
Controlling changes to attention logits

Ben Anson
University of Bristol
ben.anson@bristol.ac.uk

Laurence Aitchison
Mistral AI
laurence.aitchison@gmail.com

Abstract

Stability of neural network weights is critical when training transformer models. The query and key weights are particularly problematic, as they tend to grow large without any intervention. Applying normalization to queries and keys, known as ‘QK norm’, fixes stability issues in practice, but is not always applicable. For example, QK norm is not compatible with Multi Latent Attention (MLA) because QK norm requires full materialization of queries and keys during inference, which is not done in MLA. In this paper we suggest that controlling the changes to logits is important for stability. We show that these changes are controllable by assigning parameter-dependent learning rates to the query and key weights. We find that our cheap intervention allows us to increase the base learning rate of the network, outperform other methods in the MLA setting, and achieve performance competitive with QK norm when using Multi-head Attention.

1 Introduction

Principled scaling of transformer models is crucial for efficiently training larger and more capable architectures. Maximal Update Parametrization (μ P) (Yang et al., 2022) has emerged as a key technique in this area, enabling the transfer of optimal hyperparameters from smaller to larger models by carefully parameterizing the model. A core desideratum of μ P is to control the magnitude of activations and their updates (Dey et al., 2025), ensuring consistent training dynamics across different model widths. Regarding attention, Yang et al. (2022) address attention logits blow-up as we increase model width by proposing a static attention scaling factor. While this static scaling helps control logit magnitude across different model widths, it does not address step-to-step changes in logits during longer training runs, which can become a major source of instability, particularly at high learning rates.

Attention logits are a well-known source of training instability (Zhai et al., 2023; Bai et al., 2025), prompting the development of interventions such as QK norm (Henry et al., 2020) and QK clip/MuonClip (Bai et al., 2025) to ensure their stability. While QK norm is especially effective, it is ill-suited for Multi-head Latent attention (MLA) (Liu et al., 2024), as queries and keys are not fully materialized at inference-time for efficiency reasons (Bai et al., 2025). Other methods like QK clip require a bespoke attention to track maximum attention logits, which can complicate integration into existing codebases. Thus there is a gap for an easy-to-implement intervention that improves training stability, but is more widely applicable than QK norm.

We approach this gap with a μ P-inspired desideratum for the attention logits. Instead of constraining logit magnitudes, like QK clip, MuonClip, and QK norm, we seek to control the *change* in logits as we train. This preserves expressivity, while reining-in instability. To validate our approach, we conduct pretraining experiments on a 1B parameter model. Our results demonstrate that our method is as stable as QK norm, particularly at high base learning rates. While not quite reaching the same peak performance as QK norm in the standard Multi-Head Attention (MHA) setting, our method is computationally cheaper and is applicable to MLA. When used with MLA our method enables higher base learning rates, outperforms QK clip, highlighting its practical value for training modern, efficient transformer architectures.

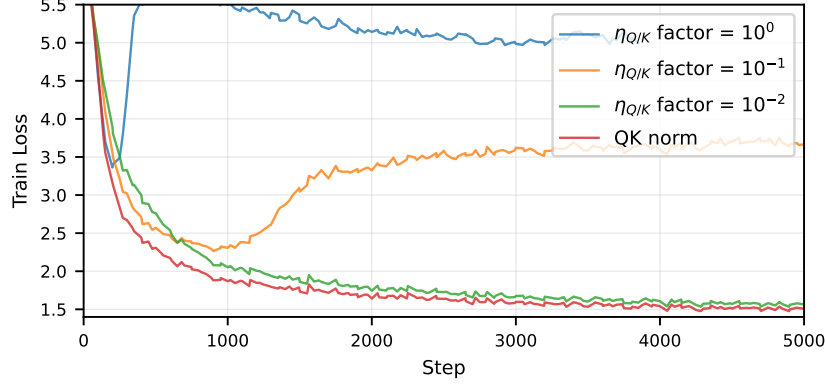


Figure 1: **Learning rate of query/key matrices is a critical factor for transformer pretraining stability.** Here, 4 models are trained with a large base learning rate of $\eta = 3e-2$ for each parameter. Decreasing the learning rates of query and key weights alone (by a factor of $\eta_{Q/K}$), fully stabilizes pretraining. QK norm is shown to illustrate a stable baseline.

In summary, our contributions are as follows,

- We propose that change in logits is an important metric to account for when training attention modules.
- We show that we can control the change in logits by modulating the learning rate of query weight based on the norms of corresponding key weight, and vice versa.
- We demonstrate that this change to the learning rate leads to better validation loss performance than alternatives when training with MLA.

2 Related Work

Many have encountered instability issues when training transformers, and it has been studied extensively (Liu et al., 2020; Dehghani et al., 2023; Henry et al., 2020; Bai et al., 2025; Wortsman et al., 2023; Qi et al., 2023; Kim et al., 2025; Zhai et al., 2023; Takase et al., 2023; Rybakov et al., 2024). Below we discuss the past literature relevant to our work.

Controlling attention logits. Training instabilities are often encountered in the attention layer itself. Attention logits may become large (Bai et al., 2025), potentially inducing collapse in attention entropy (Zhai et al., 2023), where attention distributions become highly concentrated. QK normalization (Henry et al., 2020), which applies normalization to query and key activations, has emerged as a simple and effective remedy, preventing large logits (Dehghani et al., 2023) and allowing larger learning rates (Wortsman et al., 2023). Similar methods such as logit soft-capping apply normalization to logits directly (Bello et al., 2016; Riviere et al., 2024). Other methods normalize the weights rather than activations: σ Reparam (Zhai et al., 2023) parameterizes weights into a matrix and a scalar component, with the matrix having a maximum singular value of 1, and a weights rather than the activations; QK clip (Bai et al., 2025) controls attention logits by clipping weights whenever the logits grow beyond a certain threshold.

Parameter-specific learning rates. While it is common to share the same learning rate across all parameters in a neural network, parameter-specific learning rates have been extensively examined (Milsom et al., 2025; You et al., 2017; Liu et al., 2019; Xu et al., 2019; Wang et al., 2025; Bernstein et al., 2020; Qi et al., 2025; Yang et al., 2023). Proposals often include adjusting the learning rate of a parameter according to the norm of step/gradient (Yang et al., 2023; Liu et al., 2019), as well as the parameter itself (Qi et al., 2025), such as LARS, LAMB, and Fromage (Bernstein et al., 2020; You et al., 2017, 2019).

