

Toward Understanding Why Adam Converges Faster Than SGD for Transformers

author names withheld

Under Review for OPT 2022

Abstract

While stochastic gradient descent (SGD) is still the most popular optimization algorithm in deep learning, adaptive algorithms such as Adam have established empirical advantages over SGD in some deep learning applications such as training transformers. However, it remains a question why Adam converges significantly faster than SGD in these scenarios. In this paper, we explore one explanation of why Adam converges faster than SGD using a new concept *directional sharpness*. We argue that the performance of optimization algorithms is closely related to the directional sharpness of the update steps, and show SGD has much worse directional sharpness compared to adaptive algorithms. We further observe that only a small fraction of the coordinates causes the bad sharpness and slow convergence of SGD, and propose to use coordinate-wise clipping as a solution to SGD and other optimization algorithms. We demonstrate the effect of coordinate-wise clipping in sharpness reduction and speeding up the convergence of optimization algorithms under various settings, and conclude that the sharpness reduction effect of adaptive coordinate-wise scaling is the reason for Adam’s success in practice.

1. Introduction

Stochastic gradient descent (SGD) [3, 25] is one of the most popular optimization algorithms for deep learning. Although SGD is efficient on various large-scale neural networks, in many tasks, such as training transformers [10, 30], people seek to use the adaptive variants of stochastic gradient methods. Adaptive algorithms, such as Adagrad [11], Adam [17], and AMSGrad [24], can find a “better coordinate-wise scaling” of the gradient, so the size of the update step is adaptive to the local geometry of the function. While adaptive algorithms can converge much faster than SGD in many applications [12, 15, 33], the understanding of the superior performance of Adam-type optimizers in these tasks is limited [7, 33].

In this paper, we explore one explanation of why Adam converges faster than SGD in practice, especially for transformers. We propose to study the *directional sharpness* of optimization algorithms, defined as $(x_t - x_{t+1})^\top \nabla^2 f(x_t)(x_t - x_{t-1})$, that occurs in the quadratic Taylor expansion of the objective function. We argue that directional sharpness of the update direction is an useful indicator of the performance of optimization algorithms, that high sharpness usually implies low performance. We observe through experiments that the update direction of SGD has much higher sharpness compared to other optimization algorithms. Furthermore, we propose that the imbalanced distribution of gradient across coordinates is the key contributor to SGD’s high sharpness. We observe that only a small fraction of the parameters contribute to the majority of SGD’s sharpness, causing slow convergence. Hence, we propose to use coordinate-wise clipping as a solution to the

problem of slow convergence and bad sharpness. We show that clipping can improve the directional sharpness and convergence rate of various non-coordinate-wise-scaling optimization algorithms, and propose that coordinate-wise clipping can be used as a generic component in optimization algorithms. We demonstrate our findings through two experiments under various settings and show that our observations are consistent across different tasks, models, and iterations. We conclude that the adaptive coordinate-wise scaling of Adam can effectively find a balance between optimizing gradient correlation and directional sharpness, and such ability is the key to Adam’s fast convergence in deep learning training.

2. Related Work

General Convergence Rates of Adaptive Algorithms. Adaptive algorithms have long been studied and applied in deep learning [1, 11, 17, 22, 24, 29]. Several previous work has proved convex and non-convex convergence rates for Adagrad [9, 11, 19, 31] and Adam or AMSGrad [4, 8, 13, 21, 24, 34, 35]. The best known non-convex convergence rate for Adagrad is $O(\frac{\log T}{\sqrt{T}})$ [9, 19] and $O(\frac{1}{\sqrt{T}})$ for AMSGrad [34]. While the result by Zhou et al. [34] matches the non-convex convergence rate $O(\frac{1}{\sqrt{T}})$ of SGD [14], there is no theoretical proof that Adam can converge asymptotically faster than SGD for general functions [7]. Therefore, there is still a significant gap of work between the theoretical understanding of Adam and its empirical fast performance.

Faster Convergence Rates Under Certain Settings. Another line of work focused on specific settings that Adam might work better than SGD. Adaptive algorithms can work asymptotically better when the stochastic gradients are sparse [11, 34] or when there is a sparse set of noise [2]. Zhang et al. [33] proved that global clipping methods outperforms SGD when the stochastic gradients have heavy-tail noise, argued that Adam can also deal with heavy-tail noise effectively, and designed a new algorithm based on coordinate-wise clipping. The effect of global clipping and normalization methods were also studied in [16, 18]. Our work is inspired by the use of coordinate-wise clipping in algorithm design in [33], but we propose different explanations of the effectiveness of coordinate-wise clipping with new empirical evidence.

3. Directional Sharpness of Adaptive Algorithms

In this section, we introduce a new metric *directional sharpness* that indicates the performance of optimization algorithms and argue it is closely related to the poor performance of SGD.

We first explain the motivation of our approach. In convex and non-convex optimization, a typical proof strategy is to consider the quadratic Taylor expansion of the objective function

$$f(x_{t+1}) = f(x_t) + \underbrace{\nabla f(x_t)^\top (x_t - x_{t+1})}_{\text{gradient correlation}} + \underbrace{(x_t - x_{t+1})^\top \nabla^2 f(x_t) (x_t - x_{t+1})}_{\text{directional sharpness}} + O(\eta^3) \quad (1)$$

where η is the step size. In order for $f(x_{t+1})$ to decrease monotonically, the optimization algorithm should minimize the two terms that depends on the update step. In typical convergence proof, f is assumed to be L -smooth, which is equivalent to $\|\nabla^2 f(x)\|_2 \leq L$ for all x . The local Hessian spectral norm is also called the *sharpness* [5]. Then, the second-order term is upper bounded by $L\|x_t - x_{t+1}\|_2^2$, so the loss can decrease when $\|x_t - x_{t+1}\|_2$ is sufficiently small and the correlation between the update step and the gradient is positive. However, there are disadvantages of

the smoothness assumptions in theoretical proof. The smoothness can adapt to the geometry of the trajectory and can vary significantly for different algorithms [5, 6]. Furthermore, even if the local geometry and Hessian are fixed, the update direction $x_{t+1} - x_t$ is extremely important to minimizing the second-order term. Motivated by the definition of sharpness and the above observations, we define the *directional sharpness* of a function f at x in the direction $v \in \mathbb{R}^d, \|v\|_2 = 1$ as $v^\top \nabla^2 f(x) v$. We can see this is an extension of the sharpness definition, as the sharpness at x is just the supremum of directional sharpness over all possible directions v . The directional sharpness at x_t in the update direction is extremely important to minimizing $f(x_{t+1})$. When gradient correlation is similar, the loss $f(x_{t+1})$ directly depends on the directional sharpness at x_t . Furthermore, since directional sharpness is quadratic in the step size η and gradient correlation is linear, a lower directional sharpness implies the potential to take a larger step size and possibly lead to a larger local decrement of the objective function.

