{t_n}\}$. To avoid confusion with our notation, consider the probability functions $p(x_{t_i}|\theta_{t_i})$ and $p(x_{t_{i-1}}|\theta_{t_i})$, which respectively represent the probability of observing $x_{t_i} \in \mathbf{x}_n$ and $x_{t_{i-1}} \in \mathbf{x}_n$ at time $t = t_i$. Here, the time index of θ aligns with the time index of the PPM, i.e., $p_{t_i}(X|\theta_{t_i}) \equiv p(X|\theta_{t_i})$, because the PPM is fully defined upon observing the parameters. In contrast, the time index on x denotes a specific observation within \mathbf{x}_n . To simplify our notation, the absence of a time index on x denotes a generic realization of the random variable X , and we write $p(x|\theta_{t_i})$ instead of $p(X|\theta_{t_i})$.

3.2 The parameters subsystem

The parameters of the neural network at each step of optimization represent realization of the parameters subsystem, with Θ degrees of freedom, and the stochastic trajectory $\boldsymbol{\theta}_n$. The statistical state of the parameters subsystem is given with the marginal $p(\theta_{t_i})$ at time step $t = t_i$. This marginal state represents the statistic of all possible outcome of training a PPM on specific learning objective. We can think of training an ensemble of computers on the same machine learning task. This allows us to think about the time evolution of the marginal $p(\theta_{t_i})$, and the joint distribution $p(x|\theta_{t_i})p(\theta_{t_i})$ during the learning process. we refer to this view as the *ensemble view* of learning process. In practice, however, we train the PPM only once, and we do not have access to the marginal $p(\theta_{t_i})$. Thus, our model-generated

samples are conditioned on specific observations of parameters, $\theta_{t_i} \sim \Theta_{t_i}$. This defines the *conditional view* of the learning process, that is fully described by the learning trajectory of the PPM.

In machine learning practice, it is desirable for a training process to exhibit a robust outcome, regardless of who is running the code. One way to achieve this is by imposing a *low-variance condition* on the statistics of parameters across the ensemble of all learning trials. This condition asserts that the parameters' trajectory across the ensemble is confined to a small region $D_n \subset \mathcal{R}^M$. As n grows larger, this region shrinks, and becomes associated with the area surrounding the target distribution as depicted in Figure 3.1. Under this condition, we can express:

$$\langle f(\theta) \rangle_{p(\theta_n)} \approx f(\theta^*), \quad \forall \theta^* \in D_n \quad (9)$$

The above approximation becomes exact when $p(\theta_n)$ assumes the form of a delta-Dirac function, indicating a zero-variance condition in the parameters' dynamics.

The low-variance condition proves invaluable when computing the information-theoretic measurements introduced in section 2. This is because the computation of the M-info $I_{B;\Theta}$ and L-info $I_{X;\Theta}$ necessitates averaging over the parameters' distribution. However, since we typically train our model just once, we lack direct access to the parameters' distribution throughout the learning trajectory. To overcome this challenge, we introduce the Conditional L-info:

$$I_{X;\Theta}(\theta_t) = \int dx p(x|\theta_t) \ln \frac{p(x|\theta_t)}{p_t(x)} \quad (10)$$

Subsequently, under the low-variance condition of the parameters subsystem, we can measure the conditional L-info as a proxy for the L-info: $I_{X;\Theta}(t) = \langle I_{X;\Theta}(\theta) \rangle_{p_t(\theta)} \approx I_{X;\Theta}(\theta_t)$. In section 4, we will delve deeper into the evidence supporting the low-variance condition of the subsystem Θ .

We refer to the joint (X, Θ) as the *learning system*, that embodies the thermodynamic process of learning a PPM. In this section, we will demonstrate that the thermodynamic exchange between model subsystem and parameters subsystem is the primary source producing M-info and L-info during the learning process. Before delving further, we establish two interconnected assumptions about the parameters subsystem: (1) The PPM is over-parameterized; specifically, the subsystem Θ has a much higher dimension compared to the subsystem X . (2) The parameters subsystem evolves in a quasi-static fashion (slow dynamics).

The foundation for these assumptions in machine learning is clear. Training over-parameterized models represents a pivotal achievement of machine learning algorithms, and the slow dynamics (often termed as lazy dynamics) of these over-parametrized models are well-documented [28, 29]. These characteristics underscore the significant role of the parameters subsystem in the learning process, akin to that of a heat reservoir. Over-parameterization implies a higher heat capacity for this subsystem compared to the model subsystem. Additionally, the quasi-dynamics align with the behavior of an ideal heat reservoir, which doesn't contribute to entropy production [30]. The role of the parameters subsystem as a reservoir aligns with the assumption of a low-variance condition for this subsystem. This is because we expect the stochastic dynamics of a reservoir in contact with the subsystem to be low-variance across the ensemble of all trials.

In this study, we attribute the role of an ideal heat reservoir to the parameters subsystem, with inverse temperature $\beta^{-1} = 1$. In section 4, we delve deeper into the rationale behind this assumption, by examining the stochastic dynamics of parameters under a vanilla stochastic gradient descent optimizer, and highlighting potential limitations of this assumption.

3.2.1 Lagged bipartite dynamics

We want to emphasize that the dynamics of subsystem X is not a mere conjecture or an arbitrary component in this study; rather, it's an integral part of training a generative PPM. This dynamics is inherent in the optimizer action, necessitating a fresh set of model-generated samples to compute the loss function or its gradients after each parameter

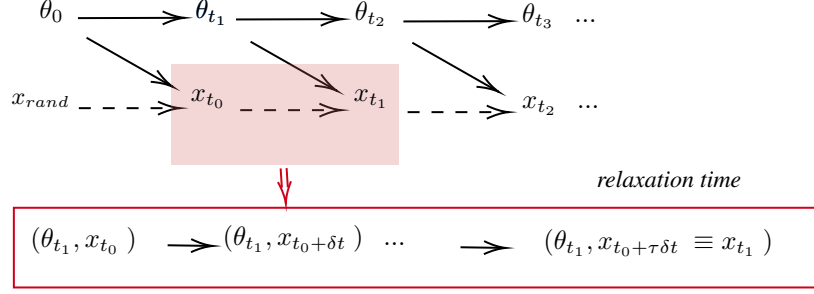


Figure 3.2: This figure shows Bayesian network for joint trajectory probability $P[\mathbf{x}_n, \boldsymbol{\theta}_n]$, based on a dual timescale bipartite dynamics.

update. For instance, in the context of EBM, a Langevin Monte Carlo (LMC) sampler can be employed to generate new samples from the model [31]. The computational cost of producing a fresh set of model-generated samples introduces a time delay in the parameter dynamics. For instance, when using an LMC sampler, the number of Monte Carlo steps dictates this lag time. Conversely, in the case of a language model, since the computation of the loss function relies on inferring subsequent tokens, the inference latency signifies the time delay. We denote the lag time parameters with τ . Here, the model subsystem evolves on the timescale δt , while the parameters subsystem evolves on the timescale $\alpha = \tau\delta t$. In the thermodynamic context, this parameter represents the *relaxation time* of the subsystem X , under fix microscopic state of subsystem Θ . Conceptually, parameter τ acts as a complexity metric, quantifying the computational resources required for each parameter optimization step. Moreover, the dynamics of the joint (X, Θ) exhibit a bipartite property. This implies that simultaneous transitions in the states of X and Θ are not allowed, given that the observation of a new set of model-generated samples occurs only after a parameter update.