Our work selects parameter-specific learning rates that control changes to attention logits. However, by considering attention logits as a whole, our parameter-specific learning rates are ‘inter-parameter’,

unlike other methods, such as LARS, which consider each parameter tensor independently. Our method is also inspired by μP (Yang et al., 2022; Dey et al., 2025); in μP , one of the desiderata is that as we make changes to our parameters in a network, the residual stream should correspondingly change in a controlled, ‘order 1-like’ manner. Our work extends this notion to logits.

3 Methods

Unlike other transformer modules, attention has quadratic structure. In particular, the attention logits are given by,

$$\mathbf{L} = \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} = \frac{\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^T\mathbf{X}^T}{\sqrt{d}}. \quad (1)$$

Inspired by μP (Yang et al., 2022; Dey et al., 2025), which (among other things) attempts to keep changes to *activations* roughly constant, we are interested in keeping the changes to *attention logits*, $\Delta\mathbf{L}$, under control. By a first-order analysis, we see that if the queries are large, then perturbations due to the keys will be amplified, and vice versa:

$$\Delta\mathbf{L} = \frac{(\mathbf{Q} + \Delta\mathbf{Q})(\mathbf{K} + \Delta\mathbf{K})^T - \mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \approx \frac{\mathbf{Q}(\Delta\mathbf{K})^T + (\Delta\mathbf{Q})\mathbf{K}^T}{\sqrt{d}}. \quad (2)$$

The main tool we have for controlling changes is the learning rate. Thus we propose to set the learning rates η_Q , η_K (for \mathbf{W}_Q , and \mathbf{W}_K respectively) such that $\mathbf{Q}(\Delta\mathbf{K})^T$ and $(\Delta\mathbf{Q})\mathbf{K}^T$ are both ‘order 1’. We formalize this notion in Lemma 1.

Lemma 1. *Let $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$ be weight matrices corresponding to a particular attention head, and consider the worst-case change in logits, for unit normed input,*

$$\max_{\|x\|_2=\|y\|_2=1} |\Delta\ell| := \max_{\|x\|_2=\|y\|_2=1} |x^\top (\mathbf{W} + \Delta\mathbf{W})y - x^\top \mathbf{W}y|,$$

where $\mathbf{W} = d_{\text{head}}^{-1/2} \mathbf{W}_Q^\top \mathbf{W}_K$. Suppose that the steps for \mathbf{W}_Q and \mathbf{W}_K are given by $\Delta\mathbf{W}_{Q/K} = -\eta_{Q/K} \mathbf{G}_{Q/K}$, where $\|\mathbf{G}_{Q/K}\| \leq D$ for some constant D (which is the case for Adam and Muon). If there is a constant c such that $0 < c \leq \|\mathbf{W}_Q\|$, $\|\mathbf{W}_K\|$, and the learning rates satisfy $\eta_Q \propto \|\mathbf{W}_K\|^{-1}$, and $\eta_K \propto \|\mathbf{W}_Q\|^{-1}$, then the worst-case change in logits is bounded above independently of the weight size.

For Lemma 1 to apply, we need a constant c such that $c \leq \|\mathbf{W}_{Q/K}\|$, but this is not unreasonable in practice. The Lemma does not specify a norm because all norms are equivalent, though in practice, we do need to pick a norm for an implementation. The most natural norm for restricting the maximum change to the logits is perhaps the spectral norm. In a preliminary experiment, we compared the performance of both Frobenius and spectral norm, with results shown in Figure 3. The benefits of using the spectral norm are very small, thus we opted to use the Frobenius norm for further experiments in Section 4.

Following the Lemma we set,

$$\eta_Q \propto \|\mathbf{W}_K\|^{-1}, \quad \eta_K \propto \|\mathbf{W}_Q\|^{-1}. \quad (3)$$

In practice we treat the constant of proportionality in Eq. (3) as a hyperparameter: at initialization, we set the learning rate for each query and key weight to be equal to $\tau\eta$ and we tune τ . Thus τ acts as a relative initial learning rate (relative to η , the base learning rate).

The above methodology applies to both the single- and multi-head (MHA) setting. In MHA, each head has its own query and key weight, so we apply Eq. (3) to each head separately. We summarize the resulting method in Algorithm 1.

We extend to MLA using a similar approach in Appendix B. A notable difference between MHA and MLA is that there are several more parameter matrices to consider; bounding the change in logits requires us to adjust the learning rate of each of these parameters. We detail exactly how to set the learning rates for MLA in Algorithm 2.

Algorithm 1 QuacK (MHA)

Require: Hyperparameter τ , base learning rate η

Make the following additions to the transformer training script:

At initialization. Calculate initial norms for query/key weights for all heads

```
for all layers  $\ell$  do
  for all heads  $h$  do
     $\mathbf{W}_Q^{\ell,h}.\text{init\_norm} \leftarrow \|\mathbf{W}_Q^{\ell,h}\|$ 
     $\mathbf{W}_K^{\ell,h}.\text{init\_norm} \leftarrow \|\mathbf{W}_K^{\ell,h}\|$ 
  end for
end for
```

During training. Prior to each optimization step, adjust learning rates

```
for all layers  $\ell$  do
  for all heads  $h$  do
     $\mathbf{W}_Q^{\ell,h}.\text{lr} \leftarrow \tau \eta \cdot \frac{\mathbf{W}_K^{\ell,h}.\text{init\_norm}}{\|\mathbf{W}_K^{\ell,h}\|}$ 
     $\mathbf{W}_K^{\ell,h}.\text{lr} \leftarrow \tau \eta \cdot \frac{\mathbf{W}_Q^{\ell,h}.\text{init\_norm}}{\|\mathbf{W}_Q^{\ell,h}\|}$ 
  end for
end for
```

4 Experiments

To evaluate our method, we trained $\sim 1\text{B}$ models based on Qwen3 (Yang et al., 2025) with both MHA (Vaswani et al., 2017) and MLA (Liu et al., 2024). All models used $d_{\text{model}} = 2048$, $d_{\text{ff}} = 4d_{\text{model}}$, $n_{\text{head}} = 32$, $n_{\text{layer}} = 14$, and were trained using gradient accumulation at 2048 context length with 96 sequences per batch (i.e. 196608 tokens per batch) for 5000 steps, using data from the Cosmopedia-V2 subset of SmolLM-corpus (Ben Allal et al., 2024). We used the GPT-2 (Radford et al., 2019) tokenizer, with vocab size 49152 and embedding/unembedding weight tying. The MHA models were trained with $d_{\text{head}} = 64$, while the MLA models were trained with $d_{\text{head}} = 128$ (with $d_{\text{rope}} = d_{\text{nope}} = 64$). MLA also used latent dimensions of $d_{\text{cq}} = 512$ for the queries and $d_{\text{ckv}} = 256$ for the keys and values. We trained using Muon (Jordan et al., 2024), with constant LR schedule and 500 warmup steps. All pretraining runs were completed on 4xGH200 nodes at bfloat16 precision.

