

Singapore Management University

## Institutional Knowledge at Singapore Management University

---

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

---

3-2024

### Towards understanding convergence and generalization of AdamW

Pan ZHOU

Singapore Management University, panzhou@smu.edu.sg

Xingyu XIE

Zhouchen LIN

Shuicheng YAN

Follow this and additional works at: [https://ink.library.smu.edu.sg/sis\\_research](https://ink.library.smu.edu.sg/sis_research)



Part of the [Graphics and Human Computer Interfaces Commons](#)

---

#### Citation

ZHOU, Pan; XIE, Xingyu; LIN, Zhouchen; and YAN, Shuicheng. Towards understanding convergence and generalization of AdamW. (2024). *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 1-8. Available at: [https://ink.library.smu.edu.sg/sis\\_research/8986](https://ink.library.smu.edu.sg/sis_research/8986)

This Journal Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email [cherylds@smu.edu.sg](mailto:cherylds@smu.edu.sg).

# Towards Understanding Convergence and Generalization of AdamW

Pan Zhou, Xingyu Xie, Zhouchen Lin, *Fellow, IEEE*, Shuicheng Yan, *Fellow, IEEE*

**Abstract**—AdamW modifies Adam by adding a decoupled weight decay to decay network weights per training iteration. For adaptive algorithms, this decoupled weight decay does not affect specific optimization steps, and differs from the widely used  $\ell_2$ -regularizer which changes optimization steps via changing the first- and second-order gradient moments. Despite its great practical success, for AdamW, its convergence behavior and generalization improvement over Adam and  $\ell_2$ -regularized Adam ( $\ell_2$ -Adam) remain absent yet. To solve this issue, we prove the convergence of AdamW and justify its generalization advantages over Adam and  $\ell_2$ -Adam. Specifically, AdamW provably converges but minimizes a dynamically regularized loss that combines vanilla loss and a dynamical regularization induced by decoupled weight decay, thus yielding different behaviors with Adam and  $\ell_2$ -Adam. Moreover, on both general nonconvex problems and PL-conditioned problems, we establish stochastic gradient complexity of AdamW to find a stationary point. Such complexity is also applicable to Adam and  $\ell_2$ -Adam, and improves their previously known complexity, especially for over-parametrized networks. Besides, we prove that AdamW enjoys smaller generalization errors than Adam and  $\ell_2$ -Adam from the Bayesian posterior aspect. This result, for the first time, explicitly reveals the benefits of decoupled weight decay in AdamW. Experimental results validate our theory.

**Index Terms**—Analysis of AdamW, Convergence of AdamW, Generalization of AdamW, Adaptive gradient algorithms

## 1 INTRODUCTION

ADAPTIVE gradient algorithms, *e.g.*, Adam [1], have become the most popular optimizers to train deep networks because of their faster convergence speed than SGD [2], with many successful applications in computer vision [3], [4] and natural language processing [5], *etc.* Similar to the precondition in the second-order algorithms [6], adaptive algorithms precondition the landscape curvature of loss objective to adjust the learning rate for each gradient coordinate. This precondition often helps these adaptive algorithms achieve faster convergence speed than their non-adaptive counterparts, *e.g.*, SGD which uses a single learning rate for all gradient coordinates. Unfortunately, this precondition also brings negative effect. That is, adaptive algorithms usually suffer from worse generalization performance than SGD [7]–[10].

As a leading adaptive gradient approach, AdamW [11] greatly improves the generalization performance of adaptive algorithms on vision transformers (ViTs) [12] and CNNs [13], [14]. The core of AdamW is a decoupled weight decay. Specifically, AdamW uses an exponential moving average to estimate the first-order moment  $\mathbf{m}_k$  and second-order moment  $\mathbf{n}_k$  like Adam, and then updates network weights  $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \mathbf{m}_k / \sqrt{\mathbf{n}_k + \delta} - \eta \lambda_k \mathbf{x}_k$  with a learning rate  $\eta$ , a weight decay parameter  $\lambda_k$ , and a small constant  $\delta$ . One can observe that AdamW decouples the weight decay from the optimization steps w.r.t. the loss function, since the weight decay is always  $-\eta \lambda_k \mathbf{x}_k$  no matter what the loss and optimization step are. This decoupled weight decay becomes  $\ell_2$ -regularization for SGD, but differs from  $\ell_2$ -regularization for adaptive algorithms. Thanks

to its effectiveness, AdamW has been widely used in network training. But there remain many mysteries about AdamW yet. Firstly, it is not clear whether AdamW can theoretically converge or not, and if yes, what convergence rate it can achieve. Moreover, for the generalization superiority of AdamW over the widely used Adam and  $\ell_2$ -regularized Adam ( $\ell_2$ -Adam), the theoretical reasons are rarely investigated though heavily desired.

**Contributions:** To resolve these issues, we provide a new viewpoint to understand the convergence and generalization behaviors of AdamW. Particularly, we theoretically prove the convergence of AdamW, and also justify its superior generalization to ( $\ell_2$ )-Adam. Our main contributions are highlighted below.

Firstly, we prove that AdamW can converge but minimizes a dynamically regularized loss that combines the vanilla loss and a dynamical regularization induced by the decoupled weight decay. Interestingly, this dynamical regularization differs from the commonly used  $\ell_2$ -regularization, and thus yields the different behaviors between AdamW and  $\ell_2$ -Adam. For convergence speed, on general nonconvex problems, AdamW finds an  $\epsilon$ -accurate first-order stationary point within stochastic gradient complexity  $\mathcal{O}(c_\infty^{2.5} \epsilon^{-4})$  when using constant learning rate and  $\mathcal{O}(c_\infty^{1.25} \epsilon^{-4} \log(\frac{1}{\epsilon}))$  with decaying learning rate, where  $c_\infty$  is the  $\ell_\infty$ -norm upper bound of stochastic gradient. When ignoring logarithm terms, both complexities match the lower complexity bound  $\mathcal{O}(\epsilon^{-4})$  in [15]. These complexities are applicable to Adam and  $\ell_2$ -Adam, and improve their previously known complexities  $\mathcal{O}(c_\infty \sqrt{d} \epsilon^{-4})$  and  $\mathcal{O}(c_\infty \sqrt{d} \epsilon^{-4} \log(\frac{1}{\epsilon}))$  when respectively using constant and decaying learning rate [16]–[18], as  $c_\infty$  is often much smaller than the network parameter dimension  $d$ . On PL-conditioned nonconvex problems, our established complexity of AdamW also enjoys similar advantages.

Next, we theoretically show the benefits of the decoupled weight decay in AdamW to the generalization performance from the Bayesian posterior aspect. Specifically, we show that a proper decoupled weight decay  $\lambda_k > 0$  helps AdamW achieve smaller

- P. Zhou is with School of Computing and Information Systems, Singapore Management University, Singapore. Shuicheng Yan is with Skywork AI. S. Yan is corresponding author. (`{panzhou3,shuicheng.yan}@gmail.com`)
- X. Xie and Z. Lin are with the National Key Lab. of General Artificial Intelligence, Peking University, China. Z. Lin is also with the Institute for Artificial Intelligence, Peking University, China, and Peng Cheng Laboratory, Shenzhen, 518055, China. (email: `{xyxie,zlin}@pku.edu.cn`)

generalization error, indicating the superiority of AdamW over vanilla Adam that corresponds to  $\lambda_k = 0$ . We further analyze  $\ell_2$ -regularized Adam, and observe that AdamW often enjoys smaller generalization error bound than  $\ell_2$ -regularized Adam. To our best knowledge, this work is the first one that explicitly shows the superiority of AdamW over Adam and its  $\ell_2$ -regularized version.

## 2 RELATED WORK

**Convergence Analysis.** Adaptive gradient algorithms, *e.g.*, Adam, have become the default optimizers in deep learning because of their fast convergence speed. Accordingly, many works investigate their convergence to deepen their understanding. On convex problems, Adam-type algorithms, *e.g.*, Adam and AMSGrad [19], enjoy the regret  $\mathcal{O}(\sqrt{T})$  under the online learning setting with training iteration number  $T$ . For nonconvex problems, Adam-type algorithms have the stochastic gradient complexity  $\mathcal{O}(c_\infty \sqrt{d} \epsilon^{-4})$  to find an  $\epsilon$ -accurate stationary point [18], [20]. RMSProp and Padam [17] are proved to have the complexity  $\mathcal{O}(\sqrt{c_\infty} d \epsilon^{-4})$  [16], and Adabelief [21] has  $\mathcal{O}(c_2^6 \epsilon^{-4})$  complexity, where  $c_2$  is the  $\ell_2$ -norm upper bound of stochastic gradient. But the convergence behaviors of AdamW remains unclear, though it is the dominant optimizer for vision transformers [12] and CNNs [13].

**Generalization Analysis.** Most works, *e.g.*, [22]–[24], analyze the generalization of an algorithm through studying its stochastic differential equations (SDEs) because of the similar convergence behaviors of an algorithm and its SDE. For instance, by formulating SGD into Brownian- or Lévy-driven SDEs, SGD always provably tends to converge to flat minima and thus enjoys good generalization [9], [24]. Recently, for weight decay, the works [25]–[27] intuitively claim that for layers followed by normalizations, *e.g.*, BatchNormalization [28], weight decay increases the effective learning rate by reducing the scale of the network weights, and higher learning rates give larger gradient noise which often acts a stochastic regularizer. But Zhou et al. [29] argued the benefits of weight decay to the layers without normalization, *e.g.*, fully-connected networks, and further empirically found the regularization effects of weight decay to the last fully-connected layer of a network. Unfortunately, none of them explicitly show the generalization benefits of weight decay in AdamW. Here we borrow the aforementioned SDE tool and PAC Bayesian framework [30] to explicitly analyze the generalization effects of decoupled weight decay of AdamW and also its superiority over  $\ell_2$ -Adam.

