

Toward Equation of Motion for Deep Neural Networks: Continuous-time Gradient Descent and Discretization Error Analysis ¹

Taiki Miyagawa ²
NEC Corporation, Japan
miyagawataik@nec.com

Abstract ³

We derive and solve an “Equation of Motion” (EoM) for deep neural networks (DNNs), a differential equation that precisely describes the discrete learning dynamics of DNNs. Differential equations are continuous but have played a prominent role even in the study of discrete optimization (gradient descent (GD) algorithms). However, there still exist gaps between differential equations and the actual learning dynamics of DNNs due to *discretization error*. In this paper, we start from gradient flow (GF) and derive a counter term that cancels the discretization error between GF and GD. As a result, we obtain *EoM*, a continuous differential equation that precisely describes the discrete learning dynamics of GD. We also derive discretization error to show to what extent EoM is precise. In addition, we apply EoM to two specific cases: scale- and translation-invariant layers. EoM highlights differences between continuous-time and discrete-time GD, indicating the importance of the counter term for a better description of the discrete learning dynamics of GD. Our experimental results support our theoretical findings. ⁴

1 Introduction ⁵

Let us first explain our primary motivation for the present paper. In *physics*, one of the fundamental goals is to predict the dynamics of matter and its fundamental constituents. Specifically, “predict” here means to construct differential equations that best describe the physical system under consideration and to solve them. Such differential equations are called *Equations of Motion* (EoM). An interesting question here may be “What is the EoM for deep neural networks (DNNs)?” That is, to what extent can we predict the discrete learning dynamics of DNNs by constructing differential equations? This is our research question. ⁶

Differential equations have played a prominent role in studying discrete optimization (gradient descent (GD) algorithms), although they are continuous [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20]. In the context of deep ⁷

learning, gradient flow (GF) and stochastic differential equations (SDEs) are used to analyze (stochastic) gradient descent ((S)GD). Research targets include: convergence [6, 7, 8, 12, 13, 9, 14, 17], ¹⁰

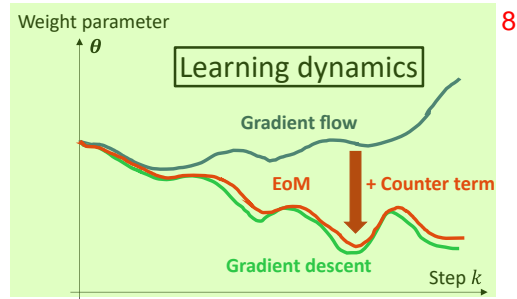


Figure 1: **Our approach.** GF fails in describing the learning dynamics of GD due to *discretization error*. Our counter term approach successfully cancels the discretization error between GF and GD and hence allows for a reliable analysis of GD. ⁹

stability of optimization [19], optimization with constraints [19], convergent states [17, 20], flatness of loss landscapes [17], empirical risk bounds [15], and online PCA [11]. Various techniques for continuous analysis have been imported to the analysis of discrete GD algorithms. 1

However, there still exist gaps between differential equations and actual learning dynamics due to *discretization error*, which is the main interest of the present paper and is often missing in the literature above. To be specific, we focus on GF $\dot{\theta}(t) = -g(\theta(t))$ as a continuous approximation of GD $\theta_{k+1} = \theta_k - \eta g(\theta_k)$, where $\theta(t) \in \mathbb{R}^d$ and $\theta_k \in \mathbb{R}^d$ are the weight parameters of a DNN at time $t \in \mathbb{R}$ and step $k \in \mathbb{Z}$, respectively, and $g \in \mathbb{R}^d$ is a gradient vector. $\eta \in \mathbb{R}$ is a learning rate and is regarded as the discretization step size when GF is discretized with the Euler method [21]: $\dot{\theta}(t = k\eta) \doteq \frac{\theta_{k+1} - \theta_k}{\eta}$. Due to this approximation, discretization error (or “continuation error”) is introduced, and thus GF cannot fully explain the dynamics of GD. For instance, we show that according to GF, the weight norm of a scale-invariant layer collapses to zero when we use weight decay, while GD does not show such behavior (Section 5.1). 2

To fill the critical gap between GF and GD, we propose modifying GF to describe the learning dynamics of GD more precisely; i.e., we add a counter term $\xi \in \mathbb{R}^d$ to the gradient g of GF that cancels the discretization error (Figure 1). This idea is motivated by backward error analysis in numerical analysis [21]. We derive a functional integral equation that determines the counter term and solve it (Section 3). As a result, we obtain a more reliable differential equation, called *EoM* here, that describes the discrete learning dynamics of GD. Using the counter term, we derive the leading order of discretization error (Section 4.1) to show to what extent GF and EoM are precise in describing GD’s dynamics. This point is often missed in the literature on the continuous approximation of discrete GD algorithms [22, 23, 24, 11, 25, 26, 27, 28]. We further derive a sufficient condition for learning rates for the discretization error to be small (Section 4.2). We show that EoM well explains empirical results. 3

Furthermore, to show the benefits of EoM, we apply it to two specific cases: scale-invariant layers [29, 30] and translation-invariant layers [31, 32] (Section 5). For scale-invariant layers, we show that a better description of GD’s discrete dynamics requires modifications to the decay rate of weight norms that is previously derived in the continuous regime (SDEs) [33]. In addition, we show that EoM successfully reproduces the limiting dynamics ($t \rightarrow \infty$) of weight norms and angular update [34] that are previously derived in the discrete regime, while GF cannot reproduce this result. For translation-invariant layers, we show that EoM rather than GF dramatically matches empirical results, indicating the importance of the counter term. To the best of our knowledge, no study analyzes the temporal evolution of translation-invariant layers except for [31] and [32], where only the sum of weights is their focus, while we derive the dynamics of the whole weights. 4

Our contribution is four-fold. Our code¹ and detailed experimental results are given as supplementary materials. 5

1. To fill the critical gap between GF and GD, we derive a counter term for GF that cancels the discretization error, and as a result, we obtain EoM, a continuous differential equation that precisely describes the discrete learning dynamics of GD. 6
2. To show to what extent GF and EoM are precise in describing discrete GD dynamics, we derive the leading order of discretization error, as is often missed in the literature on the continuous approximation of discrete GD algorithms. We further derive a sufficient condition for learning rates for the discretization error to be small.
3. We apply EoM to two specific cases: scale-invariant layers and translation-invariant layers, indicating the importance of the counter term for a better description of the discrete learning dynamics of GD.
4. Our experimental results support our theoretical findings.

Our work is the first step toward answering this research question: to what extent can we predict the discrete learning dynamics of DNNs by constructing differential equations (EoM for DNNs)? Also, our work helps researchers import continuous analysis to the discrete analysis of GD algorithms. In this sense, our work bridges discrete and continuous analyses of GD algorithms. 7

¹See Supplementary Materials at <https://openreview.net/forum?id=qq84D17BPu>.

2 Related Work ¹

The idea of approximating discrete-time stochastic algorithms with continuous equations dates back ² to stochastic approximation theory [1, 2, 3, 4, 5]. Their primary focus is convergence analysis for discrete-time algorithms, while our focus is to predict the learning dynamics (temporal evolution) of weight parameters, such as the decay rates of weight norms and effective learning rate of scale-invariant layers. Our idea of the counter term is inspired by the backward error analysis developed for numerical analysis [35]. This idea is now used to analyze discrete optimization [22, 23, 24, 11, 25, 26, 27, 28]. [18] is a pioneering work on discretization error analysis between GF and GD that is based on the numerical analysis of the Euler method [21]. They derive a sufficient condition for learning rates for the discretization error to be small. This analysis is based on a bound (inequality), while we derive an explicit relationship between learning rates and discretization error as an equality.

Neural mechanics and Noether’s learning dynamics [31, 32] provide a solution to a part of the afore- ³ mentioned problem: to what extent can we predict the learning dynamics of DNNs by constructing differential equations? They derive (the breaking of) conservation laws of weight parameters using differential equations and provide the temporal evolution of the conserved quantities. The present work is inspired by these studies but has crucial differences: 1) our focus is on the temporal evolution of all of the network parameters, not only the conserved quantities, 2) the gradient’s correction for canceling the discretization error is not limited to the first order, but all orders, and 3) the discretization error is explicitly provided in the present paper. See Appendix G for more related studies.

3 Equation of Motion for Deep Neural Networks ⁴

In the following sections, we define *EoM* by modifying GF (Section 3.1). We show that the counter ⁵ term satisfies a functional integral equation (Section 3.2), and then we solve it (Section 3.3).

3.1 Our Approach and Definitions ⁶

We begin with a simple idea: add a counter term to GF to cancel discretization error, i.e., ⁷

$$\dot{\theta}(t) = -g(\theta(t)) - \eta\xi(\theta(t)), \quad (1) \quad (8)$$

where $\theta(t) \in \mathbb{R}^d$ is the vectorized weight parameters of a DNN at time $t \in \mathbb{R}$, $d \in \mathbb{N}$ is the ⁹ dimension of the weight, and $\dot{\theta}(t)$ denotes $d\theta(t)/dt$. Gradient $g(\theta(t))$ is defined as $g(\theta(t)) := \nabla f(\theta(t)) + \lambda\theta(t)$, which consists of a loss function $f(\theta(t))$ and weight decay term $\lambda\theta(t)$, where $\lambda > 0$ controls the strength of weight decay. $\eta > 0$ is a small learning rate, and $\xi(\theta(t)) \in \mathbb{R}^d$ is the counter term. Throughout this paper, we assume all functions are sufficiently smooth. We call Equation (1) the *Equation of Motion (EoM)* for DNNs, or simply EoM.

Our aim is to find ξ that makes Equation (1) more reliable to precisely approximate GD $\theta_{k+1} =$ ¹⁰ $\theta_k - \eta g(\theta_k)$, where $\theta_k \in \mathbb{R}^d$ is the weight at step $k \in \mathbb{Z}_{\geq 0}$. To do so, we first define the *discretization error* between GF (1) and GD at step k :

$$e_k := \theta(k\eta) - \theta_k \in \mathbb{R}^d \quad (11) \quad (2)$$

and find ξ that makes e_k small. Throughout this paper, we use the standard Euler method to discretize ¹² GF: $\dot{\theta}(t) \doteq (\theta(t + \eta) - \theta(t))/\eta$ and $t = k\eta$; thus, η is identified with the discretization step size.

3.2 How to Determine Counter Term ¹³

We show that the leading order of e_k with respect to η is controlled by the counter term (Theorem ¹⁴ 3.2), and as a result, the counter term is determined via a functional integral equation (Equation (6)).

Our first theorem shows what the counter term should cancel. ¹⁵

Theorem 3.1 (Recursive formula for discretization error). *Discretization error e_k satisfies:* ¹⁶

$$e_{k+1} - e_k = -\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) + \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) - \eta^2 \xi(\theta(k\eta)) \quad (3) \quad (17)$$

$$=: -\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) + \Lambda(\theta(k\eta)). \quad (4)$$

Here, we defined $\Lambda(\theta(k\eta)) := \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) - \eta^2 \xi(\theta(k\eta)) \in \mathbb{R}^d$. The proof is based on Taylor's theorem and is given in Appendix A.1. The right-hand side of Equation (3) tells us that the counter term (third term) should cancel the first and second terms. However, the following theorem states that the first term gives only subleading contributions with respect to η .

Theorem 3.2 (Leading order of discretization error). *Suppose that $\Lambda(\theta(k\eta)) = O(\eta^\gamma)$ and $e_0 = O(\eta^\gamma)$ for some $\gamma > 0$. Then $e_k = O(\eta^\gamma)$ and $-\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) = O(\eta^{\gamma+1})$. Therefore, the first term in the right-hand side of Equation (3) is negligible compared with Λ :*

$$\begin{aligned} e_{k+1} &= e_k + \Lambda(\theta(k\eta)) - \eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) \\ &= e_k + \Lambda(\theta(k\eta)) + O(\eta^{\gamma+1}) \quad (k = 0, 1, 2, \dots) \end{aligned} \quad (5)$$

The proof is by induction and given in Appendix A.2. Therefore, the leading order of discretization error is $O(\eta^\gamma)$ and given by:

$$\Lambda(\theta(k\eta)) = O(\eta^\gamma) \iff \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) - \xi(\theta(k\eta)) = O(\eta^{\gamma-2}). \quad (6)$$

This is a functional equation of ξ because $\ddot{\theta}(t)$ contains ξ via Equation (1). A solution to Equation (6) for a large γ gives a small Λ and thus gives a small e_k via Equation (5).

3.3 Solution to Equation 6

How can we solve Equation (6)? It is not easy to find an exact solution because Equation (6) is a functional integral equation [36, 37, 38, 39, 40]; therefore, we assume a power series solution with respect to η :

$$\xi(\theta(k\eta)) = \sum_{\alpha=0}^{\infty} \eta^\alpha \xi_\alpha = \xi_0(\theta(k\eta)) + \eta \xi_1(\theta(k\eta)) + \eta^2 \xi_2(\theta(k\eta)) + \dots \quad (7)$$

In the following theorem, we successfully find a solution for *all* orders of η .

Theorem 3.3 (Solution of Equation 6). *The solution to Equation (6) of form (7) is given by*

$$\xi_\alpha(\theta) = \tilde{\xi}_\alpha(\theta) := \sum_{i=2}^{\alpha+2} \sum_{k_1+\dots+k_i=\alpha-i+2} \frac{(-1)^i}{i!} D_{k_1} \dots D_{k_{i-1}} \Xi_{k_i} \quad (8)$$

for $\alpha = 0, 1, 2, \dots$, where we use differential operators (Lie derivatives) $\mathcal{D}_\alpha := \tilde{\xi}_{\alpha-1}(\theta) \cdot \nabla$ ($\alpha = 1, 2, \dots$) and $\mathcal{D}_0 := g(\theta) \cdot \nabla$ and also defined $\Xi_\alpha(\theta) := \tilde{\xi}_{\alpha-1}(\theta)$ ($\alpha = 1, 2, \dots$) and $\Xi_0(\theta) := g(\theta)$.

The proof follows from the definition of the Lie derivative and is given in Appendix A.3. The first two orders of the solution are given by:

$$\tilde{\xi}_0(\theta) = \frac{1}{2} (g(\theta) \cdot \nabla) g(\theta) = \frac{1}{4} \nabla \|g(\theta)\|^2 \quad (9)$$

$$\tilde{\xi}_1(\theta) = \frac{1}{2} (\tilde{\xi}_0(\theta) \cdot \nabla) g(\theta) + \frac{1}{6} (g(\theta) \cdot \nabla) \tilde{\xi}_0. \quad (10)$$

Discussions. As can be inferred from Equations (8–10), $\tilde{\xi}_\alpha$ contains the $\alpha + 2_{\text{nd}}$ -order derivative of the loss function. Therefore, the higher-order counter terms cancel the higher-order smoothness of the discretization error.

Here, we note that Equation (8) can be found, e.g., in [35], as a higher-order backward error analysis. However, our derivation above has independent contributions: 1) we clarify that the counter term cancels the leading order of discretization error (Theorem 3.2), and 2) we find that the discretization error itself is also given by the counter term (Corollary 4.1 in the next section).

Equation (9) often appears in the literature on backward error analysis [21, 35] and its related topics in machine learning, e.g., [41, 23, 24, 27, 28, 31]. Typically, $\tilde{\xi}_0$ is added to gradients of continuous equations (e.g., SDE) to close the gap between continuous equations and discrete algorithms (e.g., SGD) by canceling (at least first-order) discretization error. However, higher-order discretization error is neglected in these studies. In contrast, our solution (8) cancels *all* orders of discretization error.

4 Discretization Error 1

The question here is to what extent the continuous approximation (1, 8) is precise; this point is often missed in the literature on continuous approximation [22, 23, 24, 11, 25, 26, 27, 28]. In this section, we use the counter term (8) and quantify discretization error as a function of the loss function and its derivatives (Section 4.1). We find that our result well explains empirical results. We further derive a sufficient condition for learning rates for the discretization error to be small (Section 4.2).

4.1 Counter Term Gives Leading Order of Discretization Error 3

We show that the counter term gives the leading order of discretization error between GD vs. GF and EoM. The proof follows from Theorem 3.2 and 3.3 and is given in Appendix A.4.

Corollary 4.1 (Leading order of discretization error is given by $\tilde{\xi}_\alpha$). *Suppose that we use ξ up to $O(\eta^{\gamma-1})$, i.e., $\xi = \tilde{\xi}_0 + \eta\tilde{\xi}_1 + \dots + \eta^{\gamma-1}\tilde{\xi}_{\gamma-1}$ for $\gamma \in \mathbb{Z}_{>0}$ ($\xi := \mathbf{0}$ for $\gamma = 0$). Then,*

$$e_{k+1} = e_k + \Lambda(\theta(k\eta)) + O(\eta^{\gamma+3}) = e_k + \eta^{\gamma+2}\tilde{\xi}_\gamma + O(\eta^{\gamma+3}). \quad (11)$$

First, Corollary 4.1 implies that the higher the orders of the counter term we use (large γ), the more precise EoM (1) is (small e_k). Thus, GF ($\xi = \mathbf{0}$) gives larger discretization error than EoM ($\xi \neq \mathbf{0}$). Second, Corollary 4.1 gives the *equality* of the leading order of discretization error at *arbitrary* steps. This is not a *bound* [18] nor an *asymptotic* analysis ($k \rightarrow \infty$). Third, let us give an intuition by considering $\xi = \mathbf{0}$ (GF). Then, Corollary 4.1 gives:

$$e_{k+1} = e_0 + \sum_{s=0}^k \frac{\eta^2}{2} (H(\theta(s\eta)) + \lambda I)(\nabla f(\theta(s\eta)) + \lambda \theta(s\eta)) + O(\eta^3), \quad (12)$$

where $H(\theta) \in \mathbb{R}^{d \times d}$ is the Hessian of the loss function f with respect to θ and $I \in \mathbb{R}^{d \times d}$ is the identity matrix. Equation (12) suggests that 1) large learning rates lead to a large discretization error and 2) steep loss functions (along the trajectory) lead to a large discretization error.