Empirically, we observe that the directional sharpness is much lower for adaptive algorithms than for SGD. We study the update step of different optimization algorithms under the same trajectory and local geometry using pseudo-update steps described in Appendix A, in order to rule out the impact of trajectory. We compute the directional sharpness of different optimization algorithms and visualize the optimization landscape in the update direction of a variety of optimization algorithms in Figure 3. The observation is consistent with our analysis. The update step of SGD has the best correlation with the actual gradient, so the loss decrease faster when the step size is very small, since in this case the linear term dominates the quadratic term in Equation (1). However, because of the large directional sharpness, when the step size increases the quadratic term grows faster than the linear term, so the loss reaches the local minima after a very small step size. For adaptive algorithms, the directional sharpness is much lower than SGD, so it has the potential to use a much larger step size and the optimal step will give a lower loss compared to SGD.

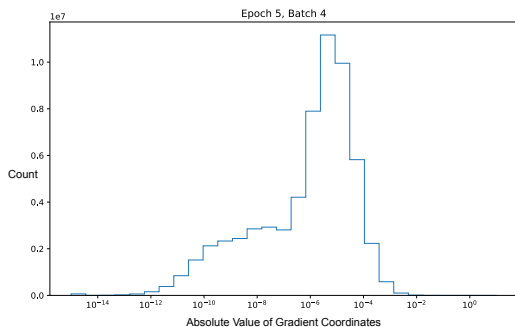


Figure 1: Histogram of momentum distribution for stochastic gradient descent on machine translation.

In order to explain the sharpness reduction effect of adaptive algorithms, since the strategy for adaptive algorithms is to find a coordinate-wise scaling of the gradient, we investigate the distribution of gradient norm across different coordinates. We visualize a histogram of the absolute value of SGD momentum coordinates in Figure 1. We observe that the gradients are distributed unevenly across the coordinates, with half of the coordinates have absolute value ranging from 10^{-12} to 10^{-6} , but also exists an innegligible portion of coordinates that can be as high as 10^{-4} to 10^{-2} , contribut-

Algorithm	Sharpness
Adam	0.16190993
SGD	31.04433435
SGD Clipping	1.77876506
Normalized SGD	0.77112307
Normalized SGD Clipping	0.38075357
Adafactor	3.1928×10^{-6}
Adafactor Clipping	2.5258×10^{-6}

Figure 2: The average sharpness of different optimization algorithms when trained on machine translation, in the same experiment and epoch as Figure 3.

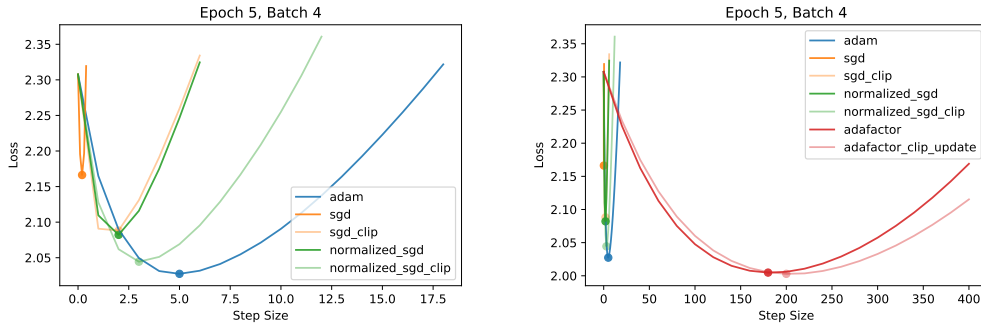


Figure 3: The loss landscape in different update directions on machine translation. The step size is the learning rate normalized by the update step ℓ_2 norm.

ing to most of the gradient norm. The histogram suggests that the gradients are concentrated on a small fraction of the coordinates, and this small fraction of coordinates can contribute to a large portion of sharpness, making optimization hard. For adaptive algorithms and normalized optimization algorithms, since they already used some forms of scaling, the imbalanced gradient distribution will not be as severe as SGD, but normalizing the large coordinates might still be beneficial.

4. Coordinate-wise Clipping

In this section, we propose to use *coordinate-wise clipping* as a solution to the aforementioned imbalanced distribution of gradient based on our experimental findings. We observe that the sharpness is also concentrated in the large coordinates in the gradient, and clipping those coordinates can significantly decrease directional sharpness. Although clipping can decrease the correlation between the update step and the true gradient, since the dependence on the clipped entry is quadratic for the second-order term and linear for the first-order term, it might not be beneficial to use these coordinates. The use of clipping in optimization algorithms is a trade-off between improving gradient correlation and reducing directional sharpness. By clipping the top coordinates in the gradient, although gradient correlation decreases, the directional sharpness can decrease even more to make up the loss.

We consider using clipping on a variety of non-coordinate-wise-scaling algorithms, including SGD, normalized SGD, and Adafactor [27]. We demonstrate that coordinate-wise clipping significantly reduces the sharpness of adaptive algorithms and speeds up the optimization process. Specifically, we compute the threshold τ for the top $k\%$ gradients in terms of the absolute value, and clip the gradient coordinates g_i to $\tilde{g}_i = \text{sgn}(g_i) \min\{|g_i|, \tau\}$ based on their sign. In practice, it is possible to simplify the procedure by setting a fixed threshold. From Figure 2, we can see that by clipping the top 5% coordinates, the directional sharpness decrease significantly. Since we normalize the update step when we compute the directional sharpness, the sharpness reduction effect of coordinate-wise clipping is not due to significant reduction of the norm of the update step, but the improved flatness of the direction. Figure 3 gives a coherent message, that clipped algorithms can find a direction that has larger maximal decrement of the loss in the local geometry.

Finally we demonstrate that clipping algorithms can converge faster than the original counterpart by directly training transformers with the clipping algorithms, with the loss curve shown

in Figure 4. According to the result, clipping algorithms can speedup training significantly. Our result suggests that clipping can be used as an generic building block in any non-coordinate-wise-scaling algorithms and speed up training. The new finding can provide insight into designing new optimization algorithms.