## 3 NOTATION AND PRELIMINARILY

**AdamW &  $\ell_2$ -Adam.** We first briefly recall the steps of AdamW, Adam and  $\ell_2$ -Adam to solve the stochastic nonconvex problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{D}} [f(\mathbf{x}, \boldsymbol{\xi})], \quad (1)$$

where loss  $f$  is differentiable and nonconvex, sample  $\boldsymbol{\xi}$  is drawn from a distribution  $\mathcal{D}$ . To solve problem (1), at the  $k$ -th iteration, AdamW estimates the current gradient  $\nabla F(\mathbf{x}_k)$  as the minibatch gradient  $\mathbf{g}_k = \frac{1}{b} \sum_{i=1}^b \nabla f(\mathbf{x}_k; \boldsymbol{\xi}_i)$ , and updates the variable  $\mathbf{x}$  with three constants  $\beta_1 \in [0, 1]$ ,  $\beta_2 \in [0, 1]$  and  $\delta > 0$ :

$$\begin{aligned} \mathbf{m}_k &= (1 - \beta_1) \mathbf{m}_k + \beta_1 \mathbf{g}_k, & \mathbf{n}_k &= (1 - \beta_2) \mathbf{n}_k + \beta_2 \mathbf{g}_k^2, \\ \mathbf{x}_{k+1} &= \mathbf{x}_k - \eta \mathbf{m}_k / \sqrt{\mathbf{n}_k + \delta} - \eta \lambda_k \mathbf{x}_k, \end{aligned} \quad (2)$$

where  $\mathbf{m}_0 = \mathbf{g}_0$ ,  $\mathbf{n}_0 = \mathbf{g}_0^2$ , and all operations (*e.g.*, product, division) involved vectors are element-wise. Here we allow  $\lambda_k$  to evolve along iteration number  $k$ , as in practice, an evolving  $\lambda_k$  often

shows better performance than a fixed one [4], [31]–[33]. See detailed AdamW in Algorithm 1 of Appendix B. AdamW differs from vanilla Adam in the third step of Eqn. (2). Specifically, AdamW decouples weight decay from the optimization steps, as weight decay is always  $-\eta \lambda_k \mathbf{x}_k$  no matter what the loss and optimization step are. But  $\ell_2$ -Adam adds a conventional weight decay  $\lambda_k \mathbf{x}_k$  into the gradient estimation  $\mathbf{g}_k = \frac{1}{b} \sum_{i=1}^b \nabla f(\mathbf{x}_k; \boldsymbol{\xi}_i) + \lambda_k \mathbf{x}_k$ , then updates  $\mathbf{m}_k$  and  $\mathbf{n}_k$  in (2), and  $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \mathbf{m}_k / \sqrt{\mathbf{n}_k + \delta}$ . The decoupled weight decay in AdamW often achieves better generalization than  $\ell_2$ -Adam on many networks, *e.g.*, [12], [14].

**Analysis Assumptions.** Here we introduce necessary assumptions for analysis, which are commonly used in [1], [8], [19], [34]–[36].

**Assumption 1 ( $L$ -smoothness).** *The function  $f(\cdot, \cdot)$  is  $L$ -smooth w.r.t. the parameter, if  $\exists L > 0$ , for  $\forall \mathbf{x}_1, \mathbf{x}_2$  and  $\boldsymbol{\xi} \sim \mathcal{D}$ , we have*

$$\|\nabla f(\mathbf{x}_1, \boldsymbol{\xi}) - \nabla f(\mathbf{x}_2, \boldsymbol{\xi})\|_2 \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

**Assumption 2 (Gradient assumption).** *The gradient estimation  $\mathbf{g}_k$  is unbiased, and its magnitude and variance are bounded:*

$$\mathbb{E}[\mathbf{g}_k] = \nabla F(\mathbf{x}_k), \quad \|\mathbf{g}_k\|_\infty \leq c_\infty, \quad \mathbb{E}[\|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|_2^2] \leq \sigma^2.$$

When a nonconvex problem satisfies Assumptions 1 and 2, the lower bound of the stochastic gradient complexity (a.k.a. IFO complexity) to find an  $\epsilon$ -accurate first-order stationary point is  $\Omega(\epsilon^{-4})$  [15]. Next, we introduce Polyak-Łojasiewicz (PL) condition which is widely used in deep network analysis, since as observed or proved in [37]–[40], deep neural networks often satisfy PL condition at least around a local minimum.

**Assumption 3 (PL Condition).** *Let  $\mathbf{x}_* \in \operatorname{argmin}_{\mathbf{x}} F(\mathbf{x})$ . We say a function  $F(\mathbf{x})$  satisfies  $\mu$ -PL condition if it satisfies  $2\mu(F(\mathbf{x}) - F(\mathbf{x}_*)) \leq \|\nabla F(\mathbf{x})\|_2^2$  ( $\forall \mathbf{x}$ ), where  $\mu$  is a universal constant.*

## 4 CONVERGENCE ANALYSIS

Here we first use a specific least square problem to compare the convergence behavior of AdamW and  $\ell_2$ -Adam. Next, we study the convergence of AdamW on general nonconvex problems and show its performance improvement on PL-conditioned problems.

### 4.1 Results on Specific Least Square Problems

Here we first use a specific least square problem (3) to analyze the different convergence performance of AdamW and  $\ell_2$ -Adam:

$$\min_{\mathbf{x} \in \mathbb{R}} F(\mathbf{x}) := \mathbb{E}_{\boldsymbol{\xi} \sim \mathcal{N}(0,1)} \frac{1}{2} \|a\mathbf{x} - \boldsymbol{\xi}\|_2^2, \quad (3)$$

where  $a \neq 0$  is a constant. Then we state our main results in Theorem 1 whose proof can be found in Appendix G.1.

**Theorem 1.** *Suppose that stochastic gradient  $\mathbf{g}_k$  is unbiased,  $\mathbb{E}[\|\mathbf{g}_k\|_2] \leq \tau$ , and  $\mathbb{E}\|\mathbf{x}_0 - \mathbf{x}_*\|_2 \leq \Delta$ . Then with learning rate  $\eta_k = \mathcal{O}(\frac{1}{k})$  and  $\lambda_k = \lambda = \mathcal{O}(\sqrt{k})$ , the sequence  $\{\mathbf{x}_k\}$  generated by AdamW obeys:*

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_*\|_2] \leq (1 - 1/\sqrt{k})^{\frac{3k}{2}} \Lambda + \frac{\tau}{k^{\frac{1}{2} + \alpha}},$$

where  $\alpha > 0$ ,  $\Lambda = \eta_0 + \Delta$ . With learning rate  $\eta_k = \mathcal{O}(\frac{1}{\sqrt{k}})$  and  $\lambda_k = \lambda = \mathcal{O}(\sqrt{k})$ , the sequence  $\{\mathbf{x}_k\}$  generated by  $\ell_2$ -Adam obeys:

$$\mathbb{E}[\|\mathbf{x}_k - \mathbf{x}_*\|_2] \leq (1 - 1/\sqrt{k})^{\frac{k}{2}} \Lambda + \frac{2\tau}{k^{\frac{1}{2}}}.$$

Theorem 1 shows that AdamW enjoys a faster convergence speed than  $\ell_2$ -Adam on the least square problem in (3). Specifically,

the first convergence term  $(1 - 1/\sqrt{k})^{\frac{3k}{2}} \Lambda$  in AdamW converges much faster than the corresponding term  $(1 - 1/\sqrt{k})^{\frac{k}{2}} \Lambda$  in  $\ell_2$ -Adam. For the second term  $\frac{\tau}{k^{\frac{1}{2} + \alpha}}$  in AdamW, it improves the corresponding term in  $\ell_2$ -Adam by a factor of  $2k^\alpha$  ( $\alpha > 0$ ). This comparison shows the superiority of AdamW over  $\ell_2$ -Adam, and thus partially explains their different convergence behaviors.

## 4.2 Results on Nonconvex Problems

Now we move on to the general and also PŁ conditioned nonconvex problems. We first define a dynamic surrogate function  $F_k(\mathbf{x})$  at the  $k$ -th iteration which is indeed the combination of the vanilla loss  $F(\mathbf{x})$  in Eqn. (1) and a dynamic regularization  $\frac{\lambda}{2} \|\mathbf{x}\|_{\mathbf{v}_k}^2$ :

$$F_k(\mathbf{x}) = F(\mathbf{x}) + \frac{\lambda_k}{2} \|\mathbf{x}\|_{\mathbf{v}_k}^2 = \mathbb{E}_\zeta[f(\boldsymbol{\theta}; \zeta)] + \frac{\lambda_k}{2} \|\mathbf{x}\|_{\mathbf{v}_k}^2, \quad (4)$$

where  $\mathbf{v}_k = \sqrt{\mathbf{n}_k + \delta}$  and  $\|\mathbf{x}\|_{\mathbf{v}_k} = \sqrt{\langle \mathbf{x}, \mathbf{v}_k \odot \mathbf{x} \rangle}$  with element-wise product  $\odot$ . To minimize (4), one can approximate vanilla loss  $F(\mathbf{x})$  by its Taylor expansion, and compute  $\mathbf{x}_{k+1}$ :

$$\begin{aligned} \mathbf{x}_{k+1} \approx \operatorname{argmin}_{\mathbf{x}} F(\mathbf{x}_k) + \langle \nabla F(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{v}_k}^2 \\ + \frac{\lambda_k}{2} \|\mathbf{x}\|_{\mathbf{v}_k}^2 = \frac{1}{1 + \lambda_k \eta} (\mathbf{x}_k - \eta \nabla F(\mathbf{x}_k) / \mathbf{v}_k). \end{aligned}$$

Then considering  $\eta$  is very small in practice, one can approximate  $\frac{1}{1 + \lambda_k \eta} \approx 1 - \lambda_k \eta$ , and the factor  $\lambda_k \eta^2$  for the term  $F(\mathbf{x}_k) / \mathbf{v}_k$  is too small and can be ignored compared with  $\eta$ . Finally, in stochastic setting, one can use the gradient estimation  $\mathbf{m}_k$  to estimate full gradient  $\nabla F(\mathbf{x}_k)$ , and thus achieves  $\mathbf{x}_{k+1} = (1 - \lambda_k \eta) \mathbf{x}_k - \eta \mathbf{m}_k / \mathbf{v}_k$  which accords with the update (2) of AdamW. From this process, one can also observe that the dynamic regularizer  $\frac{\lambda}{2} \|\mathbf{x}\|_{\mathbf{v}_k}^2$  is induced by the decoupled weight decay  $-\lambda_k \eta \mathbf{x}_k$  in AdamW. In the following, we will show that AdamW indeed minimizes the dynamic function  $F_k(\mathbf{x})$  instead of the vanilla loss  $F(\mathbf{x})$ .

