
Win: Weight-Decay-Integrated Nesterov Acceleration for Adaptive Gradient Algorithms

Pan Zhou¹ Xingyu Xie^{1,2} Shuicheng Yan¹
¹Sea AI Lab ²Peking University
{zhoupan,xyxie,yansc}@sea.com xyxie@pku.cn

Abstract

Training deep networks on increasingly large-scale datasets is computationally challenging. In this work, we explore the problem of “*how to accelerate the convergence of adaptive gradient algorithms in a general manner*”, and aim at providing practical insights to boost the training efficiency. To this end, we propose an effective Weight-decay-Integrated Nesterov acceleration (Win) for adaptive algorithms to enhance their convergence speed. Taking AdamW and Adam as examples, we minimize a dynamical loss per iteration which combines the vanilla training loss and a dynamic regularizer inspired by proximal point method (PPM) to improve the convexity of the problem. Then we respectively use the first- and second-order Taylor approximations of vanilla loss to update the variable twice while fixing the above dynamic regularization brought by PPM. In this way, we arrive at our Win acceleration (like Nesterov acceleration) for AdamW and Adam that uses a conservative step and a reckless step to update twice and then linearly combines these two updates for acceleration. Next, we extend this Win acceleration to LAMB and SGD. Our transparent acceleration derivation could provide insights for other accelerated methods and their integration into adaptive algorithms. Besides, we prove the convergence of Win-accelerated adaptive algorithms by taking AdamW and Adam as examples. Experimental results testify the faster convergence speed and superior performance of our Win-accelerated AdamW, Adam, LAMB and SGD over their vanilla counterparts on vision classification tasks and language modeling tasks with CNN and Transformer backbones.

1 Introduction

Deep neural networks (DNNs) are effective to model realistic data and have been successfully applied to various applications, *e.g.* image classification [1–10] and speech recognition [11–14]. Typically, their training models can be formulated as the following nonconvex optimization problem:

$$\min_{z \in \mathbb{R}^d} F(z) := \mathbb{E}_{\zeta \sim \mathcal{D}} [f(z, \zeta)] + \frac{\lambda}{2} \|z\|_2^2, \quad (1)$$

where z is the model parameters; sample ζ is drawn from a data distribution \mathcal{D} ; the loss f is differentiable; λ is a constant. To solve problem (1), SGD [15, 16] uses its compositional structure to efficiently estimate gradient via minibatch data, and has become a dominant algorithm to train DNNs. However, on sparse data or ill-conditioned problems, SGD suffers from slow convergence speed [17, 18], as it scales the gradient uniformly in all parameter coordinate and ignores the data or problem properties on each coordinate. To resolve this issue, recent work has proposed a variety of adaptive methods, *e.g.* Adam [17] and AdamW [19], that scale each gradient coordinate according to the current geometry curvature of the loss $F(z)$. This coordinate-wise scaling greatly accelerates the optimization convergence and helps them, *e.g.* Adam and AdamW, become much more popular in DNN training.

Unfortunately, along with the increasing scale of both datasets and models, efficient DNN training even with SGD or adaptive algorithms has become very challenging. In this work, we are particularly interested in the problem of “*how to accelerate the convergence of adaptive algorithms in a general manner*” because of their dominant popularity across many DNNs. Heavy ball acceleration [20] and Nesterov acceleration [21] are widely used in SGD but are rarely studied in adaptive algorithms.

Contributions: In this work, based on a recent Nesterov-type acceleration formulation [22] and proximal point method (PPM) [23], we propose a new *Weight-decay-Integrated Nesterov acceleration* (Win) to accelerate adaptive algorithms. By taking AdamW and Adam as examples, we follow PPM spirit and minimize a dynamically regularized loss which combines vanillas loss and a dynamical regularization, and independently approximate the vanilla loss by its first- and second-order Taylor expansions to update the variable twice while fixing the above dynamic regularization. As a result, we achieve at our Win acceleration, a Nesterov-alike acceleration, for AdamW and Adam that uses a conservative step and a reckless step to update twice and then linearly combines these two updates for acceleration. Then we extend Win acceleration to LAMB [24] and SGD. This transparent acceleration derivation may motivate other accelerations and provide examples to introduce other accelerations into adaptive algorithms. Moreover, we analyze the convergence of Win-accelerated adaptive algorithms to justify their convergence superiority by using AdamW & Adam as examples.

