New Evidence of the Two-Phase Learning Dynamics of Neural Networks

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2025

Abstract

Extensive evidence suggests that training dynamics undergo a distinct phase transition, yet our understanding of this transition still lags behind. In this paper, we introduce an interval-wise perspective that compares network states across a time window, revealing two new phenomena that illuminate the two-phase nature of deep learning. i) **The Chaos Effect.** By injecting an imperceptibly small parameter perturbation at various stages, we show that the response of the network to the perturbation exhibits a transition from chaotic to stable, suggesting there is an early critical period where the network is highly sensitive to initial conditions; ii) **The Cone Effect.** Tracking the evolution of the empirical Neural Tangent Kernel (eNTK), we find that after this transition point the model's functional trajectory is confined to a narrow cone-shaped subset: while the kernel continues to change, it gets trapped in a tight angular region. Together, these effects provide a dynamical view of how deep networks transition from sensitive exploration to stable refinement during training.

1. Introduction

Many recent studies have suggested, either implicitly or explicitly, that there is a *phase transition* point during the neural network training, where the model's properties and behaviors undergo substantial shifts before and after this time point. For example, Cohen et al. [4], Damian et al. [5], Wang et al. [14] showed that during training, the network first enters a progressive sharpening phase, and after which, the sharpness remains roughly constant for the rest of the training. Achille et al. [1] identified a critical learning period early in training, during which exposure to low-quality data can cause irreversible damage, while similar exposure later in training can be reversed.

Despite abundant evidence for the *two-phase phenomenon* [1, 4, 9, 14–16], a complete characterization and understanding of this phenomenon still lags behind. Moreover, most existing studies adopt a *point-wise* perspective: they primarily focus on examining specific properties of the network at isolated time points. This perspective, while informative, offers only a static snapshot of the model's behavior, and does not capture the temporal dynamics of learning: how a property emerges, evolves, or vanishes as training progresses.

In this paper, to gain a deeper understanding of the two-phase phenomenon, we introduce two new empirical observations that exhibit characteristics of the phase transition. Crucially, these are what we call "interval-wise" phenomena: rather than analyzing a property at specific time points, we compare the model's behavior across two different time points of training. We show that this novel approach reveals patterns that are otherwise invisible to point-wise analysis and offers new insights into the learning dynamics of neural networks. Specifically, we identify and investigate two distinct behaviors: *the Chaos Effect* and *the Cone Effect*.

The Chaos Effect. First, we observe that the learning dynamics of neural networks transition from a chaotic to a non-chaotic regime during training. Specifically, we train two networks that are initialized identically and trained with the same stochastic gradient noise. At a specific time t_0 , we apply a small perturbation to the parameters of one model, and then we compare the resulting parameters at a later time t_1 . Particularly, we observe an *inflection point* during the training process. We find that when t_0 is in the early stage of training, specifically before the inflection point, even a tiny perturbation leads to a significant divergence from the original training trajectory. This phenomenon indicates a high sensitivity of learned parameters to initial conditions, which is a hallmark of chaotic systems in physics. However, if t_0 is later in training (after the inflection point), the divergence is minimal, suggesting that the system becomes increasingly stable as training progresses.

The Cone Effect. Second, we discover that after the early training phase, the learning dynamics of neural networks remain constrained in a narrow cone in the function space. Specifically, we train a network, and starting from a chosen time point τ , we track the empirical Neural Tangent Kernel (eNTK) at later time steps and measure their deviation from the eNTK at τ . We observe that when τ is sufficiently large, the subsequent eNTKs remain confined within a narrow cone around the eNTKs at time τ . In contrast, if the τ is chosen in the early stages of training, the eNTKs experience chaotic and large unstructured changes over time, and no such confinement is observed.