The lagged bipartite dynamics described above can be represented using two time resolutions: δt and α . In the finer time resolution of δt , the Markov chain within the time interval $[t_i, t_{i+1}]$ is as follows:

$$(x_{t_i}, \theta_{t_i}) \rightarrow (x_{t_i}, \theta_{t_{i+1}}) \rightarrow (x_{t_i+\delta t}, \theta_{t_{i+1}}) \dots \rightarrow (x_{t_i+\tau\delta t}, \theta_{t_{i+1}}) \equiv (x_{t_{i+1}}, \theta_{t_{i+1}}). \quad (11)$$

We can also analyze this dynamics at a coarser time resolution of α . Within the interval $[t_0, t_n]$, the Markov Chain appears as:

$$(x_0, \theta_0) \dashrightarrow (x_{t_1}, \theta_{t_1}) \dots (x_{t_{n-1}}, \theta_{t_{n-1}}) \dashrightarrow (x_{t_n}, \theta_{t_n}). \quad (12)$$

In the above Markov chain, the dashed arrows remained us the *ignorance* of intermediate steps in the high resolution picture 11. Figure 3.2, illustrates the lagged bipartite dynamics of the learning system. An important observation is that the learning trajectory \mathcal{T} , as defined in 8, is written in the low resolution picture. Therefore, studying the learning trajectory means studying the dynamics of the system (X, Θ) in the low resolution picture.

3.3 Trajectory probabilities

To set the stage for the application of the Fluctuation Theorem (FT) to learning a PPM, we define the trajectory probability of the joint $(\mathbf{x}_n, \boldsymbol{\theta}_n)$ as the probability of observing a series of model-generated samples and parameters during the learning process:

$$P[\mathbf{x}_n, \boldsymbol{\theta}_n] := p(x_0, x_{t_1}, \dots, x_{t_n}, \theta_0, \theta_{t_1}, \dots, \theta_{t_n}) \quad (13)$$

Additionally, we can consider the time reversal of the samples' trajectory and parameters' trajectory, respectively, as $\tilde{\mathbf{x}}_n := \{x_{t_n}, x_{t_{n-1}}, \dots, x_{t_1}\}$ and $\tilde{\boldsymbol{\theta}}_n := \{\theta_{t_n}, \theta_{t_{n-1}}, \dots, \theta_{t_1}\}$. Then, the probability of observing the backward trajectory is denoted by $P[\tilde{\mathbf{x}}_n, \tilde{\boldsymbol{\theta}}_n]$.

Here, $P[\mathbf{x}_n, \boldsymbol{\theta}_n]$ represents the trajectory probability of the learning system in the ensemble view. In practice, however, we typically train our model only once, and we often lack access to the parameters' distribution. Therefore, our model is conditioned on the observation of a specific parameters' trajectory $\boldsymbol{\theta}_n$. This defines the trajectory probability in the conditional view:

$$P[\mathbf{x}_n | \boldsymbol{\theta}_n] := \frac{P[\mathbf{x}_n, \boldsymbol{\theta}_n]}{P[\boldsymbol{\theta}_n]} \quad (14)$$

where,

$$P[\boldsymbol{\theta}_n] = p(\theta_0, \theta_{t_1}, \dots, \theta_{t_n}). \quad (15)$$

Similarly, the backward conditional trajectory probability is the probability of observing the time-reversal samples' trajectory, conditioned on observation of the time-reversal parameters' trajectory: $P[\tilde{\mathbf{x}}_n | \tilde{\boldsymbol{\theta}}_n] = \frac{P[\tilde{\mathbf{x}}_n, \tilde{\boldsymbol{\theta}}_n]}{P[\tilde{\boldsymbol{\theta}}_n]}$.

We now use the Markov property in the Markov chains 11 and 12 respectively, to decompose the conditional trajectory probability and the marginal trajectory probability as follows:

$$\begin{aligned} P[\mathbf{x}_n | \boldsymbol{\theta}_n] &= p(x_{t_n} | x_{t_{n-1}}, \theta_{t_n}) \dots p(x_{t_1} | x_0, \theta_{t_1}) p(x_0 | \theta_0), \\ P[\boldsymbol{\theta}_n] &= p(\theta_{t_n} | \theta_{t_{n-1}}) \dots p(\theta_{t_1} | \theta_{t_0}) p(\theta_0), \end{aligned} \quad (16)$$

where the expressions such as $p(x_{t_n} | x_{t_{n-1}}, \theta_{t_n})$ and $p(\theta_{t_n} | \theta_{t_{n-1}})$ represent the transition probabilities that determine the probability of moving from one microscopic state to another. Additionally, we define two probability trajectories, conditioned on the initial conditions, which will be used later in the formulation of FT:

$$\begin{aligned} P[(\mathbf{x}_n | \boldsymbol{\theta}_n) | (x_0 | \theta_0)] &:= P[(\mathbf{x}_n | \boldsymbol{\theta}_n)] / p(x_0 | \theta_0), \\ P[\boldsymbol{\theta}_n | \theta_0] &:= P[\boldsymbol{\theta}_n | \theta_0] / p(\theta_0). \end{aligned} \quad (17)$$

3.4 Local Detailed Balance (LDB) for learning PPMs

The transition probabilities, represented in 16, capture the physics of the learning problem. Considering a Markov property (i.e., memoryless process) for time evolution of the model subsystem, the transition rate for this subsystem get reduced to PPM:

$$p(x_{t_i} | x_{t_{i-1}}, \theta_{t_i}) = p(x_{t_i} | \theta_{t_i}). \quad (18)$$

The above expression suggests that the transition rate between two microscopic states $x_{t_{i-1}}$ and x_{t_i} under the fixed θ_{t_i} , to be equivalent to probability of observing x_{t_i} by the PPM itself at $t = t_i$. To reiterate, this is the Markov property that suggests the element inside \mathbf{x}_n , are independently and freshly drawn from the PPM specified with given parameters along the learning trajectory \mathcal{T} . This is especially true where $\tau \gg 1$. We can generalize this observation for the backward transition probability $p(x_{t_{i-1}} | x_{t_i}, \theta_{t_i})$, that represent probability of the backward transition $(x_{t_i}, \theta_{t_i}) \rightarrow (x_{t_{i-1}}, \theta_{t_i})$ under fixed θ_{t_i} , as follows:

$$p(x_{t_{i-1}} | x_{t_i}, \theta_{t_i}) = p(x_{t_{i-1}} | \theta_{t_i}). \quad (19)$$

The above expression tells us that the probability of backward transition is equivalent with the probability of observing the sample generated at $t = t_{i-1}$ in \mathbf{x}_n with the PPM at time $t = t_i$.