Our experiments varied the attention method, $\text{attn} \in \{\text{MHA}, \text{MLA}\}$, the base learning rate $\eta \in \{3e-4, 3e-3, 3e-2\}$, as well as the attention logit interventions:

- *QK norm* applies RMS norm (with learned scaling) to both queries and keys before applying the attention operation.
- *QK clip* implements Algorithm 1 from Bai et al. (2025), where after each optimization step we multiply/clip certain weights using either $\sqrt{\gamma}$ or γ . We set $\gamma = \min\{1, \tau_{\text{QK clip}}/S_{\text{max}}^h\}$, where $\tau_{\text{QK clip}}$ is a hyperparameter denoting a threshold for the maximum logit, and S_{max}^h is the largest logit value seen by the h 'th head since the last optimization step. For MHA we use $\mathbf{W}_{Q/K}^h \leftarrow \sqrt{\gamma} \mathbf{W}_{Q/K}^h$; for MLA we use $\mathbf{W}_{\text{uq/uk}} \leftarrow \sqrt{\gamma} \mathbf{W}_{\text{uq/uk}}$ and $\mathbf{W}_{\text{qr}} \leftarrow \gamma \mathbf{W}_{\text{qr}}$. We swept over $\tau_{\text{QK clip}} \in \{30, 100\}$, the values used in the original paper.
- *Ablation* multiplies learning rates for query and key weights by a value τ , which is swept over, $\tau \in \{10^{-2}, 10^1, 10^0, 10^1\}$. For MHA we set $\eta_Q^h = \eta_K^h = \tau \cdot \eta$. For MLA, we set $\eta_{\text{uq}} = \eta_{\text{dq}} = \eta_{\text{qr}} = \eta_{\text{uk}} = \eta_{\text{dkv}} = \eta_{\text{kr}} = \tau \cdot \eta$.
- *QuacK* also multiplies learning rates for query and key weights according to a value τ , which is swept over, $\tau \in \{10^{-2}, 10^1, 10^0, 10^1\}$. For MHA we use Algorithm 1, and for MLA we use Algorithm 2. All weight norms are calculated using the Frobenius norm.

Higher learning rates are better. Figure 2, left column, shows that at the low learning rate of $\eta = 3e-4$, all logit interventions perform similarly, but with QK norm performing marginally better.

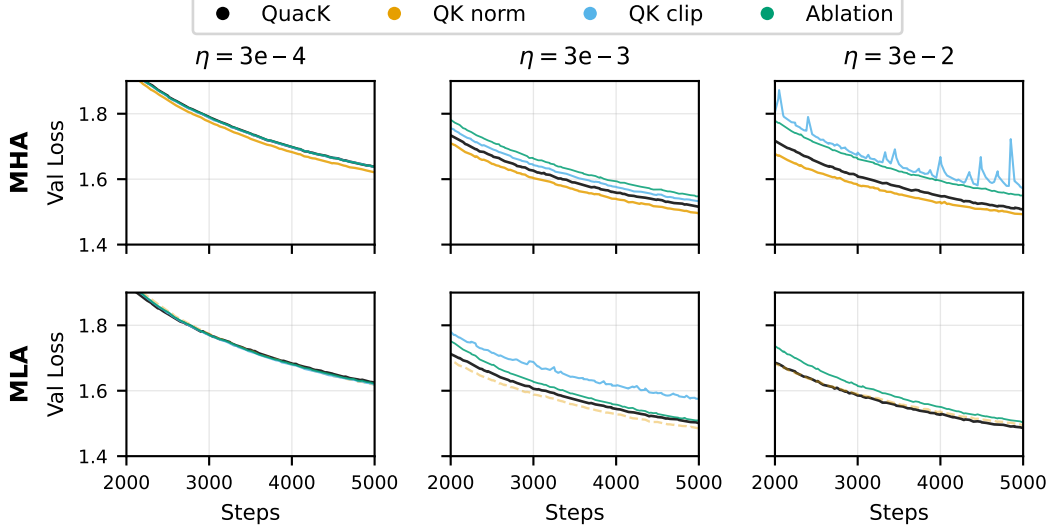


Figure 2: Validation losses when training each method with $\text{attn} \in \{\text{MHA}, \text{MLA}\}$, and learning rates, $\eta \in \{3e-4, 3e-3, 3e-2\}$. QK clip is unstable at high learning rates (it is omitted from the bottom right plot due to $\text{loss} \gg 2$). QK norm is overall the most performant, but it is not appropriate for use with MLA at inference-time for efficiency reasons (illustrated via dashed yellow line in the MLA row). QuacK is a sensible alternative, as it is stable in the high LR setting, performant, and applicable in the MLA setting.

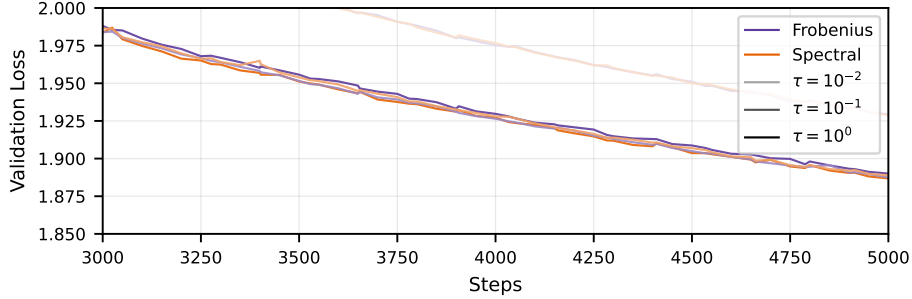


Figure 3: Performance differences when applying Algorithm 1 with different norms are small. We show validation losses when training a small model ($\sim 100\text{M}$ parameters) with Algorithm 1 to modulate the query and key weight learning rates. Different curves show results with different values of the hyperparameter τ and measuring the query and key weights with either Frobenius or spectral norm.

The lack of variety in performance is likely due to the fact the learning rate is small enough that we don’t encounter instabilities. However, performance is much improved by increasing the learning rate (column 2, 3).

