

# Training Deep Networks with Stochastic Gradient Normalized by Layerwise Adaptive Second Moments

Boris Ginsburg<sup>1</sup> Patrice Castonguay<sup>1</sup> Oleksii Hrinchuk<sup>1</sup> Oleksii Kuchaiev<sup>1</sup> Ryan Leary<sup>1</sup> Vitaly Lavrukhin<sup>1</sup>  
Jason Li<sup>1</sup> Huyen Nguyen<sup>1</sup> Yang Zhang<sup>1</sup> Jonathan M. Cohen<sup>1</sup>

## Abstract

We propose NovoGrad, an adaptive stochastic gradient descent method with layer-wise gradient normalization and decoupled weight decay. In our experiments on neural networks for image classification, speech recognition, machine translation, and language modeling, it performs on par or better than well-tuned SGD with momentum, Adam, and AdamW. Additionally, NovoGrad (1) is robust to the choice of learning rate and weight initialization, (2) works well in a large batch setting, and (3) has half the memory footprint of Adam.

## 1. Introduction

The most popular algorithms for training Neural Networks (NNs) are Stochastic Gradient Descent (SGD) with momentum (Polyak, 1964; Sutskever et al., 2013) and Adam (Kingma & Ba, 2015). SGD with momentum is the preferred algorithm for computer vision, while Adam is more commonly used for natural language processing (NLP) and speech problems. Compared to SGD, Adam is perceived as safer and more robust to weight initialization and learning rate. However, Adam has certain drawbacks. First, as noted in the original paper (Kingma & Ba, 2015), the second moment can vanish or explode, especially during the initial phase of training. Also, Adam often leads to solutions that generalize worse than SGD (Wilson et al., 2017), e.g. models for image classification trained with Adam have significantly lower accuracy than when they are trained with SGD with momentum (Loshchilov & Hutter, 2019).

Our motivation for this work was to build an algorithm which: (1) performs equally well for image classification, speech recognition, machine translation, and language modeling, (2) is robust to learning rate (LR) and weight initialization, (3) has strong regularization properties.

<sup>1</sup>NVIDIA, Santa Clara, USA. Correspondence to: Boris Ginsburg <bginsburg@nvidia.com>.

We started with Adam, and then (1) replaced the element-wise second moment with the layer-wise moment, (2) computed the first moment using gradients normalized by layer-wise second moment, (3) decoupled weight decay (WD) from normalized gradients (similar to AdamW).

The resulting algorithm, *NovoGrad*, combines SGD’s and Adam’s strengths. We applied NovoGrad to a variety of large scale problems — image classification, neural machine translation, language modeling, and speech recognition — and found that in all cases, it performs as well or better than Adam/AdamW and SGD with momentum.

## 2. Related Work

NovoGrad belongs to the family of **Stochastic Normalized Gradient Descent (SNGD)** optimizers (Nesterov, 1984; Hazan et al., 2015). SNGD uses only the direction of the stochastic gradient  $\mathbf{g}_t$  to update the weights  $\mathbf{w}_t$ :

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda_t \cdot \frac{\mathbf{g}_t}{\|\mathbf{g}_t\|}$$

The step size does not depend on the magnitude of that gradient. Hazan et al. (2015) proved that the direction of the gradient is sufficient for convergence. Ignoring the gradient magnitude makes SNGD robust to vanishing and exploding gradients, but SNGD is still sensitive to “noisy” gradients, especially during an initial training phase.

One can improve SNGD stability by gradient averaging. Adagrad (Duchi et al., 2011), RmsProp (Tieleman & Hinton, 2012), Adam (Kingma & Ba, 2015) – all these SNGD methods use the some kind of gradient averaging. Adam, the most popular one, uses moving averages  $\mathbf{m}_t, \mathbf{v}_t$ :

$$\mathbf{m}_t = \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \cdot \mathbf{g}_t \quad (1)$$

$$\mathbf{v}_t = \beta_2 \cdot \mathbf{v}_{t-1} + (1 - \beta_2) \cdot \mathbf{g}_t^2 \quad (2)$$

The weights update is computed with the first moment  $\mathbf{m}_t$  normalized by the second moment  $\mathbf{v}_t$ :

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda_t \cdot \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} \quad (3)$$

where  $\epsilon \ll 1$  is added for numerical stability. To strengthen

Adam’s robustness to “noisy” gradients, coefficients  $\beta_1$  and  $\beta_2$  are usually close to 1 (e.g.  $\beta_2 > 0.99$ ).