Finally, we write the log ratio of forward and backward transitions:

$$\ln \frac{p(x_{t_i} | x_{t_{i-1}}, \theta_{t_i})}{p(x_{t_{i-1}} | x_{t_i}, \theta_{t_i})} = \ln \frac{p(x_{t_i} | \theta_{t_i})}{p(x_{t_{i-1}} | \theta_{t_i})} = -(\phi_{\theta_{t_i}}(x_{t_i}) - \phi_{\theta_{t_i}}(x_{t_{i-1}})), \quad (20)$$

where the second equality is due to Eq. 1. The above expression resembles the celebrated Local Detailed Balance (LDB) [32] that relates the log ratio of forward and backward transition probabilities to the difference in potential energy of initial and final state in the transition. The heat reservoir that supports the legitimacy of the above LDB expression for

learning PPM is the parameters subsystem, whose temperature has been set to one, as we will discuss it in more details in section 4. We emphasize that the above LBD has emerged naturally under assumption of the Markov property and a relaxation time for learning a generic generative PPM. It is also important to note that the above LBD is only valid in the low resolution picture.

The LBD relation, presented in Eq. 20, has a profound consequence. It allows us to write the forward conditional probability trajectory, $P[\mathbf{x}_n|\boldsymbol{\theta}_n]$, and the backward conditional probability trajectory, $P[\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n]$, solely based on the series of PPMs in the learning trajectory \mathcal{T} :

$$\begin{aligned} P[\mathbf{x}_n|\boldsymbol{\theta}_n] &= p(x_{t_n}|\theta_{t_n}) \dots p(x_{t_1}|\theta_{t_1}) p(x_0|\theta_0) \\ P[\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n] &= p(x_0|\theta_{t_1}) \dots p(x_{t_{n-1}}|\theta_{t_n}) p(x_{t_n}|\theta_{t_n}) \end{aligned} \quad (21)$$

This is significant because it renders the application of the FT framework to the learning PPMs practical, as we have access to elements of the learning trajectory.

3.5 L-info from fluctuation theorem

The version of the fluctuation theorem we are about to apply to the learning PPMs is known as the Detailed Fluctuation Theorem (DFT)[33]. We also note that the machinery we are about to present for measuring information flow in PPMs has been developed to study information exchange between thermodynamic subsystems [9]. The novelty here lies merely in the application of this machinery to the learning process of a PPM. In this section, we extensively use notations presented in table 2.1. Also, note that the temperature of the parametric reservoir is set to one. Applying DFT in the conditional view, i.e., the conditional forward and backward trajectories defined in Eq. 21, results in:

$$\begin{aligned} \sigma_{\mathbf{x}_n|\boldsymbol{\theta}_n} &= \ln \frac{P[\mathbf{x}_n|\boldsymbol{\theta}_n]}{P[\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n]} \\ &= \ln \frac{P[(\mathbf{x}_n|\boldsymbol{\theta}_n)|(x_0|\theta_0)]}{P[(\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n)|(x_{t_n}|\theta_{t_n})]} + \ln \frac{p(x_0|\theta_0)}{p(x_{t_n}|\theta_{t_n})} \\ &= -q_{\mathbf{x}_n}(\boldsymbol{\theta}_n) + s[p(x_{t_n}|\theta_{t_n})] - s[p(x_{t_0}|\theta_{t_0})] \end{aligned} \quad (22)$$

The first line is due to DFT, which defines the stochastic EP to be the logarithm of the ratio of the forward and backward trajectory probabilities. The second line is due to the decomposition presented in Eq. 17. Finally, the third line is the consequence of LDB relation ??, and the definition of the stochastic heat flow $q_{\mathbf{x}_n}(\boldsymbol{\theta}_n)$, as the change in the energy of the subsystem X due to alterations in its microscopic state configuration:

$$q_{\mathbf{x}_n}(\boldsymbol{\theta}_n) := -\ln \frac{P[(\mathbf{x}_n|\boldsymbol{\theta}_n)|(x_0|\theta_0)]}{P[(\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n)|(x_{t_n}|\theta_{t_n})]} = \sum_{i=1}^n \phi_{\theta_{t_i}}(x_i) - \phi_{\theta_{t_i}}(x_{i-1}). \quad (23)$$

Note that our sing convention defines $q_{\mathbf{x}_n} > 0$ as the heat observed by the subsystem X .

The second law arises from averaging Eq.22 over the forward trajectory distribution $P_F[\mathbf{x}_n|\boldsymbol{\theta}_n]$, and recalling the non-negativity property of the KL-divergence to establish non-negativity of averaged EP: $\Sigma_{X|\Theta}(\boldsymbol{\theta}_n) := \ln \frac{P_F[\mathbf{x}_n|\boldsymbol{\theta}_n]}{P_B[\tilde{\mathbf{x}}_n|\tilde{\boldsymbol{\theta}}_n]} > P_F[\mathbf{x}_n|\boldsymbol{\theta}_n] \geq 0$. We note that the averaged EP is still conditioned on the stochastic trajectory of parameters, thus we refer to this as the conditional EP. This is indeed the consequence of working in the conditional view.

Motivated to compute L-info, in the next step, we rearrange Eq. 22 as follows:

$$\mathcal{I}[x_{t_n} : \theta_{t_n}] - \mathcal{I}[x_0 : \theta_0] = -q_{\mathbf{x}_n}(\boldsymbol{\theta}_n) + s[p(x_{t_n})] - s[p(x_{t_0})] - \sigma_{\mathbf{x}_n|\boldsymbol{\theta}_n}, \quad (24)$$

where $\mathcal{I}[x_{t_n} : \theta_{t_n}] := s[p(x_{t_n})] - s[p(x_{t_n}|\theta_{t_n})]$ is the mutual content (or stochastic mutual information) at $t = t_n$. We now arrive at the conditional L-info 10 by averaging Eq. 24 over $P_F[\mathbf{x}_n|\boldsymbol{\theta}_n]$:

$$\begin{aligned} I_{X;\Theta}(\theta_{t_n}) - I_{X;\Theta}(\theta_0) &= -Q_X(\boldsymbol{\theta}_n) + (S_X(\theta_{t_n}) - S_X(\theta_0)) - \Sigma_{X|\Theta}(\boldsymbol{\theta}_n) \\ &= \Sigma_X(\boldsymbol{\theta}_n) - \Sigma_{X|\Theta}(\boldsymbol{\theta}_n) \end{aligned} \quad (25)$$

