

Why Routers Freeze: Infinite Width Learning Dynamics for Mixture of Experts

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

Mixture-of-Experts (MoE) models scale efficiently through sparse expert activation, but their training dynamics remain poorly understood. We study MoEs in the infinite-width limit with a fixed number of experts, a regime relevant to width-based scaling for hyperparameter transfer. Using Tensor Programs, we derive the training dynamics of soft and Top- K MoEs under SGD and Adam. We show that under the Standard Parameterisation, router logits diverge after one step of feature learning, causing softmax or sigmoid gates to saturate and router gradients to vanish. In contrast, we derive μ P-MoE scaling which restores stability, but soft routing produces symmetric router dynamics: experts remain identically distributed and fail to specialise. For softmax routers, this symmetry also nullifies router gradients. We then show that Top- K routing has a qualitatively different effect: even when logits converge to a symmetric limit, finite-width fluctuations determine the selected experts, making Top- K an implicit symmetry-breaking mechanism. Experiments validate the predicted scaling laws and demonstrate hyperparameter transfer under μ P-MoE.

1. Introduction

Mixture-of-Experts (MoE) architectures partition computation across multiple expert modules. A router maps each input to expert weights, allowing different experts to specialise on different regions of the input space (Jacobs et al., 1991; Eigen et al., 2013). In sparse MoEs, only a subset of experts is activated per input, decoupling parameter count from compute (Shazeer et al., 2017). This mechanism has become central to large-scale models (Team et al., 2024; Jiang et al., 2024; Dai et al., 2024; Yang et al., 2025; Han et al., 2024; Fan et al., 2022).

Despite their empirical success, MoEs remain difficult to train. Common failure modes include router instability, load imbalance, and representation collapse, where experts learn similar functions or become under-utilised (Zoph et al., 2022; Chi et al., 2022; Pham et al., 2024; Fedus et al., 2022; Do et al., 2025). Existing fixes, such as load-balancing losses, router z -losses, router noise, and alternative gates (Fedus et al., 2022; Zoph et al., 2022; Nguyen et al., 2024; Csordás et al., 2023; Wang et al., 2024, 2025), are effective but largely heuristic. This leaves open a basic question: how should MoEs be scaled so that router learning, feature learning, and expert specialisation survive at large width?

For dense neural networks, infinite-width theory has clarified the relationship between parameterisation, stability, and feature learning (Neal, 1996; Mei et al., 2018; Yang, 2020; Bordelon and Pehlevan, 2022; Yang and Hu, 2021; Yang and Littwin, 2023). In particular, Tensor Programs provide a framework for deriving finite-time training limits and have led to parameterisations such as μ P, which enable stable feature learning and hyperparameter

transfer across widths (Yang and Hu, 2021; Yang et al., 2021). We extend this perspective to MoEs.

We study the limit where expert width $n \rightarrow \infty$ while the number of experts m remains fixed. This regime reflects practical settings where width can be scaled more readily than the number of experts, and it isolates expert-capacity scaling from expert-count scaling. Our main finding is a tension specific to MoEs: the Standard Parameterisation (SP) causes router saturation and frozen routers, whereas stable μ P-MoE scaling removes specialisation in soft-routing MoEs by enforcing symmetric router dynamics. Sparse Top- K routing avoids this latter failure by acting as an implicit symmetry-breaking mechanism.

Our contributions are:

- We derive fixed-expert infinite-width limits for soft and Top- K MoEs under SGD and Adam using Tensor Programs.
- We show that SP causes router saturation and vanishing router gradients, while stable μ P-MoE soft routing removes expert specialisation by enforcing symmetric router dynamics.
- We show that Top- K routing acts as a symmetry-breaking mechanism under μ P-MoE, enabling stable feature learning and non-trivial expert differentiation; experiments confirm the predicted scaling laws and hyperparameter transfer.

2. Setup

We consider a single MoE block with input $\boldsymbol{\xi} \in \mathbb{R}^{d_{\text{in}}}$, expert width n , and fixed number of experts m . The MoE is defined as follows.

Definition 1 *We define an MoE block as*

$$\begin{aligned} \mathbf{h}^1 &= \mathbf{W}^1 \boldsymbol{\xi} \in \mathbb{R}^n, & \mathbf{h}^2 &= \mathbf{W}^2 \mathbf{h}^1 + \mathbf{b}^2 \in \mathbb{R}^m, \\ \tilde{\mathbf{h}}^2 &= G(\mathbf{h}^2), & \mathbf{h}^{3,k} &= \mathbf{W}^{3,k} \mathbf{h}^1 \in \mathbb{R}^n, \\ \mathbf{h}^3 &= \sum_{k=1}^m \tilde{h}_k^2 \mathbf{h}^{3,k}, & \mathbf{f} &= \mathbf{W}^4 \mathbf{h}^3 \in \mathbb{R}^{d_{\text{out}}}, \end{aligned} \tag{1}$$

where G chosen to be either softmax or sigmoid.

We call \mathbf{W}^2 and \mathbf{h}^2 the router weights and logits, and $\mathbf{W}^{3,k}$ the k th expert weights. The bias \mathbf{b}^2 represents a generalisation to a learned router bias, injected router noise (Shazeer et al., 2017), or load-balancing bias (Wang et al., 2024).

For Top- K MoEs, the router logits are sparsified before applying the gate:

$$\hat{h}_k^2 = \begin{cases} h_k^2, & k \in \text{TopK}(\mathbf{h}^2), \\ -\infty, & \text{otherwise,} \end{cases} \quad \tilde{\mathbf{h}}^2 = G(\hat{\mathbf{h}}^2).$$

Only experts with non-zero gate values are executed. We use the standard straight-through estimator for Top- K gradients (Shazeer et al., 2017).

We study the Tensor Program scaling regime in which $n \rightarrow \infty$ while m , depth, batch size, data dimension, and training time remain fixed. We use a layerwise *bcd* parameterisation: weights are initialised as

$$W_{ij}^\ell \sim \mathcal{N}(0, n^{-2b_\ell}),$$

and updates are

$$\text{SGD: } \mathbf{W}^\ell \leftarrow \mathbf{W}^\ell - \eta n^{-c_\ell} \nabla_{\mathbf{W}^\ell} \mathcal{L}, \quad \text{Adam: } \mathbf{W}^\ell \leftarrow \mathbf{W}^\ell - \eta n^{-c_\ell} \frac{\widehat{m}_\ell}{\sqrt{\widehat{v}_\ell + n^{-d_\ell} \epsilon}}.$$

We ask whether a parameterisation yields four properties in the infinite-width limit: stable activations, non-vanishing feature updates, faithful Adam updates, and non-trivial expert specialisation. Formal definitions and the full scaling tables are deferred to Appendix B and Tables 1 and 2.

3. Soft-Routing MoEs

We first consider soft-routing MoEs, where every expert is active. The main conclusions are that SP is unstable, while the stable μ P-MoE limit is too symmetric to permit expert specialisation without an additional symmetry-breaking mechanism (such as a router bias).

3.1. SP causes router saturation

The Standard Parameterisation initializes $b_1 = 0$ and $b_\ell = 1/2$ for $\ell > 1$, and uses a global learning-rate scaling $c_\ell = c$ (He et al., 2015; Yang and Hu, 2021). This gives $\Theta(1)$ activations at initialization, but it does not give stable MoE training dynamics.

Proposition 2 (SP router saturation) *Consider a soft MoE with softmax or sigmoid gating under SP. For any learning-rate scaling that gives $\Theta(1)$ feature updates in the first layer after one step, the router logits diverge as $n \rightarrow \infty$. Consequently, after one step the softmax gate converges to a one-hot vector, the sigmoid gate converges to a binary vector, and the router gradient vanishes.*

The mechanism is a correlated update term. Under SP, router weights have entries of scale $\Theta(n^{-1/2})$. After one step, the update to router logits contains $\sum_{j=1}^n (\mathbf{W}_0^2)_{:,j} (\mathbf{h}_1^1 - \mathbf{h}_0^1)_j$. The feature update $(\mathbf{h}_1^1 - \mathbf{h}_0^1)_j$ is correlated with the forward weights $(\mathbf{W}_0^2)_{:,j}$, so the sum scales as $\Theta(n^{1/2})$ rather than $\Theta(1)$. Thus \mathbf{h}_1^2 diverges with width. Softmax then saturates to a one-hot gate and sigmoid to a binary gate, making the gate Jacobian vanish. The router expresses a strong preference for experts, but this preference is frozen after the first step rather than learned.

This provides a partial explanation for why auxiliary losses, especially z -losses, are useful in practice: they are explicitly designed to counteract the empirically observed large router logits (Fedus et al., 2022; Zoph et al., 2022). In Figure 1(a,b), the normalized router entropy decreases toward zero under SP, and the relative router-weight update norm vanishes with width, matching Theorem 2.