Empirical result. We find Equation (12) well explains our empirical result. We compare Equation (12) (up to $O(\eta^2)$) with the actual discretization error of GD and GF in Figure 2. First, the gap between our theoretical prediction of discretization error (orange curve) and the actual discretization error (red curve) is small because the range of *relative error* ($\|e_k\|/\|\theta_k\|$) in this plot is only 0–0.01 (see also Figure 11 in Appendix F). Second, most of the discretization error for Theory (orange curve) and Experiment (red curve) is produced within the first 100 steps. We can understand this phenomenon with the help of Equation (12). It suggests that discretization error can be enhanced when the loss function is non-smooth along the learning trajectory, which is likely to occur at the beginning of training due to random initialization. Therefore, a large part of discretization error is produced in the early stage of training. Third, we see that most of the gap between Theory (orange curve)

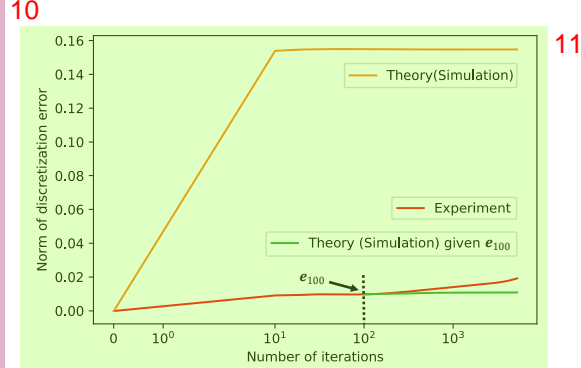


Figure 2: **Theoretical prediction of discretization error of GF and GD (Equation (12)) vs. actual discretization error of GF and GD.** The learning rate and weight decay are 10^{-2} and 10^{-2} . See Appendix F.2 for more results and details. See Section 6 for experimental settings.

and Experiment (red curve) also comes from the first 100 steps; in fact, the green curve shows that there is a much smaller enhancement of the gap after the 100th step. The source of the gap is the higher-order term $O(\eta^3)$ in Equation (12). It consists of higher-order derivatives of the loss function (Theorem 8 and Corollary 4.1) and thus can be large when the loss function is non-smooth along the learning trajectory. Therefore, by the same logic as above, the early stage of training tends to produce a gap between Theory (orange curve) and Experiment (red curve).

4.2 Discretization Error Bounds 1

We provide a sufficient condition (an upper bound for η) for GF and EoM to follow GD up to a given step k , which helps us infer desired learning rates (step sizes) for the discretization error to be small. We first consider $\xi = \mathbf{0}$ (GF). 2

Corollary 4.2 (Learning rate bound for $\xi = \mathbf{0}$). *Let $\xi = \mathbf{0}$ and assume that $e_0 = O(\eta^3)$. Let ϵ and t be arbitrary positive numbers. If the step size satisfies* 3

$$\eta < \sqrt{\frac{\epsilon}{k}} \sqrt{\frac{2}{\max_{0 \leq t' \leq t} \{ \| (H(\theta(t')) + \lambda I) g(\theta(t')) \| \}}}, \quad (13) \quad 4$$

for some $k \in \{1, 2, \dots, \lfloor \frac{t}{\eta} \rfloor\}$, then the discretization error can be arbitrarily small: 5

$$\|e_k\| < \epsilon + O(\epsilon^{\frac{3}{2}}). \quad (14) \quad 6$$

The proof follows from Equation (12) and is given in Appendix A.5. We see that 1) there is no guarantee that the discretization error is small unless the learning rate is sufficiently small, 2) we need small learning rates to keep the discretization error small for a long period, and 3) we need small learning rates to keep the discretization error small for non-smooth loss landscapes. This is consistent with our empirical results in Figure 3 and 4; in fact, 1) the discretization error blows up for a large learning rate ($\eta = 10^{-1}$ in Figure 3), 2) it increases as the number of steps increases (Figure 4), and 3) most of it is produced in the early phase of training, where the objective function tends to be non-smooth, and the gradients tend to be large. 7

We compare our bound (13) with a bound given in [18] because, to our knowledge, only [18] provides a bound for the step size with respect to discretization error in the context of deep learning. In [18], it is proved that in essence, $\eta \lesssim \epsilon / \beta_{t\epsilon} \gamma_{t\epsilon} c_t$, where $\beta_{t\epsilon}$ and $\gamma_{t\epsilon}$ measure the non-smoothness of the loss function, and c_t depends on the spectrum of the Hessian. These factors are hard to compute analytically unless the loss function and network are simple, but the qualitative behavior of this bound is the same as ours (13); i.e., both bounds become tight when the loss function is non-smooth. 8

We also derive a learning rate bound for $\xi = \tilde{\xi}_0$ (EoM) and the full statement is given in Corollary A.1 in Appendix A.6, which states that if $\eta < O(\sqrt[3]{\frac{\epsilon}{k}})$, then $\|e_k\| < \epsilon + O(\epsilon^{\frac{4}{3}})$. Therefore, larger step sizes are now allowed compared with Corollary 4.2 (GF) because of the non-zero counter term. Furthermore, we can show larger bounds for higher-order counter terms in a similar way. 9

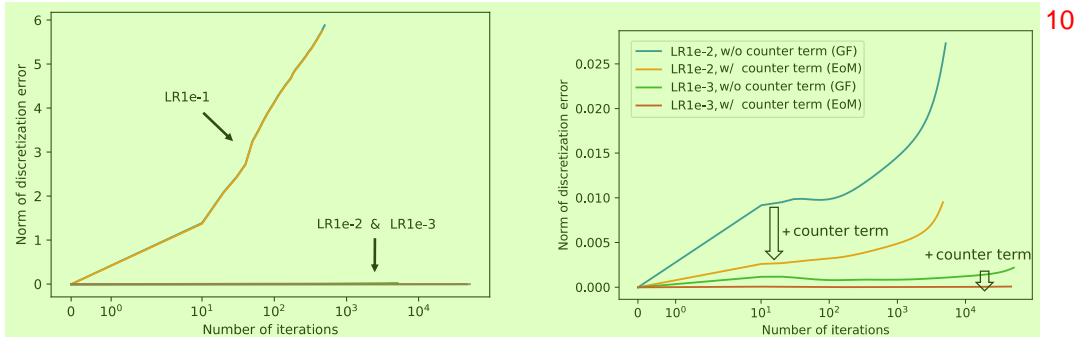


Figure 3: **Discretization error explodes for large learning rate** (10^{-1}). LR means learning rate. Weight decay is 10^{-3} . Curves include both GF and EoM. Relative discretization error is also shown in Appendix F. See Section 6 for experimental settings. 12

5 Application: Scale- and Translation-invariant Layers 13

To show the benefits of EoM, we finally apply our theory to two specific cases: scale-invariant layers [29, 30] and translation-invariant layers [31, 32]. Additionally, Appendix B provides an application 14

to broken conservation laws [31]. In the following, we simply focus on $\xi = \mathbf{0}$ and $\xi = \tilde{\xi}_0$ to analyze the differences between $\xi = \mathbf{0}$ and $\xi \neq \mathbf{0}$.

Definitions Let us first introduce our notation. A transformation ψ of $\theta \in \mathbb{R}^d$ with parameter $\alpha \in \mathbb{R}$ is said to be a *symmetry transformation* of loss function f if $f(\psi(\theta, \alpha)) = f(\theta)$. $\mathbb{1}_{\mathcal{A}} \in \{0, 1\}^d$ denotes the indicator vector of subspace $\mathcal{A} \subset \mathbb{R}^d$ (e.g., \mathcal{A} is a linear layer in the DNN). For a scalar $\alpha \in \mathbb{R}$, we define $\alpha_{\mathcal{A}} := \alpha \mathbb{1}_{\mathcal{A}} + \mathbb{1}_{\mathcal{A}^c} \in \mathbb{R}$, where \mathcal{A}^c is the complement of \mathcal{A} . For a vector $\theta \in \mathbb{R}^d$, we define $\theta_{\mathcal{A}} := \theta \odot \mathbb{1}_{\mathcal{A}} \in \mathbb{R}^d$, where \odot is the Hadamard element-wise product. For the gradient operator $\nabla = (\partial/\partial\theta_1, \dots, \partial/\partial\theta_d)^\top$, we define $\nabla_{\mathcal{A}} := \mathbb{1}_{\mathcal{A}} \odot \nabla$. We also define $r_{\mathcal{A}} := \|\theta_{\mathcal{A}}\|$ and $\hat{\theta}_{\mathcal{A}} := \theta_{\mathcal{A}}/r_{\mathcal{A}}$.

5.1 Learning Dynamics of Scale-invariant Layers

In this section, we focus on scale-invariant layers. A scale-invariant layer \mathcal{A} is defined as a subspace that is invariant under the scale transformation $\psi(\theta, \alpha) := \alpha_{\mathcal{A}}\theta = \alpha\theta_{\mathcal{A}} + \theta_{\mathcal{A}^c}$ ($\alpha > 0$). For example, a linear layer immediately before a batch normalization layer is scale-invariant. We see that for a better description of GD's discrete dynamics, we need modifications to the decay rate of $r_{\mathcal{A}}$ that is previously derived in the continuous regime [33]. In addition, we show that EoM successfully reproduces the limiting dynamics of $r_{\mathcal{A}}$ and *angular update* [34] at $t \rightarrow \infty$ that are previously derived in the discrete regime, while GF cannot. In Appendix C, we additionally show that there are crucial differences between GD and GF via the *effective learning rate* of scale-invariant layers [29, 42, 30, 43, 44, 33, 45, 34, 46, 47].

EoM for r We construct the EoM for $r_{\mathcal{A}}$ (the EoM for $\hat{\theta}_{\mathcal{A}}$ is given in Appendix C for completeness).

Theorem 5.1 (EoM for $r_{\mathcal{A}}$ and solution). *EoM (1) gives $\dot{r}_{\mathcal{A}}^2(t) = -2\lambda r_{\mathcal{A}}^2(t) - 2\eta \theta_{\mathcal{A}}(t) \cdot \xi(\theta(t))$. Specifically, this is equivalent to:*

$$\dot{r}_{\mathcal{A}}^2(t) = -2\lambda r_{\mathcal{A}}^2(t) \iff r_{\mathcal{A}}^2(t) = r_{\mathcal{A}}^2(0)e^{-2\lambda t} \quad (15)$$

for $\xi = \mathbf{0}$ (GF) and

$$\dot{r}_{\mathcal{A}}^2(t) = -2\left(\lambda + \frac{\eta\lambda^2}{2}\right)r_{\mathcal{A}}^2(t) + \frac{\eta}{r_{\mathcal{A}}^2(t)} \|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t) + \theta_{\mathcal{A}^c}(t))\|^2 \quad (16)$$

$$\iff r_{\mathcal{A}}^2(t) = r_{\mathcal{A}}^2(0)e^{-2\lambda(1+\frac{\eta\lambda}{2})t} + \eta \int_0^t d\tau e^{-2\lambda(1+\frac{\eta\lambda}{2})(t-\tau)} \frac{\|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(\tau) + \theta_{\mathcal{A}^c}(\tau))\|^2}{r_{\mathcal{A}}^2(\tau)} \quad (17)$$

for $\xi = \tilde{\xi}_0$ (EoM).

The proof is based on Equations (1, 9) and given in Appendix A.7. Equation (15) gives $r_{\mathcal{A}}^2(k\eta) = r_{\mathcal{A}}^2(0)e^{-2\eta\lambda k}$ ($k \in \mathbb{Z}_{\geq 0}$) at discretization; therefore, $\eta\lambda$ is regarded as the decay rate of $r_{\mathcal{A}}$ (*intrinsic learning rate* [33]). This is originally discussed in the continuous regime (SDE) [33]; however, we find that for a better description of the discrete dynamics of GD, the decay rate needs to be modified from $\eta\lambda$ to $\eta\lambda(1 + \frac{\eta\lambda}{2})$ (see the exponent of Equation (17)). This means that $r_{\mathcal{A}}$ in GD decays faster than expected from a naive continuous dynamics (GF (15) and SDE [33]). See Appendix G for higher-order corrections.

Limiting dynamics. We next derive the limiting dynamics ($t \rightarrow \infty$) of $r_{\mathcal{A}}$.

Corollary 5.1 ($r_{\mathcal{A}}$ at equilibrium). *When $\xi = \mathbf{0}$ (GF), $r_{\mathcal{A}}$ collapses to zero as $t \rightarrow \infty$. When $\xi = \tilde{\xi}_0$ (EoM), assume that there exist two constants $r_{\mathcal{A}*} \geq 0$ and $c_* \geq 0$ such that $r_{\mathcal{A}}(t) \xrightarrow{t \rightarrow \infty} r_{\mathcal{A}*}$ and $\|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t) + \theta_{\mathcal{A}^c}(t))\| \xrightarrow{t \rightarrow \infty} c_*$. Then $r_{\mathcal{A}*}^2 = \sqrt{\frac{\eta}{2\lambda + \eta\lambda^2} c_*}$.*

The proof follows from Theorem 5.1 and is given in Appendix A.8. The non-zero counter term successfully reproduces $r_{\mathcal{A}*}^2 \sim \sqrt{\eta/2\lambda} c_*$ [29, 34], which is originally derived in the discrete regime (SGD), although our approach is continuous (EoM (1)). Without the counter term, we cannot explain this behavior because GF gives $r_{\mathcal{A}}(t) \xrightarrow{t \rightarrow \infty} 0 (\neq \sqrt{\eta/2\lambda} c_*)$.

We next derive the limiting dynamics of *angular update* [34], which is designed to measure the temporal evolution of scale-invariant networks. It is originally defined in the discrete regime:

$\cos \Delta_k := \hat{\theta}_{\mathcal{A}k} \cdot \hat{\theta}_{\mathcal{A}k+1}$, where $\hat{\theta}_{\mathcal{A}k} := \frac{\mathbb{1}_{\mathcal{A}} \odot \theta_k}{\|\mathbb{1}_{\mathcal{A}} \odot \theta_k\|}$. That is, Δ_k represents a single-step angular change in the weight parameters of the scale-invariant layers \mathcal{A} . In the continuous regime, we can define $\cos \Delta(t) := \hat{\theta}_{\mathcal{A}}(t) \cdot \hat{\theta}_{\mathcal{A}}(t + \eta)$.

Corollary 5.2 ($\Delta(t)$ at equilibrium). *Let us use $\xi = \tilde{\xi}_0$. Suppose that the assumptions in Corollary 5.1 are satisfied. The angular update at equilibrium, denoted by Δ_* , is given by $\cos \Delta_* = \frac{1-\eta\lambda}{1-\eta^2\lambda^2/2} + O(\eta^3)$, and thus, $\Delta_* = \sqrt{2\eta\lambda} + O((\eta\lambda)^{3/2})$.*

The proof is based on Corollary 5.1 and is given in Appendix A.10. EoM successfully reproduces $\Delta_* \sim \sqrt{2\eta\lambda}$ [34], which is originally derived in the discrete regime (SGD), although EoM is continuous itself. On the other hand, GF cannot explain the limiting dynamics of $\Delta(t)$ because when $\xi = 0$, $r(t)$ goes to zero as $t \rightarrow \infty$ (Equation (15)), and thus, $\cos \Delta(t) = \frac{\theta_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t)} \cdot \frac{\theta_{\mathcal{A}}(t+\eta)}{r_{\mathcal{A}}(t+\eta)}$ is ill-defined. In summary, there are gaps between GF and GD, and our discussion above indicates that the counter term is inevitable to describe the actual dynamics of GD.

5.2 Learning Dynamics of Translation-invariant Layers

Next, we apply EoM to translation-invariant layers. To the best of our knowledge, no study analyzes the temporal evolution of translation-invariant layers except for [31] and [32], where only the sum of weights is their focus, while we derive the dynamics of the whole weights. A translation-invariant layer \mathcal{A} is defined as a layer that is invariant under the translation transformation $\psi(\theta, \alpha) := \theta + \alpha \mathbb{1}_{\mathcal{A}}$ ($\alpha \in \mathbb{R}$). For example, a linear layer immediately before the softmax layer is translation-invariant. In the following, we derive EoM and show that its theoretical prediction of decay rates dramatically matches empirical results, indicating the importance of the counter term. In Appendix D, we additionally discuss the differences between GF and GD in translation-invariant layers.

For convenience, we first decompose $\theta_{\mathcal{A}}$ to two vectors (Figure 5); $\theta_{\mathcal{A}\perp}$ is orthogonal to $\nabla f(\theta)$, and $\theta_{\mathcal{A}\parallel}$ is orthogonal to $\theta_{\mathcal{A}\perp}$. Here, note that $\nabla f(\theta)$ is orthogonal to $\mathbb{1}_{\mathcal{A}}$ because of translation invariance; in fact, differentiating both sides of $f(\theta + \alpha \mathbb{1}_{\mathcal{A}}) = f(\theta)$ with respect to α and setting $\alpha = 0$, we have $\mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta) = 0$ (see also Lemma A.7 in Appendix A.11). Formally, we define $\theta_{\mathcal{A}\perp}$, $\theta_{\mathcal{A}\parallel}$, and the projection matrix P as $\theta_{\mathcal{A}\perp} := P\theta_{\mathcal{A}} = \frac{\mathbb{1}_{\mathcal{A}} \cdot \theta_{\mathcal{A}}}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}}$, $\theta_{\mathcal{A}\parallel} := (I - P)\theta_{\mathcal{A}} = \theta_{\mathcal{A}} - \theta_{\mathcal{A}\perp}$, and $P := \frac{1}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}} \mathbb{1}_{\mathcal{A}}^{\top}$, where $d_{\mathcal{A}}$ is the dimension of \mathcal{A} .

We construct the EoM for $\theta_{\mathcal{A}\perp}$ (the EoM for $\theta_{\mathcal{A}\parallel}$ is given in Appendix D for completeness).

Theorem 5.2 (EoM for $\theta_{\mathcal{A}\perp}$). *EoM (1) gives $\dot{\theta}_{\mathcal{A}\perp}(t) = -\lambda \theta_{\mathcal{A}\perp}(t) - \eta P \xi(\theta(t))$. Specifically, this is equivalent to $\dot{\theta}_{\mathcal{A}\perp}(t) = -\lambda \theta_{\mathcal{A}\perp}(t) \iff \theta_{\mathcal{A}\perp}(t) = \theta_{\mathcal{A}\perp}(0) e^{-\lambda t}$ for $\xi = 0$ (GF) and $\dot{\theta}_{\mathcal{A}\perp}(t) = -(\lambda + \frac{\eta\lambda^2}{2}) \theta_{\mathcal{A}\perp}(t) \iff \theta_{\mathcal{A}\perp}(t) = \theta_{\mathcal{A}\perp}(0) e^{-(\lambda + \frac{\eta\lambda^2}{2})t}$ for $\xi = \tilde{\xi}_0$ (EoM).*

The proof is based on Equations (1, 9) and is given in Appendix A.11. $\theta_{\mathcal{A}\perp}$ monotonically collapses to zero as $t \rightarrow \infty$ in either case of $\xi = 0$ or $\xi \neq 0$; thus, as t increases, the dynamics is restricted onto the subspace orthogonal to $\theta_{\mathcal{A}\perp}$ (Figure 5). The decay rate is corrected by the counter term from $\eta\lambda$ to $\eta\lambda + \frac{\eta^2\lambda^2}{2}$, as is also done for $r_{\mathcal{A}}$ in Section 5.1. Therefore, the $\theta_{\mathcal{A}\perp}$ of GD decays faster than that of GF. Figure 6 and Table 1 support our findings. In particular, Table 1 shows that the decay rates predicted by EoM dramatically match those of GD, indicating the importance of the counter term.