that defines the Partially Averaged (PA) quantities,

$$\begin{aligned} Q_X(\boldsymbol{\theta}_n) &:= \sum_{i=1}^n \langle \phi_{\theta_{t_i}}(x) \rangle_{p(x|\theta_{t_i})} - \langle \phi_{\theta_{t_i}}(x) \rangle_{p(x|\theta_{t_{i-1}})} && \text{(PA Heat flow)} \\ S_{X|\Theta}(\theta_{t_i}) &:= \langle -\log(p(x|\theta_{t_i})) \rangle_{p(x|\theta_{t_i})} && \text{(PA Conditional Entropy)} \\ S_X(\theta_{t_i}) &:= \langle -\log(p(x)) \rangle_{p(x|\theta_{t_i})} && \text{(PA Marginal Entropy)} \\ \Sigma_X(\boldsymbol{\theta}_n) &:= (S_X(\theta_{t_n}) - S_X(\theta_0)) - Q_X(\boldsymbol{\theta}_n) && \text{(PA Marginal EP)} \end{aligned}$$

We note that all PA quantities are conditioned on the parameters' trajectory, i.e., the choice of $\boldsymbol{\theta}_n$ from the ensemble. This is a direct consequence of working in the conditional view. However, this also signifies that all thermodynamic quantities mentioned above are computable in the practice of machine learning, as they only require access to the time evolution of one PPM. Fortunately, thanks to the low-variance condition 9, we can use the conditional L-info as proxy to the L-info, given that: $I_{X;\Theta}(\theta_{t_n}) \approx \langle I_{X;\Theta}(\theta) \rangle_{p_t(\theta)}$, $\forall \theta_{t_n} \sim p_t(\theta)$.

Eq. 25, equates the (conditional) L-info to the difference between the Marginal EP, and the Conditional EP. We refer to this difference as the *ignorance* EP:

$$\Sigma_{ign}(\boldsymbol{\theta}_n) := \Sigma_X(\boldsymbol{\theta}_n) - \Sigma_{X|\Theta}(\boldsymbol{\theta}_n) \quad (26)$$

It is important to note that both the Marginal EP and the Conditional EP measure the EP of the same process, which is the time evolution of the subsystem X . However, the conditional EP measures this quantity with a lower time resolution of α , that is conditioned on a specific parameters' trajectory. On the other hand, the marginal EP measures this quantity with a higher time resolution of δt , including the relaxation time of the subsystem X between each parameters' update. Therefore, the term "ignorance" refers to ignorance of the full dynamic of X , and the origin of L-info is the EP between each consecutive parameters' update, i.e., the EP of generating fresh samples represented with Markov chain 11.

3.6 M-info and the role parameters subsystem

We can also apply the DFT to subsystem Θ :

$$\begin{aligned} \sigma_{\boldsymbol{\theta}_n} &= \log \frac{P[\boldsymbol{\theta}_n]}{P[\tilde{\boldsymbol{\theta}}_n]} \\ &= -q_{\boldsymbol{\theta}_n} + s[p(\theta_{t_n})] - s[p(\theta_{t_0})]. \end{aligned} \quad (27)$$

In the above expression, the second line is due to the decomposition in Eq. 17, and definition of the stochastic heat flow for parameter subsystem: $q_{\boldsymbol{\theta}_n} := \log P[\boldsymbol{\theta}_n|\theta_0]/P[\tilde{\boldsymbol{\theta}}_n|\theta_{t_n}]$.

Under the assumption that the subsystem Θ evolve quasi-statically, the EP of this subsystem is zero, as expected for an ideal reservoir. This result in $q_{\boldsymbol{\theta}_n} = \Delta_{t_n} s[p(\theta_t)]$. Furthermore, in the closed system of (X, Θ) , the heat flow of the subsystem X must be provided with an inverse flow of the subsystem Θ , i.e., $q_{\mathbf{x}_n}(\boldsymbol{\theta}_n) = -q_{\boldsymbol{\theta}_n}$. Thus, we arrive at the stochastic version of Clausius' relation for the heat reservoir:

$$\Delta_{t_n} s[p(\theta_t)] = -q_{\mathbf{x}_n}(\boldsymbol{\theta}_n) \quad (28)$$

This relation states that the heat dissipation in subsystem X ($q_{x_n}(\theta_n) < 0$) is compensated with an increase of information in subsystem Θ . We recall the definition of M-info ?? as the entropy subsystem Θ . Since heat dissipation is a source of L-info accumulation (see Eq. 25), the above Clausius' relation states that this information is stored in the parameters by increasing the entropy of this subsystem, a.k.a. the M-info, confirming the role of parameters as the memory space of the PPM.

We can also take the ensemble average of Eq. 28 (i.e., averaging over $P[x_n, \theta_n]$):

$$\Delta_{t_n} S[\Theta_t] = -Q_X(t_n), \quad (29)$$

where $Q_X(t_n) := \sum_{x_n, \theta_n} P[x_n, \theta_n] q_{x_n}(\theta_n) = \sum_{\theta_n} P[\theta_n] Q_X(\theta_n)$ is the fully averaged dissipated heat from the subsystem X . However, under the low-variance condition of learning ??, we expect $Q_X(\theta_n)$ to be independent of choice of parameters' trajectory from the ensemble of computers. Thus, we can write $Q_X(t_n) \approx Q_X(\theta_n)$.

3.7 The ideal learning process

The learning objective necessitates an increase in L-info to enhance the model's performance while simultaneously reducing M-info to minimize generalization error and prevent overfitting. As previously mentioned in Section 2, the ideal scenario is achieved when all the stored information in the parameters (M-info) matches the task-relevant information learned by the model (L-info). Now that we have studied the machinery for computing these two information-theoretic quantities through the computation of entropy production, we can formally examine this optimal learning condition.

Maximizing L-info, as described in Eq. 25, is equivalent to maximizing the marginal EP while minimizing the conditional EP. Given that the conditional EP is always non-negative, the "ideal" scenario would involve achieving a conditional EP of zero, i.e., $\Sigma_{X|\Theta}(t_n) = 0$. This condition can be realized through a quasi-static time evolution of the PPM occurring on the lower-resolution timescale α , presented in the Markov chain 12. In the context of generative models, this condition is akin to achieving perfect sampling. Under these circumstances, all EP of the subsystem X transforms into L-info, resulting in $\Delta_{t_n} I_{X;\Theta}(\theta_t) = \Sigma_X(t_n)$.

Thermodynamically, the condition of quasi-static time evolution of the PPM (and consequently zero conditional EP) can be realized by having a large relaxation parameter $\tau \gg 1$, which allows the model to reach equilibrium after each optimization step. However, a high relaxation parameter comes at the cost of requiring more computational resources and longer computation times. This introduces a fundamental trade-off between the time required to run a learning process and its efficiency - a concept central to thermodynamics and reminiscent of the Carnot cycle, representing an ideal engine that requires an infinite operation time.