### 2.1. Layer-wise Gradient Normalization

The other way to improve SNGD robustness was proposed by Yu et al. (2018), who suggested to combine Adam with **layer-wise gradient normalization**. Both moments are computed with normalized gradients  $\hat{g}_t^l = \frac{g_t^l}{\|g_t^l\|}$ , where  $g_t^l$  is the gradient for the layer  $l$  at step  $t$ :

$$\begin{aligned} m_t^l &= \beta_1 \cdot m_{t-1}^l + (1 - \beta_1) \cdot \hat{g}_t^l \\ v_t^l &= \beta_2 \cdot v_{t-1}^l + (1 - \beta_2) \cdot \|\hat{g}_t^l\|^2 \end{aligned}$$

A similar idea was used in (Singh et al., 2015) to scale up small layers gradients, while keeping large gradients unchanged:

$$\hat{g}_t^l = g_t^l \cdot \left(1 + \log\left(1 + \frac{1}{\|g_t^l\|}\right)\right)$$

### 2.2. Improving Adam Generalization

Adaptive methods like Adam generalize worse than SGD with momentum (Wilson et al., 2017). For example, Keskar & Socher (2017) proposed to use Adam during the initial stage only and then switch to SGD. Luo et al. (2019) suggested to improve Adam generalization by limiting the factor  $\frac{1}{\sqrt{v_t}}$  to a certain range: limiting from above helps to decrease the training loss while limiting from below helps to generalize better.

To improve Adam regularization, Loshchilov & Hutter (2019) proposed **AdamW**, which **decouples the weight decay**  $d \cdot w_t$  from the gradient and uses it directly in the weight update:

$$w_{t+1} = w_t - \lambda_t \cdot \left(\frac{m_t}{\sqrt{v_t} + \epsilon} + d \cdot w_t\right)$$

### 2.3. Reduction of Adam Memory Footprint

Adam needs to store the second moment, and this doubles the optimizer memory compared to SGD with momentum. This affects large models like GPT-2 – 1.5 billion parameters (Radford et al., 2019), Meena – 2.6 billion parameters (Adiwardana et al., 2020), or Megatron-LM – 8.3 billion parameters (Shoeybi et al., 2019). Shazeer & Stern (2018) proposed the **AdaFactor** algorithm, which replaced the full second moment with moving averages of the row and column sums of the squared gradients. For a layer defined by an  $n \times m$  matrix, this would reduce memory from  $\mathcal{O}(n \times m)$  to  $\mathcal{O}(n + m)$ .

## 3. Algorithm

NovoGrad combines three ideas: (1) use layer-wise second moments, (2) compute first moment with gradients normalized with layer-wise second moments, (3) decouple weight decay.

---

### Algorithm 1 NovoGrad

---

**Parameters:**

Init learning rate  $\lambda_0$ , moments  $\beta_1, \beta_2$ , weight decay  $d$ , number of steps  $T$

$t = 0$ : **weight initialization**

$w_0 \leftarrow \text{Init}()$ .

$t = 1$ : **moment initialization**

**for** each layer  $l$  **do**

$v_1^l \leftarrow \|g_1^l\|^2$ ;

$m_1^l \leftarrow \frac{g_1^l}{\sqrt{v_1^l}} + d \cdot w_1^l$ .