QuacK maintains stability and strong performance, enabling higher base learning rates. QK clip is insufficient to prevent instabilities at the highest base learning rate of $\eta = 3e-2$ (column 3, Figure 2), and it underperforms, especially in the MLA setting, when $\eta = 3e-3$ (column 2). The ablation, which sets the learning rates for query and key weights to smaller fixed values, is stable, but underperforms QuacK in both the MHA and MLA settings. QuacK has similar but slightly worse performance compared to QK norm in the MHA setting, but is the best performing method in the MLA setting (we include QK norm results in the MLA setting for comparison, but it is not viable at inference-time like the other methods).

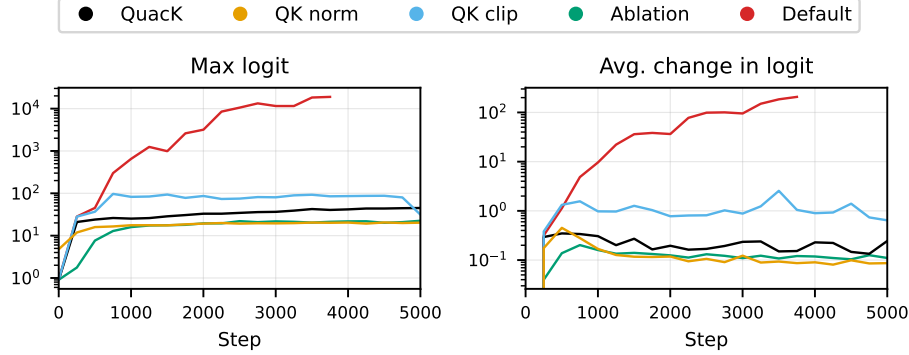


Figure 4: Max logit (left) and average absolute change in logit throughout training (right) with a base learning rate of $\eta = 3e - 3$. Here we show the middle head of the middle layer (head 16 and layer 8) while training with MLA.

QuacK controls the maximum logit as well as change in logits. During our training runs, we used a dedicated subset of our data to periodically record logit statistics. We confirm that QuacK controls the quantity we expect (average absolute change in logits) in Figure 4. We also see that either intervening with QK norm or controlling the learning rate appears sufficient to control the maximum logits. We include the ‘default’ model (i.e. using a base learning rate of $\eta = 3e - 3$ for all parameters), to illustrate behaviour of an unstable run (we killed this run prematurely, as it never got below 4.5 loss after 1.5k steps).

QuacK gives a performance boost over QK norm. QuacK yields a speedup over QK norm by removing two RMS norm computations per attention block; in practice we observed $\sim 10\%$ faster training. QuacK was also faster than QK clip, though this comparison is slightly unfair because we did not use an optimized QK clip implementation. QK clip can be efficient, but requires custom attention code to efficiently accumulate the maximum logit.

5 Limitations

Our experiments used a single model based on Qwen3 (Yang et al., 2025). Due to compute constraints, we were unable to execute longer training runs and with larger models to demonstrate that our method is widely applicable. As such, results are limited by short training durations (5k steps) and a single dataset and model architecture.

6 Conclusion

In this work, we introduced a simple and principled approach for stabilizing attention training by controlling changes in attention logits rather than their magnitude. Our analysis showed that the expected logit change can be bounded through parameter-dependent learning rates for the query and key weights, inspired by μP -style scaling principles. Empirically, our method QuacK enables the use of higher base learning rates while maintaining stability comparable to QK norm and outperforming alternative methods such as QK clip, particularly in the Multi-Latent Attention (MLA) setting where QK norm is inappropriate.

Our results demonstrate that stability in attention can be achieved without introducing additional normalization layers or specialized kernels, making QuacK a practical drop-in improvement for transformer training.

7 Acknowledgements

We thank Edward Milsom for insightful discussions.

References

- Bai, Y., Bao, Y., Chen, G., Chen, J., Chen, N., Chen, R., Chen, Y., Chen, Y., Chen, Y., et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- Bello, I., Pham, H., Le, Q. V., Norouzi, M., and Bengio, S. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- Ben Allal, L., Lozhkov, A., Penedo, G., Wolf, T., and von Werra, L. Smollm-corpus, 2024. URL <https://huggingface.co/datasets/HuggingFaceTB/smollm-corpus>.
- Bernstein, J., Vahdat, A., Yue, Y., and Liu, M.-Y. On the distance between two neural networks and the stability of learning. *Advances in Neural Information Processing Systems*, 33:21370–21381, 2020.
- Dehghani, M., Djolonga, J., Mustafa, B., Padlewski, P., Heek, J., Gilmer, J., Steiner, A. P., Caron, M., Geirhos, R., Alabdulmohsin, I., et al. Scaling vision transformers to 22 billion parameters. In *International conference on machine learning*, pp. 7480–7512. PMLR, 2023.
- Dey, N., Zhang, B. C., Noci, L., Li, M., Bordelon, B., Bergsma, S., Pehlevan, C., Hanin, B., and Hestness, J. Don’t be lazy: Completep enables compute-efficient deep transformers. *arXiv preprint arXiv:2505.01618*, 2025.
- Henry, A., Dachapally, P. R., Pawar, S., and Chen, Y. Query-key normalization for transformers. *arXiv preprint arXiv:2010.04245*, 2020.
- Ji, T., Guo, B., Wu, Y., Guo, Q., Shen, L., Chen, Z., Qiu, X., Zhang, Q., and Gui, T. Towards economical inference: Enabling deepseek’s multi-head latent attention in any transformer-based llms. *arXiv preprint arXiv:2502.14837*, 2025.
- Jordan, K., Jin, Y., Boza, V., You, J., Cesista, F., Newhouse, L., and Bernstein, J. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.
- Kim, J., Lee, B., Park, C., Oh, Y., Kim, B., Yoo, T., Shin, S., Han, D., Shin, J., and Yoo, K. M. Peri-In: Revisiting normalization layer in the transformer architecture. *arXiv preprint arXiv:2502.02732*, 2025.
- Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J., and Han, J. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019.
- Liu, L., Liu, X., Gao, J., Chen, W., and Han, J. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249*, 2020.
- Milson, E., Anson, B., and Aitchison, L. Function-space learning rates. In *Proceedings of the 42nd International Conference on Machine Learning*, 2025.
- Qi, X., Wang, J., Chen, Y., Shi, Y., and Zhang, L. Lipsformer: Introducing lipschitz continuity to vision transformers. *arXiv preprint arXiv:2304.09856*, 2023.
- Qi, X., He, Y., Ye, J., Li, C.-G., Zi, B., Dai, X., Zou, Q., and Xiao, R. Taming transformer without using learning rate warmup. *arXiv preprint arXiv:2505.21910*, 2025.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Rybakov, O., Chrzanowski, M., Dykas, P., Xue, J., and Lanir, B. Methods of improving llm training stability. *arXiv preprint arXiv:2410.16682*, 2024.