## 4.3 Results on General Nonconvex Problems

Following many works which analyze adaptive gradient algorithms [16], [18], [21], [41], [42], we first provide the convergence results of AdamW by using a constant learning rate  $\eta$ .

**Theorem 2.** *Suppose that Assumptions 1 and 2 hold. Let  $\mathbf{x}_* \in \operatorname{argmin}_{\mathbf{x}} F(\mathbf{x})$ ,  $\Delta = F(\mathbf{x}_0) - F(\mathbf{x}_*)$ ,  $\eta \leq \frac{\delta^{1.25} b \epsilon^2}{6(c_\infty^2 + \delta)^{0.75} \sigma^2 L}$ ,  $\beta_1 \leq \frac{\delta^{0.5} b \epsilon^2}{3(c_\infty^2 + \delta)^{0.5} \sigma^2}$  and  $\beta_2 \in (0, 1)$  for all iterations, and  $\lambda_k = \lambda(1 - \frac{\beta_2 c_\infty^2}{\delta})^k$  with a constant  $\lambda$ . After  $T = \mathcal{O}(\max(\frac{c_\infty^2 L \Delta \sigma^2}{\delta^{1.25} b \epsilon^4}, \frac{c_\infty^2 \sigma^4}{\delta b^2 \epsilon^4}))$  iterations, the sequence  $\{\mathbf{x}_k\}_{k=0}^T$  of AdamW in Eqn. (2) obeys*

$$\begin{aligned} \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \left[ \|\nabla F_k(\mathbf{x}_k)\|_2^2 \right] \leq \epsilon^2, \quad \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k+1}\|_{\mathbf{v}_k}^2 \right] \leq \frac{\eta^2 \epsilon^2}{4}, \\ \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \left[ \|\mathbf{m}_k - \nabla F(\mathbf{x}_k)\|_2^2 \right] \leq 8\epsilon^2. \end{aligned} \quad (5)$$

Moreover, the total stochastic gradient complexity to achieve (5) is  $\mathcal{O}(\max(\frac{c_\infty^2 L \Delta \sigma^2}{\delta^{1.25} b \epsilon^4}, \frac{c_\infty^2 \sigma^4}{\delta b^2 \epsilon^4}))$ .

See its proof in Appendix G.2. Theorem 2 shows the convergence of AdamW on the nonconvex problems. Within  $T = \mathcal{O}(\max(\frac{c_\infty^2 L \Delta \sigma^2}{\delta^{1.25} b \epsilon^4}, \frac{c_\infty^2 \sigma^4}{\delta b^2 \epsilon^4}))$  iterations, the average gradient  $\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \left[ \|\nabla F_k(\mathbf{x}_k)\|_2^2 \right]$  is smaller than  $\epsilon^2$ , indicating the convergence of AdamW. Now we show small  $\|\nabla F_k(\mathbf{x}_k)\|_2$  guarantees small  $\|\nabla F(\mathbf{x}_k)\|_2$  in Corollary 1 with proof in Appendix G.3.

**Corollary 1.** *Assume that  $\|\mathbf{v}_k\|_2 \leq \rho' \|\nabla F(\mathbf{x}_k)\|_2$  with a constant  $\rho' > 0$ , and  $1 > \lambda_k \rho' \|\mathbf{x}_k\|_\infty$ . We have  $\|\nabla F(\mathbf{x}_k)\|_2 \leq \frac{1}{1 - \lambda_k \rho' \|\mathbf{x}_k\|_\infty} \|\nabla F_k(\mathbf{x}_k)\|_2$ .*

The assumptions in Corollary 1 are mild. As  $\mathbf{n}_k$  is the moving average of stochastic square version of full gradient  $\nabla F(\mathbf{x}_k)$ , one can assume  $\|\mathbf{n}_k\|_2 \leq \rho \|\nabla F(\mathbf{x}_k)\|_2^2$ , especially for the late training phase where  $\mathbf{x}_k$  is updated very slowly. Indeed, this assumption is validated in Adam analysis works, e.g., [9]. Specifically, since  $\delta$  is extremely small in  $\mathbf{v}_k = \sqrt{\mathbf{n}_k + \delta}$ , one can find a constant  $\rho' \approx \rho$  so that  $\|\mathbf{v}_k\|_2 \leq \|\nabla F(\mathbf{x}_k)\|_2$ . For assumption  $1 > \lambda_k \rho' \|\mathbf{x}_k\|_\infty$ , it is mild, since a)  $\lambda_k$  is often very small in practice, e.g.,  $10^{-4}$ , and b) the magnitude  $\|\mathbf{x}_k\|_\infty$  of network parameter is not large as observed and proved in [43] because of the auto-adaptive tradeoff among the parameter magnitude at different layers. Also, we empirically find  $\|\mathbf{x}_k\|_\infty \approx 8.0$  in the well-trained ViT-small across different training epoch numbers. Indeed, for  $\rho'$ , Zhou et al. [9] empirically finds it around 1.0 on CNNs (see their Fig. 2).

The second inequality in Eqn. (5) guarantees the small distance between two neighboring solutions  $\mathbf{x}_k$  and  $\mathbf{x}_{k+1}$ , also showing the good convergence behaviors of AdamW. The last inequality in Eqn. (5) reveals that the exponential moving average (EMA)  $\mathbf{m}_k$  of all historical stochastic gradient is close to the full gradient  $\nabla F(\mathbf{x}_k)$  and explains the success of EMA gradient estimation.

Besides, in Theorem 2, to find an  $\epsilon$ -accurate first-order stationary point ( $\epsilon$ -ASP), the stochastic gradient complexity of AdamW is  $\mathcal{O}(c_\infty^{2.5} \epsilon^{-4})$  which matches the lower bound  $\Omega(\epsilon^{-4})$  in [15] (up to constant factors). Moreover, AdamW enjoys lower complexity than Adabelief [21] of  $\mathcal{O}(c_2^6 \epsilon^{-4})$  and LAMB [44] of  $\mathcal{O}(c_2 \sqrt{d} \epsilon^{-4})$ , especially on over-parameterized networks, where  $c_2$  upper bounds the  $\ell_2$ -norm of stochastic gradient. This is because for the  $d$ -dimensional gradient, its  $\ell_\infty$ -norm  $c_\infty$  is often much smaller than its  $\ell_2$ -norm  $c_2$ , and can be  $\sqrt{d} \times$  smaller for the best case. Appendix D discusses the proof technique differences among ours and the above works. One can extend the results in Theorem 2 to  $\ell_2$ -Adam. See the proof of Corollary 2 in Appendix G.4.

**Corollary 2.** *With the same parameter settings in Theorem 2, to achieve (5), the total stochastic gradient complexity of Adam and  $\ell_2$ -Adam is  $\mathcal{O}(\max(\frac{c_\infty^{2.5} L \Delta \sigma^2}{\delta^{1.25} b \epsilon^4}, \frac{c_\infty^2 \sigma^4}{\delta b^2 \epsilon^4}))$ .*

Corollary 2 shows that the complexity of Adam and  $\ell_2$ -Adam is  $\mathcal{O}(c_\infty^{2.5} \epsilon^{-4})$ , and is superior than the previously known complexity  $\mathcal{O}(c_\infty \sqrt{d} \epsilon^{-4})$  of Adam-type optimizers analyzed in [16]–[18], e.g., ( $\ell_2$ -)Adam, AdaGrad [34], AdaBound [8]. Though sharing the same complexity with Adam and  $\ell_2$ -Adam, AdamW separates the  $\ell_2$ -regularizer with the loss objective via the decoupled weight decay whose generalization benefits have been validated empirically in many works, e.g., [12], and theoretically in our Sec. 5.

Now we investigate the convergence performance of AdamW when using a decayed learning rate  $\eta_k$ . Compared with the constant learning rate, this decay strategy is more widely used in practice, but is rarely investigated in other optimization analysis (e.g., [16], [21], [44]) except for [18]. Theorem 2 states our main results.