Based on our experimental findings, we conjecture that there is a positive correlation between the magnitude of Hessian coordinates and gradient coordinates. The positive correlations is also mentioned in [32], but their proposed correlation is between the norm of Hessian and norm of gradient. We further suggest that there is a positive correlation between the *coordinates* of gradient and Hessian, and the success of Adam is due to the ability to scale down the bad coordinates and reduce the sharpness through coordinate-wise scaling of the gradient. Understanding of this phenomenon could be essential in proving convergence rates for Adam that are faster than SGD.

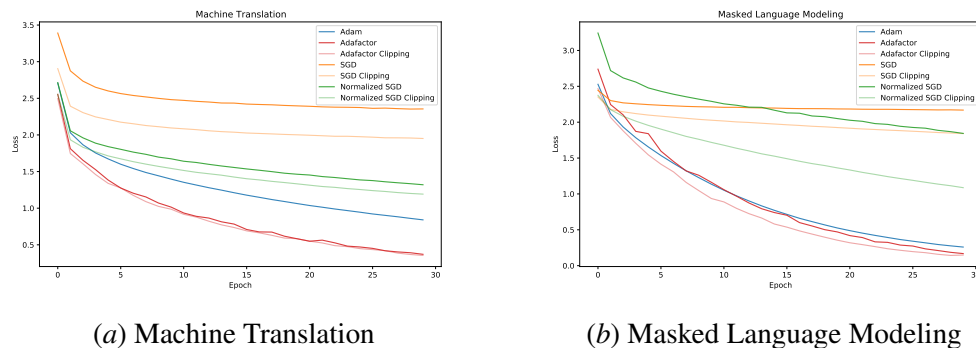


Figure 4: Clipped optimization algorithms generally converge faster than the original algorithms.

5. Experiments

In Appendices A and B, we demonstrate our findings with two types of experiments. We explore several different tasks and settings and show our results hold in various setting, including training t5 [23] architecture on machine translation datasets and DistilRoBERTa [26] on masked language modeling datasets. We compute the directional sharpness of a variety of optimization algorithms, including SGD, normalized SGD, and Adafactor [27], and visualize the corresponding loss landscape direction, under different local geometry. We show that SGD has bad sharpness under all of the settings, regardless of the task, model, or local geometry. In addition, we demonstrate clipping can always improve the directional sharpness of optimization algorithms, and often result in better local decrement of loss function. We also implement clipping algorithms and use them to train different models, and demonstrate that clipping algorithms converge faster in practice.

6. Conclusion

In summary, our work provides a new insight of why Adam converges faster than SGD in practice. In contrast to assumptions on properties of the gradient, we propose to study directional sharpness as an important indicator for the performance of optimization algorithms in deep learning. We show that adaptive algorithms and clipped optimization algorithms can generally achieve significantly better directional sharpness compared to SGD. We argue that the slow convergence of SGD is related

to the high directional sharpness, caused by a positive coordinate-wise gradient-Hessian correlation. We propose to use coordinate-wise clipping as a solution to the problem of high sharpness. We demonstrate the sharpness reduction effect of coordinate-wise clipping and show that it is possible to step into a lower loss in the update direction of clipping algorithms compared to the original algorithms. We further demonstrate the effectiveness of coordinate-wise clipping in a wide range of optimization algorithms without coordinate-wise scaling, including SGD, normalized SGD, and Adafactor. We suggest the use of coordinate-wise clipping as a generic building block in non-convex optimization algorithms. Our work provide useful explanations and conjectures about the superior performance of Adam and further understanding of the results could be useful in theoretical understanding of the empirical advantage of Adam over SGD.

References

- [1] Larry Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of mathematics*, 16(1):1–3, 1966.
- [2] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signSGD: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pages 560–569. PMLR, 2018.
- [3] Léon Bottou et al. Stochastic gradient learning in neural networks. *Proceedings of Neuro-Nimes*, 91(8):12, 1991.
- [4] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of Adam-type algorithms for non-convex optimization. In *International Conference on Learning Representations*, 2019.
- [5] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
- [6] Jeremy M Cohen, Behrooz Ghorbani, Shankar Krishnan, Naman Agarwal, Sourabh Medapati, Michal Badura, Daniel Suo, David Cardoze, Zachary Nado, George E Dahl, et al. Adaptive gradient methods at the edge of stability. *arXiv preprint arXiv:2207.14484*, 2022.
- [7] Marina Danilova, Pavel Dvurechensky, Alexander Gasnikov, Eduard Gorbunov, Sergey Guminov, Dmitry Kamzolov, and Innokentiy Shibaev. Recent theoretical advances in non-convex optimization. In *High-Dimensional Optimization and Probability*, pages 79–163. Springer, 2022.
- [8] Soham De, Anirbit Mukherjee, and Enayat Ullah. Convergence guarantees for RMSProp and Adam in non-convex optimization and an empirical comparison to nesterov acceleration. *arXiv preprint arXiv:1807.06766*, 2018.
- [9] Alexandre Défossez, Léon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof of Adam and Adagrad. *arXiv preprint arXiv:2003.02395*, 2020.

- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- [11] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [12] John Duchi, Michael I Jordan, and Brendan McMahan. Estimation, optimization, and parallelism when data is sparse. *Advances in Neural Information Processing Systems*, 26, 2013.
- [13] Biyi Fang and Diego Klabjan. Convergence analyses of online Adam algorithm in convex setting and two-layer ReLU neural network. *arXiv preprint arXiv:1905.09356*, 2019.
- [14] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [15] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [16] Elad Hazan, Kfir Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex optimization. *Advances in neural information processing systems*, 28, 2015.
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [18] Kfir Y Levy. The power of normalization: Faster evasion of saddle points. *arXiv preprint arXiv:1611.04831*, 2016.
- [19] Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 983–992. PMLR, 2019.
- [20] Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150, 2011.
- [21] Tran Thi Phuong and Le Trieu Phong. On the convergence proof of AMSGrad and a new version. *arXiv preprint arXiv:1904.03590*, 2019.
- [22] Boris T Polyak. Introduction to optimization. 1987. *Optimization Software, Inc, New York*.
- [23] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67, 2020.
- [24] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of Adam and beyond. In *International Conference on Learning Representations*, 2018.
- [25] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.