2. Preliminary and Methodology

Basic Notations. We focus on a classification task. Denote $[k] = \{1, 2, \dots, k\}$. Let $\mathcal{D} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ be the training set of size n, where $\boldsymbol{x}_i \in \mathbb{R}^{d_0}$ represents the *i*-th input and $y_i \in [c]$ represents the corresponding target. Here, c is the number of classes. Let $f : \mathcal{D} \times \mathbb{R}^p \to \mathbb{R}$ be the NN model, and thus $f(\boldsymbol{x}, \boldsymbol{\theta}) \in \mathbb{R}$ denotes the output of model f on the input \boldsymbol{x} with parameter $\boldsymbol{\theta} \in \mathbb{R}^p$. Let $\ell(f(\boldsymbol{x}_i, \boldsymbol{\theta}), y_i)$ be the loss at the *i*-th data point, simplified to $\ell_i(\boldsymbol{\theta})$. The total loss over the dataset \mathcal{D} is then denoted as $\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\boldsymbol{\theta})$. We also use $\operatorname{Err}_{\mathcal{D}}(\boldsymbol{\theta}) / \operatorname{Acc}_{\mathcal{D}}(\boldsymbol{\theta})$ to denote the classification error/accuracy of the network $f(\boldsymbol{\theta}; \cdot)$ on the training set \mathcal{D} .

Parameter Dissimilarity. Following Singh et al. [13], we first introduce the *parameter dissimilarity* to measure the directionality of the optimization process. Specifically, given the training trajectory consisting of a sequence of checkpoints $\{\theta_t\}_{t=0}^T$, we use the pairwise cosine dissimilarity to capture the directional aspect of the trajectory. For any two time points $i, j \in [T]^2$, we define

$$(\boldsymbol{C})_{i,j} := 1 - \cos\langle \operatorname{vec}(\boldsymbol{\theta}_i), \operatorname{vec}(\boldsymbol{\theta}_j) \rangle = 1 - \langle \operatorname{vec}(\boldsymbol{\theta}_i), \operatorname{vec}(\boldsymbol{\theta}_j) \rangle / (\|\operatorname{vec}(\boldsymbol{\theta}_j)\|_2 \|\operatorname{vec}(\boldsymbol{\theta}_j)\|_2), \quad (1)$$

where $vec(\theta)$ denotes the flattened parameters of the network.

Kernel Distance. Similar to Fort et al. [8], we use the *kernel distance* to quantify the evolution of neural networks in the function space. Specifically, we measure the pairwise distance between the *empirical neural tangent kernel* (eNTK) matrix $H(\theta)$ at two different time points, namely *kernel distance*. For two time points $i, j \in [T]^2$, we define:

$$(\mathbf{S})_{i,j} := 1 - \frac{\langle \mathbf{H}(\boldsymbol{\theta}_i), \mathbf{H}(\boldsymbol{\theta}_j) \rangle}{\|\mathbf{H}(\boldsymbol{\theta})_i\|_F \|\mathbf{H}(\boldsymbol{\theta}_j)\|_F}.$$
(2)



Figure 1: The sensitivity of learning dynamics to tiny perturbations. (a) The parameter dissimilarity $(C)_{t_0,t_1}$. (b) The loss barrier $(B)_{t_0,t_1}$. (c) The disagreement rate $(D)_{t_0,t_1}$. Note that the t_0 and t_1 are presented in iterations, not epochs.

Loss Barriers. In addition to the directional and functional aspect, we also investigate the geometry of the neural network's loss landscape. In particular, we examine the *loss barriers* [2, 9] between any two points along the training trajectory. For any two points $i, j \in [T]^2$, we define:

$$(\boldsymbol{B})_{i,j} := \max_{\alpha} \mathcal{L}_{\mathcal{D}'}(\alpha \boldsymbol{\theta}_i + (1-\alpha)\boldsymbol{\theta}_j) - \frac{1}{2}(\mathcal{L}_{\mathcal{D}'}(\boldsymbol{\theta}_i) + \mathcal{L}_{\mathcal{D}'}(\boldsymbol{\theta}_j)),$$
(3)

where \mathcal{D}' denotes the unseen test set.