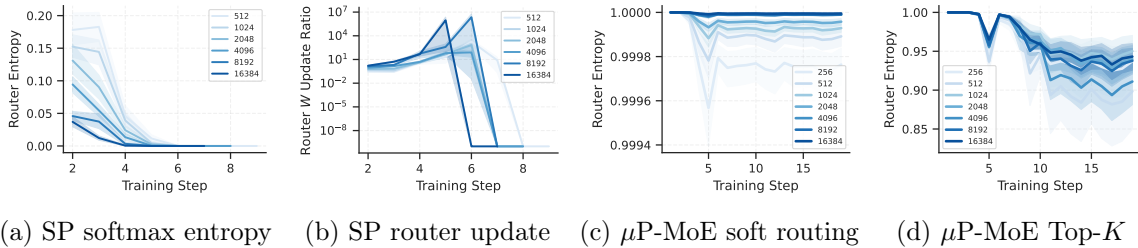


Figure 1: Empirical validation of the fixed-expert scaling predictions. Under SP, router entropy collapses and router updates vanish as width grows. Under μ P-MoE with softmax soft routing, entropy remains maximal, indicating uniform expert allocation. Under μ P-MoE with Top- K routing, entropy decreases without an explicit router bias, showing that Top- K breaks the soft-routing symmetry.

3.2. μ P-MoE (without router bias) gives stability without specialisation

Deriving the stable infinite-width limit gives a μ P-MoE scaling: the router is scaled as an output-like layer, while each expert is scaled as a matrix-like $n \times n$ layer. Full SGD and Adam scalings are given in Tables 1 and 2. This scaling resolves the SP instability, but soft routing has a different failure mode, without explicit symmetry breaking from a router bias/noise term.

Proposition 3 (μ P-MoE soft-routing symmetry) *Consider a soft MoE without router bias \mathbf{b}^2 , parameterised according to μ P-MoE. As $n \rightarrow \infty$, the router scores are identical across experts at initialization: $\tilde{\mathbf{h}}_0^2 \rightarrow m^{-1}\mathbf{1}$ for softmax and $\tilde{\mathbf{h}}_0^2 \rightarrow 2^{-1}\mathbf{1}$ for sigmoid. For every finite training step, the router remains symmetric, $\tilde{\mathbf{h}}_t^2 = C_t\mathbf{1}$, and all experts remain identically distributed. With softmax gating, the router gradient also vanishes.*

The proof is in Appendix D.2. The key point is that router feature learning requires router weights of scale $\Theta(n^{-1})$. In the fixed-expert limit, this makes the router logits deterministic and identical across experts. Since expert weights are also initialised from the same distribution, the downstream gradient into every router logit converges to the same deterministic limit. Router updates therefore remain synchronous, and the experts receive identical gradients. Softmax routing is even more restrictive: when all logits are equal, the softmax Jacobian cancels the symmetric gradient signal, so the router gradient vanishes.

Thus, μ P-MoE solves stability but not specialisation for soft-routing MoEs. Sigmoid routing permits non-zero router updates, but these updates are identical across experts; softmax routing removes even this router feature learning. A router bias or injected router noise avoids this degeneracy by breaking symmetry at initialisation. Empirically, Figure 1(c,d) shows that under μ P-MoE without router bias, softmax soft-routing entropy remains maximal and the router update norm follows the predicted vanishing width scaling.

4. Sparse MoEs

We now consider Top- K sparse routing, the standard mechanism used to reduce MoE compute cost (Shazeer et al., 2017; Jiang et al., 2024). Top- K routing inherits the SP instability:

sparsifying the gradient changes which experts are updated, but not the $\Theta(n^{1/2})$ magnitude of the correlated router-logit update. Hence the router saturation mechanism of Theorem 2 still applies under SP.

The $\mu\text{P-MoE}$ case is qualitatively different. In soft routing and without a router bias, symmetric logits at initialisation lead to symmetric gradients which persisted throughout training. In Top- K routing, however, the selected expert set depends on the ordering of logits. Even at initialisation when logits converge to the same deterministic limit, their finite-width fluctuations can determine the Top- K ordering.

Proposition 4 (Top- K symmetry breaking) *Consider the router logits $\mathbf{h}^2(n) \in \mathbb{R}^m$ under $\mu\text{P-MoE}$. At initialisation $\mathbf{h}^2(n) \rightarrow C\mathbf{1}$ almost surely and, for finite n , $\mathbf{h}^2(n) = C\mathbf{1} + \sigma_n\epsilon(n)$ with $\sigma_n \rightarrow 0$, where $\epsilon(n)$ converges in distribution to a non-degenerate Gaussian vector ϵ . Then*

$$\text{TopK}(\mathbf{h}^2(n)) \Rightarrow \text{TopK}(\epsilon).$$

Thus the Top- K index set remains random.

The proof is in Appendix E. This result shows that Top- K is not merely a computational device, but also a symmetry-breaking mechanism. Under $\mu\text{P-MoE}$, router logits can converge to a symmetric deterministic limit, but the Top- K operator acts on the rescaled finite-width fluctuations. Different experts are therefore selected, yielding expert-dependent gradients and preventing the synchronous soft-routing dynamics of Theorem 3.

We derive the infinite width limit for Top- K MoEs in Appendix F and show that $\mu\text{P-MoEs}$ satisfy stability, feature learning and faithfulness (for Adam). As the router does not remain symmetric, like in the soft case, it also allows for expert differentiation.

We also test the practical implication of the stable scaling: hyperparameter transfer. Following the usual μP methodology (Yang et al., 2021), we tune hyperparameters on the smallest model and transfer them across widths. For a Mixtral-style Top- K MoE trained with Adam, $\mu\text{P-MoE}$ transfers learning rates across widths, whereas SP does not. These results, together with additional coordinate checks and experiments with router bias/noise, are reported in Appendices H and H.4.

5. Discussion

We analysed MoE training in the fixed-expert infinite-width limit and identified two distinct degeneracies. Under SP, router logits diverge after one step of feature learning, saturating the gate and freezing router updates. Under stable $\mu\text{P-MoE}$ scaling, soft routing avoids this instability but becomes symmetric across experts: experts remain identically distributed and fail to specialise. For softmax routing, the same symmetry also nullifies the router gradient.

Top- K routing changes the picture. Even when router logits converge to a symmetric limit, the Top- K operator can act on finite-width fluctuations and produce non-deterministic expert selections. This makes Top- K an implicit symmetry-breaking mechanism, explaining a benefit of sparse routing beyond computational efficiency. Overall, stable large-width MoE training requires more than correct feature-learning scaling: it also requires controlled symmetry breaking, through Top- K , router bias, router noise, or related mechanisms.

References

- Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.
- Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. *Advances in Neural Information Processing Systems*, 35, 2022.
- Blake Bordelon, Hamza Chaudhry, and Cengiz Pehlevan. Infinite limits of multi-head transformer dynamics. *Advances in Neural Information Processing Systems*, 37, 2024a.
- Blake Bordelon, Lorenzo Noci, Mufan Bill Li, Boris Hanin, and Cengiz Pehlevan. Depthwise hyperparameter transfer in residual networks: Dynamics and scaling limit. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, et al. On the representation collapse of sparse mixture of experts. *Advances in Neural Information Processing Systems*, 35, 2022.
- Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. Approximating two-layer feedforward networks for efficient transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023.
- Damai Dai, Chengqi Deng, Chenggang Zhao, Rx Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024.
- Nolan Dey, Bin Claire Zhang, Lorenzo Noci, Mufan Li, Blake Bordelon, Shane Bergsma, Cengiz Pehlevan, Boris Hanin, and Joel Hestness. Don't be lazy: Completep enables compute-efficient deep transformers. *arXiv preprint arXiv:2505.01618*, 2025.
- Giang Do, Hung Le, and Truyen Tran. Simsmoe: Toward efficient training mixture of experts via solving representational collapse. In *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025.
- David Eigen, Marc'Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013.
- Zhiwen Fan, Rishov Sarkar, Ziyu Jiang, Tianlong Chen, Kai Zou, Yu Cheng, Cong Hao, Zhangyang Wang, et al. M³vit: Mixture-of-experts vision transformer for efficient multi-task learning with model-accelerator co-design. *Advances in Neural Information Processing Systems*, 35, 2022.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120), 2022.
- Moritz Haas, Jin Xu, Volkan Cevher, and Leena Chennuru Vankadara. $\mu\mathbf{P}^2$: Effective sharpness aware minimization requires layerwise perturbation scaling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

- Moritz Haas, Sebastian Bordt, Ulrike von Luxburg, and Leena Chennuru Vankadara. On the surprising effectiveness of large learning rates under standard width scaling. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Xing Han, Huy Nguyen, Carl Harris, Nhat Ho, and Suchi Saria. Fusemoe: Mixture-of-experts transformers for fleximodal fusion. *Advances in Neural Information Processing Systems*, 37, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, 2015.
- Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1), 1991.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Tianze Jiang, Blake Bordelon, Cengiz Pehlevan, and Boris Hanin. Hyperparameter transfer with mixture-of-expert layers. *arXiv preprint arXiv:2601.20205*, 2026.
- Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019.
- Jan Małaśnicki, Kamil Ciebiera, Mateusz Boruń, Maciej Pióro, Jan Ludziejewski, Maciej Stefaniak, Michał Krutul, Sebastian Jaszczur, Marek Cygan, Kamil Adamczewski, et al. μ -parametrization for mixture of experts. *arXiv preprint arXiv:2508.09752*, 2025.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33), 2018.
- Radford M Neal. Bayesian learning for neural networks, 1996.
- Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. Sigmoid gating is more sample efficient than softmax gating in mixture of experts. *Advances in Neural Information Processing Systems*, 37, 2024.
- Quang Pham, Giang Do, Huy Nguyen, TrungTin Nguyen, Chenghao Liu, Mina Sartipi, Binh T Nguyen, Savitha Ramasamy, Xiaoli Li, Steven Hoi, et al. Competesmoe—effective training of sparse mixture of experts via competition. *arXiv preprint arXiv:2402.02526*, 2024.