Table 1: **Decay rates of $\|\theta_{\mathcal{A}\perp}\|$.** The theoretical predictions by EoM (third column) dramatically match experimental results of GD (fourth column) much better than GF (second column), indicating the importance of the counter term. LR and WD mean learning rate and weight decay, respectively. The colors correspond to those in Figure 6. See Section 6 for experimental settings.

(LR, WD)	Theory (GF)	Theory (EoM: Ours)	Experiment (GD)
$(10^{-1}, 10^{-2})$ (blue)	10^{-3}	1.0005×10^{-3}	$1.0005003484995967 \times 10^{-3}$
$(10^{-1}, 10^{-3})$ (orange)	10^{-4}	1.00005×10^{-4}	$1.0000500182363355 \times 10^{-4}$
$(10^{-2}, 10^{-2})$ (green)	10^{-4}	1.00005×10^{-4}	$1.0000499809795858 \times 10^{-4}$
$(10^{-2}, 10^{-3})$ (red)	10^{-5}	1.000005×10^{-5}	$1.0000049776814671 \times 10^{-5}$
$(10^{-3}, 10^{-2})$ (purple)	10^{-5}	1.000005×10^{-5}	$1.0000050475312426 \times 10^{-5}$
$(10^{-3}, 10^{-3})$ (yellow)	10^{-6}	1.0000005×10^{-6}	$1.0000005475009833 \times 10^{-6}$

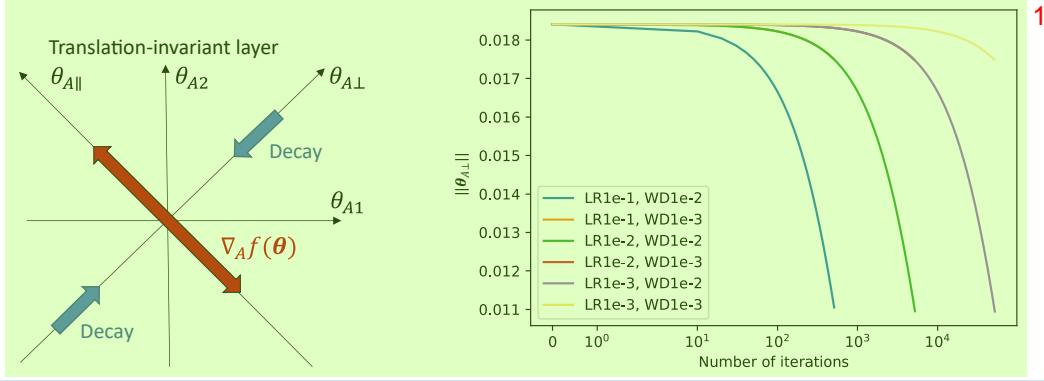


Figure 5: **Learning dynamics of translation-invariant layer.** Here, $\theta_{\mathcal{A}} = (\theta_{\mathcal{A}1}, \theta_{\mathcal{A}2})^\top$. $\theta_{\mathcal{A}\perp}$ decays to 0 (also shown in Figure 6). The decay of GD is faster than that of GF (Theorem 5.2). As t increases, the dynamics is restricted onto the subspace orthogonal to $\theta_{\mathcal{A}\perp}$. Figure 6: **Decay of $\|\theta_{\mathcal{A}\perp}\|$ (GD).** $\|\theta_{\mathcal{A}\perp}\|$ monotonically decays to zero, as suggested by Theorem 5.2. \mathcal{A} is translation-invariant layer. LR and WD mean learning rate and weight decay, respectively. Note that the orange and green curves (LR1e-1, WD1e-3 and LR1e-2, WD1e-2) and the red and purple curves (LR1e-2, WD1e-3 and LR1e-3, WD1e-2) totally overlap. The decay rates of all curves are given in Table 1. See Section 6 for experimental settings.

6 Experiment

We explain our experimental settings for Figures 2–6 and Table 1. Our network consists of a first linear layer, swish activation [48], second linear layer, batch normalization [49], third linear layer, and last softmax layer. Cross-entropy is used for the loss function. We note that the second linear layer is scale-invariant, and the last linear layer is translation invariant. The batch normalization uses fixed statistics to keep the scale invariance of the second linear layer. Swish is chosen to ensure differentiability. None of the linear layers have a bias term. The dataset is the training set of MNIST [50], and thus, the batch size is 60,000. Gradient descent is used for the optimizer. We use 64-bits of precision for all computations. To simulate GF and EoM, we use a sufficiently small learning rate (10^{-5}). The results are produced from only one random seed to save on computational costs, but we confirm that different random seeds lead to similar results. More detailed information is given in Appendix E and our code. In all experiments, we use $\xi = \xi_0$ for EoM. We do not include higher-order counter terms, such as ξ_1 , because they require third and higher order derivatives of the loss function and are thus extremely memory-consuming. We could circumvent this issue, e.g., by applying Hessian-free optimization [51], but this is out of our current scope.

7 Conclusion and Limitations

In this work, to fill the critical gap between GF and GD, we add a counter term to GF and obtain EoM, a continuous differential equation that precisely describes the discrete learning dynamics of GD. To show to what extent GF and EoM are precise in describing GD’s discrete dynamics, we derive the leading order of discretization error, as is often missed in the literature on the continuous approximation of discrete GD algorithms. We further derive a sufficient condition for learning rates for the discretization error to be small. We apply our theory to two specific cases, scale- and translation-invariant layers, indicating the importance of the counter term for a better description of the discrete learning dynamics of GD. Our experimental results support our theoretical findings.

Throughout this paper, we focus only on GD and GF to expose the ideas simply, and our study does not include stochasticity (e.g., SGD and SDE), acceleration methods (e.g., momentum and Nesterov [52]), or adaptive optimizers (e.g., Adam [53]). Nonetheless, they could be combined with our analysis, for example, using error analysis of SDEs [23, 24], continuous-time accelerated methods [7, 54, 9, 13, 14, 55, 16], and continuous-time Adam [56]. See Appendix G for more discussions. Therefore, our study could be extended to import continuous analysis to the discrete analysis of various GD algorithms. In this sense, our work bridges discrete and continuous analyses of GD algorithms.

Acknowledgment¹

We thank Shuhei M. Yoshida for his insightful comments on the dynamics of scale-invariant layers and the experimental settings. We also thank Hidenori Tanaka for his discussion that inspired us to start this study.

References³

- [1] Harold J. Kushner. Rates of convergence for sequential Monte Carlo optimization methods. *SIAM Journal on Control and Optimization*, 16(1):150–168, 1978.
- [2] Harold J. Kushner and Dean S. Clark. *Stochastic approximation methods for constrained and unconstrained systems*. Springer-Verlag New York, 1978.
- [3] Harold J. Kushner and Adam Schwartz. An invariant measure approach to the convergence of stochastic approximations with state dependent noise. *SIAM Journal on Control and Optimization*, 22(1):13–27, 1984.
- [4] L. Jung, G.C. Pflug, and H. Walk. *Stochastic Approximation and Optimization of Random Systems*. Oberwolfach Seminars. Birkhäuser Basel, 1992.
- [5] H. Kushner and G.G. Yin. *Stochastic Approximation and Recursive Algorithms and Applications*. Stochastic Modelling and Applied Probability. Springer New York, 2003.
- [6] Maxim Raginsky and Jake Brourie. Continuous-time stochastic mirror descent on a network: Variance reduction, consensus, convergence. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 6793–6800. IEEE, 2012.
- [7] Walid Krichene, Alexandre Bayen, and Peter L Bartlett. Accelerated mirror descent in continuous and discrete time. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [8] Panayotis Mertikopoulos and Mathias Staudigl. Convergence to nash equilibrium in continuous games with noisy first-order feedback. In *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*, pages 5609–5614. IEEE, 2017.
- [9] Walid Krichene and Peter L Bartlett. Acceleration and averaging in stochastic descent dynamics. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [10] Qiang Liu. Stein variational gradient descent as gradient flow. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [11] Yuanyuan Feng, Lei Li, and Jian-Guo Liu. Semigroups of stochastic gradient descent and online principal component analysis: properties and diffusion approximations. *Communications in Mathematical Sciences*, 16:777–789, 2017.
- [12] Damien Scieur, Vincent Roulet, Francis Bach, and Alexandre d'Aspremont. Integration methods and optimization algorithms. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [13] Pan Xu, Tianhao Wang, and Quanquan Gu. Accelerated stochastic mirror descent: From continuous-time dynamics to discrete-time algorithms. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1087–1096. PMLR, 09–11 Apr 2018.
- [14] Pan Xu, Tianhao Wang, and Quanquan Gu. Continuous and discrete-time accelerated stochastic mirror descent for strongly convex functions. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5492–5501. PMLR, 10–15 Jul 2018.
- [15] Alnur Ali, Edgar Dobriban, and Ryan J. Tibshirani. The implicit regularization of stochastic gradient flow for least squares. In *ICML*, pages 233–244, 2020.

- [16] Nikola B Kovachki and Andrew M Stuart. Continuous time analysis of momentum methods. *Journal of Machine Learning Research*, 22(17):1–40, 2021.
- [17] Stephan Wojtowytsch. Stochastic gradient descent with noise of machine learning type. Part II: Continuous time analysis. *arXiv preprint arXiv:2106.02588*, 2021.
- [18] Omer Elkabetz and Nadav Cohen. Continuous vs. discrete optimization of deep neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [19] Fanchen Bu and Dong Eui Chang. Feedback gradient descent: Efficient and stable optimization with orthogonality for DNNs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [20] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches zero loss? –a mathematical framework. In *International Conference on Learning Representations*, 2022.
- [21] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving Ordinary Differential Equations I (2nd Revised. Ed.): Nonstiff Problems*. Springer-Verlag, Berlin, Heidelberg, 1993.
- [22] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2101–2110. PMLR, 06–11 Aug 2017.
- [23] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of stochastic gradient algorithms I: Mathematical foundations. *Journal of Machine Learning Research*, 20(40):1–47, 2019.
- [24] Yuanyuan Feng, Tingran Gao, Lei Li, Jian-Guo Liu, and Yulong Lu. Uniform-in-time weak error analysis for stochastic gradient descent algorithms via diffusion approximation. *Communications in Mathematical Sciences*, 18(1):163–188, 2020.
- [25] Wenqing Hu, Chris Junchi Li, Lei Li, and Jian-Guo Liu. On the diffusion approximation of nonconvex stochastic gradient descent. *Annals of Mathematical Sciences and Applications*, 2019.
- [26] Jing An, Jianfeng Lu, and Lexing Ying. Stochastic modified equations for the asynchronous stochastic gradient descent. *Information and Inference: A Journal of the IMA*, 9(4):851–873, 11 2019.
- [27] David Barrett and Benoit Dherin. Implicit gradient regularization. In *International Conference on Learning Representations*, 2021.
- [28] Samuel L Smith, Benoit Dherin, David Barrett, and Soham De. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations*, 2021.
- [29] Twan Van Laarhoven. L2 regularization versus batch and weight normalization. *arXiv preprint arXiv:1706.05350*, 2017.
- [30] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse. Three mechanisms of weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [31] Daniel Kunin, Javier Sagastuy-Brena, Surya Ganguli, Daniel LK Yamins, and Hidenori Tanaka. Neural mechanics: Symmetry and broken conservation laws in deep learning dynamics. In *International Conference on Learning Representations*, 2021.
- [32] Hidenori Tanaka and Daniel Kunin. Noether’s learning dynamics: Role of symmetry breaking in neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [33] Zhiyuan Li, Kaifeng Lyu, and Sanjeev Arora. Reconciling modern deep learning with traditional optimization analyses: The intrinsic learning rate. In *NeurIPS*, 2020.
- [34] Ruosi Wan, Zhanxing Zhu, Xiangyu Zhang, and Jian Sun. Spherical motion dynamics: Learning dynamics of neural network with normalization, weight decay, and SGD, 2021.
- [35] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric Numerical Integration. Structure-Preserving Algorithms for Ordinary Differential Equations*. Springer, Berlin, 2nd ed. edition, 2006. ID: unige:12343.

- [36] Ioan A Rus. *On the problem of Darboux-Ionescu*. Universitatea Babes-Bolyai. Faculty of Mathematics, 1981.
- [37] Nicolaie Lungu and Ioan A Rus. On a functional volterra-fredholm integral equation, via picard operators. *J. math. ineq*, 3(4):519–527, 2009.
- [38] Nguyen Thanh Long et al. On a nonlinear volterra-hammerstein integral equation in two variables. *Acta Mathematica Scientia*, 33(2):484–494, 2013.
- [39] Tran Minh Thuyet, Nguyen Thanh Long, et al. A nonlinear volterra-hammerstein integral equation in three variables. *Nonlinear Functional Analysis and Applications*, 19(2):193–211, 2014.
- [40] Daniela Marian, Sorina Anamaria Ciplea, and Nicolaie Lungu. On a functional integral equation. *Symmetry*, 13(8):1321, 2021.
- [41] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2101–2110. PMLR, 06–11 Aug 2017.
- [42] Elad Hoffer, Ron Banner, Itay Golan, and Daniel Soudry. Norm matters: efficient and accurate normalization schemes in deep networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [43] Sanjeev Arora, Zhiyuan Li, and Kaifeng Lyu. Theoretical analysis of auto rate-tuning by batch normalization. In *International Conference on Learning Representations*, 2019.
- [44] Vitaliy Chiley, Ilya Sharapov, Atli Kosson, Urs Koster, Ryan Reece, Sofia Samaniego de la Fuente, Vishal Subbiah, and Michael James. Online normalization for training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [45] Zhiyuan Li and Sanjeev Arora. An exponential learning rate schedule for deep learning. In *International Conference on Learning Representations*, 2020.
- [46] Zhiyuan Li, Srinadh Bhojanapalli, Manzil Zaheer, Sashank J Reddi, and Sanjiv Kumar. Robust training of neural networks using scale invariant architectures. *arXiv preprint arXiv:2202.00980*, 2022.
- [47] Simon Roburin, Yann de Mont-Marin, Andrei Bursuc, Renaud Marlet, Patrick Pérez, and Mathieu Aubry. Spherical perspective on learning with normalization layers. *Neurocomputing*, 487:66–74, 2022.
- [48] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.
- [49] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [50] Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010. License: Creative Commons Attribution-Share Alike 3.0 license.
- [51] James Martens et al. Deep learning via hessian-free optimization. In *ICML*, volume 27, pages 735–742, 2010.
- [52] Nesterov Y. E. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. *Dokl. Akad. Nauk SSSR*, 269:543–547, 1983.
- [53] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [54] Weijie Su, Stephen Boyd, and Emmanuel J. Candès. A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights. *Journal of Machine Learning Research*, 17(153):1–43, 2016.
- [55] Jean-Francois Aujol, Charles Dossal, and Aude Rondepierre. Optimal convergence rates for Nesterov acceleration. *SIAM Journal on Optimization*, 29(4):3131–3153, 2019.
- [56] Anas Barakat and Pascal Bianchi. Convergence and dynamical behavior of the adam algorithm for nonconvex stochastic optimization. *SIAM Journal on Optimization*, 31(1):244–274, 2021.

- [57] Emmy Noether. Invariante Variationsprobleme. *Nachr. d. König. Gesellsch. d. Wiss. zu Göttingen, Math-phys. Klasse*, Seite 235-157, 1918. 1
- [58] Emmy Noether. Invariant variation problems. *Transport theory and statistical physics*, 1(3):186–207, 1971. 2
- [59] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. License: Apache License 2.0. Software available from tensorflow.org. 3
- [60] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. Del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, 09 2020. License: BSD 3-Clause "New" or "Revised" License. 4
- [61] Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009. 5
- [62] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 6
- [63] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pages 630–645, 2016. 7
- [64] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. IEEE, 2009. 8
- [65] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. 2014. 9
- [66] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 10
- [67] Jian Deng. Strong backward error analysis for Euler-Maruyama method. *Int. J. Numer. Anal. Model.*, 13:1–21, 2016. 11

Checklist 12

1. For all authors... 13

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] 14
- (b) Did you describe the limitations of your work? [Yes] See Sections 6 and 7 and Appendix G.
- (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results... 15

- (a) Did you state the full set of assumptions of all theoretical results? [Yes] 16
- (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.

3. If you ran experiments... 17

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See the code in the supplemental material. 1
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Section 6 and Appendix E.
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#) To save computational costs, we do not run experiments with multiple random seeds, but we confirm that different random seeds give similar results, as stated in Section 6.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Appendix E.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets... 2

- (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See our code. 3
- (b) Did you mention the license of the assets? [\[Yes\]](#) See our code.
- (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#) See our code.
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[N/A\]](#) We do not use such data.
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#) The data we are using do not include such information.