- Su, J., Lu, Y., Pan, S., Murtadha, A., Wen, B., and Liu, Y. Roformer: enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- Takase, S., Kiyono, S., Kobayashi, S., and Suzuki, J. Spike no more: Stabilizing the pre-training of large language models. *arXiv preprint arXiv:2312.16903*, 2023.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wang, J., Wang, M., Zhou, Z., Yan, J., Wu, L., et al. The sharpness disparity principle in transformers for accelerating language model pre-training. *arXiv preprint arXiv:2502.19002*, 2025.
- Wortsman, M., Liu, P. J., Xiao, L., Everett, K., Alemi, A., Adlam, B., Co-Reyes, J. D., Gur, I., Kumar, A., Novak, R., et al. Small-scale proxies for large-scale transformer training instabilities. *arXiv preprint arXiv:2309.14322*, 2023.
- Xu, Z., Dai, A. M., Kemp, J., and Metz, L. Learning an adaptive learning rate schedule. *arXiv preprint arXiv:1909.09712*, 2019.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Yang, G., Hu, E. J., Babuschkin, I., Sidor, S., Liu, X., Farhi, D., Ryder, N., Pachocki, J., Chen, W., and Gao, J. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. *arXiv preprint arXiv:2203.03466*, 2022.
- Yang, G., Simon, J. B., and Bernstein, J. A spectral condition for feature learning. *arXiv preprint arXiv:2310.17813*, 2023.
- You, Y., Gitman, I., and Ginsburg, B. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
- You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., and Hsieh, C.-J. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019.
- Zhai, S., Likhomanenko, T., Littwin, E., Busbridge, D., Ramapuram, J., Zhang, Y., Gu, J., and Susskind, J. M. Stabilizing transformer training by preventing attention entropy collapse. In *International Conference on Machine Learning*, pp. 40770–40803. PMLR, 2023.

A Proof of Lemma 1

Lemma 1. Let $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d_{model} \times d_{head}}$ be weight matrices corresponding to a particular attention head, and consider the worst-case change in logits, for unit normed input,

$$\max_{\|x\|_2=\|y\|_2=1} |\Delta \ell| := \max_{\|x\|_2=\|y\|_2=1} |x^\top (\mathbf{W} + \Delta \mathbf{W})y - x^\top \mathbf{W}y|,$$

where $\mathbf{W} = d_{head}^{-1/2} \mathbf{W}_Q^\top \mathbf{W}_K$. Suppose that the steps for \mathbf{W}_Q and \mathbf{W}_K are given by $\Delta \mathbf{W}_{Q/K} = -\eta_{Q/K} \mathbf{G}_{Q/K}$, where $\|\mathbf{G}_{Q/K}\| \leq D$ for some constant D (which is the case for Adam and Muon). If there is a constant c such that $0 < c \leq \|\mathbf{W}_Q\|, \|\mathbf{W}_K\|$, and the learning rates satisfy $\eta_Q \propto \|\mathbf{W}_K\|^{-1}$, and $\eta_K \propto \|\mathbf{W}_Q\|^{-1}$, then the worst-case change in logits is bounded above independently of the weight size.

Proof. The change in logits is given by,

$$\begin{aligned} d_{head}^{1/2} |\Delta \ell| &= |(q + \Delta q)^T (k + \Delta k) - q^T k| \\ &= |(\Delta q)^T k + q^T \Delta k + (\Delta q)^T \Delta k| \\ &\leq |(\Delta q)^T k| + |q^T \Delta k| + \|\Delta q\| \|\Delta k\|. \end{aligned} \quad (4)$$

where,

$$q = \mathbf{W}_Q x, \quad k = \mathbf{W}_K y \quad (5)$$

The query and key perturbations are given by,

$$\Delta q = (\mathbf{W}_Q + \Delta \mathbf{W}_Q)x - \mathbf{W}_Q x = \Delta \mathbf{W}_Q x, \quad (6)$$

$$\Delta k = (\mathbf{W}_K + \Delta \mathbf{W}_K)y - \mathbf{W}_K y = \Delta \mathbf{W}_K y. \quad (7)$$

We now bound the first order terms in Eq. (4), assuming inputs are unit normed,

$$\begin{aligned} |(\Delta q)^T k| &= |(\Delta \mathbf{W}_Q x)^T (\mathbf{W}_K y)| = |x^T \Delta \mathbf{W}_Q^T \mathbf{W}_K y| \\ &\leq \|x\| \|\Delta \mathbf{W}_Q^T \mathbf{W}_K\| \|y\| \\ &\leq \eta_Q D \|\mathbf{W}_K\|, \end{aligned} \quad (8a)$$

$$\begin{aligned} |q^T \Delta k| &= |(\mathbf{W}_Q x)^T (\Delta \mathbf{W}_K y)| = |x^T \mathbf{W}_Q^T \Delta \mathbf{W}_K y| \\ &\leq \|x\| \|\mathbf{W}_Q^T \Delta \mathbf{W}_K\| \|y\| \\ &\leq \eta_K D \|\mathbf{W}_Q\|. \end{aligned} \quad (8b)$$

For some constants τ_Q, τ_K , set,

$$\eta_Q = \tau_Q \|\mathbf{W}_K\|^{-1} \quad (9a)$$

$$\eta_K = \tau_K \|\mathbf{W}_Q\|^{-1}. \quad (9b)$$

Substituting these into Eqs. (8), we obtain the bounds,

$$|(\Delta q)^T k| \leq (\tau_Q \|\mathbf{W}_K\|^{-1}) D \|\mathbf{W}_K\| = \tau_Q D, \quad (10)$$

$$|q^T \Delta k| \leq (\tau_K \|\mathbf{W}_Q\|^{-1}) D \|\mathbf{W}_Q\| = \tau_K D. \quad (11)$$

Note that even if RoPE is applied, such that $q = R_x \mathbf{W}_Q x$, the bound remains identical as $\|R \mathbf{W}\| = \|\mathbf{W}\|$ (if the Frobenius or spectral is used).