- Joan Puigcerver, Carlos Riquelme Ruiz, Basil Mustafa, and Neil Houlsby. From sparse to soft mixtures of experts. In *The Twelfth International Conference on Learning Representations*, 2024.
- Noam Shazeer, *Azalia Mirhoseini, *Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multi-modal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Leena Chennuru Vankadara, Jin Xu, Moritz Haas, and Volkan Cevher. On feature learning in structured state space models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free load balancing strategy for mixture-of-experts. *arXiv preprint arXiv:2408.15664*, 2024.
- Ziteng Wang, Jun Zhu, and Jianfei Chen. Remoe: Fully differentiable mixture-of-experts with reLU routing. In *The Thirteenth International Conference on Learning Representations*, 2025.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Greg Yang. Wide feedforward or recurrent neural networks of any architecture are gaussian processes. *Advances in Neural Information Processing Systems*, 32, 2019.
- Greg Yang. Tensor programs iii: Neural matrix laws. *arXiv preprint arXiv:2009.10685*, 2020.
- Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*. PMLR, 2021.
- Greg Yang and Etai Littwin. Tensor programs ivb: Adaptive optimization in the infinite-width limit. *arXiv preprint arXiv:2308.01814*, 2023.
- Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: tuning large neural networks via zero-shot hyperparameter transfer. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, 2021.
- Greg Yang, Dingli Yu, Chen Zhu, and Soufiane Hayou. Tensor programs VI: Feature learning in infinite depth neural networks. In *The Twelfth International Conference on Learning Representations*, 2024.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.

Appendix Roadmap: The Logic of Router Dynamics

This appendix provides the rigorous derivations for the infinite-width limits of Mixture-of-Experts (MoE) architectures. To navigate the theoretical contributions, we structure the analysis into a narrative arc: establishing the Tensor Program machinery, demonstrating the inevitable failure of standard soft routing, and deriving the conditional limit for Top- K routing.

- **Full Related Work**
- **Desirable Asymptotics**
Desirable behaviour as width increases.
- **Tensor Program for Soft MoE Architectures (Appendix C)**
Develops the tensor-program representation of the forward/backward passes and training dynamics for soft-routed MoEs.
 - **Tensor Program for Soft MoE with SGD (Appendix C.1)**
Specifies the $\text{NETSOR}\top^+$ program for SGD training, including all primitives and scaling conventions.
 - **Corresponding infinite-width limit for Soft MoE with SGD (Appendix C.2)**
Derives the deterministic and distribution limit objects (Master Theorem recursion) induced by the SGD tensor program.
 - **Tensor Program for Soft MoE with Adam (Appendix C.3)**
Extends the program to Adam using $\text{NE}\otimes\text{ORT}$, introducing the OuterNonlin structure needed for adaptive updates.
 - **Corresponding infinite-width limit for Soft MoE with Adam (Appendix C.4)**
Computes the Adam infinite-width dynamics, including the additional limit rule for OuterNonlin and faithfulness constraints.
- **Deriving Results for Soft MoE (Appendix D)**
Uses the above limits to prove the main qualitative claims about soft routing: collapse under SP and lack of expert specialisation under μP -style stability.
 - **Router Collapse with Standard Parameterisation (Appendix D.1)**
Shows that under SP scaling the router logits diverge, causing gating saturation and vanishing router gradients.
 - **Derivation of μP -MoE and Lack of Expert Specialiaation (Appendix D.2)**
Proves that under stable μP -style scaling the router updates become identical across experts, preventing specialisation in the limit.
- **Top- K Mask Depends on Finite-Width Noise (Appendix E)**
Establishes that when logits collapse to a common value (in the first forward pass under μP -MoE), the Top- K selection remains random in the limit because it is determined by vanishing-but-ordering noise.

- **Tensor Program for Top- K MoE (Appendix F)**
Develops the tensor-program representation of the forward/backward passes and training dynamics for Top- K routed MoEs.
 - **Tensor Program for Top- K MoE (Appendix F.1)**
Defines the Top- K training program.
 - **Corresponding infinite-width limit for Top- K MoE (Appendix F.2)**
Derives the resulting infinite-width limit.
- **Additional Experiments (Appendix H)**
Collects supplementary empirical results that validate the theoretical predictions across parametrisations.

Appendix A. Full Related Work

MoEs were originally introduced to encourage specialisation by decomposing the parameter space into modular experts (Jacobs et al., 1991). Lately, their application is driven by the ability to decouple parameter count from compute cost via sparse routing (Shazeer et al., 2017; Jiang et al., 2024). Despite their efficiency, MoEs are notoriously difficult to train compared to their dense counterparts. Practitioners frequently encounter exploding router logits (Zoph et al., 2022; Puigcerver et al., 2024) and representation collapse — where experts become redundant or underutilised (Chi et al., 2022; Pham et al., 2024). Mitigation of these issues has largely come from empirical heuristics, including auxiliary losses such as the router z -loss (Zoph et al., 2022) and load-balancing regularisers (Fedus et al., 2022), and alternative gating mechanisms, such as sigmoid and ReLU routing (Nguyen et al., 2024; Csordás et al., 2023; Wang et al., 2025). While effective in practice, these solutions are often ad-hoc, and the underlying theoretical mechanisms driving these instabilities remain poorly understood.

In contrast to the heuristic fixes used in MoEs, the training dynamics of dense neural networks have been rigorously characterised through infinite width limits. Early work focused on the “lazy” regime (e.g. NNNGP, NTK), where weights stay close to initialization and feature learning is suppressed (Jacot et al., 2018; Lee et al., 2019; Neal, 1996). More recently, the Tensor Programs framework (Yang, 2020; Yang and Hu, 2021) and related mean-field approaches (Mei et al., 2018; Bordelon and Pehlevan, 2022) have characterized the “feature learning” (maximal update) limit. This theory led to the μ P, which prescribes layer-wise scaling rules to ensure stable feature learning and enables zero-shot hyperparameter transfer across model scales (Yang and Hu, 2021; Yang et al., 2021). These limits have been extended to various architectures (Bordelon et al., 2024a; Haas et al., 2024; Vankadara et al., 2024), and for scaling other quantities such as depth (Dey et al., 2025; Yang et al., 2024). An important empirical consequence of such analysis is the ability for zero shot hyperparameter transfer across width and depth (Yang et al., 2021; Bordelon et al., 2024b).

While independent and concurrent empirical work has attempted to apply μ P directly to MoEs (Małaśnicki et al., 2025), they do not derive the limiting dynamics and thus do not characterise the breakdown of expert specialisation without the router bias, or that Top- K acts as a symmetry breaker. Concurrent work by Jiang et al. (2026) studies hyperparameter transfer for transformer MoEs and proposes scaling rules for transferring learning rates and initialisation scales as width, depth, expert hidden size, and the number of experts are varied. Our work is complementary. First, the two limits are not interchangeable: Jiang et al. (2026) take a joint proportional limit under *gradient flow dynamics* in which width, depth, expert size, and expert count diverge together, whereas we study a width-only Tensor Program limits for both soft and Top- K , sigmoid and softmax routing under SGD and Adam with a fixed number of experts. Since iterated limits do not commute in general, their joint limit does not recover the fixed-expert width limit we study. Second, the two limits are qualitatively different in their predictions: in the proportional limit of Jiang et al. (2026), experts become asymptotically independent, whereas in our width limit experts remain strongly correlated. This fixed-expert regime lets us isolate router-specific degeneracies: under SP the router logits saturate and router gradients vanish, while under stable μ P-heuristic scaling soft routing remains symmetric across experts unless an explicit symmetry-breaking mechanism is present.

Appendix B. Desirable Asymptotics

Building on the asymptotic framework of Yang and Hu (2021); Yang and Littwin (2023), which formalises stability, maximal feature learning, and faithfulness, we introduce a fourth desideratum specific to MoE architectures: expert specialisation. We aim to identify the class of *bcd*-parameterizations under which network training dynamics satisfy these asymptotic properties. We denote the value of a quantity at training step t using a subscript, e.g. \mathbf{h}_t^ℓ .

We first require that activations do not diverge with increasing width.

Definition 5 (Stability) *A network is stable if, for every layer ℓ , the activations satisfy $\mathbf{h}^\ell = \mathcal{O}(1)$ at initialisation and throughout training.*

Next, to ensure feature learning is taking place, we require that the change in activations $\mathbf{h}_t - \mathbf{h}_{t-1}$ at each time step do not vanish with increasing width (Yang and Hu, 2021). Activation updates can be decomposed into two terms

$$\mathbf{h}_t^\ell - \mathbf{h}_{t-1}^\ell = \underbrace{\mathbf{W}_{t-1}^\ell (\mathbf{h}_t^{\ell-1} - \mathbf{h}_{t-1}^{\ell-1})}_{\text{Propagating Updates}} + \underbrace{(\mathbf{W}_t^\ell - \mathbf{W}_{t-1}^\ell) \mathbf{h}_t^{\ell-1}}_{\text{Effective Updates}},$$

where the propagating updates are contributions from the update to the previous layer features, and effective updates are the contribution due to the current weight updates. *Maximal* feature learning requires that both the above terms do not vanish with increasing width. For conciseness we refer to this as simply feature learning throughout.