- [26] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [27] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018.
- [28] Jörg Tiedemann. Parallel data, tools and interfaces in OPUS. In *Eight International Conference on Language Resources and Evaluation, MAY 21-27, 2012, Istanbul, Turkey*, pages 2214–2218, 2012.
- [29] Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- [30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [31] Rachel Ward, Xiaoxia Wu, and Leon Bottou. AdaGrad stepsizes: Sharp convergence over nonconvex landscapes. In *International Conference on Machine Learning*, pages 6677–6686. PMLR, 2019.
- [32] Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, 2020.
- [33] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33, 2020.
- [34] Dongruo Zhou, Jinghui Chen, Yuan Cao, Yiqi Tang, Ziyang Yang, and Quanquan Gu. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv preprint arXiv:1808.05671*, 2018.
- [35] Fangyu Zou and Li Shen. On the convergence of weighted adagrad with momentum for training deep neural networks. *arXiv preprint arXiv:1808.03408*, 2018.

Appendix A. Experimental Details

A.1. Tasks, Datasets, and Models

We run our experiments on two tasks, including machine translation and masked language modeling. The details of the dataset, training set size, and model we use are in Table 1. For each dataset, we select the first 10240 data as our training set. Since we’re mainly interested in minimizing the training loss, we do not use any test or validation sets, nor any evaluation metrics other than the cross-entropy loss. For machine translation, we use the English to French opus books dataset [28] and t5 model [23]. For masked language modeling, we use the imdb dataset [20] and DistilRoBERTa model [26]. In order to evaluate the function in a offline setting, we generate fixed masks with probability 0.15 at the beginning of the training and does not generate new masks whenever we collate the data.

Task	Dataset	Size	Model
Machine Translation	opus books [28]	10240	t5-small [23]
Masked Language Modeling	imdb [20]	10240	DistilRoBERTa-base [26]

Table 1: Details of the tasks, datasets, training set sizes, and models we use for the two different experiments.

A.2. Optimization Algorithms and Clipping Methods

We use 4 optimization algorithms, including Adam, SGD, normalized SGD, and Adafactor [27]. We use momentum for all of the optimization algorithms to rule out any potential effect of momentum. The clipped optimization algorithms are described in Algorithms 1 to 4. Notice that for Adafactor, we only clip the gradient in the nominator of the final update step, since otherwise the scaling effect could cancel out or even increase the norm. Adafactor is originally used with the relative step sizes α_t , but in certain cases we use a fixed learning rate in place of α_t . In the algorithms, we assume $\text{Clip}(g)$ calculates the clipping threshold τ for the top $k\%$ coordinates and returns \tilde{g} where $\tilde{g}_i = \text{sgn}(g_i) \min\{|g_i|, \tau\}$. We use $k = 5$ in all of our experiments.

Algorithm 1: SGD with momentum

Data: initial point x_0 , learning rate η , momentum term β

for $t \leftarrow 1, \dots, T$ **do**

$g_t \leftarrow \nabla f(x_t)$

$\hat{g}_t \leftarrow \text{Clip}(g_t)$

$m_t \leftarrow \beta m_{t-1} + (1 - \beta)\hat{g}_t$

$x_t \leftarrow x_{t-1} - \eta m_t$

end

Algorithm 2: Normalized SGD with momentum for weight matrices and vectors**Data:** initial point $x_0 \in \mathbb{R}^{m \times n}$, learning rate η , momentum term β

for $t \leftarrow 1, \dots, T$ **do**
 $g_t \leftarrow \nabla f(x_t)$
 $\hat{g}_t \leftarrow \text{Clip}(g_t)$
 $m_t \leftarrow \beta m_{t-1} + (1 - \beta) \hat{g}_t$
 $v_t \leftarrow \frac{m_t}{\|m_t\|_2} \cdot \sqrt{mn}$
 $x_t \leftarrow x_{t-1} - \eta v_t$
end

Algorithm 3: Adafactor for weight matrices [27]**Data:** initial point $x_0 \in \mathbb{R}^{m \times n}$, relative step sizes $\rho_t = \min\{10^{-2}, \frac{1}{\sqrt{t}}\}$, second moment decay $\hat{\beta}_{2t} = 1 - t^{-0.8}$, regularization constants $\epsilon_1 = 10^{-30}$ and $\epsilon_2 = 10^{-3}$, clipping threshold $d = 1$, $\text{RMS}(x) := \frac{\|x\|_F}{\sqrt{mn}}$

for $t \leftarrow 1, \dots, T$ **do**
 $\alpha_t \leftarrow \max\{\epsilon_2, \text{RMS}(x_{t-1})\} \rho_t$
 $G_t \leftarrow \nabla f(x_{t-1})$
 $\hat{G}_t \leftarrow \text{Clip}(G_t)$
 $R_t \leftarrow \hat{\beta}_{2t} R_{t-1} + (1 - \hat{\beta}_{2t})(G_t^2 + \epsilon_1) \mathbf{1}_m$
 $C_t \leftarrow \hat{\beta}_{2t} C_{t-1} + (1 - \hat{\beta}_{2t}) \mathbf{1}_n^\top (G_t^2 + \epsilon_1)$
 $\hat{V}_t \leftarrow R_t C_t / \mathbf{1}_n^\top R_t$
 $U_t \leftarrow \hat{G}_t / \sqrt{\hat{V}_t}$
 $\hat{U}_t \leftarrow U_t / \max\{1, \text{RMS}(U_t)/d\}$
 $x_t \leftarrow x_{t-1} - \alpha_t \hat{U}_t$
end

Algorithm 4: Adafactor for weight vectors [27]**Data:** initial point $x_0 \in \mathbb{R}^n$, relative step sizes $\rho_t = \min\{10^{-2}, \frac{1}{\sqrt{t}}\}$, second moment decay $\hat{\beta}_{2t} = 1 - t^{-0.8}$, regularization constants $\epsilon_1 = 10^{-30}$ and $\epsilon_2 = 10^{-3}$, clipping threshold $d = 1$, $\text{RMS}(x) := \frac{\|x\|_2}{\sqrt{n}}$

for $t \leftarrow 1, \dots, T$ **do**
 $\alpha_t \leftarrow \max\{\epsilon_2, \text{RMS}(x_{t-1})\} \rho_t$
 $G_t \leftarrow \nabla f(x_{t-1})$
 $\hat{G}_t \leftarrow \text{Clip}(G_t)$
 $\hat{V}_t \leftarrow \hat{\beta}_{2t} \hat{V}_{t-1} + (1 - \hat{\beta}_{2t})(G_t^2 + \epsilon_1)$
 $U_t \leftarrow \hat{G}_t / \sqrt{\hat{V}_t}$
 $\hat{U}_t \leftarrow U_t / \max\{1, \text{RMS}(U_t)/d\}$
 $x_t \leftarrow x_{t-1} - \alpha_t \hat{U}_t$
end