5. If you used crowdsourcing or conducted research with human subjects... 4

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#) 5
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

Appendices 1

A Proofs 2

A.1 Proof of Theorem 3.1 3

Proof. The integral form of Taylor's theorem gives 4

$$\begin{aligned}\theta(k\eta + \eta) - \theta(k\eta) &= \eta \dot{\theta}(k\eta) + \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) \\ &= -\eta g(\theta(k\eta)) - \eta^2 \xi(\theta(k\eta)) + \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s).\end{aligned}\quad (18)$$

Remember the definition of the discrete gradient descent: 6

$$\theta_{k+1} - \theta_k = -\eta g(\theta_k). \quad (19)$$

Subtracting Equation (19) from Equation (18), we have 8

$$\begin{aligned}e_{k+1} - e_k &= -\eta(g(\theta(k\eta)) - g(\theta_k)) - \eta^2 \xi(\theta(k\eta)) + \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) \\ &= -\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) - \eta^2 \xi(\theta(k\eta)) + \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s).\end{aligned}\quad (20)$$

(21) \square

A.2 Proof of Theorem 3.2 10

Proof. The proof is by induction. For $k = 0$, $e_0 = O(\eta^\gamma)$ by assumption. If $e_k = O(\eta^\gamma)$ for $k \geq 1$, Theorem 3.1 gives 11

$$e_{k+1} = e_k - \eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) + \Lambda(\theta(k\eta)) = O(\eta^\gamma) + O(\eta^{\gamma+1}) + O(\eta^\gamma) = O(\eta^\gamma). \quad (22)$$

$\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) = O(\eta^{\gamma+1})$ follows from Taylor's expansion of $g(\theta(k\eta) - e_k)$ around $\theta(k\eta)$ and from assumption $e_k = O(\eta^\gamma)$: 12

$$-\eta(g(\theta(k\eta)) - g(\theta(k\eta) - e_k)) = \eta(e_k \cdot \nabla g(\theta(k\eta)) + O(\|e_k\|^2)) = O(\eta^{\gamma+1}). \quad (23)$$

\square

A.3 Proof of Theorem 3.3 15

Proof. The proof of Theorem 3.3 consists of the following three Lemmas, all of which are proved in the following sections. 16

Lemma A.1. 17

$$\int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) = \sum_{n=0}^{\infty} \frac{\eta^n}{(n+2)!} \frac{d^{n+2}}{dt^{n+2}} \theta(k\eta) \quad (24)$$

Lemma A.2. For $n \geq 1$, 19

$$\frac{d^n}{dt^n} \theta(t) = (-1)^n \sum_{k_1, \dots, k_n=0}^{\infty} \eta^{k_1 + \dots + k_n} \mathcal{D}_{k_1} \dots \mathcal{D}_{k_{n-1}} \Xi_{k_n}, \quad (25)$$

where $\mathcal{D}_{k_1} \dots \mathcal{D}_{k_{n-1}} := 1$ for $n = 1$. 21

Lemma A.3. 22

$$\int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) = \sum_{j=0}^{\infty} \sum_{i=2}^{j+2} \sum_{k_1 + \dots + k_i = j-i+2} \frac{(-1)^i}{i!} \eta^j \mathcal{D}_{k_1} \dots \mathcal{D}_{k_{i-1}} \Xi_{k_i} \quad (26)$$

Theorem 3.3 follows by comparing both sides of Equation (6) order-by-order with using Equation (26) and the expansion of ξ (7). \square 24

A.3.1 Proof of Lemma A.1 ¹

Proof. ²

$$\begin{aligned} & \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) \\ &= \frac{1}{\eta^2} \int_{k\eta}^{k\eta+\eta} ds \ddot{\theta}(s)(k\eta + \eta - s) \end{aligned} \quad (27)$$

$$= \frac{1}{\eta^2} \int_0^\eta ds' [\ddot{\theta}(k\eta)(\eta - s') + \ddot{\theta}(k\eta)(\eta - s')s' + \frac{1}{2!} \ddot{\theta}(k\eta)(\eta - s')s'^2 + \dots] \quad (28)$$

$$= \sum_{n=0}^{\infty} \frac{\eta^n}{(n+2)!} \frac{d^{n+2}}{dt^{n+2}} \theta(k\eta) \quad (29)$$

From Line (27) to (28), we used $s' := s - k\eta$ and the Taylor expansion of $\ddot{\theta}(k\eta + s')$ around $k\eta$.⁴
 From Line (28) to (29), we used $\int_0^\eta ds' (\eta - s')s'^n = \frac{\eta^{n+2}}{(n+1)(n+2)}$ for $n \geq 0$. \square

A.3.2 Proof of Lemma A.2 ⁵

Proof. Note that given $\dot{\theta}(t) = -g(\theta(t)) - \eta\xi(\theta(t))$, we have ⁶

$$\frac{d}{dt} \left(\frac{d^{n-1}}{dt^{n-1}} \theta(t) \right) = -\mathcal{D} \left(\frac{d^{n-1}}{dt^{n-1}} \theta(t) \right) \quad (n \geq 1), \quad (30)$$

where $d^0\theta/dt^0 := \theta$. Therefore, ⁸

$$\frac{d^n}{dt^n} \theta(t) = (-1)^{n-1} \mathcal{D}^{n-1}(-g - \eta\xi) = (-1)^n \mathcal{D}^{n-1} \Xi \quad (n \geq 1). \quad (31)$$

Thus, by definition of \mathcal{D} , \mathcal{D}_α , and Ξ_α (Theorem 3.3 in Section 3.3), we have ¹⁰

$$\frac{d^n}{dt^n} \theta(t) = (-1)^n \left(\sum_{k_1=0}^{\infty} \eta^{k_1} \mathcal{D}_{k_1} \right) \cdots \left(\sum_{k_{n-1}=0}^{\infty} \eta^{k_{n-1}} \mathcal{D}_{k_{n-1}} \right) \Xi \quad (32)$$

$$= (-1)^n \sum_{k_1, \dots, k_n=0}^{\infty} \eta^{k_1 + \dots + k_n} \mathcal{D}_{k_1} \cdots \mathcal{D}_{k_{n-1}} \Xi_{k_n}. \quad (33)$$

\square

A.3.3 Proof of Lemma A.3 ¹²

Proof. From Lemma A.1 and A.2, we have ¹³

$$\begin{aligned} & \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) \\ &= \sum_{n=0}^{\infty} \frac{\eta^n}{(n+2)!} \frac{d^{n+2}}{dt^{n+2}} \theta(k\eta) \end{aligned} \quad (34)$$

$$= \sum_{n=0}^{\infty} \frac{\eta^n}{(n+2)!} (-1)^{n+2} \sum_{k_1, \dots, k_{n+2}=0}^{\infty} \eta^{k_1 + \dots + k_{n+2}} \mathcal{D}_{k_1} \cdots \mathcal{D}_{k_{n+1}} \Xi_{k_{n+2}} \quad (35)$$

$$= \sum_{n=0}^{\infty} \sum_{k_1, \dots, k_{n+2}=0}^{\infty} \frac{(-1)^n}{(n+2)!} \eta^{n+k_1 + \dots + k_{n+2}} \mathcal{D}_{k_1} \cdots \mathcal{D}_{k_{n+1}} \Xi_{k_{n+2}} \quad (36)$$

$$= \sum_{j=0}^{\infty} \sum_{i=2}^{j+2} \sum_{k_1 + \dots + k_i = j-i+2} \frac{(-1)^{i-2}}{i!} \eta^j \mathcal{D}_{k_1} \cdots \mathcal{D}_{k_{i-1}} \Xi_{k_i}. \quad (37)$$

On the last line, we replaced $n+2$ and $n+k_1 + \dots + k_{n+2}$ with i and j , respectively. ¹⁵ \square

A.4 Proof of Corollary 4.1 ¹

Proof. By assumption, we use ²

$$\xi(\theta) = \eta^2 \sum_{\alpha=0}^{\gamma-1} \eta^\alpha \tilde{\xi}_\alpha. \quad (38) \quad \text{3}$$

From Theorem 3.3, we have ⁴

$$\Lambda(\theta) = \eta^2 \int_0^1 ds \ddot{\theta}(\eta(k+s))(1-s) - \eta^2 \xi(\theta(k\eta)) \quad (39) \quad \text{5}$$

$$= \eta^2 \sum_{\alpha=0}^{\infty} \eta^\alpha \tilde{\xi}_\alpha - \eta^2 \sum_{\alpha=0}^{\gamma-1} \eta^\alpha \tilde{\xi}_\alpha \quad (40)$$

$$= \eta^2 \sum_{\alpha=\gamma}^{\infty} \eta^\alpha \tilde{\xi}_\alpha \quad (41)$$

$$= \eta^{\gamma+2} \tilde{\xi}_\gamma + O(\eta^{\gamma+3}). \quad (42)$$

Therefore, Theorem 3.2 gives ⁶

$$e_{k+1} = e_k + \Lambda(\theta(k\eta)) + O(\eta^{\gamma+3}) \quad (43) \quad \text{7}$$

$$= e_k + \eta^{\gamma+2} \tilde{\xi}_\gamma + O(\eta^{\gamma+3}) + O(\eta^{\gamma+3}) \quad (44)$$

$$= e_k + \eta^{\gamma+2} \tilde{\xi}_\gamma + O(\eta^{\gamma+3}). \quad (45)$$

□

A.5 Proof of Corollary 4.2 ⁸

Proof. From Equation (12), we have ⁹

$$e_k = e_0 + \sum_{s=0}^{k-1} \frac{\eta^2}{2} (H(\theta(s\eta)) + \lambda I) g(\theta(s\eta)) + O(\eta^3). \quad (46) \quad \text{10}$$

Because $e_0 = O(\eta^3)$ by assumption, we have ¹¹

$$e_k = \sum_{s=0}^{k-1} \frac{\eta^2}{2} (H(\theta(s\eta)) + \lambda I) g(\theta(s\eta)) + O(\eta^3) \quad (47) \quad \text{12}$$

$$\therefore \|e_k\| \leq \frac{\eta^2}{2} \sum_{s=0}^{k-1} \|(H(\theta(s\eta)) + \lambda I) g(\theta(s\eta))\| + O(\eta^3) \quad (48)$$

$$\leq \frac{\eta^2 k}{2} \max_{0 \leq s \leq k-1} \{\|(H(\theta(s\eta)) + \lambda I) g(\theta(s\eta))\|\} + O(\eta^3). \quad (49)$$

Let $t > 0$ be a given arbitrary number. Then, for $k \in \{1, 2, \dots, \lfloor \frac{t}{\eta} \rfloor\}$, ¹³

$$\|e_k\| \leq \frac{\eta^2 k}{2} \max_{0 \leq t' \leq t} \{\|(H(\theta(t')) + \lambda I) g(\theta(t'))\|\} + O(\eta^3). \quad (50) \quad \text{14}$$

Therefore, if $\eta < \sqrt{\epsilon/k} \sqrt{2/\max_{0 \leq t' \leq t} \{\|(H(\theta(t')) + \lambda I) g(\theta(t'))\|\}}$, then ¹⁵

$$\|e_k\| < \epsilon + O(\epsilon^{3/2}). \quad (51) \quad \text{16}$$

□

A.6 Proof of Corollary A.1 ¹

Corollary A.1 (Learning rate bound when $\xi = \tilde{\xi}_0$). *Let $\xi = \tilde{\xi}_0$ and assume that $e_0 = O(\eta^4)$. Let ϵ ² and t be arbitrary positive numbers. If the step size satisfies*

$$\eta < \sqrt[3]{\frac{\epsilon}{k}} \sqrt[3]{\frac{12}{\max_{0 \leq t' \leq t} \{ \|4(H(\theta(t')) + \lambda I)^2 \mathbf{g}(\theta(t')) + \mathbf{g}(\theta(t'))^\top \nabla H(\theta(t')) \mathbf{g}(t')\| \}}}, \quad (52) \quad 3$$

for some $k \in \{1, 2, \dots, \lfloor \frac{t}{\eta} \rfloor\}$, then the discretization error can be arbitrarily small: ⁴

$$\|e_k\| < \epsilon + O(\epsilon^{\frac{4}{3}}). \quad (53) \quad 5$$

Proof. From Equation (10) and Corollary 4.1 and by assumption, we have ⁶

$$e_k = e_0 + \eta^3 \sum_{s=0}^{k-1} \left\{ \frac{1}{2} (\tilde{\xi}_0(\theta(s\eta)) \cdot \nabla) \mathbf{g}(\theta(s\eta)) + \frac{1}{6} (\mathbf{g}(\theta(s\eta)) \cdot \nabla) \tilde{\xi}_0(\theta(s\eta)) \right\} + O(\eta^4). \quad (54) \quad 7$$

Because $e_0 = O(\eta^4)$ by assumption, we have ⁸

$$\begin{aligned} e_k &= \eta^3 \sum_{s=0}^{k-1} \left\{ \frac{1}{2} (\tilde{\xi}_0(\theta(s\eta)) \cdot \nabla) \mathbf{g}(\theta(s\eta)) + \frac{1}{6} (\mathbf{g}(\theta(s\eta)) \cdot \nabla) \tilde{\xi}_0(\theta(s\eta)) \right\} + O(\eta^4) \quad (55) \quad 9 \\ &= \eta^3 \sum_{s=0}^{k-1} \left\{ \frac{1}{3} (H(\theta(s\eta)) + \lambda I)^2 \mathbf{g}(\theta(s\eta)) + \frac{1}{12} \mathbf{g}^\top(\theta(s\eta)) \nabla H(\theta(s\eta)) \mathbf{g}(\theta(s\eta)) \right\} + O(\eta^4). \end{aligned} \quad (56)$$

Therefore, ¹⁰

$$\|e_k\| \leq \eta^3 \sum_{s=0}^{k-1} \left\| \frac{1}{3} (H(\theta(s\eta)) + \lambda I)^2 \mathbf{g}(\theta(s\eta)) + \frac{1}{12} \mathbf{g}^\top(\theta(s\eta)) \nabla H(\theta(s\eta)) \mathbf{g}(\theta(s\eta)) \right\| + O(\eta^4) \quad (57) \quad 11$$

$$\begin{aligned} &\leq \frac{\eta^3 k}{12} \max_{0 \leq s \leq k-1} \{ \|4(H(\theta(s\eta)) + \lambda I)^2 \mathbf{g}(\theta(s\eta)) + \mathbf{g}^\top(\theta(s\eta)) \nabla H(\theta(s\eta)) \mathbf{g}(\theta(s\eta))\| \} \\ &\quad + O(\eta^4). \end{aligned} \quad (58)$$

Let $t > 0$ be a given arbitrary number. Then, for $k \in \{1, 2, \dots, \lfloor \frac{t}{\eta} \rfloor\}$, ¹²

$$\|e_k\| \leq \frac{\eta^3 k}{12} \max_{0 \leq t' \leq t} \{ \|4(H(\theta(t')) + \lambda I)^2 \mathbf{g}(\theta(t')) + \mathbf{g}^\top(\theta(t')) \nabla H(\theta(t')) \mathbf{g}(\theta(t'))\| \} + O(\eta^4). \quad (59) \quad 13$$

Therefore, if ¹⁴

$$\eta < \sqrt[3]{\frac{\epsilon}{k}} \sqrt[3]{\frac{12}{\max_{0 \leq t' \leq t} \{ \|4(H(\theta(t')) + \lambda I)^2 \mathbf{g}(\theta(t')) + \mathbf{g}(\theta(t'))^\top \nabla H(\theta(t')) \mathbf{g}(t')\| \}}}, \quad (60) \quad 15$$

then $\|e_k\| < \epsilon + O(\epsilon^{4/3})$. ¹⁶ □

A.7 Proof of Theorem 5.1 ¹⁷

We use the following Lemmas. ¹⁸

Lemma A.4. *For scale-invariant layers \mathcal{A} , the following equations hold:*

$$\theta_{\mathcal{A}} \cdot \nabla f(\theta) = \theta_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\theta) = 0 \quad (61) \quad 19$$

$$H_{\mathcal{A}}(\theta) \theta_{\mathcal{A}} + \nabla_{\mathcal{A}} f(\theta) = 0 \quad (62)$$

$$\nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^\top f(\theta) \theta_{\mathcal{A}} = 0, \quad (63)$$

where $H_{\mathcal{A}}(\theta) := (\mathbf{1}_{\mathcal{A}} \odot \nabla)(\mathbf{1}_{\mathcal{A}} \odot \nabla)^\top f(\theta)$. ²⁰

Proof. Differentiating both sides of $f(\alpha_{\mathcal{A}} \odot \theta) = f(\theta)$ with respect to α , we have ¹

$$\theta_{\mathcal{A}} \cdot \nabla f(\theta) = \theta_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\alpha_{\mathcal{A}} \odot \theta) = 0, \quad 2 \quad (64)$$

where $\nabla_{\mathcal{A}} f(\alpha_{\mathcal{A}} \odot \theta)$ means $(\nabla_{\mathcal{A}} f(\theta))|_{\theta=\alpha_{\mathcal{A}} \odot \theta}$. For $\alpha = 1$, we have ³

$$\theta_{\mathcal{A}} \cdot \nabla f(\theta) = \theta_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\theta) = 0. \quad 4 \quad (65)$$

Applying ∇ , we have ⁵

$$(\theta_{\mathcal{A}} \cdot \nabla_{\mathcal{A}}) \nabla f(\theta) + \nabla_{\mathcal{A}} f(\theta) = 0 \quad 6 \quad (66)$$

$$\iff (\theta_{\mathcal{A}} \cdot \nabla_{\mathcal{A}}) (\nabla_{\mathcal{A}} + \nabla_{\mathcal{A}^c}) f(\theta) + \nabla_{\mathcal{A}} f(\theta) = 0 \quad (67)$$

$$\iff H_{\mathcal{A}}(\theta) \theta_{\mathcal{A}} + \nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^{\top} f(\theta) \theta_{\mathcal{A}} + \nabla_{\mathcal{A}} f(\theta) = 0. \quad (68)$$

Multiplying by $\mathbf{1}_{\mathcal{A}^c} \odot$, we have ⁷

$$\nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^{\top} f(\theta) \theta_{\mathcal{A}} = 0. \quad 8 \quad (69)$$

Therefore, ⁹

$$H_{\mathcal{A}}(\theta) \theta_{\mathcal{A}} + \nabla_{\mathcal{A}} f(\theta) = 0. \quad 10 \quad (70)$$

□

Lemma A.5. For scale-invariant layers \mathcal{A} , the following equations hold: ¹¹

$$\nabla_{\mathcal{A}} f(\theta) = \frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}), \quad 12 \quad (71)$$

where $\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}) := (\nabla_{\mathcal{A}} f(\theta))|_{\theta=\hat{\theta}_{\mathcal{A}}+\theta_{\mathcal{A}^c}}$. ¹³

Proof. Note that $f(\theta) = f(\alpha_{\mathcal{A}} \odot \theta) = f(\alpha \theta_{\mathcal{A}} + \theta_{\mathcal{A}^c})$. Differentiating both sides with respect to θ , ¹⁴ we have

$$\nabla f(\theta) \quad 15 \quad (72)$$

$$= \nabla(f(\alpha_{\mathcal{A}} \odot \theta)) \quad (73)$$

$$= (\nabla_{\mathcal{A}} + \nabla_{\mathcal{A}^c})(f(\alpha_{\mathcal{A}} \odot \theta)) \quad (74)$$

$$= \alpha \nabla_{\mathcal{A}} f(\alpha_{\mathcal{A}} \odot \theta) + \nabla_{\mathcal{A}^c} f(\alpha_{\mathcal{A}} \odot \theta). \quad (75)$$

For $\alpha = 1/r_{\mathcal{A}}$, we have ¹⁶

$$\nabla f(\theta) = \frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}) + \nabla_{\mathcal{A}^c} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}). \quad 17 \quad (76)$$