Definition 6 (Maximal Feature Learning) *A layer ℓ exhibits maximal feature learning if both effective and propagating updates after one step of training scale $\Omega(1)$ (Yang and Hu, 2021).*

For Adam updates, we require that the gradients have a non-trivial $\Theta(1)$ effect on the updates, a condition referred to as *faithfulness* (Yang and Littwin, 2023).

Definition 7 (Adam Faithfulness) *We say a *bcd*-parametrisation is faithful if the inputs to the Adam update function have scale $\Theta(1)$ for every layer's parameter updates.*

Finally, we introduce a desideratum specific to MoEs. Expert specialisation is a central tenet of MoE architectures; accordingly, we require that such specialisation does not cease at infinite width. A *necessary* condition for expert specialisation is that, for a given input, the router exhibits a preference for certain experts over others. We use this to diagnose regimes where the experts cannot specialise.

Definition 8 (Expert Specialisation) *We say that a *bcd*-parametrisation exhibits expert specialisation only if $\mathbf{h}_t^2 \neq C_t \mathbf{1}$ for any constants C_t for all steps $t > 0$.*

Appendix C. Tensor Program for Soft MoE

Here, we derive the Tensor program and the corresponding infinite width limit for soft MoEs (Theorem 1) for both softmax and sigmoid gating. We start by deriving the TP for SGD (Appendix C.1) and the infinite width limit (Appendix C.2). Then, we derive the TP for Adam (Appendix C.3) and the infinite width limit (Appendix C.4).

C.1. Tensor Program for Soft MoE with SGD

Matrices, Vectors and Scalars for the program: As we consider the limit of $n \rightarrow \infty$, we refer to quantities that have both dimensions scaled ($\in \mathbb{R}^{n \times n}$) as ‘matrices’ in the program. Similarly, if only one of the dimensions is taken to infinity, with the other being held constant (e.g. $\in \mathbb{R}^{\bullet \times n}$ for some constant dimension \bullet), then we refer to the quantity as a ‘vector’ in the program. Finally, if both dimensions of a quantity are held constant, we refer to it as a ‘scalar’ in the program. Objects like $\mathbf{W}^2 \in \mathbb{R}^{m \times n}$, are treated as multiple (m) vector objects.

Notation: For all vectors in the program \mathbf{h} , we write $\delta \mathbf{h}_t := \theta_{\delta \mathbf{h}_t}^{-1}(\mathbf{h}_t - \mathbf{h}_{t-1})$ as the normalised change in the features, where $\theta_{\delta \mathbf{h}_t}$ is the scale of $(\mathbf{h}_t - \mathbf{h}_{t-1})$. Thus, $\delta \mathbf{h}_t$ scales $\Theta(1)$. We denote scalars and vectors at training step t as $\mathbf{h}_t = \mathbf{h}_0 + (\mathbf{h}_1 - \mathbf{h}_0) + \dots + (\mathbf{h}_t - \mathbf{h}_{t-1})$. We define $\boldsymbol{\chi}_t := \frac{\partial \mathcal{L}_t}{\partial \mathbf{f}_t} \in \mathbb{R}^{d_{\text{out}} \times 1}$ which we assume is always $\Theta(1)$. This holds for common loss functions (such as mean squared error and cross entropy) as long as the function is stable (does not blow up with width). For a gradient term $\frac{\partial \mathbf{f}}{\partial \mathbf{h}}$ with scale $\theta_{\partial \mathbf{h}}$, we define the normalised gradient $d\mathbf{h} := \theta_{\partial \mathbf{h}}^{-1} \frac{\partial \mathbf{f}}{\partial \mathbf{h}}$, such that $d\mathbf{h}$ scales $\Theta(1)$. We also denote $\theta'_x := n\theta_x$. We reserve subscript for indexing training time \mathbf{h}_t , and index a component with brackets, for example $(\mathbf{h}_t)_\alpha$ denotes the α component of \mathbf{h} at step t .

Tensor Program Operations: For the SGD limit, we use NETSORT⁺ (Yang, 2020; Yang and Hu, 2021). The program contains a sequence of \mathbb{R}^n vectors and \mathbb{R} scalars generated from one of the following ways from initial vectors and scalars

- **MatMul:** Given $\mathbf{W} \in \mathbb{R}^{n \times n}$ and $\mathbf{v} \in \mathbb{R}^n$, new vectors are generated as $\mathbf{W}\mathbf{v}$ or $\mathbf{W}^T\mathbf{v}$.
- **Nonlin:** Given $\phi: \mathbb{R}^k \times \mathbb{R}^l \rightarrow \mathbb{R}$, scalars $\theta_1, \dots, \theta_l \in \mathbb{R}$ and vectors $\mathbf{v}^1, \dots, \mathbf{v}^k \in \mathbb{R}^n$, a new vector can be generated as

$$\phi(\mathbf{v}^1, \dots, \mathbf{v}^k; \theta_1, \dots, \theta_l) \in \mathbb{R}^n$$

where ϕ applies coordinate wise on \mathbf{v}_α^1 .

- **Moment:** Given the same as above, a new scalar can be generated

$$\frac{1}{n} \sum_{\alpha=1}^n \phi(\mathbf{v}_\alpha^1, \dots, \mathbf{v}_\alpha^k; \theta_1, \dots, \theta_l) \in \mathbb{R}.$$

The only assumption on ϕ operations are that they are pseudo-Lipschitz (Yang and Hu, 2021, Assumption F.4).

We first make the assumption that the parametrisation values are chosen such that the network is stable throughout training — no quantity diverges with increasing width.

Setting 9 (Stable parametrisation) *We only consider values of $\{b_\ell, c_\ell\}_\ell$ for each layer ℓ such that the network is stable throughout training.*

This implies that no quantity diverges with increasing width, leading to well-defined limits.

We will write the forward pass, backward pass, and SGD updates as tensor programs (Yang and Hu, 2021). This requires all initial matrices ($n \times n$ shape) to be samples from $\mathcal{N}(0, \frac{1}{n})$, and all initial vectors from $\mathcal{N}(0, 1)$. All θ_\bullet parameters are defined as initial scalars.

Network: We define the network as following,

$$\mathbf{h}^1 = \mathbf{W}^1 \boldsymbol{\xi}, \quad \in \mathbb{R}^{n \times 1} \quad (2)$$

$$\mathbf{h}^2 = \mathbf{W}^2 \mathbf{h}^1 + \mathbf{b}^2, \quad \in \mathbb{R}^{m \times 1} \quad (3)$$

$$\tilde{\mathbf{h}}^2 = \mathbf{G}(\mathbf{h}^2), \quad \in [0, 1]^{m \times 1} \quad (4)$$

$$\mathbf{h}^{3,k} = \mathbf{W}^{3,k} \mathbf{h}^1, \quad \in \mathbb{R}^{n \times 1} \quad (5)$$

$$\mathbf{h}^3 = \sum_{k=1}^m \mathbf{h}^{3,k} \tilde{h}_k^2, \quad \in \mathbb{R}^{n \times 1} \quad (6)$$

$$\mathbf{f} = \mathbf{W}^4 \mathbf{h}^3, \quad \in \mathbb{R}^{d_{\text{out}} \times 1} \quad (7)$$

where \mathbf{G} is the gating function (either softmax or sigmoid), we refer to \mathbf{h}^2 as the router, each $\mathbf{h}^{3,k}$ as expert k , and \mathbf{f} as the function output. We follow a simpler variation of the *abc* parametrisation (Yang and Hu, 2021). The initial scalars, vectors and matrices in the program are defined below,

$$\mathbf{W}^1 = n^{-b_1} \overline{\mathbf{W}}^1 = \theta_{W_0^1} \overline{\mathbf{W}}^1, \quad \mathbf{W}^1 \in \mathbb{R}^{n \times d_{\text{in}}} \quad \overline{W}_{ij}^1 \sim \mathcal{N}(0, 1), \quad (8)$$

$$\mathbf{W}^2 = n^{-b_2} \overline{\mathbf{W}}^2 = \theta_{W_0^2} \overline{\mathbf{W}}^2, \quad \mathbf{W}^2 \in \mathbb{R}^{m \times n} \quad \overline{W}_{ij}^2 \sim \mathcal{N}(0, 1), \quad (9)$$

$$\mathbf{W}^{3,k} = n^{\frac{1}{2}-b_3} \overline{\mathbf{W}}^{3,k} = \theta_{W_0^3} \overline{\mathbf{W}}^{3,k}, \quad \mathbf{W}^{3,k} \in \mathbb{R}^{n \times n} \quad \overline{W}_{ij}^{3,k} \sim \mathcal{N}(0, n^{-1}), \quad (10)$$

$$\mathbf{W}^4 = n^{-b_4} \overline{\mathbf{W}}^4 = \theta_{W_0^4} \overline{\mathbf{W}}^4, \quad \mathbf{W}^4 \in \mathbb{R}^{d_{\text{out}} \times n} \quad \overline{W}_{ij}^4 \sim \mathcal{N}(0, 1). \quad (11)$$

Note that we only consider the limit $n \rightarrow \infty$ with other factors held fixed, and thus $\{\mathbf{W}^{3,k}\}_{k=1}^m$ are the only initial matrices in the program ($n \times n$).