Disagreement Rate. Lastly, we concern about the similarity of the outputs at two points along the optimization trajectory. Specifically, we introduce the *disagreement rate* on the test data for any two points $i, j \in [T]^2$:

$$(\boldsymbol{D})_{i,j} := \mathbb{E}_{\boldsymbol{x} \in \mathcal{D}'} [\mathbf{1}(f(\boldsymbol{x}, \boldsymbol{\theta}_i) \neq f(\boldsymbol{x}, \boldsymbol{\theta}_j))], \tag{4}$$

where $\mathbf{1}(\cdot)$ is the indicator function and \mathcal{D}' is the test set.

Main Experimental Setup. We train the VGG-16 architecture [12] and the ResNet-20 architecture [10] on the CIFAR-10 dataset. Optimization is done using SGD with momentum (momentum set to 0.9). A weight decay of 1×10^{-4} is applied. The learning rate is initialized at 0.1 and is dropped by 10 times at 80 and 120 epochs. The total number of epochs is 160.

3. The Chaos Effect: Sensitivity of Learning Dynamics to Small Perturbations

In the section, we study the sensitivity of neural network learning dynamics to small perturbations.

Experimental Design. We train two networks with identical initializations and the same stochastic gradient noise. However, at a specific time t_0 , we introduce a small perturbation ϵ to the parameters of one network, such that $\theta'_{t_0} = \theta_{t_0} + \epsilon$. We then compare the resulting models at a later time t_1 (



Figure 2: The kernel distance between every pair of two points at the optimization trajectory $\{\theta_t\}_{t=1}^T$. Our results are reported for both VGG-16 and ResNet-20 on CIFAR-10. Note that the *i* and *j* are presented in iterations, not epochs.

 $t_1 > t_0$), with the parameters θ'_{t_1} and θ_{t_1} respectively. We consider a tiny perturbation here, where $\|\boldsymbol{\epsilon}\|_0 = 10^{-7}$. We vary the time t_0, t_1 and compare the two resulting models using different metrics.

Finding I. Optimization trajectory changes its direction at an inflection point. In Figure 1 (*a*), we present the parameter dissimilarity for any pair of t_0 and t_1 (with $t_1 \ge t_0$). Notably, there exists a specific time t_1 at which the value of $(C)_{t_0,t_1}$ remains relatively high for all choices of t_0 . Typically, a high value of $(C)_{t_0,t_1}$ indicates the directional change along the optimization trajectory. Therefore, the optimization trajectory evidently changes its direction at a fixed point, namely the *inflection point*.

Finding II. Tiny perturbations applied before the inflection point leads to significant loss barriers and disagreement rate. In Figure 1 (b), we also report the loss barrier between each pair of time points t_0, t_1 . We observe that even a tiny perturbation ($\|\epsilon\|_0 = 10^{-7}$) applied at a early time point t_0 could result in a substantial loss barrier between the resulting parameters θ_{t_1} and θ'_{t_1} in the later stage of training. This indicates that the two solutions likely reside in different, isolated "valleys" of the loss landscape. In Figure 1 (c), we further evaluate the disagreement rate between θ_{t_1} and θ'_{t_1} . The high disagreement rate validates the functional dissimilarity between θ_{t_1} and θ'_{t_1} . Together, we show that the learning dynamics pass through a chaotic regime during the early phase of training, where small perturbations might lead to substantial divergence, namely the *chaos effect*.

Conjecture I. The inflection points marks the transition from a chaotic to an non-chaotic regime. Taking VGG-16 on CIFAR-10 as an example, we observe that a significant loss barrier $(B)_{t_0,t_1}$ emerges only when $t_0 \leq 2500$ iterations and $t_1 > 2500$ iterations. Similar observations are also noted for the disagreement rate. Recall that the 2500 iteration marks the inflection point for VGG-16 on CIFAR-10. Therefore, we conjecture that *the inflection point serves as a hallmark of the transition from a chaotic to a non-chaotic training regime*. To verify, we compute the kernel distance between every two points along the optimization trajectory. Specifically, we train the neural network and obtain a sequence of checkpoints $\{\theta_t\}_{t=1}^T$. Then for any $i, j \in [T]^2$ we compute the kernel distance, i.e., $(S)_{i,j}$. In Figure 2, we observe that the eNTKs evolve significantly during the early phase of training, indicating a chaotic regime. Subsequently, the evolution of the eNTKs stabilizes, transitioning to a non-chaotic phase. Notably, the transition point in the evolution of the eNTKs aligns with the inflection points identified in earlier experiments.