A.3. Experiment for Directional Sharpness of Optimization Algorithms

Pseudo-Update Step. Since all algorithms we use has momentum part, we need to compute the momentum term in a different trajectory using “pseudo-update step.” Specifically, we compute the momentum term for all the optimization algorithms at time t using the past values of x_1, \dots, x_{t-1} , regardless of the optimization algorithm we use to perform the actual update step. The values we computed for the algorithms were only used to visualize the landscape and compare the sharpness, but not used for training. The momentum parameters are set to the default values [17, 27].

Training Optimizer. We use different training optimizers to compare our results across different local geometry and optimization trajectory. We use SGD momentum with learning rate 10^{-3} and Adam with learning rate 10^{-4} as training optimizers. The momentum parameters are set to the default values [17].

Test Batch. Since computation on the full-batch objective function is very computationally expensive, we sample a fixed random subset of size 1024 as the test dataset at the beginning of the training, and fix it during all epochs and batches, in order to speed up the landscape visualization process. The losses in all the plots are the losses on the test batch.

Landscape Visualization. To visualize the landscape, we update the weight with the desired update step and compute the loss. Afterwards, we reset the weight back to the original value before the update, and repeat the above step with a new step size.

Directional Sharpness. We utilize PyTorch’s Hessian-vector product to efficiently compute directional sharpness. We sample 5 batches from the 10 batches in the epoch. To show the effect of clipping and adaptive update steps on the sharpness compared to SGD, we calculate the mean of the ratio of sharpness versus SGD sharpness for the sampled batches in these epochs.

A.4. Experiment for Convergence of Clipped Optimization Algorithms

We demonstrate the convergence of clipped optimization algorithms using a 5% clipping threshold. We manually tune the learning rate to find the best learning rate for the experiments. The learning rate configuration of our experiment is shown in Table 2.

Task	Algorithm	Learning Rate
Machine Translation	Adam	2×10^{-3}
	Adafactor, 5% Clipping	Relative
	Adafactor	Relative
	Normalized SGD, 5% Clipping	5×10^{-4}
	Normalized SGD	5×10^{-4}
	SGD, 5% Clipping	8×10^{-1}
	SGD	1×10^{-3}
Masked Language Modeling	Adam	2×10^{-3}
	Adafactor, 5% Clipping	1×10^{-2}
	Adafactor	1×10^{-2}
	Normalized SGD, 5% Clipping	5×10^{-5}
	Normalized SGD	5×10^{-5}
	SGD, 5% Clipping	8×10^{-1}
	SGD	1×10^{-2}

Table 2: Learning rate configuration of our experiments. The relative learning rate for Adafactor is defined in Algorithms 3 and 4 and [27].

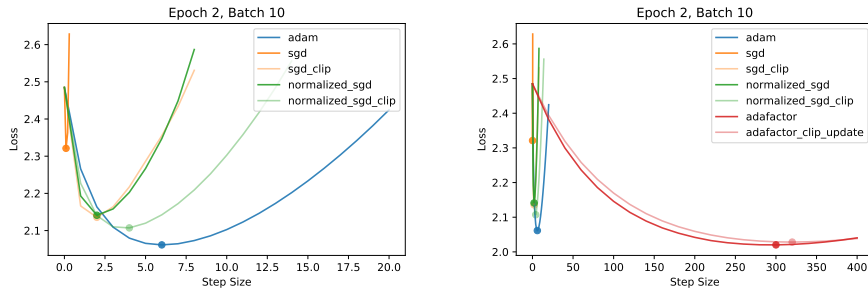
Appendix B. Directional Sharpness Results

In this section we show our experimental result for the directional sharpness of optimization algorithms. For each of the landscape visualization, we show two plots, where one of them has Adafactor and the other does not. The rest of the plots are the same with different scales.

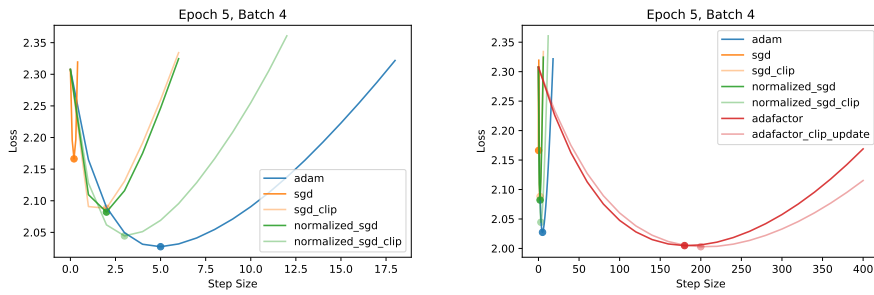
B.1. SGD Trajectory

Task	Epoch	Algorithm	Sharpness Ratio
Machine Translation	2	SGD	1
		SGD, 5% Clipping	0.43546787
		Adam	0.04595183
		Adafactor	1.6432×10^{-5}
		Adafactor, 5% Clipping	9.9343×10^{-6}
		Normalized SGD	0.20991666
		Normalized SGD, 5% Clipping	0.06885640
	5	SGD	1
		SGD, 5% Clipping	0.02766145
		Adam	0.00250487
		Adafactor	6.0198×10^{-7}
		Adafactor, 5% Clipping	4.7015×10^{-7}
		Normalized SGD	0.01427107
		Normalized SGD, 5% Clipping	0.00610456
	10	SGD	1
		SGD, 5% Clipping	0.06981113
		Adam	0.00546100
		Adafactor	1.5864×10^{-6}
		Adafactor, 5% Clipping	1.2210×10^{-6}
		Normalized SGD	0.03849150
		Normalized SGD, 5% Clipping	0.01631434
	20	SGD	1
		SGD, 5% Clipping	0.11132615
		Adam	0.00914083
Adafactor		2.4767×10^{-6}	
Adafactor, 5% Clipping		1.9563×10^{-6}	
Normalized SGD		0.06279230	
Normalized SGD, 5% Clipping		0.02604086	