Finally, experimental results on both vision classification tasks and language modeling tasks show that our Win-accelerated algorithms, i.e. accelerated AdamW, Adam, LAMB and SGD, can accelerate the convergence speed and also improve the performance of their corresponding non-accelerated counterparts by a remarkable margin on both CNN and transformer architectures.

2 Weight-decay-Integrated Nesterov Acceleration

In deterministic optimization, one widely used optimization-stabilizing and acceleration approach is proximal point method (PPM) [23, 25]. At the k -th iteration, PPM optimizes an ℓ_2 -regularized loss $F(\mathbf{z}) + \frac{1}{2\eta_k} \|\mathbf{z} - \mathbf{z}_{k-1}\|_2^2$ instead of the vanilla loss $F(\mathbf{z})$. This change enhances the convexity of the problem, accelerating and also stabilizing optimization [26, 27]. To make the ℓ_2 -regularized problem solvable iteratively, PPM approximates $F(\mathbf{z})$ by its first- or second-order Taylor expansion to get a close-form solution. At below, we borrow the idea in PPM to induce a *Weight-decay-Integrated Nesterov acceleration* (Win) for adaptive algorithms by using AdamW and Adam as examples.

Win-Accelerated AdamW and Adam. To begin with, following most adaptive gradient algorithms, e.g. Adam and AdamW, we estimate the first- and second-order moments \mathbf{m}_k and \mathbf{v}_k of gradient as

$$\mathbf{g}_k = \frac{1}{b} \sum_{i=1}^b \nabla f(\mathbf{z}_k; \zeta_i), \quad \mathbf{m}_k = (1 - \beta_1)\mathbf{m}_{k-1} + \beta_1\mathbf{g}_k, \quad \mathbf{v}_k = (1 - \beta_2)\mathbf{v}_{k-1} + \beta_2\mathbf{g}_k^2, \quad (2)$$

where $\mathbf{m}_0 = \mathbf{g}_0$, $\mathbf{v}_0 = \mathbf{g}_0^2$, $\beta_1 \in [0, 1]$ and $\beta_2 \in [0, 1]$. For brevity, with a small scaler $\nu > 0$, we define

$$\mathbf{s}_k = \sqrt{\mathbf{v}_k + \nu}, \quad \mathbf{u}_k = \mathbf{m}_k / \sqrt{\mathbf{v}_k + \nu}. \quad (3)$$

Then following PPM spirit, at the k -th iteration, we minimize a regularized loss $F(\mathbf{x}) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{s}_k}^2$. Here we use the regularizer $\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{s}_k}^2$ instead of the ℓ_2 -regularization $\|\mathbf{x} - \mathbf{x}_k\|_2^2$, since 1) this new regularization allows us to handle adaptive algorithms as shown below Eqn. (4), and 2) it also helps increase the problem convexity to speed up the convergence. To make problem solvable iteratively, we approximate $F(\mathbf{z})$ by its first-order Taylor expansion at the point \mathbf{z}_k and update \mathbf{x}_{k+1} as

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} F(\mathbf{z}_k) + \langle \mathbf{m}_k, \mathbf{x} - \mathbf{z}_k \rangle + \frac{1}{2\eta_k} \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{s}_k}^2 + \frac{\lambda}{2} \|\mathbf{x}\|_{\mathbf{s}_k}^2 = \frac{1}{1 + \lambda\eta_k} (\mathbf{x}_k - \eta_k \mathbf{u}_k), \quad (4)$$