4 The parameters' reservoir

In the formulation of the previous section, we make this assumption that the subsystem Θ behaves as an ideal reservoir. In this section, we get deeper on the premises of this assumption by studying the dynamic of the parameters subsystem. To facilitate our formulation, we adapt negative log-likelihood as a fairly general form for the loss function:

$$\ell(b_t, \theta) := -\frac{1}{|b_t|} \sum_{x \in b_t} \log(p_\theta(x)) = \phi_\theta(b_t), \quad (30)$$

Here, the loss function is computed according to the empirical average of a random mini-batch $b_t \in B$ drawn from the training dataset at time step t . The last equality is due to the PPM defined in Eq. 1, and $\phi_\theta(b_t) := \frac{1}{|b_t|} \sum_{x \in b_t} \phi_\theta(x)$, where notation $|\cdot|$ shows the size of a set. We also use a vanilla Stochastic Gradient Descent (SGD) optimizer, with the learning rate r , to take gradient steps iteratively for n steps, in the direction of loss function minimization:

$$\theta_{t+1} = \theta_t - r \nabla_\theta \phi_\theta(b_t)|_{\theta=\theta_t} \quad (31)$$

To render the dynamic of parameters in the form of a conventional overdamped Langevin dynamic, we introduce the following conservative potential, defined by the entire training dataset B :

$$U_B(\theta) := \frac{1}{|B|} \sum_{x \in B} \phi_\theta(x). \quad (32)$$

The negative gradient of this potential gives rise to a deterministic vector force. Additionally, we define the fluctuation term, that represents the source of random forces due to selection of a mini-batch at time step t_n :

$$\eta(t_n) := -\nabla_\theta \phi_{\theta_t}(b_{t_n}) + \nabla_\theta U_B(\theta_{t_n}).$$

We now reformulate the SGD optimizer 31, in the guise of overdamped Langevin dynamics, dividing it by the parameters' update timescale α to convert the learning protocol into a dynamic over time:

$$\frac{\theta_{t_{n+1}} - \theta_{t_n}}{\alpha} = -\mu \nabla_\theta U_B(\theta_{t_n}) + \mu \eta(t_n), \quad (33)$$

where $\mu := r/\alpha$ is known as the mobility constant, in the context of Brownian motion.

We note that Eq. 33 is merely a rearrangement of the standard SGD. For us to interpret it as a Langevin equation, the term $\eta(t_n)$ must represent a stationary stochastic process to serve as the *noise* term in the Langevin equation. To demonstrate this property of $\eta(t_n)$, we must examine the characteristic of its Time Correlation Function (TCF)[34]: $C_{i,j}(t, t-t') := \delta_{i,j} \langle \eta_i(t) \eta_j(t') \rangle$, where indices i, j represent different components of the vector θ , and $\delta_{i,j}$ is the Kronecker delta.

If the fluctuation term, η , satisfies the condition of the white noise (uncorrelated stationary random process), and assuming that Eq. 33 describes a motion akin to Brownian motion, we can apply the fluctuation-dissipation theorem to write:

$$\langle \eta_i(t) \eta_j(t') \rangle = \frac{2k_B T}{\mu} \delta(t - t') \delta_{i,j} \quad (34)$$

Here, $\delta(t - t')$ is a delta Dirac, and the constant T symbolizes the temperature. The constant k_B stands for the Boltzmann constant. To render our framework unitless, we treat the product of the Boltzmann factor and temperature as dimensionless. Moreover, regardless of the noise width we set $T = 1$, and henceforth it will not appear in our formulation. This is possible by adjusting the Boltzmann factor according to the noise width, i.e., $k_B = \mu \langle \eta_i(t) \eta_i(t) \rangle / 2$.