**Theorem 3.** *Suppose that Assumptions 1 and 2 hold. Let  $\eta_k = \frac{\gamma \delta^{0.75}}{2(c_\infty^2 + \delta)^{0.25} L \sqrt{k+1}}$ ,  $\beta_{1k} = \frac{\gamma}{\sqrt{k+1}}$ ,  $\beta_{2k} = \beta_2 \in (0, 1)$  with  $\gamma = \max(1, \frac{c_\infty^{0.25} L^{0.5} \Delta^{0.5}}{\delta^{0.125} \sigma})$ , and  $\lambda_k = \lambda(1 - \frac{\beta_2 c_\infty^2}{\delta})^k$  with a constant  $\lambda$  for the  $k$ -th training iteration. To achieve the results in Eqn. (5) with  $\eta$  replaced by  $\eta_1$ , the stochastic gradient complexity of AdamW in Eqn. (2) is  $\mathcal{O}(\max(\frac{c_\infty^{1.25} L^{0.5} \Delta^{0.5} \sigma}{\delta^{0.625} \epsilon^4} \log(\frac{1}{\epsilon}), \frac{c_\infty \sigma^2}{\delta^{0.5} \epsilon^4} \log(\frac{1}{\epsilon})))$ .*

See its proof in Appendix G.5. Theorem 3 shows that with decaying learning rate  $\eta_k = \frac{1}{\sqrt{k+1}}$ , AdamW converges and shares almost the same results in Theorem 2 where it uses constant learning rate. To achieve  $\epsilon$ -ASP, the complexity of AdamW with decaying learning rate is  $\mathcal{O}\left(\max\left(\frac{c_\infty^{1.25} L^{0.5} \Delta^{0.5} \sigma}{\delta^{0.625} \epsilon^4} \log\left(\frac{1}{\epsilon}\right), \frac{c_\infty \sigma^2}{\delta^{0.5} \epsilon^4} \log\left(\frac{1}{\epsilon}\right)\right)\right)$  and slightly differs from the one  $\mathcal{O}\left(\max\left(\frac{c_\infty^{2.5} L \Delta \sigma^2}{\delta^{1.25} \epsilon^4}, \frac{c_\infty^2 \sigma^4}{\delta b \epsilon^4}\right)\right)$  of AdamW using constant learning rate. By comparing each complexity term, decaying learning rate respectively improves the constant one by factors  $\frac{c_\infty^{1.25} L^{0.5} \Delta^{0.5} \sigma}{\delta^{0.625} \epsilon^4} \log^{-1}\left(\frac{1}{\epsilon}\right)$  and  $\frac{c_\infty^2 \sigma^2}{\delta^{0.5} \epsilon^4} \log^{-1}\left(\frac{1}{\epsilon}\right)$ . Consider that  $\frac{c_\infty^{1.25} L^{0.5} \Delta^{0.5} \sigma}{\delta^{0.625} \epsilon^4}$  and  $\frac{c_\infty \sigma^2}{\delta^{0.5} \epsilon^4}$  are often large than  $\log\left(\frac{1}{\epsilon}\right)$ , as the  $\ell_1$ -norm upper bound  $c_\infty$  of stochastic gradient is often not small and  $\delta$  is very small, e.g.,  $10^{-4}$  by default, decaying learning rate is superior than constant one which accords with the practical observations. When 1)  $\lambda_k=0$  or 2) the loss  $F(\mathbf{x})$  is a  $\ell_2$ -regularized loss, Theorem 3 still holds. So the stochastic complexity in Theorem 3 is applicable to  $\ell_2$ -Adam. Guo et al. [18] proved the complexity  $\mathcal{O}\left(\max\left(\frac{c_\infty^{2.5} L^2 \sigma^2}{\delta^{2.5} \epsilon^4} \log\left(\frac{1}{\epsilon}\right), \frac{c_\infty^2 \sigma^4}{\delta^2 \epsilon^4} \log\left(\frac{1}{\epsilon}\right)\right)\right)$  of Adam-type algorithms, e.g., Adam and  $\ell_2$ -Adam, with decaying learning rate, which but is inferior than the complexity in this work, since as aforementioned,  $\delta$  is often very small.

#### 4.4 Results on PL-conditioned Nonconvex Problems

In this work, we are also particularly interested in the nonconvex problems under PL condition, since as observed or proved in [37], [38], deep learning models often satisfy PL condition at least around a local minimum. For this special nonconvex problem, we follow [18], and divide the whole optimization into  $K$  stages. Specifically, for constant learning rate setting, AdamW uses learning rate  $\eta_k$  in the whole  $k$ -th stage; while for decayed learning rate setting, it uses a decayed  $\eta_{k_i}$  for the  $k$ -th stage which satisfies  $\eta_{k_i} < \eta_{k_j}$  if  $i > j$ , where  $\eta_{k_i}$  denotes the learning rate of the  $i$ -th iteration of the  $k$ -th stage. Moreover, for both learning rate settings, at the  $k$ -th stage, AdamW is allowed to run  $T_k$  iterations for achieving  $\mathbb{E}[F_k(\mathbf{x}_k) - F_k(\mathbf{x}_*)] \leq \epsilon_k$ , where  $\mathbf{x}_* \in \operatorname{argmin}_{\mathbf{x}} F(\mathbf{x})$ ,  $\mathbf{x}_k$  is the output of the  $k$ -stage and  $\epsilon_k = \frac{1}{2k}[F_0(\mathbf{x}_0) - F_0(\mathbf{x}_*)]$  denotes the optimization accuracy. See detailed Algorithm 2 in Appendix B. At below, we provide the convergence results of AdamW under both settings of constant or decayed learning rate in Theorem 4 with proof in appendix G.6.

**Theorem 4.** Suppose Assumptions 1 and 2 hold, and  $\mathbf{x}_* \in \operatorname{argmin}_{\mathbf{x}} F(\mathbf{x})$ . Assume the loss  $F_k(\mathbf{x}_k)$  in (4) and  $F_k(\mathbf{x}_*)$  satisfy the PL condition in Assumption 3.

1) For constant learning rate setting, assume a constant learning rate  $\eta_k \leq \frac{\delta^{1.25} \mu b \epsilon_k}{12(c_\infty^2 + \delta)^{0.75} \sigma^2 L}$ , constant  $\beta_{1k} \leq \frac{\delta^{0.5} \mu b \epsilon_k}{6(c_\infty^2 + \delta)^{0.5} \sigma^2}$ ,  $\beta_{2k} \in (0, 1)$  and  $\lambda_k = \lambda(1 - \frac{\beta_2 c_\infty^2}{\delta})^k$  at the  $k$ -th stage. We have:

1.1) For the  $k$ -th stage, AdamW runs at most  $T_k = \mathcal{O}\left(\max\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} b \epsilon_k}, \frac{c_\infty^2 \sigma^2}{\mu \delta b \epsilon_k}\right)\right)$  iterations to achieve  $\mathbb{E}[F_k(\mathbf{x}_k) - F_k(\mathbf{x}_*)] \leq \epsilon_k$ , where the output  $\mathbf{x}_k$  is uniformly randomly selected from the sequence  $\{\mathbf{x}_{k_i}\}_{i=1}^{T_k}$  at the  $k$ -th stage.

1.2) For  $K$  stages, the total stochastic complexity is  $\mathcal{O}\left(\max\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} \epsilon}, \frac{c_\infty^2 \sigma^2}{\mu \delta \epsilon}\right)\right)$  to achieve

$$\min_{1 \leq k \leq K} \mathbb{E}[F_k(\mathbf{x}_k) - F_k(\mathbf{x}_*)] \leq \epsilon. \quad (6)$$

2) For decaying learning rate setting, let  $\eta_{k_i} \leq \frac{\gamma \delta^{0.75}}{2(c_\infty^2 + \delta)^{0.25} L \sqrt{i+1}}$ ,  $\beta_{1k_i} \leq \frac{\gamma}{\sqrt{i+1}}$ ,  $\beta_{2k_i} = \beta_{2k} \in (0, 1)$ ,  $\lambda_{k_i} = \lambda(1 - \frac{\beta_2 c_\infty^2}{\delta})^i$  at the  $i$ -th iteration of the  $k$ -th stage with  $\gamma = \max\left(1, \frac{(c_\infty^2 + \delta)^{0.125} L^{0.5} b^{0.5} \epsilon_k^{0.5}}{\delta^{0.125} \sigma}\right)$ .

2.1) For the  $k$ -th stage, AdamW runs at most  $T_k = \mathcal{O}\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} b \epsilon}\right)$  iterations to achieve  $\mathbb{E}[F_k(\mathbf{x}_k) - F_k(\mathbf{x}_*)] \leq \epsilon_k$ , where the output  $\mathbf{x}_k$  is randomly selected from the sequence  $\{\mathbf{x}_{k_i}\}_{i=1}^{T_k}$  at the  $k$ -th stage according to the distribution  $\left\{\frac{\eta_{k_i}}{\sum_{j=1}^{T_k} \eta_{k_j}}\right\}_{i=1}^{T_k}$ .

2.2) The total complexity is  $\mathcal{O}\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} \epsilon}\right)$  to achieve (6).

Theorem 4 shows that AdamW can converge under both constant and decaying learning rate settings. Moreover, by comparison, to achieve  $\epsilon$ -ASP in Eqn. (6), the decaying learning rate has the total complexity  $\mathcal{O}\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} \epsilon}\right)$ , and could be better than the constant learning rate whose complexity is  $\mathcal{O}\left(\max\left(\frac{c_\infty^{2.5} L \sigma^2}{\mu^2 \delta^{1.25} \epsilon}, \frac{c_\infty^2 \sigma^2}{\mu \delta \epsilon}\right)\right)$ . It should be also noted that the complexity of AdamW on this special nonconvex problems (i.e. with PL condition) enjoys lower complexity than the one on the general nonconvex problems, since PL condition ensures a convexity-alike landscape of the loss objective and thus can be optimized faster.