Therefore, ¹⁸

$$\nabla_{\mathcal{A}} f(\theta) = \mathbf{1}_{\mathcal{A}} \odot \nabla f(\theta) = \mathbf{1}_{\mathcal{A}} \odot \left(\frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}) + \nabla_{\mathcal{A}^c} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}) \right) \quad 19 \quad (77)$$

$$= \frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}). \quad (78)$$

□

Lemma A.6. For scale-invariant layers \mathcal{A} , the following equations hold for all $\alpha > 0$: ²⁰

$$H(\theta) = \alpha^2 H_{\mathcal{A}}(\alpha_{\mathcal{A}} \odot \theta) + \alpha (\nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^{\top} f(\alpha_{\mathcal{A}} \odot \theta) + \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\alpha_{\mathcal{A}} \odot \theta)) + H_{\mathcal{A}^c}(\alpha_{\mathcal{A}} \odot \theta) \quad 21 \quad (79)$$

$$H(\theta) \theta_{\mathcal{A}} = \alpha^2 H_{\mathcal{A}}(\alpha_{\mathcal{A}} \odot \theta) \theta_{\mathcal{A}} \quad (80)$$

$$H(\theta) \theta_{\mathcal{A}} = H_{\mathcal{A}}(\theta) \theta_{\mathcal{A}}, \quad (81)$$

where $H_{\mathcal{A}}(\alpha_{\mathcal{A}} \odot \theta) := ((\mathbf{1}_{\mathcal{A}} \odot \nabla)(\mathbf{1}_{\mathcal{A}} \odot \nabla)^{\top} f(\theta))|_{\theta=\alpha_{\mathcal{A}} \odot \theta}$. ²²

Proof. Because $\nabla f(\boldsymbol{\theta}) = \alpha \nabla_{\mathcal{A}} f(\alpha \boldsymbol{\theta}_{\mathcal{A}}) + \nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}_{\mathcal{A}^c})$ (Equation 75), ¹

$$H(\boldsymbol{\theta}) = \nabla \nabla^\top f(\boldsymbol{\theta}) \quad (82)$$

$$= \nabla(\alpha \nabla_{\mathcal{A}}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta}) + \nabla_{\mathcal{A}^c}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})) \quad (83)$$

$$= (\nabla_{\mathcal{A}} + \nabla_{\mathcal{A}^c})(\alpha \nabla_{\mathcal{A}}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta}) + \nabla_{\mathcal{A}^c}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})) \quad (84)$$

$$= \alpha^2 H_{\mathcal{A}}(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta}) + \alpha (\nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta}) + \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})) + H_{\mathcal{A}^c}(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta}). \quad (85)$$

Therefore, ³

$$H(\boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}} = \alpha^2 H_{\mathcal{A}}(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}} + \alpha \nabla_{\mathcal{A}^c} \nabla_{\mathcal{A}}^\top f(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}} \quad (86)$$

$$= \alpha^2 H_{\mathcal{A}}(\alpha \boldsymbol{\theta}_{\mathcal{A}} \odot \boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}}. \quad (87)$$

For $\alpha = 1$, we have ⁵

$$H(\boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}} = H_{\mathcal{A}}(\boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}}. \quad (88)$$

□

We now prove Theorem 5.1. ⁷

Proof. We use Lemmas A.4, A.5, and A.6. ⁸

$$\dot{r}_{\mathcal{A}}^2(t) = 2\boldsymbol{\theta}_{\mathcal{A}}(t) \cdot \dot{\boldsymbol{\theta}}_{\mathcal{A}}(t) \quad (89)$$

$$= 2\boldsymbol{\theta}_{\mathcal{A}}(t) \cdot (-\nabla_{\mathcal{A}} f(\boldsymbol{\theta}(t)) - \lambda \boldsymbol{\theta}_{\mathcal{A}}(t) - \eta \boldsymbol{\xi}(\boldsymbol{\theta}(t))) \quad (90)$$

$$= -2\lambda r_{\mathcal{A}}^2(t) - 2\eta \boldsymbol{\theta}_{\mathcal{A}}(t) \cdot \boldsymbol{\xi}(\boldsymbol{\theta}(t)). \quad (91)$$

For $\boldsymbol{\xi} = \mathbf{0}$, ¹⁰

$$\dot{r}_{\mathcal{A}}^2(t) = -2\lambda r_{\mathcal{A}}^2(t). \quad (92)$$

For $\boldsymbol{\xi} = \tilde{\boldsymbol{\xi}}_0$, ¹²

$$\dot{r}_{\mathcal{A}}^2(t) = -2\lambda r_{\mathcal{A}}^2(t) - 2\eta \boldsymbol{\theta}_{\mathcal{A}}(t) \cdot \tilde{\boldsymbol{\xi}}_0(\boldsymbol{\theta}(t)) \quad (93)$$

$$= -2\lambda r_{\mathcal{A}}^2(t) - \eta(\lambda^2 r_{\mathcal{A}}^2(t) - \|\nabla_{\mathcal{A}} f(\boldsymbol{\theta}(t))\|^2) \quad (94)$$

$$= -2\lambda(1 + \frac{\eta\lambda}{2})r_{\mathcal{A}}^2(t) + \frac{\eta}{r_{\mathcal{A}}^2(t)}\|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}}(t))\|^2. \quad (95)$$

We used ¹⁴

$$\boldsymbol{\theta}_{\mathcal{A}} \cdot \tilde{\boldsymbol{\xi}}_{0,\mathcal{A}} = \frac{1}{2}\boldsymbol{\theta}_{\mathcal{A}} \cdot (H(\boldsymbol{\theta}) + \lambda I)(\nabla f(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}) \quad (96)$$

$$= \frac{1}{2}\boldsymbol{\theta}_{\mathcal{A}} \cdot (H(\boldsymbol{\theta})\nabla f(\boldsymbol{\theta}) + \lambda H(\boldsymbol{\theta})\boldsymbol{\theta} + \lambda \nabla f(\boldsymbol{\theta}) + \lambda^2 \boldsymbol{\theta}) \quad (97)$$

$$= \frac{1}{2}(\boldsymbol{\theta}_{\mathcal{A}}^\top H_{\mathcal{A}}(\boldsymbol{\theta})\nabla_{\mathcal{A}} f(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}_{\mathcal{A}}^\top H_{\mathcal{A}}(\boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}} + \lambda^2 r_{\mathcal{A}}^2) \quad (98)$$

$$= \frac{1}{2}(-\|\nabla_{\mathcal{A}} f(\boldsymbol{\theta})\|^2 + \lambda^2 r_{\mathcal{A}}^2). \quad (99)$$

Using $\dot{\boldsymbol{x}}(t) = -a\boldsymbol{x} + \boldsymbol{y}(t) \Leftrightarrow \boldsymbol{x}(t) = \boldsymbol{x}(0)e^{-at} + \int_0^t d\tau e^{-a(t-\tau)}\boldsymbol{y}(\tau)$, we can show the remaining equations. ¹⁶

□

A.8 Proof of Corollary 5.1 ¹⁷

Proof. When $\boldsymbol{\xi} = \mathbf{0}$, $r_{\mathcal{A}} \xrightarrow{t \rightarrow \infty} 0$ is obvious from the EoM for $r_{\mathcal{A}}$ (Theorem 5.1). When $\boldsymbol{\xi} = \tilde{\boldsymbol{\xi}}_0$, ¹⁸ EoM is given by

$$\dot{r}_{\mathcal{A}}^2(t) = -2\lambda(1 + \frac{\eta\lambda}{2})r_{\mathcal{A}}^2(t) + \frac{\eta}{r_{\mathcal{A}}^2(t)}\|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}}(t) + \boldsymbol{\theta}_{\mathcal{A}^c}(t))\|^2. \quad (100)$$

At equilibrium, $\dot{r}_{\mathcal{A}} = 0$ and $\|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\| = c_*$ by assumption; thus, we have ¹

$$0 = -2\lambda(1 + \frac{\eta\lambda}{2})r_{\mathcal{A}^*}^2 + \frac{\eta}{r_{\mathcal{A}^*}^2}c_*^2 \quad (101)$$

$$\iff r_{\mathcal{A}^*}^2 = \sqrt{\frac{\eta}{2\lambda + \eta\lambda^2}}c_* \quad (102)$$

□

A.9 Proof of Theorem C.1 ³

Proof. We use Lemmas A.4, A.5, and A.6: ⁴

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}} = \frac{d}{dt} \frac{\boldsymbol{\theta}_{\mathcal{A}}}{r_{\mathcal{A}}} \quad (103)$$

$$= -\frac{\dot{r}_{\mathcal{A}}}{r_{\mathcal{A}}^2} \boldsymbol{\theta}_{\mathcal{A}} + \frac{1}{r_{\mathcal{A}}} \dot{\boldsymbol{\theta}}_{\mathcal{A}} \quad (104)$$

$$= \frac{\boldsymbol{\theta}_{\mathcal{A}}}{r_{\mathcal{A}}^2} (\lambda r_{\mathcal{A}} + \eta \hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \boldsymbol{\xi}(\boldsymbol{\theta})) + \frac{1}{r_{\mathcal{A}}} (-\nabla_{\mathcal{A}} f(\boldsymbol{\theta}) - \lambda \boldsymbol{\theta}_{\mathcal{A}} - \eta \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta})) \quad (105)$$

$$= \frac{\eta}{r_{\mathcal{A}}} \hat{\boldsymbol{\theta}}_{\mathcal{A}} (\hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta})) - \frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} f(\boldsymbol{\theta}) - \frac{\eta}{r_{\mathcal{A}}} \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta}) \quad (106)$$

$$= -\frac{1}{r_{\mathcal{A}}^2} \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) + \frac{\eta}{r_{\mathcal{A}}} ((\hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta})) \hat{\boldsymbol{\theta}}_{\mathcal{A}} - \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta})), \quad (107)$$

where $\boldsymbol{\xi}_{\mathcal{A}} := \mathbb{1}_{\mathcal{A}} \odot \boldsymbol{\xi}$. We used $\dot{r}_{\mathcal{A}} = -\lambda r_{\mathcal{A}} - \eta \hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \boldsymbol{\xi}(\boldsymbol{\theta})$ (Theorem 5.1). Note that $\frac{\eta}{r_{\mathcal{A}}} ((\hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta})) \hat{\boldsymbol{\theta}}_{\mathcal{A}} - \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta}))$ has no $\hat{\boldsymbol{\theta}}_{\mathcal{A}}$ component; i.e., it is orthogonal to $\hat{\boldsymbol{\theta}}_{\mathcal{A}}$. When $\boldsymbol{\xi} = \mathbf{0}$, Equation (107) is equivalent to $\dot{\boldsymbol{\theta}}_{\mathcal{A}} = -\frac{1}{r_{\mathcal{A}}^2} \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}})$. When $\boldsymbol{\xi} = \tilde{\boldsymbol{\xi}}_0$, note that from Equation (99), ⁶

$$\boldsymbol{\theta}_{\mathcal{A}} \cdot \tilde{\boldsymbol{\xi}}_{0\mathcal{A}} = \frac{1}{2} (-\|\nabla_{\mathcal{A}} f(\boldsymbol{\theta})\|^2 + \lambda^2 r_{\mathcal{A}}^2). \quad (108)$$

Therefore, ⁸

$$(\hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \tilde{\boldsymbol{\xi}}_{0\mathcal{A}}) \hat{\boldsymbol{\theta}}_{\mathcal{A}} = -\frac{1}{2} \frac{1}{r_{\mathcal{A}}^3} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}} + \frac{\lambda^2}{2} \boldsymbol{\theta}_{\mathcal{A}}. \quad (109)$$

Also, ¹⁰

$$\tilde{\boldsymbol{\xi}}_{0\mathcal{A}} = \frac{1}{2} \mathbb{1}_{\mathcal{A}} \odot (H(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) + \lambda H(\boldsymbol{\theta}) \boldsymbol{\theta} + \lambda \nabla f(\boldsymbol{\theta}) + \lambda^2 \boldsymbol{\theta}) \quad (110)$$

$$= \frac{1}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) + \lambda \mathbb{1}_{\mathcal{A}} \odot H(\boldsymbol{\theta}) (\boldsymbol{\theta}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) + \lambda \nabla_{\mathcal{A}} f(\boldsymbol{\theta}) + \lambda^2 \boldsymbol{\theta}_{\mathcal{A}}) \quad (111)$$

$$= \frac{1}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) + \lambda H_{\mathcal{A}}(\boldsymbol{\theta}) \boldsymbol{\theta}_{\mathcal{A}} + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \boldsymbol{\theta}_{\mathcal{A}^c} + \lambda \nabla_{\mathcal{A}} f(\boldsymbol{\theta}) + \lambda^2 \boldsymbol{\theta}_{\mathcal{A}}) \quad (112)$$

$$= \frac{1}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \boldsymbol{\theta}_{\mathcal{A}^c} + \lambda^2 \boldsymbol{\theta}_{\mathcal{A}}). \quad (113)$$

Therefore, ¹

$$(\hat{\boldsymbol{\theta}}_{\mathcal{A}} \cdot \tilde{\boldsymbol{\xi}}_{0,\mathcal{A}}(\boldsymbol{\theta}))\hat{\boldsymbol{\theta}}_{\mathcal{A}} - \tilde{\boldsymbol{\xi}}_{0,\mathcal{A}} \quad 2 \quad (114)$$

$$= -\frac{1}{2} \frac{1}{r_{\mathcal{A}}^3} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}} + \frac{\lambda^2}{2} \boldsymbol{\theta}_{\mathcal{A}} - \frac{1}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \boldsymbol{\theta}_{\mathcal{A}^c} + \lambda^2 \boldsymbol{\theta}_{\mathcal{A}}) \quad (115)$$

$$= -\frac{1}{2} \frac{1}{r_{\mathcal{A}}^3} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}} - \frac{1}{2} \nabla_{\mathcal{A}} \nabla^{\top} f(\boldsymbol{\theta}) \nabla f(\boldsymbol{\theta}) - \frac{\lambda}{2} \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \boldsymbol{\theta}_{\mathcal{A}^c} \quad (116)$$

$$= -\frac{1}{2} \frac{1}{r_{\mathcal{A}}^3} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}} - \frac{1}{2} H_{\mathcal{A}}(\boldsymbol{\theta}) \nabla_{\mathcal{A}} f(\boldsymbol{\theta}) - \frac{1}{2} \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}) - \frac{1}{2} \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\boldsymbol{\theta}) \lambda \boldsymbol{\theta}_{\mathcal{A}^c} \quad (117)$$

$$= -\frac{1}{2} \frac{1}{r_{\mathcal{A}}^3} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}} - \frac{1}{2} \frac{1}{r_{\mathcal{A}}} H_{\mathcal{A}} \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) - \frac{1}{2} \frac{1}{r_{\mathcal{A}}} \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) (\nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}_{\mathcal{A}^c}). \quad (118)$$

Hence, ³

$$\begin{aligned} \dot{\hat{\boldsymbol{\theta}}}_{\mathcal{A}} &= -\frac{1}{r_{\mathcal{A}}^2} \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) - \frac{\eta}{2r_{\mathcal{A}}^2} (H_{\mathcal{A}}(\boldsymbol{\theta}) \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) \\ &\quad + \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) (\nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}) \lambda \boldsymbol{\theta}_{\mathcal{A}^c}) + \frac{1}{r_{\mathcal{A}}^2} \|\nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})\|^2 \hat{\boldsymbol{\theta}}_{\mathcal{A}}) \end{aligned} \quad 4 \quad (119)$$

$$= -\frac{1}{r_{\mathcal{A}}^2} (I + \frac{\eta}{2} H_{\mathcal{A}}(\boldsymbol{\theta}) + \frac{\eta}{2} (\nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}_{\mathcal{A}^c}) \cdot \nabla_{\mathcal{A}^c} + \frac{\eta}{2} \frac{1}{r_{\mathcal{A}}^2} \hat{\boldsymbol{\theta}}_{\mathcal{A}} \nabla_{\mathcal{A}}^{\top} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c})) \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}) \quad (120)$$

$$= -\frac{1}{r_{\mathcal{A}}^2} (I + \frac{\eta}{2} H_{\mathcal{A}}(\boldsymbol{\theta}) + \frac{\eta}{2} (\nabla_{\mathcal{A}^c} f(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}_{\mathcal{A}^c}) \cdot \nabla_{\mathcal{A}^c} + \frac{\eta}{2} \hat{\boldsymbol{\theta}}_{\mathcal{A}} \nabla_{\mathcal{A}}^{\top} f(\boldsymbol{\theta})) \nabla_{\mathcal{A}} f(\hat{\boldsymbol{\theta}}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c}). \quad (121)$$

□

A.10 Proof of Corollary 5.2 ⁵

Proof. We use Lemmas A.4 and A.5. The angular update is defined as ⁶

$$\cos \Delta(t) = \frac{\boldsymbol{\theta}_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t)} \cdot \frac{\boldsymbol{\theta}_{\mathcal{A}}(t + \eta)}{r_{\mathcal{A}}(t + \eta)}. \quad 7 \quad (122)$$

We evaluate the higher order terms in $\boldsymbol{\theta}_{\mathcal{A}}(t + \eta)$ and $r_{\mathcal{A}}(t + \eta)$. First, ⁸

$$\begin{aligned} \boldsymbol{\theta}_{\mathcal{A}}(t + \eta) &= \boldsymbol{\theta}_{\mathcal{A}}(t) + \eta \dot{\boldsymbol{\theta}}_{\mathcal{A}}(t) + \frac{\eta^2}{2} \ddot{\boldsymbol{\theta}}_{\mathcal{A}}(t) + O(\eta^3) \\ &= \boldsymbol{\theta}_{\mathcal{A}}(t) - \eta \nabla f(\boldsymbol{\theta}_{\mathcal{A}}(t)) - \eta \lambda \boldsymbol{\theta}_{\mathcal{A}}(t) - \eta^2 \boldsymbol{\xi}_{\mathcal{A}}(\boldsymbol{\theta}(t)) + \frac{\eta^2}{2} \ddot{\boldsymbol{\theta}}_{\mathcal{A}}(t) + O(\eta^3). \end{aligned} \quad 9 \quad (123)$$

The second derivative $\ddot{\theta}(t)$ is given by ¹

$$\ddot{\theta}_{\mathcal{A}} = \frac{d}{dt} \dot{\theta}_{\mathcal{A}} \quad (124)$$

$$= \frac{d}{dt} (-\nabla_{\mathcal{A}} f(\theta) - \lambda \theta_{\mathcal{A}}) + O(\eta) \quad (125)$$

$$= -(\dot{\theta} \cdot \nabla) \nabla_{\mathcal{A}} f(\theta) - \lambda \dot{\theta}_{\mathcal{A}} + O(\eta) \quad (126)$$

$$= \nabla_{\mathcal{A}} \nabla^{\top} f(\theta) (\nabla f(\theta) + \lambda \theta) + \lambda (\nabla_{\mathcal{A}} f(\theta) + \lambda \theta_{\mathcal{A}}) + O(\eta) \quad (127)$$

$$= \mathbb{1}_{\mathcal{A}} \odot H(\theta) \nabla f(\theta) + \lambda H_{\mathcal{A}}(\theta) \theta_{\mathcal{A}} + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\theta) \theta_{\mathcal{A}^c} + \lambda \nabla_{\mathcal{A}} f(\theta) + \lambda^2 \theta_{\mathcal{A}} + O(\eta) \quad (128)$$

$$= \mathbb{1}_{\mathcal{A}} \odot H(\theta) \nabla f(\theta) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\theta) \theta_{\mathcal{A}^c} + \lambda^2 \theta_{\mathcal{A}} + O(\eta). \quad (129)$$