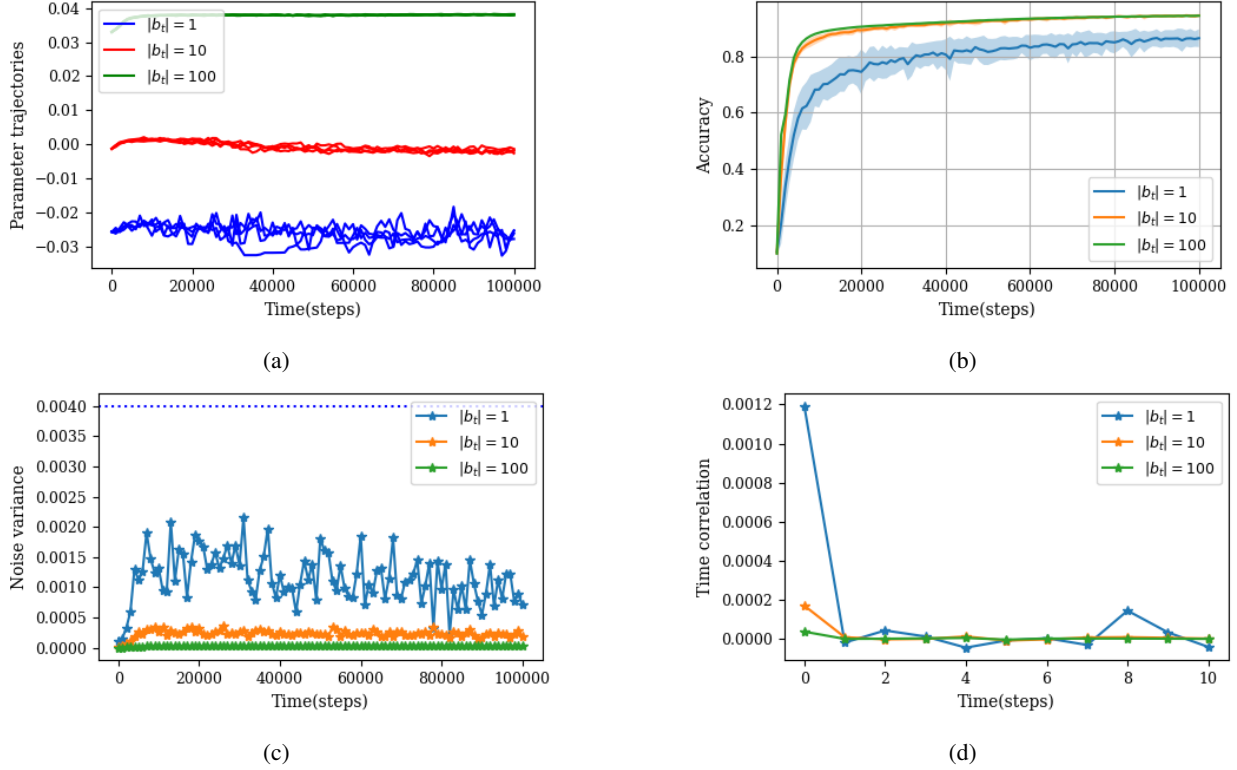


Figure 4.1: This experiment contrasts the parameter dynamics with three different mini-batch sizes: $|b_t| = 1$, $|b_t| = 10$ and $|b_t| = 100$. The model under consideration is a four-layer feedforward neural network with a uniform width of 200 neurons. It was trained on the MNIST classification task using a vanilla SGD optimizer. The experiment was replicated over 50 trials to generate an ensemble of parameters. a) One random parameter from the model’s last layer is chosen for each batch size scenario, and four of its dynamic realizations are depicted. b) Illustrates both the average accuracy (solid line) and the variance of accuracy within the ensemble (shaded area), emphasizing the low-variance condition, which asserts that macroscopic quantities such as accuracy have low variance statistics across the ensemble. c) Displays the noise variance averaged over all parameters, i.e., $\frac{1}{\dim(\theta)} \sum_{i=0}^{\dim(\theta)} C_{i,i}(t, 0)$, for each mini-batch size scenario, underscoring the stationary nature of η . This part also highlights the role of mini-batch size in determining the noise width, i.e., the temperature of the environment. The horizontal dashed line indicates the maximum absolute value observed from $\nabla_{\theta} U_B(\theta_{t_n})$, serving as a reference point for the magnitude of the noise. d) Exhibits the autocorrelation of the term η averaged over all parameters. For instance, computing this quantity at step 1000 reads: $\frac{1}{\dim(\theta)} \sum_{i=0}^{\dim(\theta)} C_{i,i}(t = 1000, t' - t)$. The rapid decline in autocorrelation with time lag indicating the white noise characteristic of η .

We still need to investigate if the fluctuation term indeed describes an uncorrelated stationary random process, as presented in Eq. 34. To this end, we conducted an experiment by training an ensemble of 50 models for the classification of the MNIST dataset. To induce different level of stochastic behavior, i.e., different "temperatures", we consider three different mini-batch sizes. A smaller mini-batch size leads to a bigger deviation in the fluctuation term, consequently amplifying the influence of random forces. Results are presented in Fig. 4.1. The plot 4.1c represents the TCF function at no time lag $t = t'$, i.e., variance of $\eta(t)$, as a function of time. The constant value of variance suggests the stationary property of $\eta(t)$. Moreover, Fig. 4.1d illustrates the autocorrelation of $\eta(t)$ at different time lags, indicating white noise characteristic for this term.

However, it would be naive to draw a generic conclusion regarding the nature of the fluctuation term as an uncorrelated stationary random process solely based on a simple experiment. Indeed, research has demonstrated that the noise term can be influenced by the Hessian matrix of the loss function [35]. This observation aligns with our definition of the fluctuation term presented in Eq. 33, where η is defined in relation to the gradient of the loss itself. Consequently, as

the optimizer explores the landscape of the loss function, the characteristics of the fluctuation term η can vary. We can grasp this concept in the context of Brownian motion by envisioning a Brownian particle transitioning from one medium to another, each with distinct characteristics. This implies that there could be intervals during training where η stays independent of the loss function and exhibits a stationary behavior.

Moreover, we overlooked the fact that $\eta(t)$ is also a function of θ itself. This could potentially jeopardize its stationary property. To address this issue, we refer to the slow dynamic (lazy dynamic) [28, 29] of over-parameterized models under SGD optimization. This slow dynamic allows us to write the Taylor expansion¹ of the loss function around a microscopic state θ^* , sampled from its current state $p_t(\theta)$:

$$\phi_{\theta_t}(b_t) = \phi_{\theta^*}(b_t) + (\theta_t - \theta^*) \nabla_{\theta} \phi_{\theta^*}(b_t) \quad (35)$$

As a result, the gradient of the loss $\nabla_{\theta} \phi_{\theta_t}(b_t) = \nabla_{\theta} \phi_{\theta^*}(b_t)$, signifying an independent behavior from the specific value of the parameters θ_t at a given time t . We can extend this concept to the deterministic force $-\nabla_{\theta} U_B(\theta_t) = -\nabla_{\theta} U_B(\theta^*) = F(\theta^*)$, which indicates a conservative force in lazy dynamics regime, denoted as $F(\theta^*)$. The key point here is that the value of this force is not dependent on the microscopic state of θ_t , but rather on any typical sample, θ^* , from Θ_t . In Appendix 5, we illustrate how the condition of lazy dynamics leads to a thermodynamically reversible dynamic of the subsystem Θ .

4.0.1 Naive parametric reservoir

The stationary state of subsystem Θ , under the dynamic of Eq. 33, satisfying the fluctuation-dissipation relation in Eq. 34, corresponds to the thermal equilibrium state (the canonical state):

$$p^{eq} = e^{-U_B(\theta) + F_{\Theta}} \quad (36)$$

where $F_{\Theta} := -\log(\int d\theta e^{-U_B(\theta)})$ is the free energy of the subsystem θ . Recall that, the temperature has been set to one. This state, also, satisfies the detailed balance condition, that define the log ratio between forward and backward transition probability as follows:

$$\log \frac{p(\theta_{t_i} | \theta_{t_{i-1}})}{p(\theta_{t_{i-1}} | \theta_{t_i})} = -\left(U_B(\theta_{t_i}) - U_B(\theta_{t_{i-1}})\right) \quad (37)$$

The standard plot of the loss function versus optimization steps in machine learning practice can help us to visualize the dynamics of the subsystem Θ . A rapid decline in the loss function signals a swift relaxation of the subsystem Θ to its equilibrium state. It is important to note that this *self-equilibrating property* is determined by the training dataset B through the definition of the potential function $U_B(\theta)$. These swift and self-equilibrating properties mirror the characteristics of a heat reservoir in thermodynamics [30]. Hence, we refer to the subsystem Θ as the *parametric reservoir*. After a swift decline, a gradual reduction of the loss function, can be sign of a quasi-statistic process, when subsystem Θ evolve from one equilibrium state to another. This can be due to the lazy dynamic condition, as discussed in Appendix 5. Additionally, the requirement of a high heat capacity for the reservoir, represented as $\dim(\Theta) \gg \dim(X)$, offers a thermodynamic justification for the use of over-parameterized models in machine learning.

4.0.2 Realistic parametric reservoir

We refer to the assumption of the parametric reservoir with an equilibrium state expressed in Eq. 36 as the "naive assumption" due to several issues that were previously sidestepped. The first issue stems from the assumption that all components of the parameter vector θ are subject to the same temperature, i.e., $\langle \eta_i(t) \eta_i(t) \rangle = \frac{2k_B T}{\mu}$ for all index i . In practice, we might find different values of noise width, particularly with respect to different layers of a deep neural

¹Similar to what has been done in Neural tangent kernel theory [36], but with a different purpose.

network. Furthermore, the weights or biases within a specific layer might experience different amounts of fluctuation. This scenario is entirely acceptable, if we consider each group of parameters as a subsystem that contributes to the formation of the parametric reservoir Θ . Consequently, each subsystem possesses different environmental temperatures and distinct stationary states. This observation may explain, in thermodynamic terms, why a deep neural network can offer a richer model. As it encompasses multiple heat reservoirs at varying temperatures, it presents a perfect paradigm for the emergence of non-equilibrium thermodynamic properties.