**end for**

**while**  $t \leq T$  **do**

compute global learning rate  $\lambda_t \leftarrow LR(\lambda_0, t, T)$

**for** each layer  $l$  **do**

$g_t^l \leftarrow \nabla_l L(w_t)$

$v_t^l \leftarrow \beta_2 \cdot v_{t-1}^l + (1 - \beta_2) \cdot \|g_t^l\|^2$

$m_t^l \leftarrow \beta_1 \cdot m_{t-1}^l + \left(\frac{g_t^l}{\sqrt{v_t^l} + \epsilon} + d \cdot w_t^l\right)$

$w_{t+1}^l \leftarrow w_t^l - \lambda_t \cdot m_t^l$

**end for**

**end while**

---

Let  $g_t^l$  be the stochastic gradient for layer  $l$  at step  $t$ . First, we compute the layer-wise second moment  $v_t^l$  using  $\|g_t^l\|$ :

$$v_t^l = \beta_2 \cdot v_{t-1}^l + (1 - \beta_2) \cdot \|g_t^l\|^2 \quad (4)$$

where  $0 \leq \beta_2 \leq 1$ . We use much smaller  $\beta_2$  than in Adam,<sup>1</sup> The moment  $v_t^l$  is used to normalize the gradient  $g_t^l$  before calculating the first moment  $m_t^l$ .

$$m_t^l = \beta_1 \cdot m_{t-1}^l + \frac{g_t^l}{\sqrt{v_t^l} + \epsilon} w_t^l \quad (5)$$

Last, we decouple weight decay  $d \cdot w_t$  from the stochastic gradient similarly to AdamW, but we add it to normalized gradient before computing moment  $m_t^l$ :<sup>2</sup>

$$m_t^l = \beta_1 \cdot m_{t-1}^l + \left(\frac{g_t^l}{\sqrt{v_t^l} + \epsilon} + d \cdot w_t^l\right) \quad (6)$$

where  $0 < \beta_1 < 1$  is the momentum, typically in the same range as in SGD or Adam [0.9 – 0.95]. The first moment

<sup>1</sup>The default  $\beta_2 = 0.25$  which we used in all experiments except Imagenet classification where  $\beta_2 = 0.98$  was used.

<sup>2</sup>One can also use decoupled weight decay in the weights update, as in AdamW. We didn’t find a significant difference between these two options.

can be also computed with an exponential moving average in Adam-like style:

$$\mathbf{m}_t^l = \beta_1 \cdot \mathbf{m}_{t-1}^l + (1 - \beta_1) \cdot \left( \frac{\mathbf{g}_t^l}{\sqrt{v_t^l + \epsilon}} + d \cdot \mathbf{w}_t^l \right)$$

We use the following moments initialization to remove bias:

$$v_1^l = \|\mathbf{g}_1^l\|^2; \quad \mathbf{m}_1^l = \frac{\mathbf{g}_1^l}{\|\mathbf{g}_1^l\|} + d \cdot \mathbf{w}_1^l$$

Finally, weights are updated the same way as in SGD with momentum:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda_t \cdot \mathbf{m}_t$$

Similar to Adam, one can construct a counter-example for NovoGrad in the stochastic convex optimization settings (Wilson et al., 2017). However, the ‘‘AMS-Grad’’ fix (Reddi et al., 2018) for Adam can also be applied in this case to make sure that  $\frac{\lambda_t}{\sqrt{v_t^l}}$  is monotonically decreasing:

$$\begin{aligned} v_t^l &= \beta_2 \cdot v_{t-1}^l + (1 - \beta_2) \cdot \|\mathbf{g}_t^l\|^2 \\ \hat{v}_t^l &= \max(\hat{v}_{t-1}^l, v_t^l) \\ \mathbf{m}_t^l &= \beta_1 \cdot \mathbf{m}_{t-1}^l + \left( \frac{\mathbf{g}_t^l}{\sqrt{\hat{v}_t^l + \epsilon}} + d \cdot \mathbf{w}_t^l \right) \end{aligned}$$

Notes.

1. If we set  $\beta_2 = 0$ , then  $v_t^l = \|\mathbf{g}_t^l\|^2$ , and NovoGrad becomes layer-wise NGD.
2. We use gradient normalization in the first moment computation instead of moment normalization in weights update to improve the algorithm robustness against very large ‘‘outliers’’ gradients.
3. NovoGrad has half the memory footprint of Adam.