where $\|\mathbf{x}\|_{\mathbf{s}_k} = \sqrt{\langle \mathbf{x}, \mathbf{s}_k * \mathbf{x} \rangle}$ with element-wise product $*$, \mathbf{m}_k is used to approximate the full gradient $\nabla F(\mathbf{z}_k)$ for Taylor expansion. We add a small regularization $\frac{\lambda}{2} \|\mathbf{x}\|_{\mathbf{s}_k}^2$, as 1) it can improve the generalization in practice [19, 28]; 2) it allows us to derive Adam ($\lambda = 0$) and AdamW ($\lambda > 0$). If $\lambda = 0$, the updating (4) becomes the exact Adam algorithm. If $\lambda > 0$, the updating (4) can approximate the updating rule $\mathbf{x}_{k+1} = (1 - \lambda\eta_k)\mathbf{x}_k - \eta_k\mathbf{u}_k$ of AdamW. This is because consider $\lambda\eta_k$ is small in practice, we approximate $(1 + \lambda\eta_k)^{-1} = 1 - \lambda\eta_k + \mathcal{O}(\lambda^2\eta_k^2)$ and thus $\frac{1}{1 + \lambda\eta_k} (\mathbf{x}_k - \eta_k\mathbf{u}_k) = [1 - \lambda\eta_k + \mathcal{O}(\lambda^2\eta_k^2)]\mathbf{x}_k - [\eta_k - \mathcal{O}(\lambda\eta_k^2) + \mathcal{O}(\lambda^3\eta_k^3)]\mathbf{u}_k$ which becomes AdamW

Algorithm 1: Win-Accelerated AdamW, Adam and LAMB

Input: initialization $\mathbf{x}_0 = \mathbf{z}_0 = \mathbf{0}$, step size $\{(\eta_k, \bar{\eta}_k)\}_{k=0}^T$, moment parameters $\{\beta_1, \beta_2\}$.

```

1 while  $k < T$  do
2    $\mathbf{g}_k = \frac{1}{b} \sum_{i=1}^b \nabla f(\mathbf{z}_k; \zeta_i)$ 
3    $\mathbf{m}_k = (1 - \beta_1)\mathbf{m}_{k-1} + \beta_1\mathbf{g}_k$  /*  $\mathbf{m}_0 = \mathbf{g}_0$  */
4    $\mathbf{v}_k = (1 - \beta_2)\mathbf{v}_{k-1} + \beta_2\mathbf{g}_k^2$  /*  $\mathbf{v}_0 = \mathbf{g}_0^2$  */
5    $\mathbf{u}_k = \frac{\mathbf{m}_k}{\sqrt{\mathbf{v}_k + \nu}}$  for AdamW and Adam,  $\mathbf{u}_k = \frac{\|\mathbf{x}_k\|_2}{\|\mathbf{m}_k/\sqrt{\mathbf{v}_k + \nu}\|_2} \frac{\mathbf{m}_k}{\sqrt{\mathbf{v}_k + \nu}}$  for LAMB
6    $\mathbf{x}_{k+1} = \frac{1}{1 + \lambda\eta_k} (\mathbf{x}_k - \eta_k\mathbf{u}_k)$ 
7    $\mathbf{z}_{k+1} = \bar{\eta}_k\tau_k\mathbf{x}_{k+1} + \eta_k\tau_k(\mathbf{z}_k - \bar{\eta}_k\mathbf{u}_k)$  with  $\tau_k = \frac{1}{\eta_k + \bar{\eta}_k + \lambda\eta_k\bar{\eta}_k}$ 
8 end while

```

by ignoring $\mathcal{O}(\eta_k^2)$ and $\mathcal{O}(\eta_k^3)$. This is one reason that we adopt the regularizer $\|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{s}_k}^2$ in (4) instead of the ℓ_2 -regularization in PPM, as we can flexibly derive Adam and AdamW.