Figure 3: Constrained learning dynamics in the second phase. (a) The kernel distance $(S)_{\tau,t}$ v.s. training iteration t. (b) The kernel distance $(S)_{t,t+dt}$ vs. training iteration t. (c) The visualization of the changes of the eNTK matrices $H(\theta_t)$.

4. The Cone Effect: Constrained Learning Dynamics in the Second Phase

We have seen that the learning dynamics of neural networks undergo a transition from a highly chaotic to a more stable, non-chaotic phase. In this section, we dig deeper into the second "stable" phase. Surprisingly, we find that, contrary to the typical assumption of the *lazy training regime* [3, 6, 7, 11, 17], the neural network continues to evolve. However, this evolution is confined within a narrow, "cone"-like region in function space, namely the *cone effect*.

Beyond the Lazy Regime: The Cone Effect. First, we compute the kernel distance between two adjacent points $\theta(t)$ and $\theta(t + dt)$. In Figure 3 (b), we observe that the kernel distance $(S)_{t,t+dt}$ is significant in the early training phase and then drops quickly to a low but non-negligible value. However, surprisingly, we note that in the later training phase, the values of $(S)_{t,t+dt}$ are upperbounded by the same value for different dt. One possible explanation for this phenomenon is that during the second phase, the eNTK matrix evolves in a constrained space.

To validate this, we further measure how the distance between the kernel matrices at the current iterate θ_t and a referent point θ_{τ} changes during training. In Figure 3 (*a*), for different referent points τ , the kernel distance $S(\theta_t, \theta_{\tau}$ first increases and then keeps nearly constant in training. This result suggests that during the second phase, beyond the lazy regime, the model operates in a constrained function space. The visualization in Figure 3 (*c*) further confirms the existence of the cone effect, where a clear "cone" pattern is observed during the evolution of eNTK matrices.

5. Conclusion

In this paper, we introduced an interval-wise perspective on neural network training dynamics. Through this lens, we identified two novel empirical phenomena that characterize a two-phase transition in deep learning: the chaos effect and the cone effect. Together, these findings suggest a transition from an exploratory, unstable phase to a more stable, refinement-oriented phase during training.

References

- [1] Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep neural networks. *arXiv preprint arXiv:1711.08856*, 2017.
- [2] Samuel K. Ainsworth, Jonathan Hayase, and Siddhartha S. Srinivasa. Git re-basin: Merging models modulo permutation symmetries. In *ICLR*. OpenReview.net, 2023.
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings* of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pages 242–252. PMLR, 09–15 Jun 2019.
- [4] Jeremy M Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. *arXiv preprint arXiv:2103.00065*, 2021.
- [5] Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=nhKHA59gXz.
- [6] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019.
- [7] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
- [8] Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy, and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. *Advances in Neural Information Processing Systems*, 33:5850–5861, 2020.
- [9] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR, 2020.
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- [11] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

- [12] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http://arxiv.org/abs/1409.1556.
- [13] Sidak Pal Singh, Bobby He, Thomas Hofmann, and Bernhard Schölkopf. The directionality of optimization trajectories in neural networks. In *The Thirteenth International Conference* on Learning Representations, 2025. URL https://openreview.net/forum?id= JY6P45sFDS.
- [14] Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along gd trajectory: Progressive sharpening and edge of stability. *Advances in Neural Information Processing Systems*, 35: 9983–9994, 2022.
- [15] Yongyi Yang, Core Francisco Park, Ekdeep Singh Lubana, Maya Okawa, Wei Hu, and Hidenori Tanaka. Swing-by dynamics in concept learning and compositional generalization. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [16] Zhanpeng Zhou, Yongyi Yang, Xiaojiang Yang, Junchi Yan, and Wei Hu. Going beyond linear mode connectivity: The layerwise linear feature connectivity. *Advances in neural information* processing systems, 36:60853–60877, 2023.
- [17] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes overparameterized deep relu networks. *Machine learning*, 109:467–492, 2020.