## 5 GENERALIZATION ANALYSIS

### 5.1 Generalization Results

**Analysis on hypothesis posterior.** As shown in the classical PACBayesian framework [30], [45] there is strong relations between the generalization error bound and the hypothesis posterior learned by an algorithm. So we first analyze the hypothesis posterior learned by AdamW, and then investigate the generalization error of AdamW. Specifically, following [9], [22]–[24], [46], we study the corresponding stochastic differential equations (SDEs) of an algorithm to investigate its posterior and generalization behaviors because of the similar convergence behaviors of an algorithm and its SDE. Firstly, the updating rule of AdamW can be formulated as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{Q}_t \nabla F(\mathbf{x}_t) - \eta \lambda \mathbf{x}_t + \eta \mathbf{Q}_t \mathbf{u}_t, \quad (7)$$

where  $\mathbf{u}_t = \nabla F(\mathbf{x}_t) - \mathbf{m}_t$  is gradient noise,  $\mathbf{Q}_t = \operatorname{diag}\left(\mathbf{n}_t^{-\frac{1}{2}}\right)$  is a diagonal matrix. In Eqn. (7), the small  $\delta$  in Eqn. (2) is ignored for convenience which does not affect our following results. Then following [23], [47], [48], we assume the gradient noise  $\mathbf{u}_t$  obeys Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{x}_t})$  because of the Central Limit theory. Accordingly, one can write the SDE of AdamW as

$$d\mathbf{x}_t = -\mathbf{Q}_t \nabla F(\mathbf{x}_t) dt - \lambda \mathbf{x}_t dt + \mathbf{Q}_t (2\Sigma_t)^{\frac{1}{2}} d\zeta_t,$$

where  $d\zeta_t \sim \mathcal{N}(0, \mathbf{Id}t)$  and  $\Sigma_t = \frac{\eta}{2} \mathbf{C}_{\mathbf{x}_t}$ . Here  $\mathbf{C}_{\mathbf{x}_t}$  is defined as  $\mathbf{C}_{\mathbf{x}_t} = \frac{1}{b} \left[ \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{x}_t; \zeta_i) \nabla f(\mathbf{x}_t; \zeta_i)^\top - \nabla F(\mathbf{x}_t) \nabla F(\mathbf{x}_t)^\top \right]$ , where  $n$  is the training sample number, and  $b$  is minibatch size. For analysis, we make some necessary assumptions.

**Assumption 4.** a) Assume  $\mathbf{C}_{\mathbf{x}_t}$  can approximate the Fisher matrix  $\mathbf{F}_{\mathbf{x}_t} = \frac{1}{n} \sum_{i=1}^n \nabla F(\mathbf{x}_t; \zeta_i) \nabla F(\mathbf{x}_t; \zeta_i)^\top$ , i.e.,  $\mathbf{C}_{\mathbf{x}_t} \approx \mathbf{F}_{\mathbf{x}_t}$ . b) Assume  $\mathbf{F}_{\mathbf{x}_t}$  can approximate the Hessian matrix  $\mathbf{H}_{\mathbf{x}_t}$  near a minimum, i.e.,  $\mathbf{F}_{\mathbf{x}_t} \approx \mathbf{H}_{\mathbf{x}_t}$ . c) Suppose  $\mathbf{n}'_{t+1} = (1 - \beta_2) \mathbf{n}'_t + \beta_2 \mathbf{g}_t \mathbf{g}_t^\top$  (virtual sequence) with  $\mathbf{n}'_0 = \mathbf{g}_0 \mathbf{g}_0^\top$  is a good estimation to  $\mathbf{F}_{\mathbf{x}_t}$ , i.e.,  $\mathbf{n}'_{t+1} \approx \mathbf{F}_{\mathbf{x}_t}$ .

Assumption 4 is widely used. Specifically, we follow [23], [47], [48], and approximate  $\mathbf{C}_{\mathbf{x}_t} \approx \mathbf{F}_{\mathbf{x}_t}$ , since we analyze the local convergence around an optimum, leading to 1)  $\nabla F(\mathbf{x}_t) \approx 0$  and 2) a dominated variance of gradient noise. Assumption 4 b) is used in [24], [49] for analysis, and holds when  $\mathbf{x}_t$  is around a minimum. Since most works analyze the generalization performance of an algorithm around a local minimum, e.g., [9], [23], [24], [46], [47],

[47], [48], [50], Assumption 4 b) holds in their setting and thus is mild. For Assumption 4 c), Staib et al. [51] proved that the matrix-based second-order moment  $\mathbf{n}'_t$  is a good estimation to the Fisher matrix  $\mathbf{F}_{\mathbf{x}_t}$  after running a certain iteration number. Please refer to the theoretical details of Assumption 4 in Appendix E. Then we can derive the hypothesis posterior learnt by AdamW.

**Lemma 5.** *Assume the loss can be approximated by a second-order Taylor approximation, i.e.,  $F(\mathbf{x}) \approx F(\mathbf{x}_*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^\top \mathbf{H}_*(\mathbf{x} - \mathbf{x}_*)$  where  $\mathbf{H}_*$  is systemic. With Assumption 4, the solution  $\mathbf{x}_t$  of AdamW obeys a Gaussian distribution  $\mathcal{N}(\mathbf{x}_*, \mathbf{M}_{AdamW})$  where the covariance matrix  $\mathbf{M}_{AdamW} = \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top]$  is defined as*

$$\mathbf{M}_{AdamW} = \frac{\eta}{2b} (\mathbf{Q}\mathbf{H}_* + \lambda\mathbf{I})^{-1} \mathbf{Q}\mathbf{H}_*\mathbf{Q}.$$

where  $\mathbf{Q} = \text{diag}[\mathbf{H}_{*(11)}^{-\frac{1}{2}}, \mathbf{H}_{*(22)}^{-\frac{1}{2}}, \dots, \mathbf{H}_{*(dd)}^{-\frac{1}{2}}]$  is diagonal matrix.

See its proof in Appendix H.1. Lemma 5 tells that AdamW can converge to a solution which concentrates around the minimum  $\mathbf{x}_*$ . This also guarantees the good convergence behaviors of AdamW but from an SDE aspect. From the covariance matrix  $\mathbf{M}_{AdamW}$ , one can see that all singular values of  $\mathbf{M}_{AdamW}$  become smaller when increases and is large enough to ensure  $\mathbf{Q}\mathbf{H}_* + \lambda\mathbf{I} \succeq \mathbf{0}$ . This indicates that proper weight decay in AdamW can stabilize the algorithm, and benefits its convergence to the minimizer  $\mathbf{x}_*$ .

**Generalization analysis.** Based on the above posterior analysis, we employ the PAC Bayesian framework [30] to explicitly analyze the generalization performance of AdamW. Given an algorithm  $\mathcal{A}$  and a training dataset  $\mathcal{D}_{tr}$  whose samples  $\xi$  are drawn from an unknown distribution  $\mathcal{D}$ , one often trains a model to obtain a posterior hypothesis  $\mathbf{x}$  drawn from a hypothesis distribution  $\mathcal{P} \sim \mathcal{N}(\mathbf{x}_*, \mathbf{M}_{AdamW})$  in Lemma 5. Then we denote the expected risk w.r.t. the hypothesis distribution  $\mathcal{P}$  as  $\mathbb{E}_{\xi \sim \mathcal{D}, \mathbf{x} \sim \mathcal{P}}[f(\mathbf{x}, \xi)]$  and the empirical risk w.r.t. the distribution  $\mathcal{P}$  as  $\mathbb{E}_{\xi \in \mathcal{D}_{tr}, \mathbf{x} \sim \mathcal{P}}[f(\mathbf{x}, \xi)]$ . In practice, one often assumes that the prior hypothesis satisfies Gaussian distribution  $\mathcal{P}_{pre} \sim \mathcal{N}(\mathbf{0}, \rho\mathbf{I})$  [13], [50], [52], since we do not know any information on the posterior hypothesis. Based on Lemma 5, we can derive the generalization error bound of AdamW.

**Theorem 6.** *Assume that  $\mathbf{x}_0$  satisfies  $\mathcal{P}_{pre} \sim \mathcal{N}(\mathbf{0}, \rho\mathbf{I})$ . Then with at least probability  $1 - \tau$  ( $\tau \in (0, 1)$ ), the expected risk for the posterior hypothesis  $\mathbf{x} \sim \mathcal{P}$  of AdamW learned on training dataset  $\mathcal{D}_{tr} \sim \mathcal{D}$  with  $n$  samples holds*

$$\mathbb{E}_{\xi \sim \mathcal{D}, \mathbf{x} \sim \mathcal{P}}[f(\mathbf{x}, \xi)] - \mathbb{E}_{\xi \in \mathcal{D}_{tr}, \mathbf{x} \sim \mathcal{P}}[f(\mathbf{x}, \xi)] \leq \Phi_{AdamW},$$

where  $\Phi_{AdamW} = \frac{\sqrt{8}}{\sqrt{n}} (\bar{\text{err}}_{AdamW} + c_0)^{\frac{1}{2}}$  with  $\bar{\text{err}}_{AdamW} = -\log \det(\mathbf{M}_{AdamW}) + \frac{\eta}{2\rho b} \text{Tr}(\mathbf{M}_{AdamW}) + d \log \frac{2b\rho}{\eta}$ ,  $c_0 = \frac{1}{2\rho} \|\mathbf{x}_*\|^2 - \frac{d}{2} + 2 \ln \left( \frac{2n}{\tau} \right)$ . Here  $\det(M)$  and  $\text{tr}(M)$  denote the determinant and trace of matrix  $M$  respectively.

See its proof in Appendix H.2. Theorem 6 shows that the generalization error of AdamW is upper bounded by  $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$  (up to other factors) which matches the error bound in [53]–[56] derived from the PAC theory or stability aspects. When  $\lambda$  is large, the first term  $-\log \det(\mathbf{M}_{AdamW})$  in  $\mathbf{M}_{AdamW}$  becomes larger since the singular values of  $\mathbf{M}_{AdamW}$  become small, and leads to small  $\det(\mathbf{M}_{AdamW})$ , while the second term  $\frac{\eta}{2\rho b} \text{Tr}(\mathbf{M}_{AdamW})$  is small. But for small  $\lambda$ , the first term  $-\log \det(\mathbf{M}_{AdamW})$  is small, while the second term becomes large. Though it is hard to precisely decide the best  $\lambda$ , from the above discussion, at least we know that tuning  $\lambda$  can yield smaller generalization error, partly explaining the better performance of AdamW over vanilla Adam ( $\lambda = 0$ ).