Similarly, we minimize a regularized loss $F(\mathbf{z}) + \frac{1}{2\eta_k}\|\mathbf{z} - \mathbf{x}_{t+1}\|_{\mathbf{s}_k}^2$, and further approximate $F(\mathbf{z})$ by its second-order approximation $F(\mathbf{z}_k) + \langle \mathbf{m}_k, \mathbf{z} - \mathbf{z}_k \rangle + \frac{1}{2\bar{\eta}_k}\|\mathbf{z} - \mathbf{z}_k\|_{\mathbf{s}_k}^2$:

$$\begin{aligned} \mathbf{z}_{k+1} &= \operatorname{argmin}_{\mathbf{z}} F(\mathbf{z}_k) + \langle \mathbf{m}_k, \mathbf{z} - \mathbf{z}_k \rangle + \frac{1}{2\bar{\eta}_k}\|\mathbf{z} - \mathbf{z}_k\|_{\mathbf{s}_k}^2 + \frac{1}{2\eta_k}\|\mathbf{z} - \mathbf{x}_{k+1}\|_{\mathbf{s}_k}^2 + \frac{\lambda}{2}\|\mathbf{z}\|_{\mathbf{s}_k}^2 \\ &= \bar{\eta}_k\tau_k\mathbf{x}_{k+1} + \eta_k\tau_k(\mathbf{z}_k - \bar{\eta}_k\mathbf{u}_k), \end{aligned} \quad (5)$$

where $\tau_k = \frac{1}{\eta_k + \bar{\eta}_k + \lambda\eta_k\bar{\eta}_k}$, \mathbf{m}_k can approximate $\nabla F(\mathbf{x}_k)$ as shown in Theorem 1 in Appendix B.

For more clear, we introduce a virtual sequence $\mathbf{y}_{k+1} = \mathbf{z}_k - \bar{\eta}_k\mathbf{u}_k$ in Win, and rewrite (5) as

$$\mathbf{x}_{k+1} = (1 + \lambda\eta_k)^{-1}(\mathbf{x}_k - \eta_k\mathbf{u}_k), \quad \mathbf{y}_{k+1} = \mathbf{z}_k - \bar{\eta}_k\mathbf{u}_k, \quad \mathbf{z}_{k+1} = \bar{\eta}_k\tau_k\mathbf{x}_{k+1} + \eta_k\tau_k\mathbf{y}_{k+1}. \quad (6)$$

See detailed steps in Algorithm 1. Interestingly, Win acceleration is similar to Nesterov-type acceleration, since they both use a conservative step η_k and a reckless step $\bar{\eta}_k$ to update \mathbf{x}_{k+1} and \mathbf{y}_{k+1} respectively, and then linearly combine \mathbf{x}_{k+1} and \mathbf{y}_{k+1} to obtain \mathbf{z}_{k+1} .

Our Win-acceleration is quit simple and efficient, as our accelerated AdamW/Adam only adds an extra simple algorithmic step, i.e. the 7th step in Algorithm 1, on vanilla AdamW/Adam. Moreover, for the only extra hyper-parameter, the reckless step $\bar{\eta}_k$, in Algorithm 1 over AdamW/Adam, we always set it as $\bar{\eta}_k = 2\eta_k$, which works well in our all experiments.

Extension to LAMB and SGD. Here we generalize Win acceleration to LAMB [24] and SGD [15]. For LAMB, it scales the update \mathbf{u}_k of AdamW in Eqn. (3) so that \mathbf{u}_k is at the same magnitude of the network weight \mathbf{x}_k . That is, it changes the update rule $\mathbf{x}_{k+1} = (1 - \lambda\eta_k)\mathbf{x}_k - \eta_k\mathbf{m}_k/\mathbf{s}_k$ in AdamW to $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \frac{\|\mathbf{x}_k\|_2}{\|\mathbf{r}_k + \lambda\mathbf{x}_k\|_2}(\mathbf{r}_k + \lambda\mathbf{x}_k)$ where $\mathbf{r}_k = \mathbf{m}_k/\mathbf{s}_k$. This modification is to avoid too large or small update, improving optimization efficiency. To extend Win acceleration to LAMB, we inherit this scaling spirit, and scale the update \mathbf{u}_k in (3) to the following one:

$$\mathbf{u}_k = (\|\mathbf{x}_k\|_2/\|\mathbf{m}_k/\mathbf{s}_k\|_2) \cdot (\mathbf{m}_k/\mathbf{s}_k). \quad (7)$$

Next, we can follow Eqn. (4) and (5) to update, and summarize detailed steps in Algorithm 1.