Finally, we consider the quadratic term $\|\Delta q\| \|\Delta k\|$,

$$\begin{aligned} \|\Delta q\| \|\Delta k\| &\leq \|\Delta \mathbf{W}_Q\| \|\Delta \mathbf{W}_K\| \\ &= \frac{\tau_Q \tau_K \|\mathbf{G}_Q\| \|\mathbf{G}_K\|}{\|\mathbf{W}_Q\| \|\mathbf{W}_K\|} \\ &\leq \frac{\tau_Q \tau_K D^2}{c^2}. \end{aligned} \quad (12)$$

Thus the change in logits is bounded by a constant. \square

Algorithm 2 QuacK (MLA)

Require: Hyperparameter τ , base learning rate η

Make the following additions to the transformer training script:

```
function compute_lr_factors()
  for all layers  $\ell$  do
    for all heads  $h$  do
       $\mathbf{W}_{\text{uq}}^{\ell,h}.\text{factor} \leftarrow (\|\mathbf{W}_{\text{dq}}^{\ell}\| \|\mathbf{W}_{\text{uk}}^{\ell,h}\| \|\mathbf{W}_{\text{dkv}}^{\ell}\|)^{-1}$ 
       $\mathbf{W}_{\text{uk}}^{\ell,h}.\text{factor} \leftarrow (\|\mathbf{W}_{\text{uq}}^{\ell,h}\| \|\mathbf{W}_{\text{dq}}^{\ell}\| \|\mathbf{W}_{\text{dkv}}^{\ell}\|)^{-1}$ 
       $\mathbf{W}_{\text{qr}}^{\ell,h}.\text{factor} \leftarrow (\|\mathbf{W}_{\text{dq}}^{\ell}\| \|\mathbf{W}_{\text{kr}}^{\ell}\|)^{-1}$ 
    end for
     $\mathbf{W}_{\text{dq}}^{\ell}.\text{factor} \leftarrow \min \left\{ (\max_h \|\mathbf{W}_{\text{uq}}^{\ell,h}\| \|\mathbf{W}_{\text{uk}}^{\ell,h}\| \|\mathbf{W}_{\text{dkv}}^{\ell}\|)^{-1}, (\max_h \|\mathbf{W}_{\text{qr}}^{\ell,h}\| \|\mathbf{W}_{\text{kr}}^{\ell}\|)^{-1} \right\}$ 
     $\mathbf{W}_{\text{dkv}}^{\ell}.\text{factor} \leftarrow (\max_h \|\mathbf{W}_{\text{uq}}^{\ell,h}\| \|\mathbf{W}_{\text{dq}}^{\ell}\| \|\mathbf{W}_{\text{uk}}^{\ell,h}\|)^{-1}$ 
     $\mathbf{W}_{\text{kr}}^{\ell}.\text{factor} \leftarrow (\max_h \|\mathbf{W}_{\text{qr}}^{\ell,h}\| \|\mathbf{W}_{\text{dq}}^{\ell}\|)^{-1}$ 
  end for
end function

 $\{\text{attention\_weights}\} \leftarrow \{\mathbf{W}_{\text{uq}}^{\ell,h}, \mathbf{W}_{\text{uk}}^{\ell,h}, \mathbf{W}_{\text{qr}}^{\ell,h}, \mathbf{W}_{\text{dq}}^{\ell}, \mathbf{W}_{\text{dkv}}^{\ell}, \mathbf{W}_{\text{kr}}^{\ell} \text{ for all layers } \ell \text{ for all heads } h\}$ 

# At initialization. Compute initial learning rate factors for all attention weights
compute_lr_factors()
for all  $\mathbf{W}$  in  $\{\text{attention\_weights}\}$  do
   $\mathbf{W}.\text{init\_factor} \leftarrow \mathbf{W}.\text{factor}$ 
end for

# During training. Prior each optimization step, adjust learning rates
compute_lr_factors()
for all  $\mathbf{W}$  in  $\{\text{attention\_weights}\}$  do
   $\mathbf{W}.\text{lr} \leftarrow \tau \eta \cdot \frac{\mathbf{W}.\text{factor}}{\mathbf{W}.\text{init\_factor}}$ 
end for
```

B Extension to MLA

In this section we motivate Algorithm 2, specifically the factors associated with each weight

We use a similar approach to Section 3 / Appendix A when extending to MLA (Ji et al., 2025; Liu et al., 2024). For now, assume the single-head setting. MLA tells us to calculate queries and keys as follows,

$$q = \text{Concat}(q_{\text{nope}}, q_{\text{rope}}) \quad (13a)$$

$$k = \text{Concat}(k_{\text{nope}}, k_{\text{rope}}) \quad (13b)$$

$$q_{\text{nope}} = \mathbf{W}_{\text{uq}} \mathbf{W}_{\text{dq}} x \quad (13c)$$

$$q_{\text{rope}} = R_x(\mathbf{W}_{\text{qr}} \mathbf{W}_{\text{dq}} x) \quad (13d)$$

$$c_{\text{kv}} = \mathbf{W}_{\text{dkv}} y \quad (13e)$$

$$k_{\text{nope}} = \mathbf{W}_{\text{uk}} c_{\text{kv}} = \mathbf{W}_{\text{uk}} \mathbf{W}_{\text{dkv}} y \quad (13f)$$

$$k_{\text{rope}} = R_y(\mathbf{W}_{\text{kr}} y). \quad (13g)$$

Here, x and y are two token embeddings. The ‘down’ matrices, \mathbf{W}_{dq} and \mathbf{W}_{dkv} , project queries and keys/values respectively down to a lower dimensional latent space. This enables efficient caching of c_{kv} . The ‘up’ matrices \mathbf{W}_{uq} , \mathbf{W}_{uk} project these latents up to a higher dimensional space for attention calculations on each head. The \mathbf{W}_{qr} and \mathbf{W}_{kr} matrices are used to produce decoupled queries and keys for RoPE (Su et al., 2021) embeddings, with the position embedding applied via the rotation matrices R_x and R_y .

The change in logits is given by,

$$d_{\text{head}}^{1/2} |\Delta \ell| = |(q + \Delta q)^T (k + \Delta k) - q^T k| = |(\Delta q)^T k + q^T \Delta k + (\Delta q)^T \Delta k|, \quad (14)$$

and we can bound the change,

$$\begin{aligned} d_{\text{head}}^{1/2} |\Delta \ell| &\leq |(\Delta q)^T k| + |q^T \Delta k| + \|\Delta q\| \|\Delta k\| \\ &\leq |(\Delta q_{\text{nope}})^T k_{\text{nope}}| + |(\Delta q_{\text{rope}})^T k_{\text{rope}}| + |q_{\text{nope}}^T \Delta k_{\text{nope}}| + |q_{\text{rope}}^T \Delta k_{\text{rope}}| + \|\Delta q\| \|\Delta k\|. \end{aligned} \quad (15)$$