## 5.2 Comparison with $\ell_2$ -regularized Adam

Now we compare AdamW with  $\ell_2$ -Adam. To diminish the effects of historical gradient to the current optimization and also analyze the effects of current gradient to the behaviors of adaptive algorithms, many works, e.g., [57], [58], set  $\beta_1 = 1$  in (2) to focus on concurrent optimization process of adaptive algorithms. Here we follow this setting to investigate  $\ell_2$ -Adam with updating rule:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{Q}_t (\nabla F(\mathbf{x}_t) + \lambda \mathbf{x}_t) + \eta \mathbf{Q}_t \mathbf{u}_t,$$

where  $\mathbf{u}_t = \nabla F(\mathbf{x}_t) - \mathbf{m}_t$  and  $\mathbf{Q}_t = \text{diag}(\mathbf{n}_t^{-\frac{1}{2}})$  have the same meanings in Eqn. (7). Then one can write the SDE of  $\ell_2$ -Adam:

$$d\mathbf{x}_t = -\mathbf{Q}_t (\nabla F(\mathbf{x}_t) + \lambda \mathbf{x}_t) dt + \mathbf{Q}_t (2\Sigma_t)^{\frac{1}{2}} d\zeta_t,$$

where  $d\zeta_t \sim \mathcal{N}(0, \mathbf{I} dt)$ ,  $\Sigma_t = \frac{\eta}{2} \mathbf{C}_{\mathbf{x}_t}$  and  $\mathbf{C}_{\mathbf{x}_t}$  is given above.

**Theorem 7.** *Assume  $\mathbf{x}_0$  satisfies  $\mathcal{P}_{pre} \sim \mathcal{N}(\mathbf{0}, \rho\mathbf{I})$ . With at least probability  $1 - \tau$  and a constant  $c_0$  in Theorem 6, the expected risk for the posterior hypothesis  $\mathbf{x} \sim \mathcal{P}_{\ell_2\text{-Adam}}$  of  $\ell_2$ -Adam learned on training dataset  $\mathcal{D}_{tr} \sim \mathcal{D}$  with  $n$  samples can be upper bounded:*

$$\mathbb{E}_{\xi \sim \mathcal{D}, \mathbf{x} \sim \mathcal{P}_{\ell_2\text{-Adam}}}[f(\mathbf{x}, \xi)] - \mathbb{E}_{\xi \in \mathcal{D}_{tr}, \mathbf{x} \sim \mathcal{P}}[f(\mathbf{x}, \xi)] \leq \Phi_{\ell_2\text{-Adam}},$$

where  $\Phi_{\ell_2\text{-Adam}} = \frac{\sqrt{8}}{\sqrt{n}} (\bar{\text{err}}_{\ell_2\text{-Adam}} + c_0)^{\frac{1}{2}}$  with  $\bar{\text{err}}_{\ell_2\text{-Adam}} = -\log \det(\mathbf{M}_{AdamW}) + \frac{\eta}{2\rho b} \text{Tr}(\mathbf{M}_{\ell_2\text{-Adam}}) + d \log \frac{2b\rho}{\eta}$ .

See its proof in Appendix H.3. Theorem 7 shows the generalization error bound  $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$  of  $\ell_2$ -Adam. Moreover, when  $\lambda = 0$ , AdamW and  $\ell_2$ -Adam are exactly the same, and their error bounds are also the same as shown in Theorems 6 and 7.

Next, we compare the generalization error bounds of AdamW and  $\ell_2$ -Adam. To this end, we follow the similar spirit in [9] and approximate  $\mathbf{Q} \approx \mathbf{H}_*^{-\frac{1}{2}}$  to simplify  $\Phi_{AdamW}$  and  $\Phi_{\ell_2\text{-Adam}}$  in the Corollary 3 whose proof can be found in Appendix H.4.

**Corollary 3.** *Assume  $\mathbf{Q} \approx \mathbf{H}_*^{-\frac{1}{2}}$ . Then we have*

$$\Phi_{AdamW} \approx \frac{\sqrt{8}}{\sqrt{n}} (\text{err}_{AdamW} + c_0)^{\frac{1}{2}}, \quad \Phi_{\ell_2\text{-Adam}} \approx \frac{\sqrt{8}}{\sqrt{n}} (\text{err}_{\ell_2\text{-Adam}} + c_0)^{\frac{1}{2}},$$

where  $\text{err}_{AdamW} = \sum_{i=1}^d h(x_{AdamW}^{(i)})$  with  $x_{AdamW}^{(i)} = 2\eta^{-1} \rho b (\sigma_i^{\frac{1}{2}} + \lambda)$ ,  $\text{err}_{\ell_2\text{-Adam}} = \sum_{i=1}^d h(x_{\ell_2\text{-Adam}}^{(i)})$  with  $x_{\ell_2\text{-Adam}}^{(i)} = 2\eta^{-1} \rho b (\sigma_i^{\frac{1}{2}} + \lambda \sigma_i^{-\frac{1}{2}})$ . Here  $h(x) = \log x + \frac{1}{x}$ .

Then we only need to compare the different terms, i.e.,  $\text{err}_{AdamW}$  and  $\text{err}_{\ell_2\text{-Adam}}$ . For  $h(x)$ , since  $h'(x) = \frac{x-1}{x^2}$ ,  $h(x)$  will increase when  $x \in (1, +\infty)$ . Meanwhile, generally, we have  $x_{\ell_2\text{-Adam}}^{(i)} > x_{AdamW}^{(i)} > 1$  for most  $i \in [d]$  due to three reasons. 1) Most of the singular values  $\{\sigma_i\}_{i=1}^d$  of Hessian matrix in deep networks are much smaller than one which is well observed in many works, e.g., fully connected networks, AlexNet, VGG and ResNet [49], [59]–[61] and our experimental results on ResNet50 and ViT-small in Fig. 1. 2) The learning rate when reaching the minimum is set to be very small in practice. 3) The minibatch size  $b$  is often thousand to train a modern network, and the variance  $\rho$  for the initialization distribution  $\mathcal{P}_{pre} \sim \mathcal{N}(\mathbf{0}, \rho\mathbf{I})$  is often of the order  $\mathcal{O}(1/\sqrt{d_i})$  [62], where  $d_i$  is input dimension. These factors indicate  $x_{\ell_2\text{-Adam}}^{(i)} > x_{AdamW}^{(i)} > 1$ . So the generalization error term  $\text{err}_{AdamW}$  is smaller than  $\text{err}_{\ell_2\text{-Adam}}$ , testified by our experimental results on ResNet50 and ViT-small in Sec. 6. So AdamW often enjoys better generalization performance than  $\ell_2$ -Adam, also validated in Sec. 6. Appendix C intuitively discusses the generalization benefits of coordinate-adaptive regularization in AdamW.

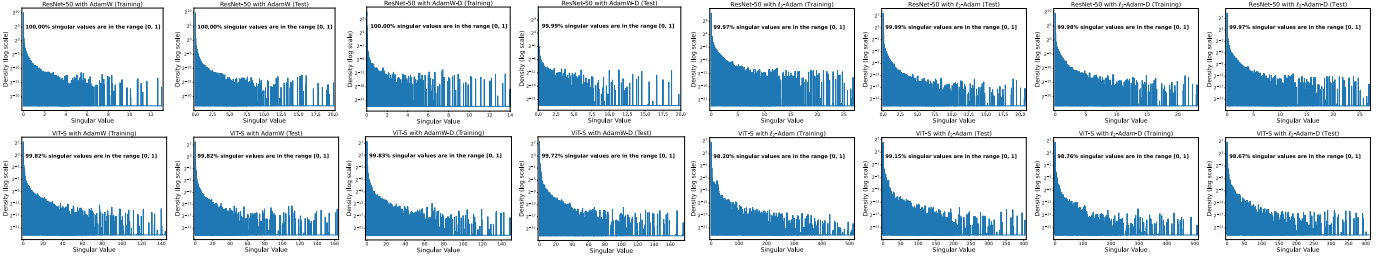


Fig. 1: Visualization of singular values in ResNet50 and ViT-small trained by AdamW (constant weight decay), AdamW-D (decreasing weight decay),  $\ell_2$ -Adam (constant weight decay) and  $\ell_2$ -Adam-D (decreasing weight decay). See more visualization results, e.g., ResNet18, in Fig. 3 of Appendix A.

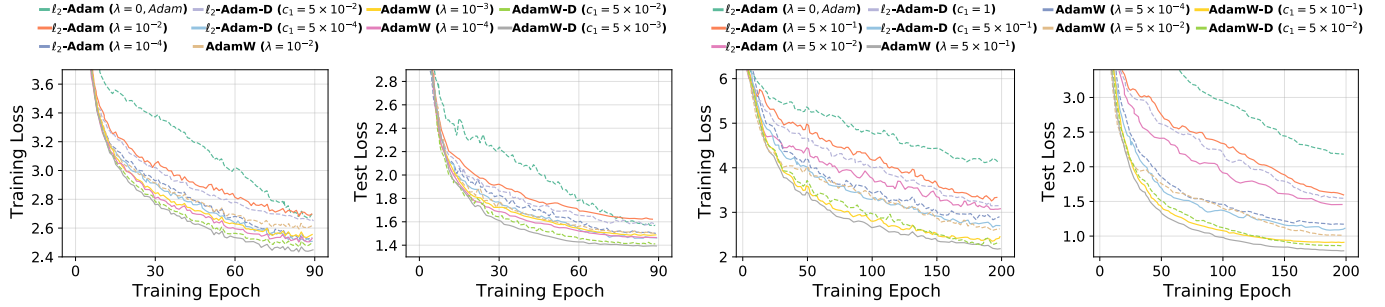


Fig. 2: Training and test curves of  $\ell_2$ -Adam,  $\ell_2$ -Adam-D, AdamW and AdamW-D on ImageNet. See more results in Appendix A.

model train epoch optimizer	ResNet18		ResNet50		ViT-small		ViT-small		ViT-small	
	90	90	100	100	100	100	200	200	300	300
err in bound	AdamW/D	$\ell_2$ -Adam/D	AdamW/D	$\ell_2$ -Adam/D	AdamW/D	$\ell_2$ -Adam/D	AdamW/D	$\ell_2$ -Adam/D	AdamW/D	$\ell_2$ -Adam/D
	3.43 / 3.40	3.85 / 3.82	3.42 / 3.41	3.78 / 3.77	3.62 / 3.63	3.75 / 3.76	3.58 / 3.57	3.72 / 3.71	3.47 / 3.45	3.70 / 3.69
test acc. (%)	67.9 / 70.1	67.2 / 67.4	77.0 / 77.1	76.5 / 76.4	76.1 / 75.9	75.3 / 75.4	79.2 / 79.3	77.6 / 77.7	79.8 / 80.0	78.5 / 78.6

TABLE 1: Generalization of AdamW (constant weight decay), AdamW-D (decaying weight decay),  $\ell_2$ -Adam (constant weight decay) and  $\ell_2$ -Adam-D (decreasing weight decay) on ImageNet. AdamW/D denotes AdamW/AdamW-D;  $\ell_2$ -Adam/D has the same meaning.