For SGD, applying Win acceleration to it is quite direct. Specifically, the only algorithmic difference between SGD and AdamW on the ℓ_2 -regularized problems is that SGD has no second-order moment \mathbf{v}_k while AdamW has. So we can borrow the acceleration framework of AdamW to accelerate SGD by setting $\mathbf{s}_k = \mathbf{1} \in \mathbb{R}^d$ in Eqn. (3), (4) and (5), and obtain WIN-accelerated SGD:

$$\mathbf{m}_k = \beta_1\mathbf{m}_{k-1} + \beta_1'\mathbf{g}_k, \quad \mathbf{x}_{k+1} = \frac{1}{1 + \lambda\eta_k}(\mathbf{x}_k - \eta_k\mathbf{m}_k), \quad \mathbf{z}_{k+1} = \bar{\eta}_k\tau_k\mathbf{x}_{k+1} + \eta_k\tau_k(\mathbf{z}_k - \bar{\eta}_k\mathbf{m}_k), \quad (8)$$

where $\beta_1' \in [0, 1]$ is dampening parameter. Here we slightly modify the moment \mathbf{m}_k to accord with the one used in Nesterov-accelerated SGD (e.g. SGD-M in Pytorch).

Convergence Analysis. Theorem 1 in Appendix B analyzes the convergence of Win-accelerated adaptive algorithms to justify their convergence superiority by using AdamW & Adam as examples.

Table 1: ImageNet top-1 accuracy (%) of ResNet50&101 whose official optimizer is LAMB due to the stronger data augmentation for better performance. * is reported in [29].

Epoch	ResNet50				ResNet101			
	100	200	300	avg.	100	200	300	avg.
SAM	77.3	78.7	79.4	78.5	79.5	81.1	81.6	80.7
SGD-H	75.3	76.9	77.2	76.5	77.7	78.6	78.8	78.4
SGD-M	77.0	78.6	79.3	78.3	79.3	81.0	81.4	80.6
SGD-Win	78.0	79.2	79.7	79.0 _{+0.7}	80.1	81.2	81.6	81.0 _{+0.4}
Adam	76.9	78.4	78.8	78.1	78.4	80.2	80.6	79.7
Adam-Win	77.8	78.8	79.3	78.7 _{+0.6}	79.2	80.6	81.0	80.3 _{+0.6}
AdamW	77.0	78.9	79.3	78.4	78.9	79.9	80.4	79.7
AdamW-Win	78.0	79.3	79.9	79.1 _{+0.7}	80.2	81.1	81.3	80.9 _{+1.2}
LAMB	77.0	79.2	79.8*	78.7	79.4	81.1	81.3*	80.6
LAMB-Win	78.4	79.7	80.1	79.4 _{+0.7}	80.6	81.5	81.7	81.3 _{+0.7}

Table 2: ImageNet top-1 accuracy (%) of ViT and PoolFormer whose default optimizers are both AdamW. * and \diamond are respectively reported in [28] and [30].

Epoch	ViT-S			ViT-B			PoolFormer-S12		
	150	300	avg.	150	300	avg.	150	300	avg.
SGD-M	77.4	79.4	78.4	79.6	80.0	79.8	69.7	74.3	72.0
SGD-Win	78.1	80.1	79.1 _{+0.7}	80.4	80.8	80.6 _{+0.8}	71.1	74.5	72.8 _{+0.8}
Adam	77.3	79.3	78.3	79.0	79.7	79.4	74.3	76.3	75.3
Adam-Win	78.6	80.2	79.4 _{+1.1}	80	80.5	80.3 _{+0.9}	75.6	77.1	76.4 _{+1.1}
AdamW	78.3	79.8*	79.1	79.5	81.8*	80.7	75.2	77.1*	76.2
AdamW-Win	79.3	80.8	80.1 _{+1.0}	81.0	82.2	81.6 _{+0.9}	76.7	77.6	77.2 _{+1.0}
LAMB	78.0	79.6	78.8	80.3	80.8	80.6	75.4	77.4	76.4
LAMB-Win	79.3	80.6	80.0 _{+1.2}	81.0	81.4	81.2 _{+0.6}	76.7	78.0	77.4 _{+1.0}

3 Experiments

For clarity, we call our accelerated algorithm “X-Win”, where “X” denotes vanilla optimizers. In all experiments, our accelerated algorithms, *e.g.* AdamW-Win, always use the default hyper-parameters of vanilla optimizers, *e.g.* moment parameters β_1 and β_2 in AdamW; and set $\bar{\eta}_k = 2\eta_k$.