Therefore, ³

$$\begin{aligned} \theta_{\mathcal{A}}(t + \eta) &= \theta_{\mathcal{A}}(t) - \eta \nabla_{\mathcal{A}} f(\theta(t)) - \eta \lambda \theta_{\mathcal{A}}(t) - \eta^2 \xi_{\mathcal{A}}(\theta(t)) \\ &\quad + \frac{\eta^2}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\theta(t)) \nabla f(\theta(t)) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\theta(t)) \theta_{\mathcal{A}^c}(t) + \lambda^2 \theta_{\mathcal{A}}(t)) + O(\eta^3). \end{aligned} \quad (130)$$

Next, ⁵

$$r_{\mathcal{A}}(t + \eta) = r_{\mathcal{A}}(t) + \dot{r}_{\mathcal{A}}(t) \eta + \frac{\eta^2}{2} \ddot{r}_{\mathcal{A}}(t) + O(\eta^3). \quad (131)$$

Because $\dot{r}_{\mathcal{A}} = -\lambda r_{\mathcal{A}} - \eta \hat{\theta}_{\mathcal{A}} \cdot \xi$ (use Equation (91) and $\dot{r}_{\mathcal{A}} = 2r_{\mathcal{A}} \dot{r}_{\mathcal{A}}$), ⁷

$$r_{\mathcal{A}}(t + \eta) = r_{\mathcal{A}}(t) - \eta \lambda r_{\mathcal{A}}(t) - \eta^2 \hat{\theta}_{\mathcal{A}}(t) \cdot \xi_{\mathcal{A}}(\theta(t)) + \frac{\eta^2}{2} \ddot{r}_{\mathcal{A}}(t) + O(\eta^3). \quad (132)$$

In addition, because $\ddot{r}_{\mathcal{A}} = -\lambda \dot{r}_{\mathcal{A}} + O(\eta) = \lambda^2 r_{\mathcal{A}} + O(\eta)$, ⁹

$$r_{\mathcal{A}}(t + \eta) = r_{\mathcal{A}}(t) - \eta \lambda r_{\mathcal{A}}(t) - \eta^2 \hat{\theta}_{\mathcal{A}}(t) \cdot \xi(\theta(t)) + \frac{\eta^2}{2} \lambda^2 r_{\mathcal{A}}(t) + O(\eta^3). \quad (133)$$

Therefore, ¹¹

$$\begin{aligned} \cos \Delta(t) &= \frac{\theta_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t)} \cdot \frac{\theta_{\mathcal{A}}(t + \eta)}{r_{\mathcal{A}}(t + \eta)} \\ &= \frac{\theta_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t)} \cdot \left(\theta_{\mathcal{A}}(t) - \eta \nabla_{\mathcal{A}} f(\theta(t)) - \eta \lambda \theta_{\mathcal{A}}(t) - \eta^2 \xi_{\mathcal{A}}(\theta(t)) \right. \\ &\quad \left. + \frac{\eta^2}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\theta(t)) \nabla f(\theta(t)) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\theta(t)) \theta_{\mathcal{A}^c}(t) + \lambda^2 \theta_{\mathcal{A}}(t)) \right) / \left(r_{\mathcal{A}}(t) - \eta \lambda r_{\mathcal{A}}(t) \right. \\ &\quad \left. - \eta^2 \hat{\theta}_{\mathcal{A}}(t) \cdot \xi(\theta(t)) + \frac{\eta^2}{2} \lambda^2 r_{\mathcal{A}}(t) \right) \\ &\quad + O(\eta^3). \end{aligned} \quad (134)$$

Substituting $\xi_{\mathcal{A}} = \tilde{\xi}_{0,\mathcal{A}}$, and using ¹³

$$\tilde{\xi}_{0,\mathcal{A}} = \frac{1}{2} (\mathbb{1}_{\mathcal{A}} \odot H(\theta) \nabla f(\theta) + \lambda \nabla_{\mathcal{A}} \nabla_{\mathcal{A}^c}^{\top} f(\theta) \theta_{\mathcal{A}^c} + \lambda^2 \theta_{\mathcal{A}}) \quad (\text{Equation 113}) \quad (136)$$

$$\theta_{\mathcal{A}} \cdot \tilde{\xi}_{0,\mathcal{A}} = \frac{1}{2} (-\|\nabla_{\mathcal{A}} f(\theta)\|^2 + \lambda^2 r_{\mathcal{A}}^2) \quad (\text{Equation 108}), \quad (137)$$

we have ¹⁵

$$\cos \Delta(t) = \frac{\theta_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t)} \cdot \frac{\theta_{\mathcal{A}}(t) - \eta \nabla_{\mathcal{A}} f(\theta(t)) - \eta \lambda \theta_{\mathcal{A}}(t)}{r_{\mathcal{A}}(t) - \eta \lambda r_{\mathcal{A}}(t) - \frac{\eta^2}{2r_{\mathcal{A}}(t)} (-\|\nabla_{\mathcal{A}} f(\theta)\|^2 + \lambda^2 r_{\mathcal{A}}^2(t)) + \frac{\eta^2}{2} \lambda^2 r_{\mathcal{A}}(t)} + O(\eta^3) \quad (138)$$

$$= \frac{(1 - \eta \lambda) r_{\mathcal{A}}^2(t)}{(1 - \eta \lambda) r_{\mathcal{A}}^2(t) + \frac{\eta^2}{2r_{\mathcal{A}}(t)} \|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c})\|^2} + O(\eta^3). \quad (139)$$

At equilibrium, we have $r_{\mathcal{A}}^2 \xrightarrow{t \rightarrow \infty} r_{\mathcal{A}^*}^2 = \sqrt{\frac{\eta}{2\lambda + \eta\lambda^2}} c_*$ and $\|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t) + \theta_{\mathcal{A}^c}(t))\| \xrightarrow{t \rightarrow \infty} c_*$ because of Corollary 5.1. Thus, ¹

$$\cos \Delta_* = \frac{(1 - \eta\lambda)r_{\mathcal{A}^*}^2}{(1 - \eta\lambda)r_{\mathcal{A}^*}^2 + \frac{\eta^2}{2r_{\mathcal{A}^*}^2} c_*^2} + O(\eta^3) \quad (140)$$

$$= \frac{1 - \eta\lambda}{1 - \eta^2\lambda^2/2} + O(\eta^3), \quad (141)$$

and we have shown the first statement of the theorem. ³

The second statement follows from Equation (141). By definition of cosine and tangent, we have ⁴

$$\tan \Delta_* = \frac{\sqrt{(1 - \eta^2\lambda^2/2)^2 - (1 - \eta\lambda)^2}}{1 - \eta\lambda} + O(\eta^3) = \frac{\sqrt{2\eta\lambda - 2\eta^2\lambda^2 + \eta^4\lambda^4/4}}{1 - \eta} + O(\eta^3). \quad (142)$$

Therefore, using Taylor's series of the tangent function, we have ⁶

$$\Delta_* = \tan \Delta_* - \frac{1}{3}\Delta_*^3 - \frac{2}{15}\Delta_*^5 - \dots = \sqrt{2\eta\lambda} + O((\eta\lambda)^{3/2}). \quad (143)$$

This concludes the proof. ⁸ □

A.11 Proof of Theorem 5.2 ⁹

We use the following Lemma: ¹⁰

Lemma A.7. *For translation-invariant layers \mathcal{A} , the following equations hold:* ¹¹

$$\theta_{\mathcal{A}\perp} \cdot \theta_{\mathcal{A}\parallel} = 0 \quad (144)$$

$$\mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta) = \mathbb{1}_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\theta) = 0 \quad (145)$$

$$P\nabla f(\theta) = P\nabla_{\mathcal{A}} f(\theta) = 0 \quad (146)$$

$$\theta_{\mathcal{A}\perp} \cdot \nabla f(\theta) = \theta_{\mathcal{A}\perp} \cdot \nabla_{\mathcal{A}} f(\theta) = 0 \quad (147)$$

$$H(\theta)\mathbb{1}_{\mathcal{A}} = 0 \quad (148)$$

$$PH(\theta) = 0 \quad (149)$$

$$H(\theta)\theta_{\mathcal{A}\perp} = 0 \quad (150)$$

$$\nabla f(\theta) = \nabla f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) \quad (151)$$

$$H(\theta) = H(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}). \quad (152)$$

Proof. Note that $P^\top = P$, $P^2 = P$, and thus, $P^\top(I - P) = P(I - P) = P - P = 0$. Therefore, ¹³

$$\theta_{\mathcal{A}\perp} \cdot \theta_{\mathcal{A}\parallel} = \theta_{\mathcal{A}}^\top P^\top (I - P) \theta_{\mathcal{A}} = 0. \quad (153)$$

Next, differentiating $f(\theta) = f(\theta + \alpha \mathbb{1}_{\mathcal{A}})$ with respect to α , we have ¹⁵

$$\mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta + \alpha \mathbb{1}_{\mathcal{A}}) = 0. \quad (154)$$

For $\alpha = 0$, we have ¹⁷

$$\mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta) = \mathbb{1}_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\theta) = 0. \quad (155)$$

Therefore, ¹⁹

$$P\nabla f(\theta) = P\nabla_{\mathcal{A}} f(\theta) = (\mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta)) \frac{1}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}} = 0 \quad (156)$$

and ²¹

$$\theta_{\mathcal{A}\perp} \cdot \nabla f(\theta) = \frac{\mathbb{1}_{\mathcal{A}} \cdot \theta_{\mathcal{A}}}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}} \cdot \nabla f(\theta) = \frac{\mathbb{1}_{\mathcal{A}} \cdot \theta_{\mathcal{A}}}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f(\theta) = 0. \quad (157)$$

Next, differentiating Equation 155 with respect to θ , we have ²³

$$H(\theta)\mathbb{1}_{\mathcal{A}} = 0. \quad (158)$$

Therefore, ¹

$$PH(\boldsymbol{\theta}) = \frac{\mathbb{1}_{\mathcal{A}}}{d_{\mathcal{A}}} \mathbb{1}_{\mathcal{A}}^{\top} H(\boldsymbol{\theta}) = 0 \quad ^2 \quad (159)$$

and ³

$$H(\boldsymbol{\theta})\boldsymbol{\theta}_{\mathcal{A}\perp} = \frac{\mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}}{d_{\mathcal{A}}} H(\boldsymbol{\theta})\mathbb{1}_{\mathcal{A}} = 0. \quad ^4 \quad (160)$$

Next, differentiating $f(\boldsymbol{\theta}) = f(\boldsymbol{\theta} + \alpha \mathbb{1}_{\mathcal{A}})$ with respect to $\boldsymbol{\theta}$, we have ⁵

$$\nabla f(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta} + \alpha \mathbb{1}_{\mathcal{A}}) \quad ^6 \quad (161)$$

and ⁷

$$H(\boldsymbol{\theta}) = H(\boldsymbol{\theta} + \alpha \mathbb{1}_{\mathcal{A}}). \quad ^8 \quad (162)$$

For $\alpha = -\frac{\mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}}{d_{\mathcal{A}}}$, we have ⁹

$$\nabla f(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta} - P\boldsymbol{\theta}_{\mathcal{A}}) = \nabla f(\boldsymbol{\theta}_{\mathcal{A}} + \boldsymbol{\theta}_{\mathcal{A}^c} - P\boldsymbol{\theta}_{\mathcal{A}}) = \nabla f(\boldsymbol{\theta}_{\mathcal{A}\parallel} + \boldsymbol{\theta}_{\mathcal{A}^c}) \quad ^{10} \quad (163)$$

and ¹¹

$$H(\boldsymbol{\theta}) = H(\boldsymbol{\theta}_{\mathcal{A}\parallel} + \boldsymbol{\theta}_{\mathcal{A}^c}). \quad ^{12} \quad (164)$$

□

We begin the proof of Theorem 5.2. ¹³

Proof. We use Lemma A.7. ¹⁴

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}\perp} = P\dot{\boldsymbol{\theta}}_{\mathcal{A}} = P(-\nabla_{\mathcal{A}}f(\boldsymbol{\theta}) - \lambda\boldsymbol{\theta}_{\mathcal{A}} - \eta\xi_{\mathcal{A}}) = -\lambda\boldsymbol{\theta}_{\mathcal{A}\perp} - \eta P\xi_{\mathcal{A}}. \quad ^{15} \quad (165)$$

When $\boldsymbol{\theta} = \mathbf{0}$, EoM is ¹⁶

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}\perp}(t) = -\lambda\boldsymbol{\theta}_{\mathcal{A}}(t). \quad ^{17} \quad (166)$$

When $\xi = \tilde{\xi}_{0,\mathcal{A}}$, note that ¹⁸

$$\tilde{\xi}_0 = \frac{1}{2}(H(\boldsymbol{\theta})\nabla f(\boldsymbol{\theta}) + \lambda\nabla f(\boldsymbol{\theta}) + \lambda H(\boldsymbol{\theta})\boldsymbol{\theta} + \lambda^2\boldsymbol{\theta}) \quad ^{19} \quad (167)$$

and ²⁰

$$\tilde{\xi}_0 \cdot \mathbb{1}_{\mathcal{A}} = \tilde{\xi}_{0,\mathcal{A}} \cdot \mathbb{1}_{\mathcal{A}} = \frac{\lambda^2}{2} \mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}. \quad ^{21} \quad (168)$$

Thus, ²²

$$P\tilde{\xi}_0 = \frac{\lambda^2}{2} \boldsymbol{\theta}_{\mathcal{A}\perp}. \quad ^{23} \quad (169)$$

Therefore, ²⁴

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}\perp} = -\lambda\boldsymbol{\theta}_{\mathcal{A}\perp} - \eta P\tilde{\xi}_{0,\mathcal{A}} = -\lambda\boldsymbol{\theta}_{\mathcal{A}\perp} - \eta \frac{\lambda^2}{2} \boldsymbol{\theta}_{\mathcal{A}\perp} = -(\lambda + \frac{\eta\lambda^2}{2})\boldsymbol{\theta}_{\mathcal{A}\perp}. \quad ^{25} \quad (170)$$

Using $\dot{\boldsymbol{v}}(t) = -a\boldsymbol{v}(t) \Leftrightarrow \boldsymbol{v}(t) = \boldsymbol{v}(0)e^{-at}$, we can show the remaining equations. ²⁶ □

A.12 Proof of Theorem D.1 ²⁷

Proof. We use Lemma A.7. First, note that ²⁸

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}\parallel} = \dot{\boldsymbol{\theta}}_{\mathcal{A}} - \dot{\boldsymbol{\theta}}_{\mathcal{A}\perp}. \quad ^{29} \quad (171)$$

Because ³⁰

$$\dot{\boldsymbol{\theta}}_{\mathcal{A}} = -\nabla_{\mathcal{A}}f(\boldsymbol{\theta}) - \lambda\boldsymbol{\theta}_{\mathcal{A}} - \eta\xi_{\mathcal{A}} \quad ^{31} \quad (172)$$

and ¹

$$\dot{\theta}_{\mathcal{A}\perp} = -\lambda\theta_{\mathcal{A}\perp} - \eta P\xi_{\mathcal{A}}, \quad (173)$$

we have ³

$$\dot{\theta}_{\mathcal{A}\parallel} = \dot{\theta}_{\mathcal{A}} - \dot{\theta}_{\mathcal{A}\perp} = -\nabla_{\mathcal{A}}f(\theta) - \lambda\theta_{\mathcal{A}\parallel} - \eta(I - P)\xi_{\mathcal{A}} \quad (174)$$

$$= -\nabla_{\mathcal{A}}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) - \lambda\theta_{\mathcal{A}\parallel} - \eta(I - P)\xi_{\mathcal{A}}. \quad (175)$$

Note that $\dot{\theta}_{\mathcal{A}\parallel}$ is orthogonal to $\theta_{\mathcal{A}\parallel}$ because $\theta_{\mathcal{A}\perp} \cdot \dot{\theta}_{\mathcal{A}\parallel} = -\theta_{\mathcal{A}\perp} \cdot \nabla_{\mathcal{A}}f(\theta) - \lambda\theta_{\mathcal{A}\perp} \cdot \theta_{\mathcal{A}\parallel} - \eta\theta_{\mathcal{A}\perp}^{\top}(I - P)\xi_{0,\mathcal{A}} = 0 - 0 - 0 = 0$ (we used $\theta_{\mathcal{A}\perp}^{\top}(I - P) = \theta_{\mathcal{A}}^{\top}P^{\top}(I - P) = \theta_{\mathcal{A}}^{\top}(P - P) = 0$).