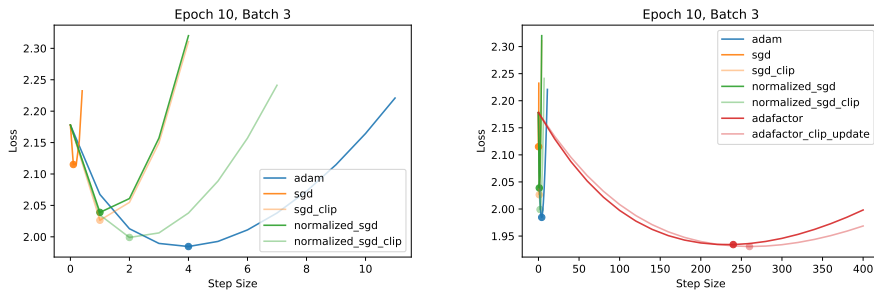
Table 3: Average ratio of directional sharpness of optimization algorithms with respect to SGD on the machine translation task in SGD trajectory.



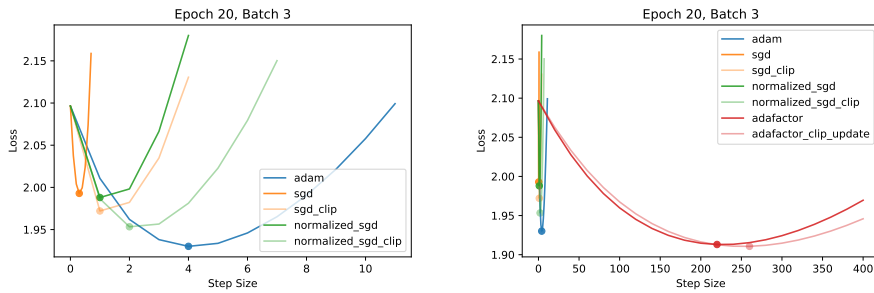
(a) Epoch 2, Batch 10



(b) Epoch 5, Batch 4



(c) Epoch 10, Batch 3

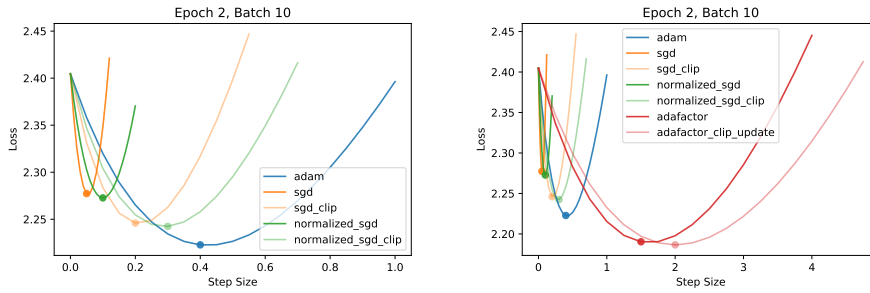


(d) Epoch 20, Batch 3

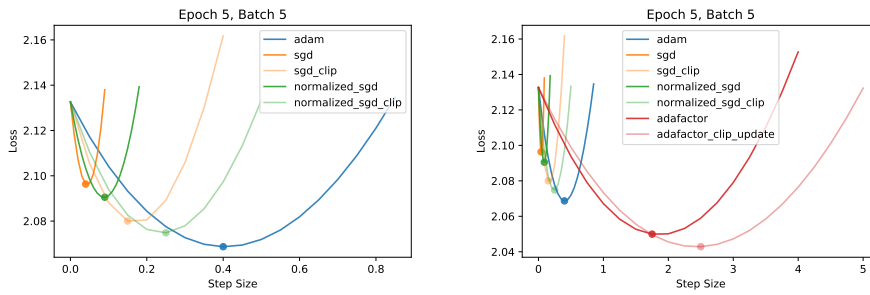
Figure 5: Landscape visualization of machine translation in SGD trajectory.

Task	Epoch	Algorithm	Sharpness Ratio
Masked Language Modeling	2	SGD	1
		SGD, 5% Clipping	0.10431360
		Adam	0.03353230
		Adafactor	0.00282802
		Adafactor, 5% Clipping	0.00196985
		Normalized SGD	0.28767978
		Normalized SGD, 5% Clipping	0.05043281
	5	SGD	1
		SGD, 5% Clipping	0.12449840
		Adam	0.02682579
		Adafactor	0.00191930
		Adafactor, 5% Clipping	0.00123799
		Normalized SGD	0.26880846
		Normalized SGD, 5% Clipping	0.05914444
	10	SGD	1
		SGD, 5% Clipping	0.11993282
		Adam	0.01577251
		Adafactor	0.00441034
		Adafactor, 5% Clipping	0.00278653
		Normalized SGD	0.30607491
		Normalized SGD, 5% Clipping	0.05791050
	20	SGD	1
		SGD, 5% Clipping	0.09336649
		Adam	0.01320078
Adafactor		0.00893414	
Adafactor, 5% Clipping		0.00542521	
Normalized SGD		0.28826713	
Normalized SGD, 5% Clipping		0.04681089	

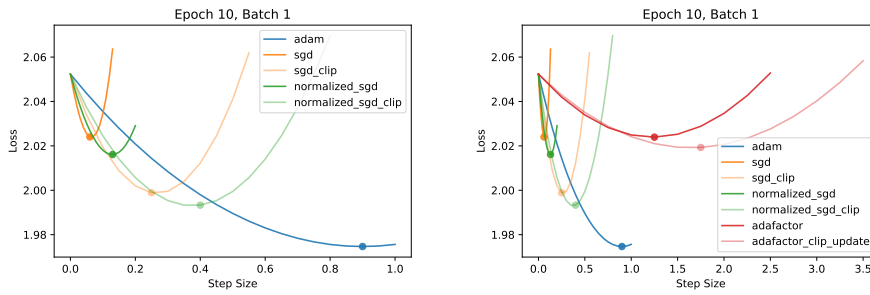
Table 4: Average ratio of directional sharpness of optimization algorithms with respect to SGD on the masked language modeling task in SGD trajectory.