Second, the fluctuation term η may exhibit an autocorrelation property that characterizes colored noise, as presented in Ref [37]. While this introduces additional richness to the problem, potentially displaying non-Markovian properties, it does not impede us from deriving the equilibrium state of the subsystem Θ , as demonstrated in [38].

We also overlooked the irregular behavior of the loss function, such as spikes or step-like patterns. These irregularities are considered abnormal as we typically expect the loss function to exhibit a monotonous decline, but in practice, such behaviors are quite common. These anomalies may be associated with a more intricate process, such as a phase transition or a shock, experienced by the reservoir. Nevertheless, we can still uphold the basic parametric reservoir assumption during the time intervals between these irregular behaviors.

The mentioned issues are attributed to a richer and more complex dynamic of subsystem Θ , and do not fundamentally contradict the potential role of subsystem Θ as a reservoir. Examples of these richer dynamics can be found in a recent study [39], that shows the limitation of Langevin formulation of SGD, and Ref. [40] that investigates exotic non-equilibrium characteristic of parameters' dynamics under SGD optimization.

Before closing this section, it is worth mentioning that the experimental results presented in Figure 4.1 support the assumption of a low-variance condition for the stochastic dynamics of the subsystem Θ . For instance, panel (a) shows that even in the high noise regime ($|b_t| = 1$), the dynamics of parameters remain confined to a small region across the ensemble. Furthermore, panel (b) demonstrates the low-variance characteristics of the model's performance accuracy. Finally, the large magnitude of deterministic force (dashed line in panel (c)) to random force, is an evidence of low-variance dynamics.

5 Discussion

In this study, we delved into the thermodynamic aspects of machine learning algorithms. Our approach involved first formulating the learning problem as the time evolution of a PPM. Consequently, the learning process naturally emerged as a thermodynamic process. This process is driven by the work of the optimizer, which can be considered as a thermodynamic work since parameters' optimization constantly change the system's energy landscape through $\phi_\theta(x)$. The optimizer action is fueled by the input trajectory, a series of samples drawn from the ground truth system. The work and heat exchange of the subsystem X can be computed practically along the learning trajectory \mathcal{T} , as outlined below:

$$W_X(\theta_n) = \sum_{t=1}^n \langle \phi_{\theta_t}(x) \rangle_{p(x|\theta_{t-1})} - \langle \phi_{\theta_{t-1}}(x) \rangle_{p(x|\theta_{t-1})} \quad (38)$$

$$Q_X(\theta_n) = \sum_{t=1}^n \langle \phi_{\theta_t}(x) \rangle_{p(x|\theta_t)} - \langle \phi_{\theta_t}(x) \rangle_{p(x|\theta_{t-1})} \quad (39)$$

We use the term "practically" because, when running a machine learning algorithm, we have access to the function $\phi_\theta(x)$ at each training instance. We also note that these quantities are conditioned on specific parameters' trajectory, as the result of working in the conditional view. Finally, the learning process can be summarized as follows: *The model learns by dissipating heat, and the dissipated heat increases the entropy of parameters, which act as the heat reservoir (a memory space) for learned information.* This means the learning process must be irreversible, this is the only way to increase the mutual information between the two subsystem X and Θ [10].

It is important to note that despite the wealth of literature highlighting the significance of information content in parameters [22, 23, 26], calculating these quantities remains difficult due to the lack of access to the parameter distribution. In contrast, the thermodynamic approach compute the information-theoretic metrics indirectly, as heat and work of the process. Moreover, the mysterious success of over-parametrized models can be explained within the thermodynamic framework, where over-parameterization plays a crucial role in allowing the parameter subsystem to function as a heat reservoir.

At the same time, we are aware of the strong assumptions made during this study. Addressing each of these assumptions or providing justifications for them represents a direction for future research. For instance, we assumed slow dynamics of parameters for the over-parameterized regime under the SGD optimizer. This formed the basis for treating the parameters' degrees of freedom as an ideal heat reservoir, evolving in a thermodynamically reversible manner. Breaking this assumption due to rapid changes in parameter values would violate this assumption. Exploring these more complex scenarios would only serve to enrich the thermodynamics of this problem.

We have also sidestepped the role of changes in the marginal entropy of the model's subsystem, $\Delta_{t_n} S_X(t)$. This term can be estimated by computing the entropy of the empirical distribution of generated samples. For a model initialized randomly, this term is always negative, as the initial model produces uncorrelated patterns with maximum entropy. Then, the negative value of this term must converge when the entropy of the generated patterns reaches the entropy of the training dataset. However, if we look at Eq. 25 as an optimization objective to maximize L-info, then an increase in the model's generated samples, $S_X(t)$, is favorable. This might act as a regularization term to improve the generalization power of the model by forcing it to avoid easy replication of the dataset.

Appendix A: Reversibility under lazy dynamic regime

In this appendix, we establish the thermodynamic reversibility of parameter evaluation as a consequence of training an over-parameterized model with lazy dynamics. The forward action of the optimizer can be summarized as follows: $\mathbf{b}_n \rightarrow \boldsymbol{\theta}_n$, where the optimizer samples an i.i.d trajectory of inputs from the training dataset $\mathbf{b}_n = \{b_{t_1}, b_{t_2}, \dots, b_{t_n}\}$ to generate a trajectory of updated parameters $\mathbf{b}_n = \{\theta_0, \theta_{t_1}, \dots, \theta_{t_n}\}$.

The backward (time-reversal) action of the optimizer is defined as: $\tilde{\mathbf{b}}_n \rightarrow \boldsymbol{\theta}_n^\dagger$, where $\tilde{\mathbf{b}}_n = \{b_{t_n}, b_{t_{n-1}}, \dots, b_{t_1}\}$ represents the time-reversal of the input trajectory, and gradient descent is reversed to gradient ascent, resulting in a new parameters' trajectory $\boldsymbol{\theta}_n^\dagger$.

In general, the backward action of SGD does not yield the time-reversal of forward parameters' trajectory:

$$\boldsymbol{\theta}_n^\dagger \neq \tilde{\boldsymbol{\theta}}_n = \{\theta_{t_n}, \theta_{t_{n-1}}, \dots, \theta_{t_0}\}$$

To illustrate this, let's examine a single forward and backward action of the optimizer:

$$\begin{aligned} \theta_{t+1} &= \theta_t - r \nabla_{\theta} \phi_{\theta_t}(b_t) & (\text{Forward step}) \\ \theta_t^\dagger &= \theta_{t+1} + r \nabla_{\theta} \phi_{\theta_{t+1}}(b_t) & (\text{Backward step}) \end{aligned}$$

This discrepancy arises due to the gradient step's dependence on the current value of parameters in both the forward and backward optimizations, i.e., $\nabla_{\theta} \phi_{\theta_t}(b_t) \neq \nabla_{\theta} \phi_{\theta_{t+1}}(b_t)$.