## 6 EXPERIMENTS

**Investigation on singular values of Hessian.** We respectively use AdamW and  $\ell_2$ -Adam to train two popular networks on ImageNet [63], i.e. ResNet50 [13] and vision transformer small (ViT-small) [3] for both 100 epochs. Then we adopt the method in [64] to estimate the singular values of these two trained networks. AdamW/ $\ell_2$ -Adam uses constant weight decay  $\lambda_k$ , while AdamW-D/ $\ell_2$ -Adam-D adopts exponentially-decaying weight decay  $\lambda_k = c_1 \cdot \lambda^k$  with two constants  $c_1 > 0$  and  $\lambda \in (0, 1)$ . Fig. 1 plots the spectral density of these singular values on training/test data of ImageNet, and shows that there more than 99% singular values are in the range  $[0, 1]$  and are much smaller than one. This accords with the observations on AlexNet, VGG and ResNet in [49], [59]–[61]. All these observations support the results in Sec. 5.2.

**Investigation on generalization.** To compute the key generalization error terms i.e.,  $\bar{\text{err}}_{\text{AdamW}}$  and  $\bar{\text{err}}_{\ell_2\text{-Adam}}$  in Theorems 6 and 7, one needs to compute the full Hessian for matrix multiplication that however is prohibitively computable. So we compute their approximations  $\text{err}_{\text{AdamW}}$  and  $\text{err}_{\ell_2\text{-Adam}}$  in Corollary 3 to compare the generalization error bounds of AdamW and  $\ell_2$ -Adam. For comprehension, we also compute  $\text{err}_{\text{AdamW-D}}$  of AdamW-D and  $\text{err}_{\ell_2\text{-Adam-D}}$  of  $\ell_2$ -Adam-D which respectively share the same formulation with  $\text{err}_{\text{AdamW}}$  and  $\text{err}_{\ell_2\text{-Adam}}$  but performs computation on the models respectively trained by AdamW-D and  $\ell_2$ -Adam-D with the above exponentially-decaying weight decay  $\lambda_k$ .

Then we respectively use AdamW, AdamW-D,  $\ell_2$ -Adam and  $\ell_2$ -Adam-D to train three models, i.e., ResNet18, ResNet50 and ViT-small, on ImageNet, and well tune their hyper-parameters, e.g., learning rate and weight decay parameter  $\lambda_k$ . Note,  $\ell_2$ -

Adam includes Adam by setting  $\lambda_k = 0$ . Next, we compute  $\text{err}_{\text{AdamW}}$ ,  $\text{err}_{\text{AdamW-D}}$ ,  $\text{err}_{\ell_2\text{-Adam}}$  and  $\text{err}_{\ell_2\text{-Adam-D}}$  on the test dataset of ImageNet, as test data can better reveal the generalization ability of an algorithm. Table 1 shows that on all test cases,  $\text{err}_{\text{AdamW}}$  and  $\text{err}_{\text{AdamW-D}}$  are smaller than  $\text{err}_{\ell_2\text{-Adam}}$  and  $\text{err}_{\ell_2\text{-Adam-D}}$  by a remarkable margin.  $\text{err}_{\text{AdamW-D}}$  and  $\text{err}_{\ell_2\text{-Adam-D}}$  respectively enjoy similar values with their corresponding  $\text{err}_{\text{AdamW}}$  and  $\text{err}_{\ell_2\text{-Adam}}$ . These results empirically support the superior generalization error of AdamW over  $\ell_2$ -Adam. Moreover, Table 1 also reveals that 1) AdamW and AdamW-D have higher test accuracy than  $\ell_2$ -Adam and  $\ell_2$ -Adam-D; 2) AdamW-D ( $\ell_2$ -Adam-D) enjoys very similar performance as AdamW ( $\ell_2$ -Adam). All these results accord with our theoretical results in Sec. 5.2.

**Investigation on convergence.** We plot the training/test curves of AdamW, AdamW-D,  $\ell_2$ -Adam and  $\ell_2$ -Adam-D on ImageNet in Fig. 2. For AdamW-D and  $\ell_2$ -Adam-D, we fix  $\lambda = 0.99999$  and tune  $c_1$  to compute its weight decay  $\lambda_k$ . One can find that on ResNet50 and ViT-small, 1) AdamW and AdamW-D show faster convergence speed than  $\ell_2$ -Adam (including Adam via  $\lambda=0$ ) and  $\ell_2$ -Adam-D when their weight decay parameter are well-tuned, e.g.,  $\lambda = 5 \times 10^{-1}$  for AdamW and  $\ell_2$ -Adam,  $c_1 = 5 \times 10^{-2}$  for AdamW-D on ViT-small; 2) AdamW and AdamW-D share similar convergence behaviors; 3) weight decay parameter greatly affects the convergence speed of the three optimizers. So under the same training cost, the faster convergence of AdamW could also partially explain its better generalization performance over  $\ell_2$ -Adam.

## 7 CONCLUSION

In this work, we first prove the convergence of AdamW using both constant and decaying learning rates on the general nonconvex

problems and PL-conditioned problems. Moreover, we find that AdamW provably minimizes a dynamically regularized loss that combines a vanilla loss and a dynamical regularization, and thus its behaviors differ from those in Adam and  $\ell_2$ -Adam. Besides, for the first time, we quantitatively justify the generalization superiority of AdamW over both Adam and  $\ell_2$ -Adam. Finally, experimental results validate the implications of our theory.

## ACKNOWLEDGEMENTS

Pan Zhou is supported by the Singapore Ministry of Education (MOE) Academic Research Fund (AcRF) Tier 1 grant. Zhouchen Lin is supported by the NSF China (No. 62276004) and the major key project of PCL, China (No. PCL2021A12).