When $\xi = 0$, we have

$$\dot{\theta}_{\mathcal{A}\parallel} = -\nabla_{\mathcal{A}}f(\theta) - \lambda\theta_{\mathcal{A}\parallel} = -\nabla_{\mathcal{A}}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) - \lambda\theta_{\mathcal{A}\parallel}. \quad (176)$$

When $\xi = \tilde{\xi}_0$, we have ⁷

$$\dot{\theta}_{\mathcal{A}\parallel} = -\nabla_{\mathcal{A}}f(\theta) - \lambda\theta_{\mathcal{A}\parallel} \quad (177)$$

$$- \eta\left(\frac{1}{2}(\mathbb{1}_{\mathcal{A}}H(\theta)\nabla f(\theta) + \lambda\nabla_{\mathcal{A}}f(\theta) + \lambda\mathbb{1}_{\mathcal{A}} \odot H(\theta)\theta + \lambda^2\theta_{\mathcal{A}}) - \frac{\lambda^2}{2}\theta_{\mathcal{A}\perp}\right) \quad (178)$$

$$= -\nabla_{\mathcal{A}}f(\theta) - \lambda\theta_{\mathcal{A}\parallel} - \eta\left(\frac{1}{2}\mathbb{1}_{\mathcal{A}} \odot H(\theta)\nabla f(\theta) + \frac{\lambda}{2}\nabla_{\mathcal{A}}f(\theta) + \frac{\lambda}{2}\mathbb{1}_{\mathcal{A}} \odot H(\theta)\theta + \lambda\lambda^2\theta_{\mathcal{A}\parallel}\right) \quad (179)$$

$$= -\lambda\theta_{\mathcal{A}\parallel} - \frac{\eta\lambda^2}{2}\theta_{\mathcal{A}\parallel} - \nabla_{\mathcal{A}}f(\theta) - \frac{\eta\lambda}{2}\nabla_{\mathcal{A}}f(\theta) - \frac{\eta}{2}\mathbb{1}_{\mathcal{A}} \odot H(\theta)\nabla f(\theta) - \frac{\eta\lambda}{2}\mathbb{1}_{\mathcal{A}} \odot H(\theta)\theta \quad (180)$$

$$= -(1 + \frac{\eta\lambda}{2})(\nabla_{\mathcal{A}}f(\theta) + \lambda\theta_{\mathcal{A}\parallel}) - \frac{\eta}{2}H_{\mathcal{A}}(\theta)\nabla_{\mathcal{A}}f(\theta) - \frac{\eta}{2}\nabla_{\mathcal{A}}\nabla_{\mathcal{A}^c}^{\top}f(\theta)\nabla_{\mathcal{A}^c}f(\theta) - \frac{\eta\lambda}{2}H_{\mathcal{A}}(\theta)\theta_{\mathcal{A}} - \frac{\eta\lambda}{2}\nabla_{\mathcal{A}}\nabla_{\mathcal{A}^c}^{\top}f(\theta)\theta_{\mathcal{A}^c} \quad (181)$$

$$= -(I + \frac{\eta\lambda}{2}I + \frac{\eta}{2}H_{\mathcal{A}}(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}))(\nabla_{\mathcal{A}}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) + \lambda\theta_{\mathcal{A}\parallel}) - \frac{\eta}{2}\nabla_{\mathcal{A}}\nabla_{\mathcal{A}^c}^{\top}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c})(\nabla_{\mathcal{A}^c}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) + \lambda\theta_{\mathcal{A}^c}) \quad (182)$$

$$= -\lambda(I + \frac{\eta\lambda}{2}I + \frac{\eta}{2}H_{\mathcal{A}}(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}))\theta_{\mathcal{A}\parallel} - (I + \frac{\eta\lambda}{2}I + \frac{\eta}{2}H_{\mathcal{A}}(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}))\theta_{\mathcal{A}^c} - \frac{\eta}{2}I((\nabla_{\mathcal{A}^c}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}) + \lambda\theta_{\mathcal{A}^c}) \cdot \nabla_{\mathcal{A}^c})\nabla_{\mathcal{A}}f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^c}). \quad (183)$$

□

A.13 Proof of Theorem B.1 ¹⁰

Proof. First, note that ¹¹

$$\nabla f(\theta) \cdot G(\theta, \alpha) = 0, \quad (184)$$

which can be shown by differentiating $f(\theta) = f(G(\theta, \alpha))$ with respect to α . Thus, assuming $\theta \cdot ((\nabla f(\theta) \cdot \nabla)G(\theta, \alpha))$ and using $\dot{\theta}(t) = -\nabla f(\theta(t)) - \lambda\theta(t) - \eta\xi(\theta(t))$, we have ¹³

$$\frac{d}{dt}(\theta(t) \cdot G(\theta(t), \alpha)) \quad (185)$$

$$= \dot{\theta} \cdot G(\theta, \alpha) + \theta \cdot (\dot{\theta} \cdot \nabla G(\theta, \alpha)) \quad (186)$$

$$= -\nabla f(\theta) \cdot G(\theta, \alpha) - \lambda\theta \cdot G(\theta, \alpha) - \eta\xi(\theta) \cdot G(\theta, \alpha) + \theta \cdot (-(\nabla f(\theta) \cdot \nabla) - \lambda(\theta \cdot \nabla) - \eta\xi(\theta) \cdot \nabla)G(\theta, \alpha) \quad (187)$$

$$= -\lambda(\theta \cdot G(\theta, \alpha) + \theta \cdot ((\theta \cdot \nabla)G(\theta, \alpha))) - \eta\xi(\theta) \cdot G(\theta, \alpha) - \theta \cdot ((\eta\xi(\theta) \cdot \nabla)G(\theta, \alpha)). \quad (188)$$

B Learning Dynamics Induced by Symmetry Breaking: Neural Mechanics

To show the benefits of the counter term, we apply it to broken conservation laws [31]. In [31],¹ the authors build relationships between the symmetries of weights and conserved quantities (i.e., Noether’s theorem [57, 58] for DNNs), and they also investigate the dynamics of DNNs under symmetry breaking. We address three shortcomings of their analysis: 1) it includes a counter term only up to order one, 2) a discretization error analysis is missing, and 3) their experiment makes too optimistic an assumption on gradients.

First, we generalize broken conservation laws (Equations (18–20) in [31]) by adding all orders of the counter term. Let $\mathbf{G}(\boldsymbol{\theta}, \alpha) := \partial_\alpha \psi(\boldsymbol{\theta}, \alpha)$, which is called the generator of symmetry transformation ψ .

Theorem B.1 (Generalized broken conservation law). *Let f be symmetric under transformation ψ .³ Assume that \mathbf{G} satisfies $\boldsymbol{\theta}(t) \cdot \{(\nabla f(\boldsymbol{\theta}(t)) \cdot \nabla) \mathbf{G}(\boldsymbol{\theta}(t), \alpha)\} = 0$. Then,*

$$\begin{aligned} \frac{d}{dt}(\boldsymbol{\theta}(t) \cdot \mathbf{G}(\boldsymbol{\theta}(t), \alpha)) = \\ -\lambda \boldsymbol{\theta}(t) \cdot \mathbf{G}(\boldsymbol{\theta}(t), \alpha) - \lambda \boldsymbol{\theta}(t) \cdot \{(\boldsymbol{\theta}(t) \cdot \nabla) \mathbf{G}(\boldsymbol{\theta}(t), \alpha)\} - \eta(\boldsymbol{\xi}(\boldsymbol{\theta}(t)) \cdot \nabla) \cdot (\boldsymbol{\xi}(\boldsymbol{\theta}(t)) \cdot \mathbf{G}(\boldsymbol{\theta}(t), \alpha)). \end{aligned} \quad (192)$$

Note that the assumption holds for translation, scale, and rescale transformation [31]. Furthermore,⁵ Equation (192) can be formally solved:

$$\begin{aligned} \boldsymbol{\theta}(t) \cdot \mathbf{G}(\boldsymbol{\theta}(t), \alpha) &= \boldsymbol{\theta}(0) \cdot \mathbf{G}(\boldsymbol{\theta}(0), \alpha) e^{-\lambda t} \\ &- \lambda \int_0^t e^{-\lambda(t-\tau)} \boldsymbol{\theta}(\tau) \cdot \{(\boldsymbol{\theta}(\tau) \cdot \nabla) \mathbf{G}(\boldsymbol{\theta}(\tau), \alpha)\} d\tau \\ &- \eta \int_0^t e^{-\lambda(t-\tau)} (\boldsymbol{\xi}(\boldsymbol{\theta}(\tau)) \cdot \nabla) (\boldsymbol{\xi}(\boldsymbol{\theta}(\tau)) \cdot \mathbf{G}(\boldsymbol{\theta}(\tau), \alpha)) d\tau. \end{aligned} \quad (193)$$

The proof is given in Appendix A.13. Now, Equation (193) includes all orders of the counter term $\boldsymbol{\xi} = \sum_{\alpha=0}^\infty \boldsymbol{\xi}_\alpha$. We can reproduce [31] by setting $\boldsymbol{\xi} = \boldsymbol{\xi}_0$. In addition, we already know the discretization error (Corollary 4.1), which is lacking in [31]. We also provide empirical results on Equation (193) in the following sections.

B.1 Scale-invariant Layers⁸

For scale transformation, $\mathbf{G}(\boldsymbol{\theta}, \alpha) = \alpha_{\mathcal{A}} \boldsymbol{\theta}$, and thus, the left hand side of Equation (193) becomes $\|\boldsymbol{\theta}_{\mathcal{A}}\|^2$. Therefore, Equation 193 describes the temporal evolution of the weight norm of scale-invariant layers. Figure 7 shows the temporal evolution of $\|\boldsymbol{\theta}_{\mathcal{A}}\|^2$ for the network explained in Section 6. Figure 8 shows the gap of $\|\boldsymbol{\theta}_{\mathcal{A}}\|^2$ between GD and its theoretical predictions (GF and EoM) (Equation 193). We see that the counter term reduces the gap. There is an improvement in the experimental settings compared with [31]. As described in [31], they substitute the gradients computed in GD for the gradients used for GF’s simulation instead of using small learning rates to simulate continuous trajectories of GF. This approximation reduces computational costs, but it causes an additional gap between the surrogate gradients and the true gradients of GF along the continuous trajectories. Therefore, we avoid this approximation; we use a small learning rate ($\eta = 10^{-5}$) to simulate GF and EoM, as explained in Section 6.

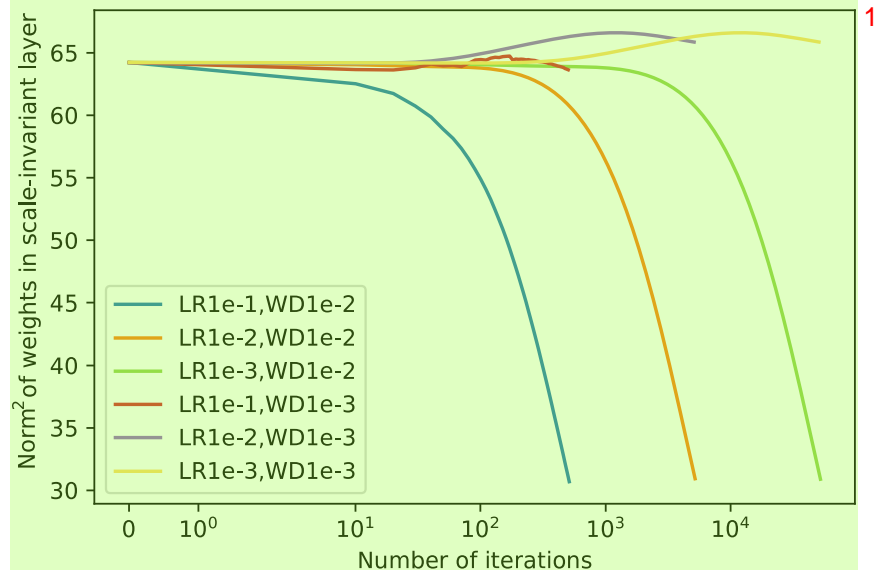


Figure 7: **Dynamics of squared weight norm of scale-invariant layer.** LR and WD mean learning 2
rate and weight decay, respectively. See Section 6 for experimental settings.

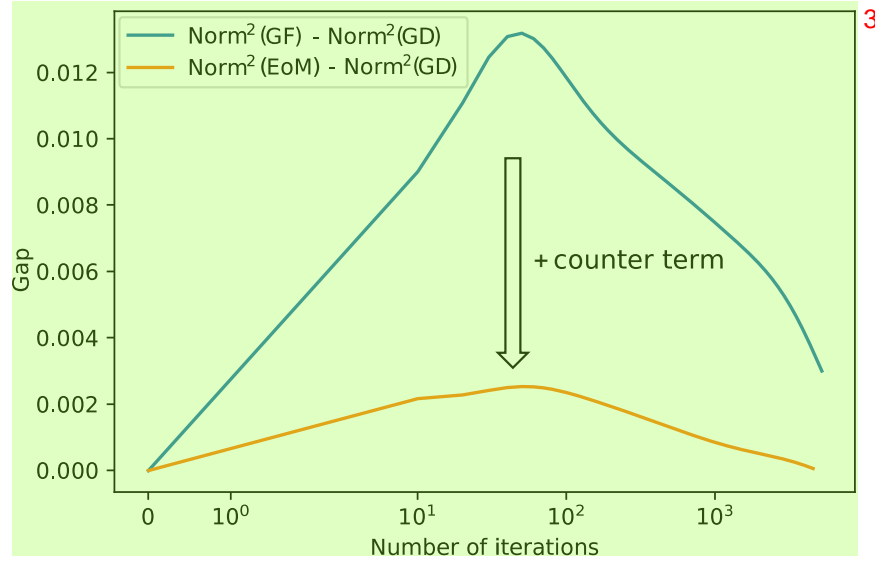


Figure 8: **Discrepancy between actual dynamics of GD and its theoretical prediction (GF and 4
EoM) of squared weight norm of scale-invariant layer.** We see that our counter term reduces the
gap between the actual dynamics of GD and its theoretical prediction. See Section 6 for experimental
settings.

B.2 Translation-invariant Layers 1

We also provide an empirical result for translation-invariant layers. For translation transformation, $\mathbf{G}(\boldsymbol{\theta}, \alpha) = \alpha \mathbb{1}_{\mathcal{A}}$ and thus the left hand side of Equation (193) becomes $\mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}$ (sum of weights). Therefore, Equation (193) describes the temporal evolution of the sum of weights of translation-invariant layers. Figure 9 shows the temporal evolution of $\mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}$ for the network described in Section 6. Figure 10 shows the gap of $\mathbb{1}_{\mathcal{A}} \cdot \boldsymbol{\theta}_{\mathcal{A}}$ between GD and its theoretical predictions (GF and EoM) (Equation 193). We see that the counter term reduces the gap.

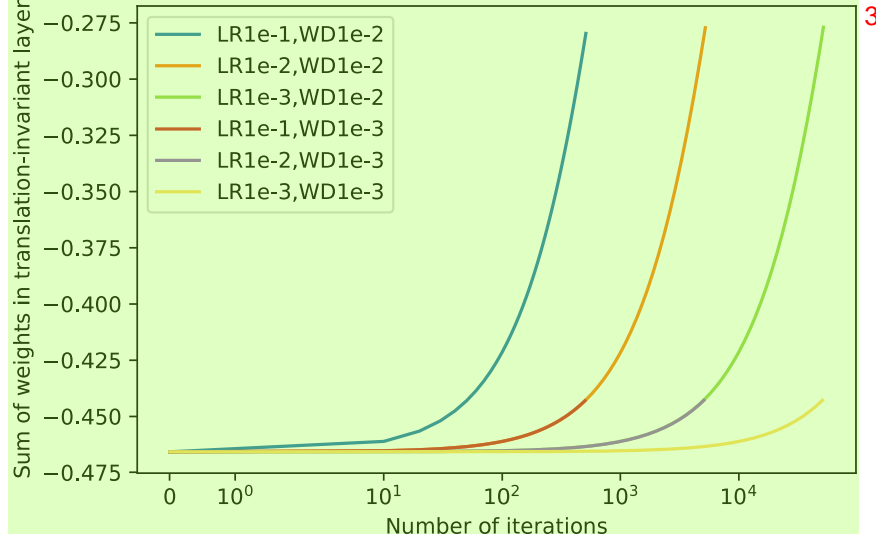


Figure 9: **Sum of weights of translation-invariant layer.** LR and WD mean learning rate and weight decay, respectively. See Section 6 for experimental settings.

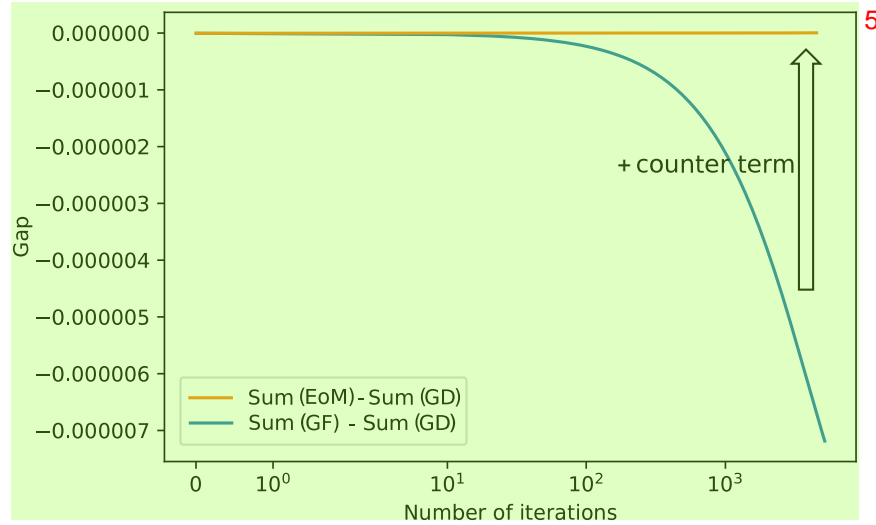


Figure 10: **Discrepancy between actual dynamics of GD and its theoretical prediction (GF and EoM) of sum of weights of translation-invariant layer.** We see that our counter term reduces the gap. See Section 6 for experimental settings.

C Equation of Motion for $\hat{\theta}_{\mathcal{A}}$ ¹

For completeness, we construct the EoM for $\hat{\theta}_{\mathcal{A}}$ for scale-invariant layers \mathcal{A} . See Section 5.1 for the EoM for $r_{\mathcal{A}}$. ²

Theorem C.1 (EoM for $\hat{\theta}_{\mathcal{A}}$). *EoM (I) gives $\dot{\hat{\theta}}_{\mathcal{A}}(t) = -\frac{1}{r_{\mathcal{A}}^2(t)} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t)) + \frac{\eta}{r_{\mathcal{A}}(t)} ((\hat{\theta}_{\mathcal{A}}(t) \cdot \xi(\theta(t))) \hat{\theta}_{\mathcal{A}}(t) - \xi(\theta(t)))$. Specifically, this is equivalent to:* ³

$$\dot{\hat{\theta}}_{\mathcal{A}}(t) = -\frac{1}{r_{\mathcal{A}}^2(t)} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t)) \quad (194) \quad \text{⁴}$$

for $\xi = 0$ (GF) and ⁵

$$\dot{\hat{\theta}}_{\mathcal{A}} = -\frac{1}{r_{\mathcal{A}}^2} \left(I + \frac{\eta}{2} H_{\mathcal{A}}(\theta) + \frac{\eta}{2} I((\nabla_{\mathcal{A}^c} f(\theta) + \lambda \theta_{\mathcal{A}^c}) \cdot \nabla_{\mathcal{A}^c}) + \frac{\eta}{2} \hat{\theta}_{\mathcal{A}} \nabla_{\mathcal{A}}^{\top} f(\theta) \right) \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c}) \quad (195) \quad \text{⁶}$$

for $\xi = \tilde{\xi}_0$ (EoM), where $H_{\mathcal{A}}(\hat{\theta}_{\mathcal{A}}) := (\mathbf{1}_{\mathcal{A}} \odot \nabla)(\mathbf{1}_{\mathcal{A}} \odot \nabla)^{\top} f(\theta)|_{\theta=\hat{\theta}_{\mathcal{A}}}$.