Expanding further, for the queries, we have,

$$\begin{aligned} \Delta q_{\text{nope}} &= (\mathbf{W}_{\text{uq}} + \Delta \mathbf{W}_{\text{uq}})(\mathbf{W}_{\text{dq}} + \Delta \mathbf{W}_{\text{dq}})x - \mathbf{W}_{\text{uq}} \mathbf{W}_{\text{dq}} x \\ &= \Delta \mathbf{W}_{\text{uq}} \mathbf{W}_{\text{dq}} x + \mathbf{W}_{\text{uq}} \Delta \mathbf{W}_{\text{dq}} x + \Delta \mathbf{W}_{\text{uq}} \Delta \mathbf{W}_{\text{dq}} x, \end{aligned} \quad (16a)$$

$$\begin{aligned} \Delta q_{\text{rope}} &= R_x [(\mathbf{W}_{\text{qr}} + \Delta \mathbf{W}_{\text{qr}})(\mathbf{W}_{\text{dq}} + \Delta \mathbf{W}_{\text{dq}})x - \mathbf{W}_{\text{qr}} \mathbf{W}_{\text{dq}} x] \\ &= R_x [\Delta \mathbf{W}_{\text{qr}} \mathbf{W}_{\text{dq}} x + \mathbf{W}_{\text{qr}} \Delta \mathbf{W}_{\text{dq}} x + \Delta \mathbf{W}_{\text{qr}} \Delta \mathbf{W}_{\text{dq}} x], \end{aligned} \quad (16b)$$

and for the keys,

$$\begin{aligned} \Delta k_{\text{nope}} &= (\mathbf{W}_{\text{uk}} + \Delta \mathbf{W}_{\text{uk}})(\mathbf{W}_{\text{dkv}} + \Delta \mathbf{W}_{\text{dkv}})y - \mathbf{W}_{\text{uk}} \mathbf{W}_{\text{dkv}} y \\ &= \Delta \mathbf{W}_{\text{uk}} \mathbf{W}_{\text{dkv}} y + \mathbf{W}_{\text{uk}} \Delta \mathbf{W}_{\text{dkv}} y + \Delta \mathbf{W}_{\text{uk}} \Delta \mathbf{W}_{\text{dkv}} y \end{aligned} \quad (17a)$$

$$\Delta k_{\text{rope}} = R_y [(\mathbf{W}_{\text{kr}} + \Delta \mathbf{W}_{\text{kr}})y - \mathbf{W}_{\text{kr}} y] = R_y \Delta \mathbf{W}_{\text{kr}} y. \quad (17b)$$

We now use these expressions, and the expressions for q_{nope} , k_{nope} , q_{rope} , k_{rope} , to bound each of the terms in Eq. (15). We will make some assumptions (similar to Lemma 1,

- the inputs x and y are unit normed;
- we use the Frobenius norm;
- conditioned gradients are bounded by a constant, i.e. $\Delta \mathbf{W}_x = -\eta_x \mathbf{G}_x$ where $\|\mathbf{G}_x\| \leq D$ (valid for Muon and Adam);
- the weight norms are lower bounded by a constant c .

We consider the first order terms. We have,

$$\begin{aligned} |\Delta q_{\text{nope}}^T k_{\text{nope}}| &= |(\Delta \mathbf{W}_{\text{uq}} \mathbf{W}_{\text{dq}} x + \mathbf{W}_{\text{uq}} \Delta \mathbf{W}_{\text{dq}} x + \Delta \mathbf{W}_{\text{uq}} \Delta \mathbf{W}_{\text{dq}} x)^T k_{\text{nope}}| \\ &\leq \|\Delta \mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|k_{\text{nope}}\| + \|\mathbf{W}_{\text{uq}}\| \|\Delta \mathbf{W}_{\text{dq}}\| \|k_{\text{nope}}\| + \|\Delta \mathbf{W}_{\text{uq}}\| \|\Delta \mathbf{W}_{\text{dq}}\| \|k_{\text{nope}}\| \\ &\leq \eta_{\text{uq}} D \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\| + \eta_{\text{dq}} D \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\| + O(\eta_{\text{uq}} \eta_{\text{dq}} \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|) \end{aligned} \quad (18a)$$

$$\begin{aligned} |(\Delta q_{\text{rope}})^T k_{\text{rope}}| &= |(R_x [\Delta \mathbf{W}_{\text{qr}} \mathbf{W}_{\text{dq}} x + \mathbf{W}_{\text{qr}} \Delta \mathbf{W}_{\text{dq}} x + \Delta \mathbf{W}_{\text{qr}} \Delta \mathbf{W}_{\text{dq}} x])^T k_{\text{rope}}| \\ &\leq (\|\Delta \mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{dq}}\| + \|\mathbf{W}_{\text{qr}}\| \|\Delta \mathbf{W}_{\text{dq}}\| + \|\Delta \mathbf{W}_{\text{qr}}\| \|\Delta \mathbf{W}_{\text{dq}}\|) \|k_{\text{rope}}\| \\ &\leq \eta_{\text{qr}} D \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{kr}}\| + \eta_{\text{dq}} D \|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{kr}}\| + O(\eta_{\text{dq}} \eta_{\text{qr}} \|\mathbf{W}_{\text{kr}}\|), \end{aligned} \quad (18b)$$

$$\begin{aligned} |q_{\text{nope}}^T \Delta k_{\text{nope}}| &= |q_{\text{nope}}^T (\Delta \mathbf{W}_{\text{uk}} \mathbf{W}_{\text{dkv}} y + \mathbf{W}_{\text{uk}} \Delta \mathbf{W}_{\text{dkv}} y + \Delta \mathbf{W}_{\text{uk}} \Delta \mathbf{W}_{\text{dkv}} y)| \\ &\leq \|q_{\text{nope}}\| \|\Delta \mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\| + \|q_{\text{nope}}\| \|\mathbf{W}_{\text{uk}}\| \|\Delta \mathbf{W}_{\text{dkv}}\| + \|q_{\text{nope}}\| \|\Delta \mathbf{W}_{\text{uk}}\| \|\Delta \mathbf{W}_{\text{dkv}}\| \\ &\leq \eta_{\text{uk}} D \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{dkv}}\| + \eta_{\text{dkv}} D \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\| + O(\eta_{\text{uk}} \eta_{\text{dkv}} \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\|), \end{aligned} \quad (18c)$$

$$|q_{\text{rope}}^T \Delta k_{\text{rope}}| = |q_{\text{rope}}^T (R_y \Delta \mathbf{W}_{\text{kr}} y)| \leq \|q_{\text{rope}}\| \|\Delta \mathbf{W}_{\text{kr}}\| \leq \eta_{\text{kr}} D \|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{dq}}\|. \quad (18d)$$

We used the fact that for rotation matrices R , $\|R\mathbf{W}\| = \|R\mathbf{W}\| = \|\mathbf{W}\|$.