### 3.1. Training with NovoGrad

NovoGrad has initial learning rates different than both SGD and Adam. In the initial phase, normalized gradients have larger magnitudes than non-normalized gradients used by SGD, so NovoGrad uses smaller learning rate than SGD. In Adam, on the other side, normalization is done by element-wise second moments, which are significantly smaller than per-layer second moments used in NovoGrad. So for NovoGrad, safe learning rates are somewhere between those of SGD and Adam, as the gradients are normalized by the per-layer gradient norm.

## 4. Experiments

We train four deep models: ResNet-50 (He et al., 2016) — for ImageNet classification, Transformer-big (Vaswani et al.,

2017) — for WMT 2014 translation, Jasper (Li et al., 2019) — for LibriSpeech speech recognition, and Transformer-XL (Dai et al., 2019) — for WikiText-103 word-level language modeling, with NovoGrad, SGD with momentum, and Adam/AdamW. Each model was trained on a single DGX-1 with 8 NVIDIA V100 GPUs with gradient accumulation used for large batch training. In all the experiments, NovoGrad performed on par or better than other algorithms.

### 4.1. Image Classification

We used ResNet-50v2 model (He et al., 2016) for ImageNet classification task (Russakovsky et al., 2015). We trained the model with three optimizers: SGD with momentum (SGD), AdamW, and NovoGrad using the batch of 1024 for 100 epochs. We used quadratic LR decay for SGD with momentum and cosine decay (Loshchilov & Hutter, 2016) for AdamW and NovoGrad. We could not find any training recipe for ResNet-50 with AdamW, so we report the best accuracy we achieved after extensive hyper-parameter search. We used only standard data augmentation methods: resize, flip, random crop, and did not employ any additional training tricks (He et al., 2018). The single-crop validation accuracy for each algorithm is reported in Table 1.

Table 1: ImageNet: ResNet-50v2 trained with batch 1024, top-1/top-5 accuracy (%).

Optimizer	LR	WD	Epoch	Top-1/Top-5, %
SGD	0.400	0.0001	100	76.38/93.08
			200	76.33/92.96
AdamW	0.002	0.120	100	76.36/93.01
			200	76.48/92.94
NovoGrad	0.007	0.002	100	76.94/93.41
			200	77.74/93.70
			300	77.65/93.62

NovoGrad outperformed both AdamW and SGD with the top-1 accuracy of 76.94% after 100 epochs. SGD and Adam accuracy remained under 76.5% if we trained for 200 epochs instead, while NovoGrad accuracy improved to 77.74%. NovoGrad demonstrated powerful regularization capabilities: training for 100 additional epochs kept top-1=77.65% without overfitting. Note that this is ‘‘vanilla’’ ResNet-50, without sophisticated data augmentation or model tweaking.

### 4.2. Large Batch Training for Image Classification

Hazan et al. (2015) showed that large batch size is beneficial for SNGD convergence, which motivated us to explore NovoGrad for large batch training. We trained ResNet-50 v2 with batch sizes of 8K and 32K. To compare with the previous methods, we train the model for 90 epochs us-

Table 2: ImageNet, large batch training comparison: ResNet-50v2, top-1 accuracy(%).

Optimizer	Reference	Bag of Tricks	Epochs	B=8K	B=32K
SGD	(Goyal et al., 2017)	warmup	90	76.26	72.45
SGD	(You et al., 2018)	warmup, LARS	90	75.30	75.40
SGD	(Codreanu et al., 2017)	warmup, multi-step WD	100	76.60	75.31
NovoGrad		—	90	76.64	75.78
		warmup	90	—	75.99

ing cosine LR decay. To emulate a large batch, we used a mini-batch of 128 per GPU and accumulated gradients from several mini-batches before each weight update.

To establish the baseline for NovoGrad training with batch 32K we first used the method similar to proposed in Goyal et al. (2017): scaling the learning rate linearly with the batch size and using LR warmup. This method gives top-1=75.09% and top-5=92.27%. We found that we get better results when we increase both the learning rate  $\lambda$  and the weight decay  $d$  to improve the regularization (see Table 3).

Table 3: ImageNet, large batch training with NovoGrad: ResNet-50v2, 90 epochs. top-1/top-5 accuracy (%).