The proof is given in Appendix A.9. ⁷

Effective learning rate. This result highlights the differences between GD and GF on scale-invariant layers. The factor $\frac{1}{r_{\mathcal{A}}^2}$ (Equation (194)), which is $\frac{\eta}{r_{\mathcal{A}}^2}$ at discretization, is called the *effective learning rate* [29, 42, 30, 43, 44, 33, 45, 34, 46, 47]. The dynamics of $\hat{\theta}_{\mathcal{A}}$ is induced by $\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c})$ with the effective learning rate $\frac{\eta}{r_{\mathcal{A}}^2}$, not η . We find that the counter term corrects the effective learning rate to a matrix operator form (Equation (195)). Let us see the meaning of each correction in order. First, I (identity matrix) corresponds to the original effective learning rate. Second, $\frac{\eta}{2} H_{\mathcal{A}}$ directs the gradient $\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c})$ toward the maximum eigenvector of $H_{\mathcal{A}}$, i.e., a flat direction. Therefore, GD tends to go through flatter regions than GF. Third, $\frac{\eta}{2} I((\nabla_{\mathcal{A}^c} f(\theta) + \lambda \theta_{\mathcal{A}^c}) \cdot \nabla_{\mathcal{A}^c})$ involves $\nabla_{\mathcal{A}^c} f$ into the learning dynamics of \mathcal{A} ; therefore, \mathcal{A} is explicitly affected by \mathcal{A}^c in GD, unlike in GF. This point is often missing in the literature on scale-invariant networks because it is often assumed that the whole network is scale-invariant. Fourth, $\frac{\eta}{2} \hat{\theta}_{\mathcal{A}} \nabla_{\mathcal{A}}^{\top} f(\theta)$ cancels the $\hat{\theta}_{\mathcal{A}}$ component of the right hand side of Equation (195), which may not seem obvious but can be seen from the proof of Theorem C.1 (see Appendix A.9), and thus, $\dot{\hat{\theta}}_{\mathcal{A}}$ is orthogonal to $\hat{\theta}_{\mathcal{A}}$, which should be satisfied anyway because $\|\hat{\theta}_{\mathcal{A}}\|^2 \equiv 1 \implies 2\dot{\hat{\theta}}_{\mathcal{A}} \cdot \hat{\theta}_{\mathcal{A}} = 0$. ⁸

D Equation of Motion for $\theta_{\mathcal{A}\parallel}$ ¹

For completeness, we provide the EoM for $\theta_{\mathcal{A}\parallel}$. The proof is given in Appendix A.12. ²

Theorem D.1 (EoM for $\theta_{\mathcal{A}\parallel}$). *EoM (I) gives*

$$\dot{\theta}_{\mathcal{A}\parallel}(t) = -\lambda\theta_{\mathcal{A}\parallel}(t) - \nabla f(\theta_{\mathcal{A}\parallel}(t)) - \eta(I - P)\xi(\theta(t)). \quad \text{³ (196)}$$

Specifically, this is equivalent to: ⁴

$$\dot{\theta}_{\mathcal{A}\parallel}(t) = -\lambda\theta_{\mathcal{A}\parallel}(t) - \nabla f(\theta_{\mathcal{A}\parallel}(t) + \theta_{\mathcal{A}^\perp}) \quad \text{⁵ (197)}$$

for $\xi = 0$ (GF) and ⁶

$$\begin{aligned} \dot{\theta}_{\mathcal{A}\parallel}(t) = & -\lambda\left(I + \frac{\eta\lambda}{2}I + \frac{\eta}{2}H_{\mathcal{A}}(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^\perp})\right)\theta_{\mathcal{A}\parallel} \\ & - \left(I + \frac{\eta\lambda}{2}I + \frac{\eta}{2}H_{\mathcal{A}}(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^\perp}) + \frac{\eta}{2}I((\nabla_{\mathcal{A}^\perp} f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^\perp}) + \lambda\theta_{\mathcal{A}^\perp}) \cdot \nabla_{\mathcal{A}^\perp})\right) \nabla_{\mathcal{A}} f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^\perp}) \end{aligned} \quad \text{⁷ (198)}$$

for $\xi = \tilde{\xi}_0$ (EoM). ⁸

This result highlights the differences between the dynamics of GD and GF. The two factors $\frac{\eta\lambda}{2}I$ in Equation (198) mean that the existence of weight decay increases the learning rate (increases the velocity $\dot{\theta}_{\mathcal{A}\parallel}$). The factor $\frac{\eta}{2}H$ means that, as mentioned in Appendix C, GD tends to go along sharper paths than GF. Note that velocity $\dot{\theta}_{\mathcal{A}\parallel}$ is orthogonal to $\theta_{\mathcal{A}\perp}$ because ∇f , $\theta_{\mathcal{A}\parallel}$, and $H(\nabla f + \lambda\theta_{\mathcal{A}\parallel})$ are orthogonal to $\theta_{\mathcal{A}\perp}$. $H(\nabla f + \lambda\theta_{\mathcal{A}\parallel}) \perp \theta_{\mathcal{A}\perp}$ follows because $Hv \perp \theta_{\mathcal{A}\perp}$ for arbitrary non-zero vector $v \in \mathbb{R}^d$ ($\cdot: H\mathbb{1}_{\mathcal{A}} = H\theta_{\mathcal{A}\perp} = 0$) (see Lemma A.7). $\frac{\eta}{2}I((\nabla_{\mathcal{A}^\perp} f(\theta_{\mathcal{A}\parallel} + \theta_{\mathcal{A}^\perp}) + \lambda\theta_{\mathcal{A}^\perp}) \cdot \nabla_{\mathcal{A}^\perp})$ involves $\nabla_{\mathcal{A}^\perp} f$ into the learning dynamics of \mathcal{A} . We see that the dynamics of $\theta_{\mathcal{A}\parallel}$ is also independent of that of $\theta_{\mathcal{A}\perp}$, and thus, they are completely separable. A summary of Theorems 5.2 and D.1 is given in Figure 5. ⁹

E Details of Experiment ¹

We provide detailed experimental settings (see also Section 6). Our computational infrastructure is a DGX-1 server. The fundamental libraries used in the experiment are TensorFlow 2.3 [59], Numpy 1.18 [60], and Python 3.6.8 [61]. The random seeds used for TensorFlow and Numpy are both 7. The input image is first divided by 127.5 and subtracted by 1. The maximum total number of iterations is 5 million steps for GF and EoM. The total runtime is approximately a month. We use least square fitting (`np.polyfit`) to calculate the decay rates in Table 1. More information and detailed experimental results can be found in our code. ²

In Figures 2 and 12, the theoretical prediction of discretization error is defined as $\|e_k\| = \frac{\eta^2}{2} \|\sum_{s=0}^{k-1} (H(\theta(s\eta)) + \lambda I)\mathbf{g}(\theta(s\eta))\|$ (Equation (12)). To reduce computational costs, we approximate the r.h.s.: $(H(\theta(t)) + \lambda I)\mathbf{g}(\theta(t)) \sim \frac{\mathbf{g}(\theta(t)+\epsilon\mathbf{g}(\theta(t))) - \mathbf{g}(\theta(t)-\epsilon\mathbf{g}(\theta(t)))}{2\epsilon}$, where ϵ is set to 10^{-7} . The green curve in Figure 2 is defined as $e_k = \tilde{e}_{100} + \frac{\eta^2}{2} \sum_{s=100}^{k-1} (H(\theta(s\eta)) + \lambda I)\mathbf{g}(\theta(s\eta))$ (compare this with Equation (12)), where \tilde{e}_{100} is the actual discretization error at the 100th step that is obtained from GD. Therefore, the green curve represents the theoretical prediction of discretization error after the 100th step, given \tilde{e}_{100} . ³

F Supplementary Experiment 1

F.1 Relative Discretization Error 2

We provide the relative discretization error, which is defined as $\|e_k\|/\|\theta_k\|$ ($k \in \mathbb{Z}_{\geq 0}$). See Figure 11. We can see that a large learning rate ($\eta = 10^{-1}$) leads to a large discretization error (Figure 11 (a) and (c)). We also see that the counter term reduces the discretization error as expected (Figure 11 (b) and (d)).

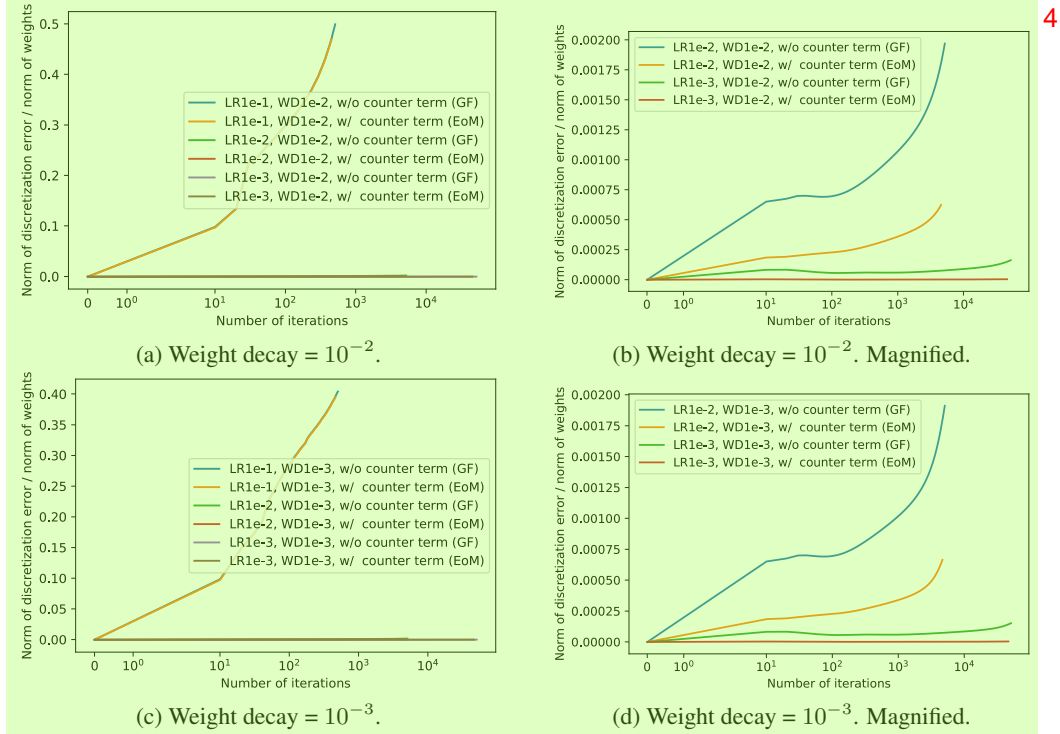


Figure 11: **Relative discretization error.** In (a) and (c), the LR1e-1 curves overlap each other, and the LR1e-2 and LR1e-3 curves collapse in the lower region of the figure. The LR1e-2 and LR1e-3 are magnified and shown in (c) and (d). See Section 6 and Appendix E for experimental settings.

F.2 Theoretical Prediction Vs. Experimental Result of Discretization Error 1

We compare the theoretical prediction of discretization error between GF and GD (Equation (12)) 2 with the actual discretization error obtained in the experiment. The green curve is defined as $e_k = \tilde{e}_{100} + \frac{\eta^2}{2} \sum_{s=100}^{k-1} (H(\theta(s\eta)) + \lambda I) \mathbf{g}(\theta(s\eta)) + O(\eta^3)$ (compare this with Equation (12)), where \tilde{e}_{100} is the actual discretization error at the 100th step. Therefore, the green curve represents the theoretical prediction of discretization error after the 100th step given \tilde{e}_{100} .

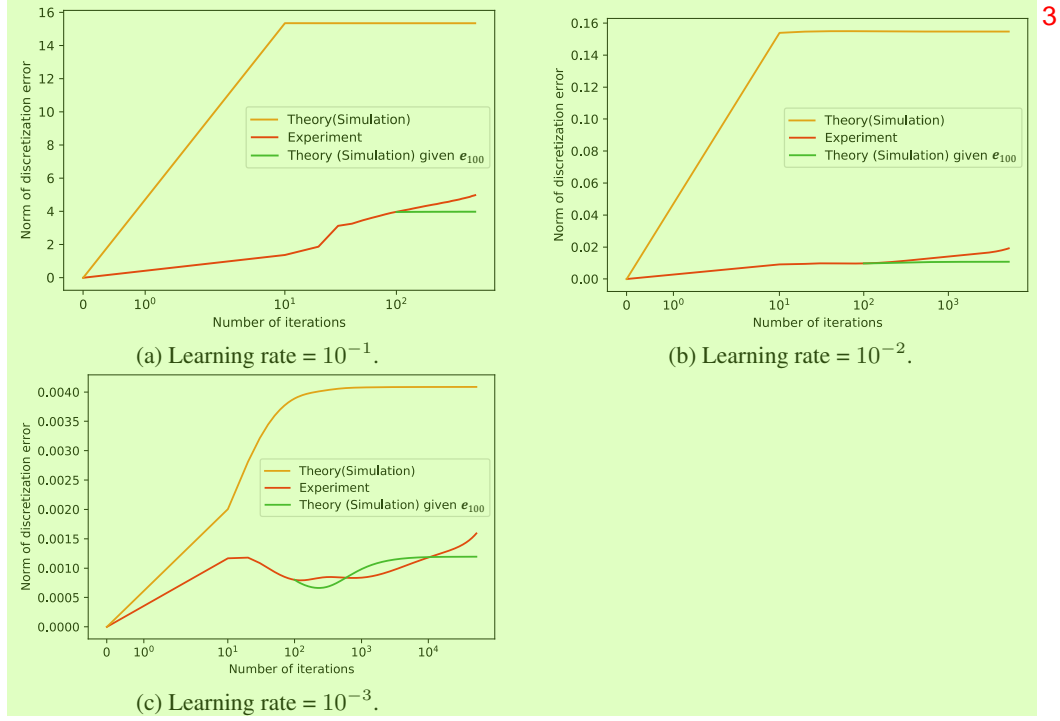


Figure 12: Theoretical prediction (Equation (12)) vs. experimental result of discretization error 4 between GF and GD. The weight decay is 10^{-2} . See Section 6 and Appendix E for experimental settings.

G Supplementary Discussion ¹

Supplementary related work (Section 2). To show the benefits of EoM, we focus on scale-invariant layers [29, 42, 30, 43, 44, 33, 45, 34, 46, 47] and translation-invariant layers [31, 32] in Section 5. To carry over the stability of a continuous optimization algorithm to a discretized system, the authors of [19] add a feedback term to the optimization, and after that, they apply a discretization method to it. The authors’ primary motivation is to keep the orthogonality of the weight parameters of DNNs, which is different from ours.

Convergence of ξ (Section 3.3). Note that the expansion of ξ in terms of η is not necessarily convergent, as is also pointed out in [35]. Thus, we have to truncate the expansion at a suitable order. The discretization error at the truncation is given in Theorem 4.1.

Beyond leading order of discretization error (Theorem 3.2 and Section 4.1). In this work, we analyze the leading order of discretization error. However, higher-order terms cannot always be negligible. We discuss in Section 4.1 that the higher-order terms are important at the beginning of training.

Existence of \mathcal{A}^c (Section 5). In our theoretical analysis of scale- and translation-invariant layers, the network contains both invariant (\mathcal{A}) and non-invariant layers (\mathcal{A}^c), while previous works assume the whole network is invariant for simplicity [29, 42, 30, 43, 44, 33, 45, 34, 46, 47]. We avoid this assumption and show that such mixed networks require appropriate modifications to analyses of invariant networks. For example, $\nabla f(\theta) = \frac{1}{\|\theta\|} \nabla f(\hat{\theta})$ for invariant networks, while $\nabla_{\mathcal{A}} f(\theta) = \frac{1}{\|\theta_{\mathcal{A}}\|} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}} + \theta_{\mathcal{A}^c})$ for mixed networks (Lemma A.5), not $\frac{1}{\|\theta_{\mathcal{A}}\|} \nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}})$. Such a naive replacement is not allowed.

Higher-order corrections to decay rate of $r_{\mathcal{A}}$ (Section 5.1). We can compute more corrections to the decay rate of $r_{\mathcal{A}}$ (\mathcal{A} is a scale-invariant layer), using more counter terms. For example, a long algebra gives decay rate $\eta\lambda(1 + \frac{\eta\lambda}{2} + \frac{\eta^2\lambda^2}{3})$ for $\xi = \tilde{\xi}_0 + \eta\tilde{\xi}_1$. The proof is similar to Appendix A.7.

On equilibrium assumptions in Corollaries 5.1 and 5.2 (Section 5.1). We make assumptions in Corollaries 5.1 and 5.2; there exist two constants $r_{\mathcal{A}*} \geq 0$ and $c_* \geq 0$ such that $r_{\mathcal{A}}(t) \xrightarrow{t \rightarrow \infty} r_{\mathcal{A}*}$ and $\|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t) + \theta_{\mathcal{A}^c}(t))\| \xrightarrow{t \rightarrow \infty} c_*$. These assumptions are similar to those given in previous studies [29, 34]. However, whether the assumptions are valid in the actual learning dynamics of DNNs is of independent interest. In fact, the equilibrium assumption ($r_{\mathcal{A}*}(t)$ and $\|\nabla_{\mathcal{A}} f(\hat{\theta}_{\mathcal{A}}(t) + \theta_{\mathcal{A}^c}(t))\| \xrightarrow{t \rightarrow \infty} \text{constant}$) could not be satisfied even at one million steps of GD, and potentially because of it, $r_{\mathcal{A}*}$ and Δ_* have a large discrepancy between the empirical results and theoretical predictions. Deeper analyses on this point are needed. Under what conditions are the equilibrium assumptions valid? Can we relax the equilibrium assumptions and obtain realistic limiting dynamics of scale-invariant layers? This is exciting future work.

In contrast to our empirical result mentioned above, in [34], their experiments dramatically match their theoretical prediction. This is potentially because of differences in experimental settings; in [34], SGD is used (ours is GD) and variance is induced, ResNet-50 [62, 63] is used (ours is a fully-connected network with three layers), ImageNet [64, 65] and MSCOCO [66] are used (ours is MNIST [50]), and large learning rates ($\sim 10^{-1}$) and small weight decays ($\sim 10^{-4}$) are used (ours are given in Appendix E).

Extension of EoM to general settings (Section 7). While we focus on GD and GF for simplicity, our counter-term-based approach and discretization error analysis can be extended to more general settings, such as SGD, acceleration methods (e.g., momentum SGD), and adaptive optimizers (e.g., Adam [53]). First, to extend our analysis to SGD, discretization error analysis of the Euler-Maruyama method, e.g., [67], can be used. SDE’s error analysis [23, 24] is also relevant. Second, we can extend our counter-term-based approach and discretization error analysis to acceleration methods by modifying the analysis for different differential equations from GF and different discretization schemes from the Euler method, as is discussed in [7, 14, 12]. Third, [56] is the first work that

provides a continuous approximation of Adam. However, its counter term and discretization error are open questions. 1