Initial Matrices: $\overline{\mathbf{W}}^{3,k}$.

Initial Vectors: $\overline{\mathbf{W}}^1, \overline{\mathbf{W}}^2, \overline{\mathbf{W}}^4$, where $\overline{\mathbf{W}}^1$ is a set of d_{in} initial vectors, $\overline{\mathbf{W}}^2$ is a set of m initial vector and $\overline{\mathbf{W}}^4$ is a set of d_{out} initial vectors.

Initial Scalars: All θ_{\bullet} parameters. The router bias \mathbf{b}^2 is a set of m initial scalars sampled from some density such that $\mathbb{P}(b_i = b_j) = 0$. In practice, this is satisfied by sampling from $\mathcal{N}(0, \sigma_b^2)$ where σ_b^2 is some hyperparameter.

Noise and Learnable Bias: We consider two separate cases here:

1. Case 1: A different noise is added at each time step $\{\mathbf{b}_t^2\}_{t=1}^T$ (Shazeer et al., 2017).
2. Case 2: The term \mathbf{b}_0^2 is learned, either through gradients of the loss or load balancing considerations (Wang et al., 2024).

Note that the updates below are written only for learnable biases from gradient of the loss. For the noise case, a new noise is sampled every forward pass.

Conditioned on router bias/noise: The standard Tensor Programs Master Theorem assumes initial scalars converge almost surely to deterministic limits. A router bias/noise vector sampled from a continuous distribution does not satisfy this assumption unconditionally. Instead, we use a conditional formulation: we first sample the finite collection of router bias/noise scalars and condition on their realised values. Conditional on these values, they are treated as fixed scalar constants in the tensor program. Thus, for any fixed finite time

horizon T , we condition on $\{\mathbf{b}_t^2\}_{t=0}^T$ in the noise case, or on the initialisation of \mathbf{b}_0^2 in the learnable-bias case. The resulting infinite-width limits are deterministic conditional on the realised bias/noise values, but random unconditionally through their dependence on those values. Practically, this corresponds to coupling widths by using the same sampled router bias/noise values across widths.

First Forward Pass:

$$\begin{aligned}
 \mathbf{h}_0^1 &= \theta_{W_0^1} \overline{\mathbf{W}}_0^1 \boldsymbol{\xi}, & (\mathbf{h}_0^1)_\alpha &= \phi \left((\overline{\mathbf{W}}_0^1)_{\alpha 1}, \dots, (\overline{\mathbf{W}}_0^1)_{\alpha d_{\text{in}}}; \xi_1, \dots, \xi_{d_{\text{in}}}, \theta_{W_0^1} \right) = \sum_{j=1}^{d_{\text{in}}} \theta_{W_0^1} (\overline{\mathbf{W}}_0^1)_{\alpha j} \xi_j, & \text{(Nonl)} \\
 \mathbf{h}_0^2 &= \theta'_{W_0^2} \frac{1}{n} \overline{\mathbf{W}}_0^2 \mathbf{h}_0^1 + \mathbf{b}_0^2, & (\mathbf{h}_0^2)_k &= \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}}_0^2)_{k\alpha}, (\mathbf{h}_0^1)_\alpha; \theta'_{W_0^2}, \mathbf{b}_0^2 \right) = \frac{1}{n} \theta'_{W_0^2} \sum_{\alpha=1}^n (\overline{\mathbf{W}}_0^2)_{k\alpha} (\mathbf{h}_0^1)_\alpha + \mathbf{b}_0^2, & \text{(Mome)} \\
 \tilde{\mathbf{h}}_0^2 &= \mathbf{G}(\mathbf{h}_0^2), & & \text{(Moment, See G)} \\
 \hat{\mathbf{h}}_0^{3,k} &= \overline{\mathbf{W}}^{3,k} \mathbf{h}_0^1, & & \text{(MatMul)} \\
 \mathbf{h}_0^{3,k} &= \theta_{W_0^3} \hat{\mathbf{h}}_0^{3,k}, & (\mathbf{h}_0^{3,k})_\alpha &= \phi \left((\hat{\mathbf{h}}_0^{3,k})_\alpha; \theta_{W_0^3} \right) = \theta_{W_0^3} (\hat{\mathbf{h}}_0^{3,k})_\alpha & \text{(Nonl)} \\
 \mathbf{h}_0^3 &= \sum_{k=1}^m \mathbf{h}_0^{3,k} \tilde{h}_{0,k}^2, & (\mathbf{h}_0^3)_\alpha &= \phi \left((\mathbf{h}_0^{3,1})_\alpha, \dots, (\mathbf{h}_0^{3,m})_\alpha; \tilde{h}_{0,1}^2, \dots, \tilde{h}_{0,m}^2 \right) = \sum_{k=1}^m (\mathbf{h}_0^{3,k})_\alpha \tilde{h}_{0,k}^2 & \text{(Nonl)} \\
 \mathbf{f}_0 &= \theta'_{W_0^4} \frac{1}{n} \overline{\mathbf{W}}_0^4 \mathbf{h}_0^3, & (\mathbf{f}_0)_\beta &= \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}}_0^4)_{\beta\alpha}, (\mathbf{h}_0^3)_\alpha; \theta'_{W_0^4} \right) = \frac{1}{n} \theta'_{W_0^4} \sum_{\alpha=1}^n (\overline{\mathbf{W}}_0^4)_{\beta\alpha} (\mathbf{h}_0^3)_\alpha. & \text{(Mome)}
 \end{aligned}$$

Stability at Initialisation: Satisfying stability at initialisation (Theorem 5) fixes the values of b_1, b_2, b_3, b_4 , and the corresponding θ parameters. We find the values that satisfy stability at initialisation and fix them for the rest of the program. \mathbf{h}_0^1 is trivially a Gaussian with scale $\theta_{W_0^1} = n^{-b_1} \|\boldsymbol{\xi}\|$. Hence, we set $b_1 = 0$ to satisfy stability. Using the Master theorem (Yang and Hu, 2021), $\mathbb{E} \left[(Z^{h_0^2})^2 \right] = n^{(1-2b_2)} \mathbb{E} \left[(Z^{h_0^1})^2 \right] = n^{(1-2b_2)} \|\boldsymbol{\xi}\|^2$. Stability is satisfied for $b_2 \geq \frac{1}{2}$. We can use the Master theorem to find the distribution of \mathbf{h}_0^3 as $Z^{h_0^3} = \mathcal{N}(0, n^{(1-2b_3)} \mathbb{E} \left[(Z^{h_0^1})^2 \right])$. We thus set $b_3 = \frac{1}{2}$. Finally, using similar reasoning as for \mathbf{h}_0^2 , setting $b_4 \geq \frac{1}{2}$ ensures that \mathbf{f}_0 is stable. Note that $\tilde{\mathbf{h}}_0^2$ is always stable due to the bounded softmax operation.

Remark 10 (Feature Learning considerations to constrain output-like initialisation)

Stability constrains the output-like layers to initialise with $b_2 \geq 1/2$ and $b_4 \geq 1/2$. For ease of exposition, we use the Tensor Program derived below to constrain these values. Looking at the first feature updates of the router

$$\mathbf{h}_1^2 - \mathbf{h}_0^2 = \theta_{\delta h_1^1} \theta'_{W_0^2} \frac{\overline{\mathbf{W}}_0^2 \delta \mathbf{h}_1^1}{n} - \theta'_{\delta W_1^2} \eta d \mathbf{h}_0^2 \chi_0 \frac{(\mathbf{h}_0^1)^T \mathbf{h}_1^1}{n}.$$

By the Master Theorem, $\frac{\overline{\mathbf{W}}_0^2 \delta \mathbf{h}_1^1}{n}$ which is a Moment operation converges to an expectation (Appendix C.2). If we have feature learning for $\mathbf{h}_1^1 - \mathbf{h}_0^1$, we require that $\theta_{\delta h_1^1} = 1$. Thus, for the first term to scale $\Theta(1)$, we require $\theta'_{W_0^2} = n \cdot n^{-b_2} = 1$, which implies setting $b_2 = 1$. A similar argument holds for setting $b_4 = 1$.

The final values for stability at initialisation (and feature learning for ‘output-like’ layers) are

$$b_1 = 0, \quad b_2 = 1, \quad b_3 = \frac{1}{2}, \quad b_4 = 1, \implies \theta_{W_0^1} = 1, \quad \theta_{W_0^2} = n^{-1}, \quad \theta_{W_0^3} = 1, \quad \theta_{W_0^4} = n^{-1}.$$

Setting 11 (Stability at initialisation) *For stability at initialisation, and feature learning considerations for ‘output-like’ layers (Theorem 10), we choose*

$$b_1 = 0, \quad b_2 = 1, \quad b_3 = \frac{1}{2}, \quad b_4 = 1,$$

which gives

$$\theta_{W_0^1} = 1, \quad \theta_{W_0^2} = n^{-1}, \quad \theta_{W_0^3} = 1, \quad \theta_{W_0^4} = n^{-1}.$$

We set these values in the program for ease of exposition.