Batch	LR	WD	Top-1/Top-5, %
1K	0.070	0.002	76.86/93.31
8K	0.016	0.006	76.64/93.14
32K	0.026	0.010	75.78/92.54

For comparison, we took three methods, which (1) use fixed batch size during training and (2) do not modify the original model. All three methods employ SGD with momentum. The first method (Goyal et al. (2017)) scales LR linearly with batch size and uses the LR warmup to stabilize the initial training phase. The second method (You et al. (2018)) combines warmup with Layer-wise Adaptive Rate Scaling (LARS) (You et al., 2017). The last method (Codreanu et al. (2017)) uses warmup and dynamic weight decay (WD). NovoGrad outperformed other methods without any additional techniques like LR warmup (Goyal et al., 2017), dynamic weight decay, special batch norm, etc. Using warmup (500 steps) improves top-1 accuracy to 75.99%.

### 4.3. Speech Recognition

We conducted experiments with Jasper DR 10x5 (Li et al. (2019)), a deep convolutional neural acoustic model, on the LibriSpeech speech recognition task (Panayotov et al., 2015). Jasper was trained with SGD with momentum (SGD), Adam and NovoGrad for 400 epochs with a batch of 256, polynomial LR decay, and Layerwise Adaptive Rate Clipping (LARC). We found that NovoGrad yields lower Word

Error Rates (WER) comparing to SGD, especially for the long runs. The model and training parameters are described in Li et al. (2019).

Table 4: Speech Recognition: Jasper-10x5 trained on LibriSpeech for 400 epochs, greedy WER(%).

Optimizer	Dev		Test	
	clean	other	clean	other
Adam	5.06	15.76	5.27	15.94
SGD	4.08	13.23	4.22	13.15
NovoGrad	3.71	11.75	3.71	11.85

### 4.4. Large Batch Training for Speech Recognition

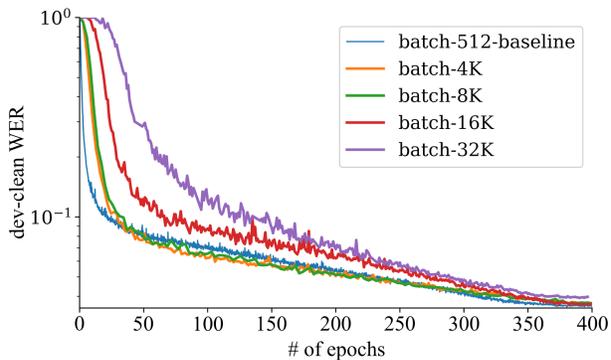


Figure 1: Speech Recognition, large batch training. Jasper-10x5 trained with NovoGrad on LibriSpeech, WER(%).

We trained Jasper DR 10x5 with batch sizes of 512, 4K, 8K, 16K and 32K on LibriSpeech. In all cases, we trained the model for 400 epochs. For batch size up to 8K, we scaled LR linearly with the batch size and used LR warmup. To scale batch to 16K and 32K we also increased weight decay (see Table 5). The batch 16K leads to WER comparable to the baseline. Batch 32K has higher WER due to the smaller number of training steps (9 weights updates per epoch). Figure 1 shows WER on dev-clean during training for different batch sizes.

Table 5: Speech Recognition, large batch training. Jasper-10x5 trained on LibriSpeech for 400 epochs, WER(%).

Batch	LR	Warmup	WD	Dev		Test	
				clean	other	clean	other
512	0.015	-	0.001	3.58	11.30	3.85	11.29
4K	0.03	0.05	0.001	3.66	11.67	3.92	11.68
8K	0.06	0.05	0.001	3.69	11.76	3.96	11.75
16K	0.06	0.05	0.003	3.67	11.03	3.94	11.19
32K	0.06	0.08	0.004	4.01	11.73	4.14	11.89

#### 4.5. Neural Machine Translation

We trained Transformer (Vaswani et al., 2017) on WMT 2014 English-to-German benchmark. For all the experiments, we used a 12-layer Transformer-big model with 185M parameters ( $d_{\text{model}} = 1024$ ,  $d_{\text{ff}} = 4096$ ,  $h = 16$ ) with the vocabulary of 8192 tokens based on joint source-target byte-pair-encodings (Sennrich et al., 2015). For Adam and AdamW we used dropout of  $P_{\text{drop}} = 0.3$  and for NovoGrad we used  $P_{\text{drop}} = 0.2$ . We trained all algorithms with mixed-precision (Micikevicius et al., 2017) for 100K steps (approximately 150 epochs) with a 4K steps warmup on batches of up to 490K source and target tokens obtained via gradient accumulation (Ott et al., 2018) with cosine learning rate annealing. We did not use checkpoint averaging, all the results are reported for the last checkpoint in the corresponding run.

