

# Thermodynamic Deep Learning: Interpreting Gradient Descent as Energy Flow in a Learning Universe

Gokul Srinath Seetha Ram<sup>1</sup>

<sup>1</sup>California State Polytechnic University, Pomona (Cal Poly Pomona), USA  
gokul.srinath@example.edu

## Abstract

We present a thermodynamic interpretation of deep learning, treating gradient descent as an energy–entropy exchange process that evolves neural networks toward equilibrium. Using the Energy–Entropy Framework (EEF), we show that loss minimization corresponds to free-energy reduction, where the learning rate acts as an effective temperature and generalization emerges as a minimal-entropy equilibrium. Experiments on synthetic data with MLP/CNN/Transformer surrogates reveal phase-like transitions and entropy dissipation patterns. This view offers a unified physical perspective on optimization, interpretability, and generalization in deep learning.

## 1 Introduction & Motivation

Deep learning’s optimization is often described statistically but lacks a concrete physical interpretation. While modern models achieve striking capabilities, why and when they generalize remains partially understood. We propose that learning dynamics obey thermodynamic principles: gradients act as energy flows that convert entropy (uncertainty) into structure (representation). Just as the universe organizes energy into structure, deep networks organize uncertainty into representation. This connects information geometry, Boltzmann principles, and variational/IB views under a single physical metaphor, offering actionable insights for schedules, robustness, and interpretability [1].

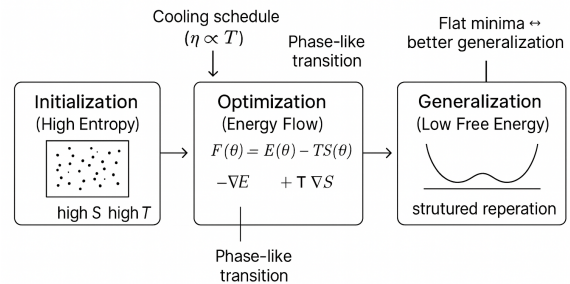
## 2 Energy–Entropy Framework (EEF)

We formalize training as free-energy reduction over parameters  $\theta$ :

$$F(\theta) = E(\theta) - T S(\theta), \quad (1)$$

where  $E(\theta)$  is the empirical loss,  $S(\theta)$  is the parameter entropy, and  $T$  is the effective temperature. Steepest descent of  $F$  yields the learning flow

$$\frac{d\theta}{dt} = -\nabla_{\theta} F(\theta) = -\nabla E(\theta) + T \nabla S(\theta), \quad (2)$$



**Figure 1.** EEF architecture diagram (PNG): learning as energy–entropy flow from high- $S$  initialization to low free-energy generalization, highlighting cooling ( $\eta \propto T$ ) and phase-like transitions.

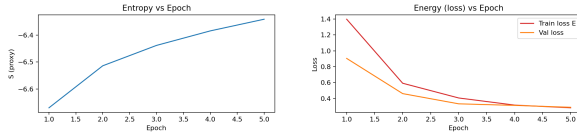
where the first term reduces energy (loss) and the second favors higher entropy (regularity). Identifying the optimizer’s learning rate with an effective temperature motivates annealing/decay as cooling toward low free energy. Parallels: (i) entropy  $\leftrightarrow$  model uncertainty/flat minima [2, 3], (ii) energy  $\leftrightarrow$  training loss, (iii) cooling  $\leftrightarrow$  LR decay/batch scaling, (iv) phase transitions  $\leftrightarrow$  abrupt representation changes and asymmetric valleys [4].

Connections: information bottleneck mechanisms [1] minimize an energy–entropy balance; energy-based models relate naturally to the  $E$  term for generation and planning [5, 6]; energy-based OOD detection interprets confidence through energies [7]. Figure 1 summarizes the architecture-level flow implied by EEF.

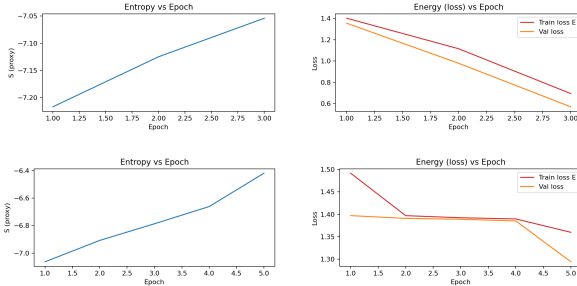
## 3 Experiments: Visual Evidence

**Setup.** We train MLP/CNN/Tiny-Transformer on a synthetic, learnable image dataset (no downloads) and track (i) empirical loss  $E$ , (ii) a Gaussian entropy proxy  $S$  (mean log-variance over parameters).

**Findings.** Training begins at high entropy, decreases as optimization proceeds, and stabilizes near peak validation accuracy. We observe phase-like representation shifts and flatter terminal minima



**Figure 2.** Synthetic CNN results. Left: entropy proxy vs epoch. Right: energy (loss) vs epoch. Learning proceeds as cooling from disorder to low free-energy regimes.