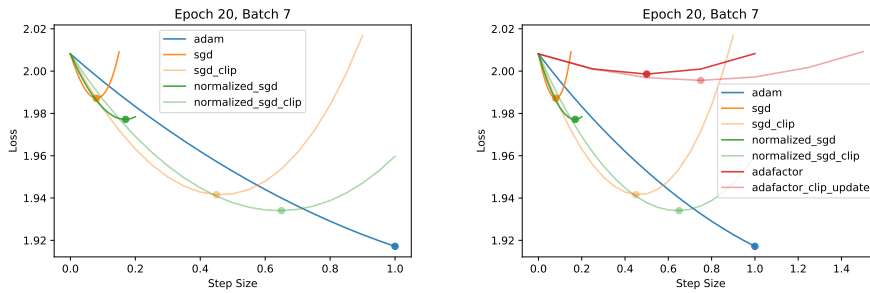
(a) Epoch 2, Batch 10



(b) Epoch 5, Batch 5



(c) Epoch 10, Batch 1



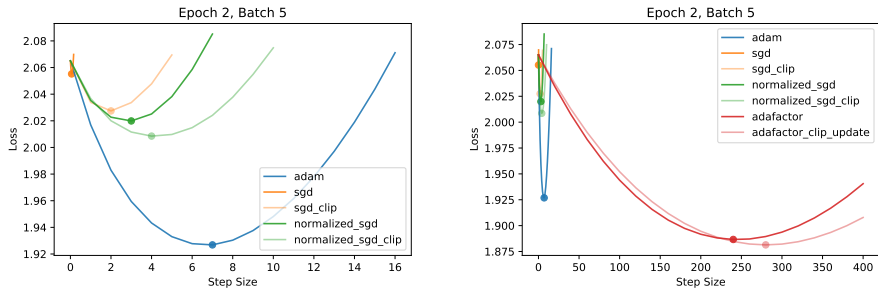
(d) Epoch 20, Batch 7

Figure 6: Landscape visualization of masked language modeling in SGD trajectory.

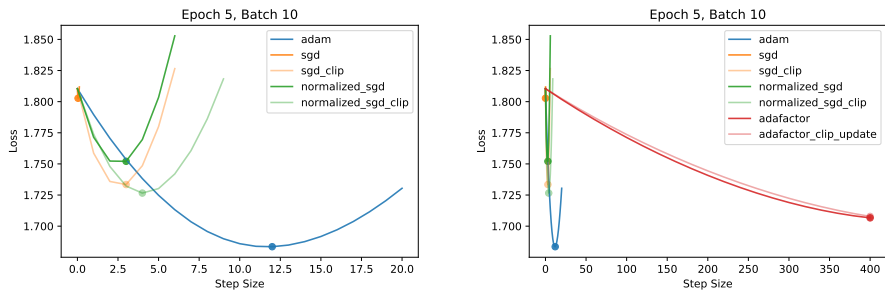
B.2. Adam Trajectory

Task	Epoch	Algorithm	Sharpness Ratio
Machine Translation	2	SGD	1
		SGD, 5% Clipping	0.01189776
		Adam	0.00014977
		Adafactor	6.2772×10^{-7}
		Adafactor, 5% Clipping	4.5537×10^{-7}
		Normalized SGD	0.00938728
		Normalized SGD, 5% Clipping	0.00334502
	5	SGD	1
		SGD, 5% Clipping	0.00341957
		Adam	0.00014246
		Adafactor	4.4903×10^{-8}
		Adafactor, 5% Clipping	3.8484×10^{-8}
		Normalized SGD	0.00173206
		Normalized SGD, 5% Clipping	0.00085046
	10	SGD	1
		SGD, 5% Clipping	0.01228746
		Adam	0.00064037
		Adafactor	3.4227×10^{-6}
		Adafactor, 5% Clipping	2.8321×10^{-6}
		Normalized SGD	0.00975973
		Normalized SGD, 5% Clipping	0.00412512
	20	SGD	1
		SGD, 5% Clipping	0.01194036
		Adam	0.00034823
Adafactor		4.0147×10^{-7}	
Adafactor, 5% Clipping		3.2190×10^{-7}	
Normalized SGD		0.00983759	
Normalized SGD, 5% Clipping		0.00338311	

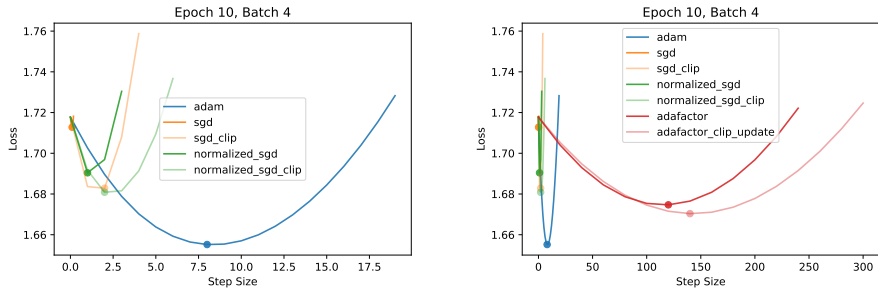
Table 5: Average ratio of directional sharpness of optimization algorithms with respect to SGD on the machine translation task in Adam trajectory.