However, the key observation here is that under the lazy dynamic regime (as described in Eq. 35), this dependency vanishes, and we have $\nabla_{\theta} \phi_{\theta^*}(b_t) \neq \nabla_{\theta} \phi_{\theta^*}(b_t)$, where θ^* is a typical sample from the stationary state (or slowly varying state) of parameters. Under such conditions, the backward action of SGD (running the learning protocol backward) results in a time-reversal of the parameters' trajectory: $\boldsymbol{\theta}_n^\dagger = \tilde{\boldsymbol{\theta}}_n$, signifying the thermodynamic reversibility of the parameters' subsystem under lazy dynamic conditions. See Ref. [41] on distinction between logical and thermodynamic reversibility.

As discussed in the paper, the lazy dynamics lead to a quasi-static evolution of the parameter subsystem, meaning that the subsystem Θ itself does not contribute to entropy production and acts as an ideal heat reservoir. Furthermore, the independence of the gradient step from the exact microscopic state of parameters aligns with path-independent forces in physics, which do not lead to dissipation and entropy production. This provides an alternative explanation for the reversibility of the parameter subsystem from a different perspective.

References

- [1] R. Landauer. Irreversibility and heat generation in the computing process. *IBM Journal of Research and Development*, 5(3):183–191, 1961.
- [2] L. Szilard. On the decrease of entropy in a thermodynamic system by the intervention of intelligent beings. *Z. Phys.* 53 (1929), 840.
- [3] Charles H. Bennett. The thermodynamics of computation—a review. *International Journal of Theoretical Physics*, 21(12):905–940, 1982.
- [4] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, Cambridge, 2010.
- [5] Ahmed Almheiri, Thomas Hartman, Juan Maldacena, Edgar Shaghoulian, and Amirhossein Tajdini. The entropy of hawking radiation. *Rev. Mod. Phys.*, 93:035002, Jul 2021.
- [6] T. Sagawa J. Parrondo, J. Horowitz. Thermodynamics of information. 2015.
- [7] Luca Peliti and Simone Pigolotti. *Stochastic Thermodynamics: An Introduction*. 2021.
- [8] Susanne Still, David A. Sivak, Anthony J. Bell, and Gavin E. Crooks. Thermodynamics of prediction. *Phys. Rev. Lett.*, 109:120604, 2012.
- [9] Takahiro Sagawa and Masahito Ueda. Fluctuation theorem with information exchange: Role of correlations in stochastic thermodynamics. *Physical review letters*, 109(18):180602, 2012.
- [10] Massimiliano Esposito, Katja Lindenberg, and Christian Van den Broeck. Entropy production as correlation between system and reservoir, 2009.
- [11] Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint arXiv:2101.03288*, 2021.
- [12] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [13] Hong Jun Jeon, Yifan Zhu, and Benjamin Van Roy. An information-theoretic framework for supervised learning. *arXiv preprint arXiv:2203.00246*, 2022.
- [14] Jirong Yi, Qiaosheng Zhang, Zhen Chen, Qiao Liu, and Wei Shao. Mutual information learned classifiers: an information-theoretic viewpoint of training deep learning classification systems. *arXiv preprint arXiv:2210.01000*, 2022.
- [15] Ravid Shwartz-Ziv and Yann LeCun. To compress or not to compress- self-supervised learning and information theory: A review, 2023.
- [16] Shujian Yu, Luis G. Sánchez Giraldo, and José C. Príncipe. Information-theoretic methods in deep neural networks: Recent advances and emerging opportunities. 2021.
- [17] Bernhard C Geiger. On information plane analyses of neural network classifiers—a review. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [18] Alessandro Achille, Giovanni Paolini, and Stefano Soatto. Where is the information in a deep neural network?, 2019.
- [19] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- [20] Andrew M. Saxe et al. On the information bottleneck theory of deep learning. 2018.
- [21] Bernhard C. Geiger. On information plane analyses of neural network classifiers—a review. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

- [22] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pages 5–13, 1993.
- [23] Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep representations. *The Journal of Machine Learning Research*, 19(1):1947–1980, 2018.
- [24] Tomas M. Cover and Joy A . Thomas. Elements of information theory. *John Wiley and Sons, New York, NY*, 1991.
- [25] Jorma Rissanen. Stochastic complexity and modeling. *Annals of Statistics*, 14:1080–1100, 1986.
- [26] Yuheng Bu, Shaofeng Zou, and Venugopal V Veeravalli. Tightening mutual information-based bounds on generalization error. *IEEE Journal on Selected Areas in Information Theory*, 1(1):121–130, 2020.
- [27] Christian Van den Broeck and Massimiliano Esposito. Ensemble and trajectory thermodynamics: A brief introduction. *Physica A: Statistical Mechanics and its Applications*, 418:6–16, 2015.
- [28] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks, 2018.
- [29] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in neural information processing systems*, 31, 2018.
- [30] Sebastian Deffner and Christopher Jarzynski. Information processing and the second law of thermodynamics: An inclusive, hamiltonian approach. *Physical Review X*, 3(4):041003, 2013.
- [31] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models, 2019.
- [32] Christian Maes. Local detailed balance. *SciPost Physics Lecture Notes*, Jul 2021.
- [33] Riccardo Rao and Massimiliano Esposito. Detailed fluctuation theorems: A unifying perspective. *Entropy*, 20(9):635, 2018.
- [34] Robert Zwanzig. *Nonequilibrium Statistical Mechanics*. 2001.
- [35] Mingwei Wei and David J Schwab. How noise affects the hessian spectrum in overparameterized neural networks. *arXiv preprint arXiv:1910.00195*, 2019.
- [36] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [37] Marcel Kühn and Bernd Rosenow. Correlated noise in epoch-based stochastic gradient descent: Implications for weight variances, 2023.
- [38] Michele Ceriotti, Giovanni Bussi, and Michele Parrinello. Langevin equation with colored noise for constant-temperature molecular dynamics simulations. *Physical Review Letters*, 102(2), Jan 2009.
- [39] Liu Ziyin, Hongchao Li, and Masahito Ueda. Law of balance and stationary distribution of stochastic gradient descent, 2023.
- [40] Shishir Adhikari, Alkan Kabakçioğlu, Alexander Strang, Deniz Yuret, and Michael Hinczewski. Machine learning in and out of equilibrium, 2023.
- [41] Takahiro Sagawa. Thermodynamic and logical reversibilities revisited. *Journal of Statistical Mechanics: Theory and Experiment*, 2014(3):P03025, Mar 2014.