Ultimately, Eqs. (18) suggest to set the learning rates for each attention weight parameter as,

$$\eta_{\text{uq}} = \tau (\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|)^{-1} \quad (19a)$$

$$\eta_{\text{dq}} = \tau \min \{ (\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|)^{-1}, (\|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{kr}}\|)^{-1} \} \quad (19b)$$

$$\eta_{\text{qr}} = \tau (\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{kr}}\|)^{-1} \quad (19c)$$

$$\eta_{\text{uk}} = \tau (\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{dkv}}\|)^{-1} \quad (19d)$$

$$\eta_{\text{dkv}} = \tau (\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\|)^{-1} \quad (19e)$$

$$\eta_{\text{kr}} = \tau (\|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{dq}}\|)^{-1}. \quad (19f)$$

We then substitute these learning rates into Eqs. (18), to see that the bounds are given by,

$$\begin{aligned} |(\Delta q_{\text{nope}})^T k_{\text{nope}}| &\leq \tau D + \tau D + O\left(\frac{\tau^2 \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|}{\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\| \cdot \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|}\right) \\ &= 2\tau D + O\left(\frac{\tau^2}{\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|}\right), \end{aligned} \quad (20a)$$

$$\begin{aligned} |(\Delta q_{\text{rope}})^T k_{\text{rope}}| &\leq \tau D + \tau D + O\left(\frac{\tau^2 \|\mathbf{W}_{\text{kr}}\|}{\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{kr}}\| \cdot \|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{kr}}\|}\right) \\ &= 2\tau D + O\left(\frac{\tau^2}{\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{kr}}\|}\right), \end{aligned} \quad (20b)$$

$$\begin{aligned} |q_{\text{nope}}^T \Delta k_{\text{nope}}| &\leq \tau D + \tau D + O\left(\frac{\tau^2 \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\|}{\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{dkv}}\| \cdot \|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\|}\right) \\ &= 2\tau D + O\left(\frac{\tau^2}{\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|}\right), \end{aligned} \quad (20c)$$

$$|q_{\text{rope}}^T \Delta k_{\text{rope}}| \leq \tau D. \quad (20d)$$

It is reasonable to assume in practice that the weights are not arbitrarily small (i.e. their norm is lower bounded), and thus that these terms are bounded by a constant.

The only remaining term to bound in Eq. (15) is the quadratic term, $\|\Delta q\| \|\Delta k\|$. We can show that this is bounded by showing that the individual parts are bounded,

$$\begin{aligned} \|\Delta q_{\text{nope}}\| &\leq \eta_{\text{uq}} D \|\mathbf{W}_{\text{dq}}\| + \eta_{\text{dq}} D \|\mathbf{W}_{\text{uq}}\| + \eta_{\text{uq}} \eta_{\text{dq}} D^2 \\ &\leq \frac{2\tau D}{\|\mathbf{W}_{\text{uk}}\| \|\mathbf{W}_{\text{dkv}}\|} + \eta_{\text{uq}} \eta_{\text{dq}} D^2, \end{aligned} \quad (21a)$$

$$\begin{aligned} \|\Delta q_{\text{rope}}\| &\leq \eta_{\text{qr}} D \|\mathbf{W}_{\text{dq}}\| + \eta_{\text{dq}} D \|\mathbf{W}_{\text{qr}}\| + \eta_{\text{qr}} \eta_{\text{dq}} D^2 \\ &\leq \frac{2\tau D}{\|\mathbf{W}_{\text{kr}}\|} + \eta_{\text{qr}} \eta_{\text{dq}} D^2, \end{aligned} \quad (21b)$$

$$\begin{aligned} \|\Delta k_{\text{nope}}\| &\leq \eta_{\text{uk}} D \|\mathbf{W}_{\text{dkv}}\| + \eta_{\text{dkv}} D \|\mathbf{W}_{\text{uk}}\| + \eta_{\text{uk}} \eta_{\text{dkv}} D^2 \\ &\leq \frac{2\tau D}{\|\mathbf{W}_{\text{uq}}\| \|\mathbf{W}_{\text{dq}}\|} + \eta_{\text{uk}} \eta_{\text{dkv}} D^2, \end{aligned} \quad (21c)$$

$$\|\Delta k_{\text{rope}}\| \leq \eta_{\text{kr}} D \leq \frac{\tau D}{\|\mathbf{W}_{\text{qr}}\| \|\mathbf{W}_{\text{dq}}\|}. \quad (21d)$$

To extend further to the multi-head setting, we add head indices to the necessary matrices, \mathbf{W}_{uq}^h , \mathbf{W}_{uk}^h , and \mathbf{W}_{kr}^h , and their corresponding learning rates, η_{uq}^h , η_{uk}^h , η_{qr}^h . The key used for RoPE, k_{rope} is shared between all heads, therefore \mathbf{W}_{kr} surprisingly does not have a head index. The down matrices project to a latent space, so also do not have head indices. Plugging these into Eqs. (19) we have,

$$\eta_{\text{uq}}^h = \tau (\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}^h\| \|\mathbf{W}_{\text{dkv}}\|)^{-1} \quad (22a)$$

$$\eta_{\text{dq}} = \tau \min \left\{ \left(\max_h \|\mathbf{W}_{\text{uq}}^h\| \|\mathbf{W}_{\text{uk}}^h\| \|\mathbf{W}_{\text{dkv}}\| \right)^{-1}, \left(\max_h \|\mathbf{W}_{\text{qr}}^h\| \|\mathbf{W}_{\text{kr}}\| \right)^{-1} \right\} \quad (22b)$$

$$\eta_{\text{qr}}^h = \tau (\|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{kr}}\|)^{-1} \quad (22c)$$

$$\eta_{\text{uk}}^h = \tau (\|\mathbf{W}_{\text{uq}}^h\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{dkv}}\|)^{-1} \quad (22d)$$

$$\eta_{\text{dkv}} = \tau (\max_h \|\mathbf{W}_{\text{uq}}^h\| \|\mathbf{W}_{\text{dq}}\| \|\mathbf{W}_{\text{uk}}^h\|)^{-1} \quad (22e)$$

$$\eta_{\text{kr}} = \tau (\max_h \|\mathbf{W}_{\text{qr}}^h\| \|\mathbf{W}_{\text{dq}}\|)^{-1}. \quad (22f)$$

The use of $\max_h(\cdot)$ comes from the requirement that we want logit changes to be bounded for all heads.