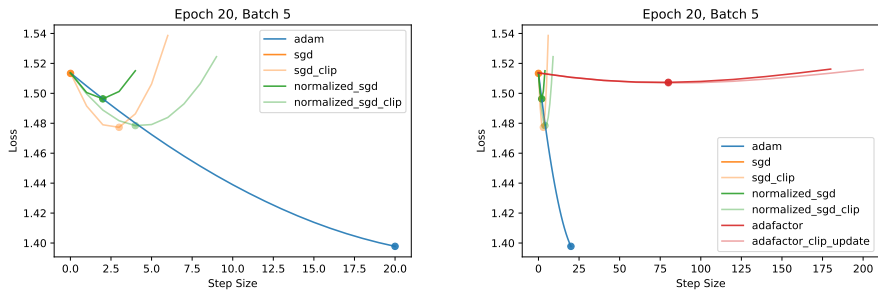
(a) Epoch 2, Batch 5



(b) Epoch 5, Batch 10



(c) Epoch 10, Batch 4

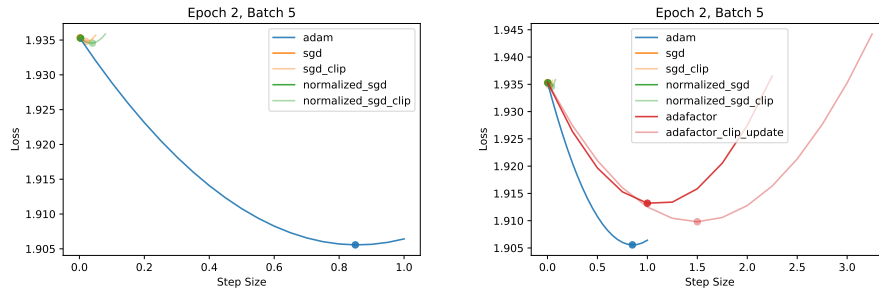


(d) Epoch 20, Batch 5

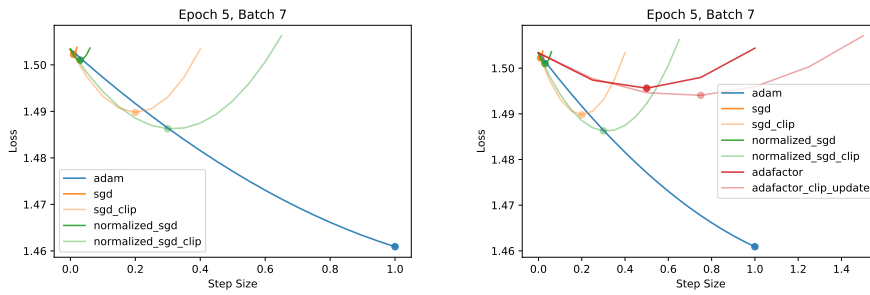
Figure 7: Landscape visualization of machine translation in Adam trajectory.

Task	Epoch	Algorithm	Sharpness Ratio
Masked Language Modeling	2	SGD	1
		SGD, 5% Clipping	0.05873187
		Adam	0.00206153
		Adafactor	0.00144918
		Adafactor, 5% Clipping	0.00088563
		Normalized SGD	0.25169848
		Normalized SGD, 5% Clipping	0.02890596
	5	SGD	1
		SGD, 5% Clipping	0.02350328
		Adam	0.00143950
		Adafactor	0.00240724
		Adafactor, 5% Clipping	0.00141641
		Normalized SGD	0.21113349
		Normalized SGD, 5% Clipping	0.01096425
	10	SGD	1
		SGD, 5% Clipping	0.02781447
		Adam	0.00266029
		Adafactor	0.00387052
		Adafactor, 5% Clipping	0.00222226
		Normalized SGD	0.20546310
		Normalized SGD, 5% Clipping	0.01329525
	20	SGD	1
		SGD, 5% Clipping	0.07415355
		Adam	0.02341449
Adafactor		0.02283569	
Adafactor, 5% Clipping		0.01298895	
Normalized SGD		0.22893490	
Normalized SGD, 5% Clipping		0.03024933	

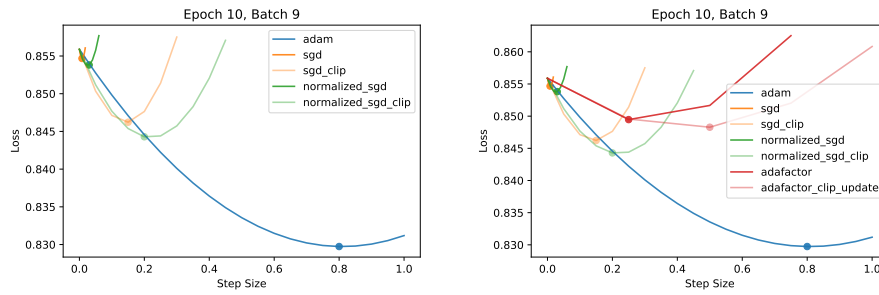
Table 6: Average ratio of directional sharpness of optimization algorithms with respect to SGD on the masked language modeling task in Adam trajectory.



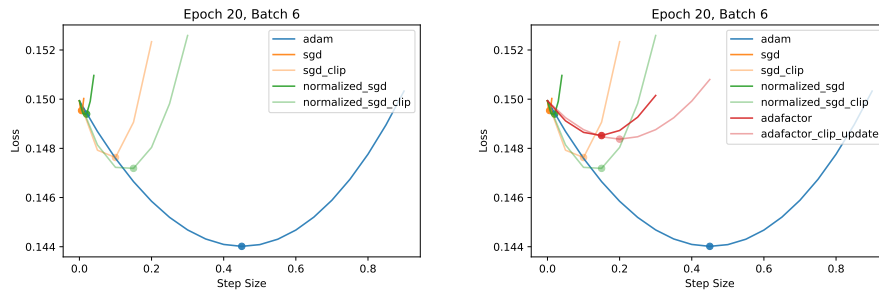
(a) Epoch 2, Batch 5



(b) Epoch 5, Batch 7



(c) Epoch 10, Batch 9



(d) Epoch 20, Batch 6

Figure 8: Landscape visualization of masked language modeling in Adam trajectory.

B.3. Discussion

As we can observe, our observation is very coherent across different tasks, model architectures, iterations, and local geometry. The directional sharpness is relatively stable for the same task across iterations, and coordinate-wise clipping always improve the sharpness of the direction and find a better direction to optimize.

Trade-off Between Directional Sharpness and Gradient Correlation. While we want the directional sharpness of our optimization algorithm to be small in order to decrease loss faster, having as small sharpness as possible does not necessarily lead to fast loss decrement. Adafactor almost always has the lowest directional sharpness across all tasks, iterations, and local geometry, but Adafactor does not always find a good direction to optimize. In many cases, the loss does not decrease significantly even for the optimal step size, and the direction can be even worse than SGD. This shows that merely minimizing the directional sharpness is not enough for an optimization algorithm to work well. As discussed in Section 4, gradient correlation is also important in the convergence of optimization algorithms. However, we can conclude that high sharpness will lead to bad performance, as demonstrated by the performance of SGD, even if SGD has good gradient correlation.

Effect of Trajectory. It is well known that different optimization algorithms can follow different trajectory and converge to different in deep learning. Landscape visualizations show that Adafactor performs well in SGD trajectory but not Adam trajectory on the machine translation task. This shows that different optimization algorithms has local geometry with different properties. The effect of trajectory is therefore an interesting problem to study. However, we point out that almost in all cases, Adam has good performance and significantly outperforms SGD, so trajectory is not necessarily related to the explanation for Adam’s excellent performance in practice.