Note that under the assumptions made in Theorem 9 and Theorem 11 certain vectors have scale $\Theta(1)$ throughout training.

Remark 12 *Under the assumption of stable parametrisations and the values selected for stability at initialisation, the vectors $\mathbf{h}_t^1, \{\mathbf{h}_t^{3,k}\}_{k=1}^m, \mathbf{h}_t^3$ have scale $\Theta(1)$ throughout training. As the softmax and sigmoid operations are bounded, the softmax output $\tilde{\mathbf{h}}_t^2$ has scale $\Theta(1)$ throughout training.*

First Backward Pass:

$$d\mathbf{h}_0^3 = (\overline{\mathbf{W}_0^4})^T, \quad \in \mathbb{R}^{n \times d_{\text{out}}}. \quad (12)$$

$$n \cdot d\tilde{h}_{0,k}^2 = (\mathbf{h}_0^{3,k})^T d\mathbf{h}_0^3, \quad \in \mathbb{R}^{1 \times d_{\text{out}}}, \quad \forall k \in [m] \quad (13)$$

$$(d\tilde{h}_{0,k}^2)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\mathbf{h}_0^{3,k})_\alpha, (d\mathbf{h}_0^3)_{\alpha\beta}; - \right), \quad (\text{Moment}) \quad (14)$$

$$\phi(\dots) := (\mathbf{h}_0^{3,k})_\alpha (d\mathbf{h}_0^3)_{\alpha\beta}, \quad (15)$$

where we repeat the operation for each expert to give $d\tilde{h}_{0,1}^2, \dots, d\tilde{h}_{0,m}^2$.

$$dh_{0,k}^2 = \sum_{j=1}^m J(\tilde{\mathbf{h}}_0^2)_{jk} d\tilde{h}_{0,j}^2, \quad \in \mathbb{R}^{1 \times d_{\text{out}}}, \quad \forall k \in [m] \quad (16)$$

$$(dh_{0,k}^2)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(-; \{J(\tilde{\mathbf{h}}_0^2)_{jk}\}_{j=1}^m, \{(d\tilde{h}_{0,j}^2)_\beta\}_{j=1}^m \right), \quad (\text{Moment}) \quad (17)$$

$$\phi(\dots) := \sum_{j=1}^m J(\tilde{\mathbf{h}}_0^2)_{jk} \overline{(d\tilde{h}_{0,j}^2)_\beta}, \quad (18)$$

where $J(\tilde{\mathbf{h}}_0^2)_{ij} = \frac{\partial(\tilde{\mathbf{h}}_0^2)_i}{\partial(\tilde{\mathbf{h}}_0^2)_j}$.

Again, we repeat this operation for each expert to give $dh_{0,1}^2, \dots, dh_{0,m}^2$. We also define

$$\mathbf{dh}_0^2 = \begin{pmatrix} dh_{0,1}^2 \\ \vdots \\ dh_{0,m}^2 \end{pmatrix} \in \mathbb{R}^{m \times d_{\text{out}}}. \quad (19)$$

$$\mathbf{dh}_0^{3,k} = \tilde{h}_{0,k}^2 \mathbf{dh}_0^3, \quad \in \mathbb{R}^{n \times d_{\text{out}}}, \quad \forall k \in [m] \quad (20)$$

$$(\mathbf{dh}_0^{3,k})_{\alpha\beta} = \phi\left((\mathbf{dh}_0^3)_{\alpha\beta}; \tilde{h}_{0,k}^2\right), \quad (\text{Nonlin}) \quad (21)$$

$$\phi(\dots) := \tilde{h}_{0,k}^2 \overline{\mathbf{dh}_0^3}, \quad (22)$$

We repeat this for each expert to give $\mathbf{dh}_0^{3,1}, \dots, \mathbf{dh}_0^{3,m}$.

$$\mathbf{dh}_0^1 = \sum_{k=1}^m (\overline{\mathbf{W}_0^2})_{:,k}^T dh_{0,k}^2 + \sum_{k=1}^m (\mathbf{W}_0^{3,k})^T \mathbf{dh}_0^{3,k}, \quad \in \mathbb{R}^{n \times d_{\text{out}}}, \quad (23)$$

$$(\mathbf{dh}_0^1)_{\alpha\beta} = \phi\left(\{(\overline{\mathbf{W}_0^2})_{\alpha k}^T, (\mathbf{W}_0^k)_{\alpha\beta}\}_{k=1}^m; \{(dh_{0,k}^2)_{\beta}\}_{k=1}^m\right), \quad (\text{Nonlin}) \quad (24)$$

$$(\mathbf{V}_0^k)_{:, \beta} := (\mathbf{W}_0^{3,k})^T (\mathbf{dh}_0^{3,k})_{:, \beta}, \quad (\text{MatMul}) \quad (25)$$

$$\phi(\dots) := \sum_{k=1}^m (\overline{\mathbf{W}_0^2})_{:,k}^T dh_{0,k}^2 + \sum_{k=1}^m \mathbf{V}_0^k. \quad (26)$$

Note that the final gradient has two branches. The scale of this term thus will be the maximum of the the scale of the twp branches. There is a $\theta'_{W_0^4} \theta_{W_0^2}$ scale coming from the router and a $\theta_{W_0^4}$ coming from the experts. The scale from the branches is the same here.

First Weight Updates: Recall that we define the gradient of the loss with respect to the function as χ_t . For the vectors $\mathbf{W} \in \{\mathbf{W}^1, \overline{\mathbf{W}^2}, \mathbf{W}^4\}$ we write $\delta \mathbf{W}_t := \theta_{\delta W_t}^{-1} (\mathbf{W}_t - \mathbf{W}_{t-1})$, and thus the weights at step t as $\mathbf{W}_t = \theta_{W_0} \overline{\mathbf{W}_0} + \dots + \theta_{\delta W_t} \delta W_t$.

$$\mathbf{W}_1^4 - \mathbf{W}_0^4 = -n^{-c_4} \eta \chi_0 (\mathbf{h}_0^3)^T \quad (27)$$

$$\theta_{\delta W_1^4} := n^{-c_4}. \quad (28)$$

$$\mathbf{W}_1^{3,k} - \mathbf{W}_0^{3,k} = -n^{-c_3} \theta_{W_0^4} \eta \mathbf{dh}_0^{3,k} \chi_0 (\mathbf{h}_0^1)^T \quad (29)$$

$$\theta_{\delta W_1^3} := n^{-(c_3+1)}. \quad (30)$$

$$\mathbf{W}_1^2 - \mathbf{W}_0^2 = -n^{-c_2} n \theta_{W_0^4} \eta \mathbf{dh}_0^2 \chi_0 (\mathbf{h}_0^1)^T \quad (31)$$

$$\theta_{\delta W_1^2} := n^{-c_2}. \quad (32)$$

$$\mathbf{W}_1^1 - \mathbf{W}_0^1 = -n^{-c_1} n \theta_{W_0^4} \theta_{W_0^2} \eta \mathbf{dh}_0^1 \chi_0 (\boldsymbol{\xi})^T \quad (33)$$

$$\theta_{\delta W_1^1} := n^{-(c_1+1)}. \quad (34)$$

The router bias in the learnable case updates as

$$\mathbf{b}_1^2 - \mathbf{b}_0^2 = -\eta \theta'_{W_0^4} \mathbf{dh}_0^2 \chi_0, \quad (35)$$

where due to stability at initialisation (Theorem 11) $\theta'_{W_0^4} = 1$.