Results on ResNets and ViTs. Table 1 reports accuracy of ResNets under the setting in [29], and Table 2 gives the performance of ViT [2] and PoolFormer [30]. Our accelerated algorithms always outperform their corresponding non-accelerated version. On ResNet, LAMB-Win achieves remarkable improvement over the official optimizer LAMB for this setting; SGD-Win also surpasses heavy-ball accelerated SGD (SGD-H) and Nesterov accelerated SGD (SGD-M). On ViTs, our accelerated algorithms consistently outperform the corresponding non-accelerated counterparts. Fig. 1 shows the faster convergence behaviors of our accelerated algorithms over non-accelerated counterparts which could benefit their better performance under the same computational cost.

Results on Transformer-XL. Table 3 shows that under different training steps on WikiText-103 dataset, our accelerated Adam-Win always achieves lower test PPL than the official Adam optimizer of Transformer-XL-base, and improves 1.5 average test PPL over Adam.

Table 3: Test PPL of Transformer-XL-B. * is officially reported.

Transformer-XL	Training Steps			
	50k	100k	200k	avg.
Adam	28.5	25.5	24.2*	26.7
Adam-Win	26.7	25.0	24.0	25.2 _{+1.5}

4 Conclusion

In this work, we adopt proximal point method to derive a weight-decay-integrated Nesterov acceleration for AdamW and Adam, and extend it to LAMB and SGD. Moreover, we prove the convergence of our accelerated algorithms, *i.e.* accelerated AdamW, Adam and SGD, and observe the superiority of the accelerated Adam-type algorithm over the vanilla ones in terms of stochastic gradient complexity. Finally, experimental results validate the advantages of our accelerated algorithms.

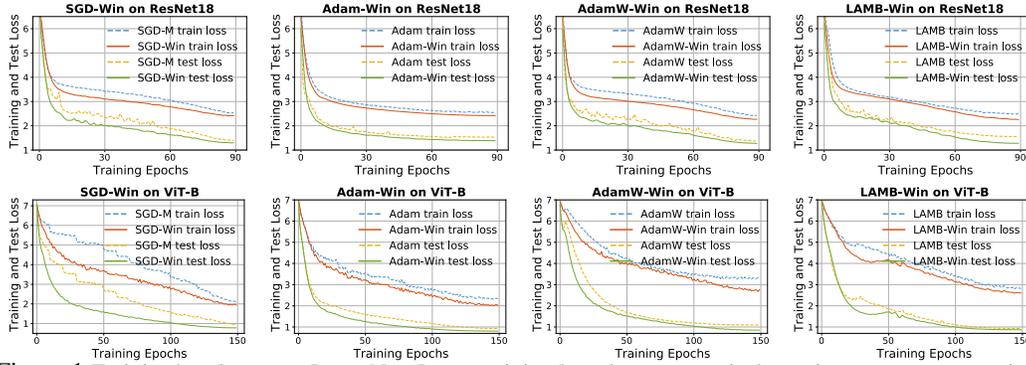


Figure 1: Training/test losses on ImageNet. Larger training loss than test one is due to its strong augmentation.