Table 6: WMT’14 English-to-German translation, Transformer-big, batch 490K tokens, 150 epochs, no checkpoint averaging. Detokenized SacreBLEU and Tokenized BLEU on WMT’14 (newstest14).

Optimizer	LR	WD	Sacre/Token-BLEU
Adam	0.0006	-	28.26/28.71
AdamW	0.0006	0.005	28.24/28.72
NovoGrad	0.04	0.0001	28.80/29.35

#### 4.6. Language Modeling

We trained Transformer-XL (Dai et al., 2019), the state-of-the-art language model on the word-level WikiText-103 (Merity et al., 2016) benchmark. For all the experiments we used a 16-layer model with 191M parameters ( $d_{\text{model}} = 512$ ,  $d_{\text{ff}} = 2048$ ,  $h = 8$ ,  $P_{\text{drop}} = 0.15$ ). All other hyper-parameters were taken from the original Transformer-XL paper, the source code was based on a publicly available implementation. Each configuration was trained for 12 billion tokens which is approximately 117 epochs and 366K iterations.

Table 7: Language Modeling. Transformer-XL trained on WikiText-103 with batch size 256, sequence length 512, 12B tokens. Validation and Test Perplexity(PPL).

Optimizer	LR	WD	Val/Test-PPL
Adam	0.00025	-	23.84/25.40
AdamW	0.00025	0.001	23.64/25.06
NovoGrad	0.01	0	20.53/21.26

NovoGrad significantly outperformed both Adam and AdamW (Table 7). NovoGrad exhibits a much smaller gap between training and validation perplexity compared to Adam (Figure 2). Even longer training for 20B tokens does not lead to overfitting, as the validation and test perplexities improve even further.

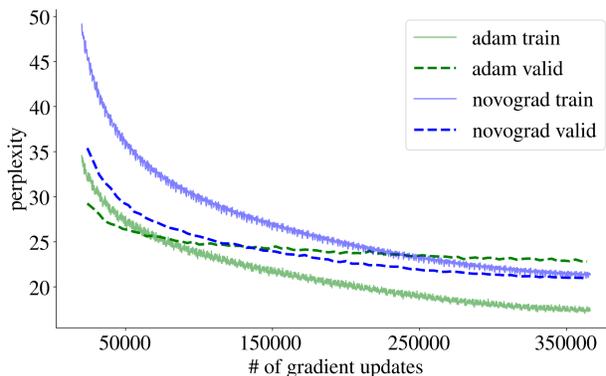


Figure 2: Language Modeling. Transformer-XL trained with Adam and NovoGrad on WikiText-103.

## 5. Conclusion

We propose NovoGrad – an adaptive SGD method with gradients normalized by the layer-wise second moment and with decoupled weight decay. We tested NovoGrad on deep models for image classification, speech recognition, translation, and language modeling.

In all experiments, NovoGrad performed equally or better than SGD and Adam/AdamW. NovoGrad is more robust to the initial learning rate and weight initialization than other methods. For example, NovoGrad works well without LR warm-up, while other methods require it. NovoGrad performs exceptionally well for large batch training, e.g. it outperforms other methods for ResNet-50 for all batches up to 32K. In addition, NovoGrad requires half the memory compared to Adam.

NovoGrad and all models in this paper are open sourced.