First Feature Updates: The feature updates are affected by the change in the previous layers features (*propagating updates*) and the update in the weights at the current layer (*effective updates*). We write the feature updates for the vectors $\mathbf{h} \in \{\mathbf{h}^1, \{\mathbf{h}^{3,k}\}_{k=1}^m, \mathbf{h}^3\}$ as $\delta\mathbf{h}_t := \theta_{\delta h_t}^{-1}(\mathbf{h}_t - \mathbf{h}_{t-1})$. This gives each features for the vectors at step t as $\mathbf{h}_t = \mathbf{h}_0 + \theta_{\delta h_1} \delta\mathbf{h}_1 + \dots + \theta_{\delta h_t} \delta\mathbf{h}_t$. For the scalars $\mathbf{h} \in \{\mathbf{h}^2, \tilde{\mathbf{h}}^2, \mathbf{f}\}$ we simply write $\delta\mathbf{h}_t := (\mathbf{h}_t - \mathbf{h}_{t-1})$ and $\mathbf{h}_t = \mathbf{h}_0 + \dots + \delta\mathbf{h}_t$.

$$\mathbf{h}_1^1 - \mathbf{h}_0^1 = (\mathbf{W}_1^1 - \mathbf{W}_0^1)\boldsymbol{\xi} \quad (36)$$

$$= -\theta_{\delta W_1^1} \eta d\mathbf{h}_0^1 \boldsymbol{\chi}_0 \boldsymbol{\xi}^T \boldsymbol{\xi} \quad (37)$$

$$(\delta\mathbf{h}_1^1)_\alpha = \phi\left(\{(\mathbf{d}\mathbf{h}_0^1)_{\alpha\beta}\}_{\beta=1}^{d_{\text{out}}}; -\eta, \{(\boldsymbol{\chi}_0 \boldsymbol{\xi}^T \boldsymbol{\xi})_\beta\}_{\beta=1}^{d_{\text{out}}}\right) \quad (\text{Nonlin}) \quad (38)$$

$$\phi(\dots) := -\eta \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_0^1)_{\alpha\beta} (\boldsymbol{\chi}_0 \boldsymbol{\xi}^T \boldsymbol{\xi})_\beta, \quad (39)$$

$$\theta_{\delta h_1^1} := \theta_{\delta W_1^1}, \quad (40)$$

$$\mathbf{h}_1^2 - \mathbf{h}_0^2 = \mathbf{W}_0^2(\mathbf{h}_1^1 - \mathbf{h}_0^1) + (\mathbf{W}_1^2 - \mathbf{W}_0^2)\mathbf{h}_1^1 + (\mathbf{b}_1^2 - \mathbf{b}_0^2) \quad (41)$$

$$= \theta_{\delta h_1^1} \theta'_{W_0^2} \frac{\overline{\mathbf{W}_0^2} \delta\mathbf{h}_1^1}{n} - \theta'_{\delta W_1^2} \eta d\mathbf{h}_0^2 \boldsymbol{\chi}_0 \frac{(\mathbf{h}_0^1)^T \mathbf{h}_1^1}{n} - \eta d\mathbf{h}_0^2 \boldsymbol{\chi}_0, \quad (42)$$

$$\theta_{\delta h_1^2} (\delta\mathbf{h}_1^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi\left(\overline{\mathbf{W}_0^2}{}_{0k\alpha}, (\delta\mathbf{h}_1^1)_\alpha, (\mathbf{h}_0^1)_\alpha, (\mathbf{h}_1^1)_\alpha; \theta'_{W_0^2}, \theta_{\delta h_1^1}, \theta'_{W_0^2}, \theta'_{\delta W_1^2}, \eta, \{(\mathbf{d}\mathbf{h}_0^2)_{k\beta}, (\boldsymbol{\chi}_0)_\beta\}_{\beta=1}^{d_{\text{out}}}\right) \quad (\text{Moment}) \quad (43)$$

$$\phi(\dots) := \theta_{\delta h_1^1} \theta'_{W_0^2} \overline{\mathbf{W}_0^2}{}_{0k\alpha} (\delta\mathbf{h}_1^1)_\alpha - \theta'_{\delta W_1^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_0^2)_{k\beta} (\boldsymbol{\chi}_0)_\beta (\mathbf{h}_0^1)_\alpha (\mathbf{h}_1^1)_\alpha - \eta \theta'_{W_0^2} \sum_{\beta=1}^{d_{\text{out}}} d(\mathbf{h}_0^2)_{k\beta} (\boldsymbol{\chi}_0)_\beta, \quad (44)$$

$$\theta_{\delta h_1^2} := \max(\theta_{\delta h_1^1} \theta'_{W_0^2}, \theta'_{\delta W_1^2}), \quad (45)$$

$$\tilde{\mathbf{h}}_1^2 - \tilde{\mathbf{h}}_0^2 = \mathbf{G}(\mathbf{h}_1^2) - \mathbf{G}(\mathbf{h}_0^2) \quad (46)$$

$$= \mathbf{G}\left(\mathbf{h}_0^2 + \theta_{\delta h_1^2} \delta\mathbf{h}_1^2\right) - \mathbf{G}(\mathbf{h}_0^2) \quad (47)$$

$$\theta_{\delta \tilde{h}_1^2} (\delta\tilde{\mathbf{h}}_1^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi_k\left(-; \{(\mathbf{h}_0^2)_k, (\delta\mathbf{h}_1^2)_k\}_{k=1}^m, \theta_{\delta h_1^2}\right), \quad (\text{Moment}), \quad (48)$$

$$\theta_{\delta \tilde{h}_1^2} := \theta_{\delta h_1^2}. \quad (49)$$

where G is the router gating function (either softmax or sigmoid). Both softmax and sigmoid operations can be represented as a Moment (see Appendix G).

$$\mathbf{h}_1^{3,k} - \mathbf{h}_0^{3,k} = \mathbf{W}_0^{3,k}(\mathbf{h}_1^1 - \mathbf{h}_0^1) + (\mathbf{W}_1^{3,k} - \mathbf{W}_0^{3,k})\mathbf{h}_1^1, \quad (50)$$

$$= \theta_{\delta h_1^1} \mathbf{W}_0^{3,k} \delta \mathbf{h}_1^1 - \theta'_{\delta W_1^3} \eta d\mathbf{h}_0^{3,k} \boldsymbol{\chi}_0 \frac{(\mathbf{h}_0^1)^T (\mathbf{h}_1^1)}{n}, \quad (51)$$

$$(\delta \mathbf{h}_1^{3,k})_\alpha = \phi \left(\Upsilon_\alpha, \{(\mathbf{d}\mathbf{h}_0^{3,k})_{\alpha l}\}_{l=1}^{d_{\text{out}}}; \theta_{\delta h_1^1 / \delta h_1^{3,k}}, \theta'_{\delta W_1^3 / \delta h_1^{3,k}}, \{(\boldsymbol{\chi}_0)_l\}_{l=1}^{d_{\text{out}}}, \eta, c_0, \right), \quad (\text{Nonlin}) \quad (52)$$

$$\Upsilon := \mathbf{W}_0^{3,k} \delta \mathbf{h}_1^1, \quad (\text{MatMul}) \quad (53)$$

$$c_0 := \frac{1}{n} \sum_{\alpha=1}^n \psi \left((\mathbf{h}_0^1)_\alpha, (\mathbf{h}_1^1)_\alpha; - \right), \quad (\text{Moment}) \quad (54)$$

$$\phi(\dots) := \theta_{\delta h_1^1 / \delta h_1^{3,k}} \Upsilon_\alpha - \theta'_{\delta W_1^3 / \delta h_1^{3,k}} \eta c_0 \sum_{l=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_0^{3,k})_{\alpha l} (\boldsymbol{\chi}_0)_l, \quad (55)$$

$$\theta_{\delta h_1^{3,k}} := \max(\theta_{\delta h_1^1}, \theta'_{\delta W_1^3}), \quad (56)$$

$$\theta_{\delta h_1^1 / \delta h_1^{3,k}} := \theta_{\delta h_1^1} / \theta_{\delta h_1^{3,k}}, \quad (57)$$

$$\theta'_{\delta W_1^3 / \delta h_1^{3,k}} := \theta'_{\delta W_1^3} / \theta_{\delta h_1^{3,k}} \quad (58)$$

We repeat the above process for each expert k to give $\{\delta \mathbf{h}_1^{3,k}\}_{k=1}^m$.

$$\mathbf{h}_1^3 - \mathbf{h}_0^3 = \sum_{k=1}^m \left(\tilde{h}_{0,k}^2 (\mathbf{h}_1^{3,k} - \mathbf{h}_0^{3,k}) + (\tilde{h}_{1,k}^2 - \tilde{h}_{0,k}^2) \mathbf{h}_1^{3,k} \right) \quad (59)$$

$$= \sum_{k=1}^m \left(\theta_{\delta h_1^{3,k}} \tilde{h}_{0,k}^2 \delta \mathbf{h}_1^{3,k} + \theta_{\delta \tilde{h}_1^2} \delta \tilde{h}_{1,k}^2 \mathbf{h}_1^{3,k} \right), \quad (60)$$

$$(\delta \mathbf{h}_1^3)_\alpha = \phi \left(\{(\delta \mathbf{h}_1^{3,k})_\alpha\}_{k=1}^m, \{(\mathbf{h}_1^{3,k})_\alpha\}_{k=1}^m; \theta_{\delta h_1^{3,k} / \delta h_1^3}, \theta_{\delta \tilde{h}_1^2 / \delta h_1^3}, \{\tilde{h}_{0,k}^2\}_{k=1}^m, \{\delta \tilde{h}_{1,k}^2\}_{k=1}^m \right), \quad (\text{Nonlin}) \quad (61)$$

$$\phi(\dots) := \sum_{k=1}^m \left(\theta_{\delta h_1^{3,k} / \delta h_1^3} \tilde{h}_{0,k}^2 (\delta \mathbf{h}_1^{3,k})_\alpha + \theta_{\delta \tilde{h}_1^2 / \delta h_1^3} \delta \tilde{h}_{1,k}^2 (\mathbf{h}_1^{3,k})_\alpha \right), \quad (62)$$

$$\theta_{\delta h_1^3} := \max(\theta_{\delta h_1^{3,k}}, \theta_{\delta \tilde{h}_1^2}), \quad (63)$$

$$\theta_{\delta h_1^{3,k} / \delta h_1^3} := \theta_{\delta h_1^{3,k}} / \theta_{\delta h_1^3}, \quad (64)$$