References

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [2] 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. *arXiv preprint arXiv:2010.11929*, 2020.
- [3] Pan Zhou, Xiaotong Yuan, Huan Xu, Shuicheng Yan, and Jiashi Feng. Efficient meta learning via mini-batch proximal update. In *Neural Information Processing Systems*, 2019.
- [4] Pan Zhou, Caiming Xiong, Richard Socher, and Steven Hoi. Theory-inspired path-regularized differential network architecture search. In *Neural Information Processing Systems*, 2020.
- [5] Pan Zhou, Caiming Xiong, Xiaotong Yuan, and Steven Hoi. A theory-driven self-labeling refinement method for contrastive representation learning. In *Neural Information Processing Systems*, 2021.
- [6] Yuxuan Liang, Pan Zhou, Roger Zimmermann, and Shuicheng Yan. Dualformer: Local-global stratified transformer for efficient video recognition. In *European Conference on Computer Vision*, pages 577–595. Springer, 2022.
- [7] Pan Zhou, Yichen Zhou, Chenyang Si, Weihao Yu, Teck Khim Ng, and Shuicheng Yan. Mugs: A multi-granular self-supervised learning framework. *arXiv preprint arXiv:2203.14415*, 2022.
- [8] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10819–10829, 2022.
- [9] Chenyang Si, Weihao Yu, Pan Zhou, Yichen Zhou, Xinchao Wang, and Shuicheng Yan. Inception transformer. *arXiv preprint arXiv:2205.12956*, 2022.
- [10] Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xinchao Wang. Metaformer baselines for vision. *arXiv preprint arXiv:2210.13452*, 2022.
- [11] T. Sainath, A. Mohamed, B. Kingsbury, and B. Ramabhadran. Deep convolutional neural networks for LVCSR. In *ICASSP*, pages 8614–8618. IEEE, 2013.
- [12] O. Abdel-Hamid, A. Mohamed, H. Jiang, L. Deng, G. Penn, and D. Yu. Convolutional neural networks for speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(10):1533–1545, 2014.
- [13] Yubei Xiao, Ke Gong, Pan Zhou, Guolin Zheng, Xiaodan Liang, and Liang Lin. Adversarial meta sampling for multilingual low-resource speech recognition. In *Association for the Advancement of Artificial Intelligence*, 2021.
- [14] Guolin Zheng, Yubei Xiao, Ke Gong, Pan Zhou, Xiaodan Liang, and Liang Lin. Wav-bert: Cooperative acoustic and linguistic representation learning for low-resource speech recognition. *arXiv preprint arXiv:2109.09161*, 2021.
- [15] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.

- [16] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.
- [17] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *Int’l Conf. Learning Representations*, 2014.
- [18] Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley-benchmarking deep learning optimizers. In *International Conference on Machine Learning*, pages 9367–9376. PMLR, 2021.
- [19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- [20] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- [21] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- [22] Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.
- [23] Jean-Jacques Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.
- [24] 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 *International Conference on Learning Representations*, 2019.
- [25] R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM journal on control and optimization*, 14(5):877–898, 1976.
- [26] Junhyung Lyle Kim, Panos Toulis, and Anastasios Kyrillidis. Convergence and stability of the stochastic proximal point algorithm with momentum. In *Learning for Dynamics and Control Conference*, pages 1034–1047. PMLR, 2022.
- [27] Pan Zhou, Hanshu Yan, Xiaotong Yuan, Jiashi Feng, and Shuicheng Yan. Towards understanding why lookahead generalizes better than sgd and beyond. *Advances in Neural Information Processing Systems*, 34:27290–27304, 2021.
- [28] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [29] Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm. *arXiv preprint arXiv:2110.00476*, 2021.
- [30] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. *arXiv preprint arXiv:2111.11418*, 2021.
- [31] 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 *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3267–3275, 2021.
- [32] 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. *Advances in Neural Information Processing Systems*, 33:18795–18806, 2020.
- [33] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International Conference on Learning Representations*, 2019.
- [34] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 248–255, 2009.
- [35] Jürgen Schmidhuber, Sepp Hochreiter, et al. Long short-term memory. *Neural Comput*, 9(8):1735–1780, 1997.
- [36] Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Using Large Corpora*, 273, 1994.

- [37] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
- [38] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6023–6032, 2019.
- [39] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020.
- [40] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- [41] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- [42] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.
- [43] Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov momentum algorithm for faster optimizing both cnns and vits. *Axriv*, 2022.