## References

- Adiwardana, D. D. F., Luong, M.-T., So, D. R., Hall, J., Fiedel, N., Thoppilan, R., Yang, Z., Kulshreshtha, A., Nemade, G., Lu, Y., and Le, Q. V. Towards a human-like open-domain chatbot. *ArXiv*, abs/2001.09977, 2020.
- Codreanu, V., Podareanu, D., and Saletore, V. Scale out for large minibatch sgd: Residual network training on imagenet-1k with improved accuracy and reduced time to train. *arXiv:1711.04291*, 2017.
- Dai, Z., Yang, Z., Yang, Y., Cohen, W. W., Carbonell, J., Le, Q. V., and Salakhutdinov, R. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv:1901.02860*, 2019.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, pp. 21212159, 2011.
- Goyal, P., Dollár, P., Girshick, R. B., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: training imagenet in 1 hour. *arXiv:1706.02677*, 2017.
- Hazan, E., Levy, K., and Shalev-Shwartz, S. Beyond convexity: Stochastic quasi-convex optimization. In *NIPS*, pp. 15851593, 2015.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. *arXiv:11603.05027*, 2016.
- He, T., Zhang, Z., Zhang, H., Zhang, Z., Xie, J., and Li, M. Bag of tricks for image classification with convolutional neural networks. *arXiv:1812.00187*, 2018.
- Keskar, N. S. and Socher, R. Improving generalization performance by switching from adam to SGD. *arXiv:1712.07628*, 2017.
- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Li, J., Lavrukhin, V., Ginsburg, B., Leary, R., Kuchaiev, O., Cohen, J. M., Nguyen, H., and Gaddei, R. T. Jasper: An end-to-end convolutional neural acoustic model. *arXiv:1904.03288*, 2019.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *ICLR*, 2016.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *ICLR*, 2019.
- Luo, L., Xiong, Y., Liu, Y., and Sun, X. Adaptive gradient methods with dynamic bound of learning rate. In *ICLR*, 2019.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. *arXiv:1609.07843*, 2016.
- Micikevicius, P., Narang, S., Alben, J., Diamos, G. F., Elsen, E., García, D., Ginsburg, B., Houston, M., Kuchaiev, O., Venkatesh, G., and Wu, H. Mixed precision training. *ICLR*, 2017.
- Nesterov, Y. E. Minimization methods for nonsmooth convex and quasiconvex functions. *Matekon*, 29:519531, 1984.
- Ott, M., Edunov, S., Grangier, D., and Auli, M. Scaling neural machine translation. *arXiv preprint arXiv:1806.00187*, 2018.
- Panayotov, V., Chen, G., Povey, D., and Khudanpur, S. Librispeech: an asr corpus based on public domain audio books. In *ICASSP*, pp. 5206–5210. IEEE, 2015.
- Polyak, B. Some methods of speeding up the convergence of iteration methods. In *USSR Computational Mathematics and Mathematical Physics*, 1964.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 2019.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of adam and beyond. In *ICLR*, 2018.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- Sennrich, R., Haddow, B., and Birch, A. Neural machine translation of rare words with subword units. *arXiv:1508.07909*, 2015.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. *arXiv:1804.04235*, 2018.

- Shoeybi, M., Patwary, M. A., Puri, R., LeGresley, P., Casper, J., and Catanzaro, B. Megatron-lm: Training multi-billion parameter language models using model parallelism. *ArXiv*, abs/1909.08053, 2019.
- Singh, B., De, S., Zhang, Y., Goldstein, T., and Taylor, G. Layer-specific adaptive learning rates for deep networks. In *ICML*, 2015.
- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In *ICML*, 2013.
- Tieleman, T. and Hinton, G. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *arXiv: 1706.03762*, 2017.
- Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., and Recht, B. The marginal value of adaptive gradient methods in machine learning. In *NIPS*, 2017.
- You, Y., Gitman, I., and Ginsburg, B. Large batch training of convolutional networks. *arXiv:1708.03888*, 2017.
- You, Y., Zhang, Z., Hsieh, C., and Demmel, J. 100-epoch imagenet training with alexnet in 24 minutes. *arXiv:1709.05011*, 2018.
- Yu, A. W., Lin, Q., Salakhutdinov, R., and Carbonell, J. Block-normalized gradient method: An empirical study for training deep neural network. *arXiv:1707.04822*, 2018.