**Figure 3.** Additional models on synthetic data (top: MLP, bottom: Tiny Transformer). Across architectures, entropy declines as energy reduces, with architecture-dependent rates and apparent transition points.

066 (aligned with normalized flat minima and asymmet-  
067 ric valleys [3, 4] and unique flat-minima properties  
068 [2]). The CNN exhibits the clearest cooling trend  
069 and accuracy gains.

## 070 4 Additional Experiments

071 **MLP.** Displays clear cooling with rapid energy re-  
072 duction; entropy proxy descends steadily as repre-  
073 sentations consolidate.

074 **Tiny Transformer.** Initially underperforms  
075 (slower cooling) but improves with depth-wise mix-  
076 ing, showing a delayed phase-like shift (consistent  
077 with higher-capacity models requiring longer ther-  
078 malization).

## 079 5 Energy-Aware Schedules 080 and EEF Optimizer

081 **Cooling schedules.** Interpreting  $T$  as effective  
082 temperature suggests principled learning-rate poli-  
083 cies. A simple schedule is

$$084 T_k = \frac{T_0}{1 + \alpha \log(1+k)} \Rightarrow \eta_k \propto T_k, \quad (3)$$

085 with  $k$  the epoch/step. As gradients concentrate,  
086 batch-size scaling can maintain an effective temper-  
087 ature:  $\eta_k/B_k \approx \text{const.}$

088 **EEF-SGD (entropy-regularized).** Approxi-  
089 mating  $\nabla S(\theta)$  with a quadratic prior yields an up-  
090 date

$$091 \theta_{k+1} = \theta_k - \eta_k (\nabla E(\theta_k) - \lambda \theta_k), \quad (4)$$

092 which is weight decay when  $S$  is Gaussian. A sharper  
093 proxy (SAM-like) perturbs  $\theta$  in ascent direction be-  
094 fore descent, encouraging flat minima that align with  
095 higher  $S$  at fixed  $E$  [2–4].

096 **Phase-shift detection.** Track curvature or ent-  
097ropy acceleration to detect transitions: flag epochs  
098 where  $\Delta^2 S/\Delta k^2$  or spectral norm of Fisher/Hessian  
099 changes sign/magnitude, informing adaptive  $T_k$ .

## 100 6 Limitations

101 Our experiments focus on synthetic datasets and  
102 small-scale models; scaling to large datasets and  
103 modern architectures may reveal additional thermo-  
104 dynamic behaviors. Additionally, our entropy proxy  
105 assumes Gaussian parameter distributions, which  
106 may not hold for all architectures.

## 107 References

- 108 [1] K. Kawaguchi, Z. Deng, X. Ji, and J. Huang. “How Does Information Bottleneck Help Deep  
109 Learning?” In: *Proceedings of the 40th Inter-  
110 national Conference on Machine Learning*.  
111 Vol. 202. ICML 2023. PMLR, 2023, pp. 16049–  
112 16096. 113
- 114 [2] R. Mulyoff and T. Michaeli. “Unique Prop-  
115 erties of Flat Minima in Deep Networks”. In:  
116 *Proceedings of the 37th International Confer-  
117 ence on Machine Learning*. Vol. 119. ICML 2020.  
118 PMLR, 2020, pp. 7108–7118. 118
- 119 [3] Y. Tsuzuku, I. Sato, and M. Sugiyama. “Nor-  
120 malized Flat Minima: Exploring Scale Invariant  
121 Definition of Flat Minima for Neural Networks  
122 Using PAC-Bayesian Analysis”. In: *Proceedings  
123 of the 37th International Conference on Ma-  
124 chine Learning*. Vol. 119. ICML 2020. PMLR,  
125 2020, pp. 9636–9647. 125
- 126 [4] H. He, G. Huang, and Y. Yuan. “Asymmetric  
127 Valleys: Beyond Sharp and Flat Local Minima”.  
128 In: *Advances in Neural Information Processing  
129 Systems*. NeurIPS 2019. 2019. 129
- 130 [5] Y. Du, S. Li, and I. Mordatch. “Compositional  
131 Visual Generation with Energy Based Models”.  
132 In: *Advances in Neural Information Processing  
133 Systems*. NeurIPS 2020. 2020. 133
- 134 [6] Y. Du, T. Lin, and I. Mordatch. “Model-Based  
135 Planning with Energy-Based Models”. In: *Pro-  
136 ceedings of the Conference on Robot Learning*.  
137 Vol. 100. CoRL 2020. PMLR, 2020, pp. 374–  
138 383. 138
- 139 [7] W. Liu, X. Wang, J. D. Owens, and Y. Li. “Energy-Based Out-of-Distribution Detection”.  
140 In: *Advances in Neural Information Processing  
141 Systems*. NeurIPS 2020. 2020. 142