## REFERENCES

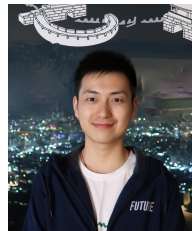
- [1] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *Int'l Conf. Learning Representations*, 2015.
- [2] Herbert Robbins and Sutton Monro, "A stochastic approximation method," *The Annals of Mathematical Statistics*, vol. 22, no. 3, pp. 400–407, 1951.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Int'l Conf. Learning Representations*, 2020.
- [4] Pan Zhou, Yichen Zhou, Chenyang Si, Weihao Yu, Teck Khim Ng, and Shuicheng Yan, "Mugs: A multi-granular self-supervised learning framework," in *arXiv preprint arXiv:2203.14415*, 2022.
- [5] Abdel-rahman Mohamed Hui Jiang Li Deng Gerald Penn Abdel-Hamid, Ossama and Dong Yu, "Convolutional neural networks for speech recognition," *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, vol. 22, no. 10, pp. 1533–1545, 2014.
- [6] Endre Süli and David F Mayers, *An introduction to numerical analysis*, Cambridge university press, 2003.
- [7] Nitish Shirish Keskar and Richard Socher, "Improving generalization performance by switching from Adam to SGD," in *Int'l Conf. Learning Representations*, 2018.
- [8] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun, "Adaptive gradient methods with dynamic bound of learning rate," in *Int'l Conf. Learning Representations*, 2018.
- [9] Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al., "Towards theoretically understanding why SGD generalizes better than adam in deep learning," in *Proc. Conf. Neural Information Processing Systems*, 2020, vol. 33, pp. 21285–21296.
- [10] Pan Zhou, Xingyu Xie, and YAN Shuicheng, "Win: Weight-decay-integrated Nesterov acceleration for adaptive gradient algorithms," in *Int'l Conf. Learning Representations*, 2022.
- [11] Ilya Loshchilov and Frank Hutter, "Decoupled weight decay regularization," in *Int'l Conf. Learning Representations*, 2018.
- [12] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou, "Training data-efficient image transformers & distillation through attention," in *Proc. Int'l Conf. Machine Learning*, 2021, pp. 10347–10357.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [14] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie, "A ConvNet for the 2020s," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2022, pp. 11966–11976.
- [15] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth, "Lower bounds for non-convex stochastic optimization," *Mathematical Programming*, vol. 199, no. 1-2, pp. 165–214, 2023.
- [16] Dongruo Zhou, Jinghui Chen, Yuan Cao, Yiqi Tang, Ziyang Yang, and Quanquan Gu, "On the convergence of adaptive gradient methods for nonconvex optimization," in *Workshop on Optimization for Machine Learning*, 2020.
- [17] Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu, "Closing the generalization gap of adaptive gradient methods in training deep neural networks," in *Proc. Int'l Joint Conf. Artificial Intelligence*, 2021, pp. 3267–3275.
- [18] Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang, "A novel convergence analysis for algorithms of the adam family," in *Workshop on Optimization for Machine Learning*, 2023.
- [19] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar, "On the convergence of Adam and beyond," in *Int'l Conf. Learning Representations*, 2019.
- [20] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong, "On the convergence of a class of adam-type algorithms for non-convex optimization," in *Int'l Conf. Learning Representations*, 2018.
- [21] Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C Tatikonda, Nicha Dvornek, Xenophon Papademetris, and James Duncan, "Adabelief optimizer: Adapting stepsizes by the belief in observed gradients," *Proc. Conf. Neural Information Processing Systems*, vol. 33, pp. 18795–18806, 2020.
- [22] Stephan Mandt, Matthew Hoffman, and David Blei, "A variational analysis of stochastic gradient algorithms," in *Proc. Int'l Conf. Machine Learning*, 2016, pp. 354–363.
- [23] Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma, "The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects," in *Proc. Int'l Conf. Machine Learning*, 2018.
- [24] Stanislaw Jastrzebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey, "Three factors influencing minima in SGD," in *Int'l Conf. Learning Representations*, 2017.
- [25] Twan Van Laarhoven, "L2 regularization versus batch and weight normalization," *arXiv preprint arXiv:1706.05350*, 2017.
- [26] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse, "Three mechanisms of weight decay regularization," in *Int'l Conf. Learning Representations*, 2018.
- [27] Elad Hoffer, Ron Banner, Itay Golan, and Daniel Soudry, "Norm matters: efficient and accurate normalization schemes in deep networks," in *Proc. Conf. Neural Information Processing Systems*, 2018, vol. 31.
- [28] Sergey Ioffe and Christian Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int'l Conf. Machine Learning*, 2015, pp. 448–456.
- [29] Yucong Zhou, Yunxiao Sun, and Zhao Zhong, "Fixnorm: Dissecting weight decay for training deep neural networks," *arXiv preprint arXiv:2103.15345*, 2021.
- [30] David A McAllester, "Some PCA-Bayesian theorems," *Machine Learning*, vol. 37, no. 3, pp. 355–363, 1999.
- [31] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin, "Emerging properties in self-supervised vision transformers," in *Int'l Conf. on Computer Vision*, 2021, pp. 9650–9660.
- [32] Shang-Hua Gao, Yong-Qiang Tan, Ming-Ming Cheng, Chengze Lu, Yunpeng Chen, and Shuicheng Yan, "Highly efficient salient object detection with 100k parameters," in *Proc. European Conf. Computer Vision*, 2020, pp. 702–721.
- [33] Johan Bjorck, Kilian Q Weinberger, and Carla Gomes, "Understanding decoupled and early weight decay," in *AAAI Conf. Artificial Intelligence*, 2021, vol. 35, pp. 6777–6785.
- [34] John Duchi, Elad Hazan, and Yoram Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *J. of Machine Learning Research*, vol. 12, no. 7, 2011.
- [35] Pan Zhou, Xiaotong Yuan, and Jiashi Feng, "Faster first-order methods for stochastic non-convex optimization on riemannian manifolds," in *Int'l Conf. Artificial Intelligence and Statistics*, 2019.
- [36] Pan Zhou, Xiaotong Yuan, Zhouchen Lin, and Steven Hoi, "A hybrid stochastic-deterministic minibatch proximal gradient method for efficient optimization and generalization," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 5933–5946, 2022.
- [37] Tengyu Ma Moritz Hardt, "Identity matters in deep learning," in *Int'l Conf. Learning Representations*, 2023.
- [38] Bo Xie, Yingyu Liang, and Le Song, "Diverse neural network learns true target functions," in *Artificial Intelligence and Statistics*, 2017, pp. 1216–1224.
- [39] Zachary Charles and Dimitris Papailiopoulos, "Stability and generalization of learning algorithms that converge to global optima," in *Proc. Int'l Conf. Machine Learning*, 2018, pp. 745–754.
- [40] Pan Zhou, Hanshu Yan, Xiaotong Yuan, Jiashi Feng, and Shuicheng Yan, "Towards understanding why lookahead generalizes better than sgd and beyond," in *Neural Information Processing Systems*, 2021.
- [41] Tieleman Tijmen and Hinton Geoffrey, "Lecture 6.5-rmsprop: Divide the gradient by a run- ning average of its recent magnitude," *COURSERA: Neural Networks for Machine Learning*, vol. 4, 2012.



- [42] Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan, "Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models," *arXiv preprint arXiv:2208.06677*, 2022.
- [43] Simon S Du, Wei Hu, and Jason D Lee, "Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced," in *Proc. Conf. Neural Information Processing Systems*, 2018.
- [44] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh, "Large batch optimization for deep learning: Training BERT in 76 minutes," in *Int'l Conf. Learning Representations*, 2019.
- [45] Tongliang Liu, Dacheng Tao, Mingli Song, and Stephen J Maybank, "Algorithm-dependent generalization bounds for multi-task learning," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 39, no. 2, pp. 227–241, 2016.
- [46] Zeke Xie, Li Yuan, Zhanxing Zhu, and Masashi Sugiyama, "Positive-negative momentum: Manipulating stochastic gradient noise to improve generalization," in *Proc. Int'l Conf. Machine Learning*, 2021, pp. 11448–11458.
- [47] Mandt Stephan, Matthew D Hoffman, David M Blei, et al., "Stochastic gradient descent as approximate bayesian inference," *J. of Machine Learning Research*, vol. 18, no. 134, pp. 1–35, 2017.
- [48] Samuel L Smith and Quoc V Le, "A Bayesian perspective on generalization and stochastic gradient descent," in *Int'l Conf. Learning Representations*, 2018.
- [49] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao, "An investigation into neural net optimization via hessian eigenvalue density," in *Proc. Int'l Conf. Machine Learning*, 2019, pp. 2232–2241.
- [50] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [51] Matthew Staib, Sashank Reddi, Satyen Kale, Sanjiv Kumar, and Suvrit Sra, "Escaping saddle points with adaptive gradient methods," in *Proc. Int'l Conf. Machine Learning*, 2019, pp. 5956–5965.
- [52] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [53] Vladimir Vapnik, *Estimation of dependences based on empirical data*, Springer Science & Business Media, 2006.
- [54] Moritz Hardt, Ben Recht, and Yoram Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in *Proc. Int'l Conf. Machine Learning*, 2016, pp. 1225–1234.
- [55] Pan Zhou and Jiashi Feng, "Empirical risk landscape analysis for understanding deep neural networks," in *Int'l Conf. Learning Representations*, 2018.
- [56] Pan Zhou and Jiashi Feng, "Understanding generalization and optimization performance of deep CNNs," in *Proc. Int'l Conf. Machine Learning*, 2018, pp. 5960–5969.
- [57] Kaifeng Lyu, Zhiyuan Li, and Sanjeev Arora, "Understanding the generalization benefit of normalization layers: Sharpness reduction," in *Proc. Conf. Neural Information Processing Systems*, 2022, vol. 35, pp. 34689–34708.
- [58] Sadhika Malladi, Kaifeng Lyu, Abhishek Panigrahi, and Sanjeev Arora, "On the SDEs and scaling rules for adaptive gradient algorithms," in *Proc. Conf. Neural Information Processing Systems*, 2022, vol. 35, pp. 7697–7711.
- [59] Levent Sagun, Leon Bottou, and Yann LeCun, "Eigenvalues of the hessian in deep learning: Singularity and beyond," *arXiv preprint arXiv:1611.07476*, 2016.
- [60] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou, "Empirical analysis of the Hessian of over-parametrized neural networks," *arXiv preprint arXiv:1706.04454*, 2017.
- [61] Adepur Ravi Sankar, Yash Khasbage, Rahul Vigneswaran, and Vineeth N Balasubramanian, "A deeper look at the hessian eigenspectrum of deep neural networks and its applications to regularization," in *AAAI Conf. Artificial Intelligence*, 2021, vol. 35, pp. 9481–9488.
- [62] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Int'l Conf. on Computer Vision*, 2015, pp. 1026–1034.
- [63] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [64] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney, "Pyhessian: Neural networks through the lens of the hessian," in *Int'l conf. Big Data*, 2020, pp. 581–590.



**Pan Zhou** received Master Degree in computer science from Peking University in 2016 and obtained Ph.D. Degree in computer science from National University of Singapore in 2019. Now he is an assistant professor at Singapore Management University, Singapore. Before he also worked as a research scientist at Salesforce and Sea AI Lab, Singapore. His research interests include computer vision, machine learning, and optimization. He was the winner of the Microsoft Research Asia Fellowship 2018.



**Xingyu Xie** is a Ph.D. candidate in the School of Intelligence Science and Technology, at Peking University. His research interests include machine learning and optimization.



**Zhouchen Lin** (M'00–SM'08–F'18) received the Ph.D. degree in applied mathematics from Peking University in 2000. He is currently a Boya Special Professor with the National Key Laboratory of General Artificial Intelligence, School of Intelligence Science and Technology, Peking University. His research interests include machine learning and numerical optimization. He has published over 270 papers, collecting more than 27800 Google Scholar citations. He is a Fellow of the IAPR, the IEEE, and the CSIG.



**Shuicheng Yan** is currently director of Sea AI Lab (SAIL) and group chief scientist of Sea. He is a Fellow of the Academy of Engineering, Singapore, AAAI Fellow, ACM Fellow, IEEE Fellow, and IAPR Fellow. His research areas include computer vision, machine learning, and multimedia analysis. Till now, he has published over 600 papers in top international journals and conferences, with H-index 120+. His team won over 10 best paper or best student paper prizes and especially, a grand slam in ACM MM, the top conference in multimedia, including Best Paper Award three times, Best Student Paper Award twice, and Best Demo Award once.