$$\theta_{\delta \tilde{h}_1^2 / \delta h_1^3} := \theta_{\delta \tilde{h}_1^2} / \theta_{\delta h_1^3}. \quad (65)$$

$$\mathbf{f}_1 - \mathbf{f}_0 = \mathbf{W}_0^4(\mathbf{h}_1^3 - \mathbf{h}_0^3) + (\mathbf{W}_1^4 - \mathbf{W}_0^4)\mathbf{h}_1^3 \quad (66)$$

$$= \theta_{\delta h_1^3} \frac{\overline{\mathbf{W}}_0^4 \delta \mathbf{h}_1^3}{n} - \theta'_{\delta W_1^4} \eta \chi_0 \frac{(\mathbf{h}_0^3)^T (\mathbf{h}_1^3)}{n}, \quad (67)$$

$$\theta_{\delta f_1}(\delta \mathbf{f}_1)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}}_{0\beta\alpha}^4, (\delta \mathbf{h}_1^3)_\alpha, (\mathbf{h}_0^3)_\alpha, (\mathbf{h}_1^3)_\alpha; \theta_{\delta h_1^3}, \theta'_{\delta W_1^4}, \eta, (\chi_0)_\beta \right), \quad (\text{Moment}) \quad (68)$$

$$\phi(\dots) := \theta_{\delta h_1^3} (\overline{\mathbf{W}}_0^4)_{\beta\alpha} (\delta \mathbf{h}_1^3)_\alpha - \theta'_{\delta W_1^4} \eta (\chi_0)_\beta (\mathbf{h}_0^3)_\alpha (\mathbf{h}_1^3)_\alpha, \quad (69)$$

$$\theta_{\delta f_1} := \max(\theta_{\delta h_1^3}, \theta'_{\delta W_1^4}). \quad (70)$$

t th Backward Pass: Recall that we write the weights at step t as $\mathbf{W}_t = \theta_{W_0} \overline{\mathbf{W}}_0 + \dots + \theta_{\delta W_t} \delta W_t$.

$$d\mathbf{h}_t^3 := (\overline{\mathbf{W}}_t^4)^T, \quad \in \mathbb{R}^{n \times d_{\text{out}}}. \quad (71)$$

$$n \cdot d\tilde{h}_{t,k}^2 = (\mathbf{h}_t^{3,k})^T d\mathbf{h}_t^3, \quad \in \mathbb{R}^{1 \times d_{\text{out}}}, \quad \forall k \in [m] \quad (72)$$

$$(d\tilde{h}_{t,k}^2)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\mathbf{h}_t^{3,k})_\alpha, (d\mathbf{h}_t^3)_{\alpha\beta}; - \right), \quad (\text{Moment}) \quad (73)$$

$$\phi(\dots) := (\mathbf{h}_t^{3,k})_\alpha (d\mathbf{h}_t^3)_{\alpha\beta}, \quad (74)$$

where we repeat this operation for each expert $d\tilde{h}_{t,1}^2, \dots, d\tilde{h}_{t,m}^2$.

$$dh_{t,k}^2 = \sum_{j=1}^m J(\tilde{\mathbf{h}}_t^2)_{jk} d\tilde{h}_{t,j}^2, \quad \in \mathbb{R}^{1 \times d_{\text{out}}}, \quad \forall k \in [m] \quad (75)$$

$$(dh_{t,k}^2)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(-; \{J(\tilde{\mathbf{h}}_t^2)_{jk}\}_{j=1}^m, \{(d\tilde{h}_{t,j}^2)_\beta\}_{j=1}^m \right), \quad (\text{Moment}) \quad (76)$$

$$\phi(\dots) := \sum_{j=1}^m J(\tilde{\mathbf{h}}_t^2)_{jk} (d\tilde{h}_{t,j}^2)_\beta, \quad (77)$$

where $J(\tilde{\mathbf{h}}_t^2)_{ij} = \frac{\partial (\tilde{\mathbf{h}}_t^2)_i}{\partial (\tilde{\mathbf{h}}_t^2)_j}$. It is useful to consider this terms for all the experts,

$$(d\mathbf{h}_t^2)_\beta = \begin{pmatrix} (dh_{t,1}^2)_\beta \\ \vdots \\ (dh_{t,m}^2)_\beta \end{pmatrix} \in \mathbb{R}^{m \times d_{\text{out}}}. \quad (78)$$

$$d\mathbf{h}_t^{3,k} = \tilde{h}_{t,k}^2 d\mathbf{h}_t^3, \quad \in \mathbb{R}^{n \times d_{\text{out}}}, \quad \forall k \in [m] \quad (79)$$

$$(d\mathbf{h}_t^{3,k})_{\alpha\beta} = \phi \left((d\mathbf{h}_t^3)_{\alpha\beta}; \tilde{h}_{t,k}^2 \right), \quad (\text{Nonlin}) \quad (80)$$

$$\phi(\dots) := \tilde{h}_{t,k}^2 (\overline{d\mathbf{h}_t^3})_{\alpha\beta}, \quad (81)$$

where we repeat this for each expert to give $\mathbf{dh}_t^{3,1}, \dots, \mathbf{dh}_t^{3,m}$.

$$\mathbf{dh}_t^1 = \sum_{k=1}^m (\overline{\mathbf{W}}_t^2)_{:,k}^T dh_{t,k}^2 + \sum_{k=1}^m (\mathbf{W}_0^{3,k})^T \mathbf{dh}_t^{3,k} - \sum_{k=1}^m \sum_{s=0}^{t-1} \eta \mathbf{h}_s^1 (\boldsymbol{\chi}_s)^T \frac{(\mathbf{dh}_s^{3,k})^T \mathbf{dh}_t^{3,k}}{n}, \quad \in \mathbb{R}^{n \times d_{\text{out}}}, \quad (82)$$

$$(\mathbf{dh}_t^1)_{\alpha\beta} = \phi \left(\{(\overline{\mathbf{W}}_t^2)_{\alpha k}^T, (\mathbf{V}_t^k)_{\alpha\beta}\}_{k=1}^m, \{(\mathbf{h}_s^1)_\alpha\}_{s=1}^{t-1}; \{(dh_{t,k}^2)_\beta\}_{k=1}^m, \{(\boldsymbol{\chi}_s)_\beta, \{C_{s;\gamma\beta}^k\}_{k=1}^m\}_{s=1}^{t-1} \right), \quad (\text{Nonlin}) \quad (83)$$

$$(\mathbf{V}_t^k)_{:, \beta} := (\mathbf{W}_0^{3,k})^T (\mathbf{dh}_t^{3,k})_{:, \beta}, \quad (\text{MatMul}) \quad (84)$$

$$C_{s;\gamma\beta}^k := \frac{1}{n} \sum_{\alpha=1}^n \psi \left((\mathbf{dh}_s^{3,k})_{\alpha\gamma}, (\mathbf{dh}_t^{3,k})_{\alpha\beta}; - \right), \quad (\text{Moment}) \quad (85)$$

$$\phi(\dots) := \sum_{k=1}^m (\overline{\mathbf{W}}_t^2)_{\alpha k}^T (dh_{t,k}^2)_\beta + \sum_{k=1}^m (\mathbf{V}_t^k)_{\alpha\beta} - \sum_{k=1}^m \sum_{s=0}^{t-1} \sum_{\gamma=1}^{d_{\text{out}}} \eta (\mathbf{h}_s^1)_\alpha (\boldsymbol{\chi}_s)_\gamma C_{s;\gamma\beta}^k. \quad (86)$$

t th Weight Updates: Recall that we define the gradient of the loss with respect to the function as $\boldsymbol{\chi}_t \in \mathbb{R}^{d_{\text{out}} \times 1}$.

$$\mathbf{W}_t^4 - \mathbf{W}_{t-1}^4 = -n^{-c_4} \eta \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^3)^T \quad (87)$$

$$\theta_{\delta W_t^4} := n^{-c_4}. \quad (88)$$

$$\mathbf{W}_t^{3,k} - \mathbf{W}_{t-1}^{3,k} = -n^{-c_3} \theta_{W_{t-1}^4} \eta \mathbf{dh}_{t-1}^{3,k} \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^1)^T \quad (89)$$

$$\theta_{\delta W_t^3} := n^{-c_3} \theta_{W_{t-1}^4} \quad (90)$$

$$\mathbf{W}_t^2 - \mathbf{W}_{t-1}^2 = -n^{-c_2} \theta'_{W_{t-1}^4} \eta \mathbf{dh}_{t-1}^2 \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^1)^T \quad (91)$$

$$\theta_{\delta W_t^2} := n^{-c_2} \theta'_{W_{t-1}^4}. \quad (92)$$

$$\mathbf{W}_t^1 - \mathbf{W}_{t-1}^1 = -n^{-c_1} \theta'_{W_{t-1}^2} \theta_{W_{t-1}^4} \eta \mathbf{dh}_{t-1}^1 \boldsymbol{\chi}_{t-1} (\boldsymbol{\xi})^T \quad (93)$$

$$\theta_{\delta W_t^1} := n^{-c_1} \theta'_{W_{t-1}^2} \theta_{W_{t-1}^4}. \quad (94)$$

For the learnable bias we have,

$$\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2 = -\eta \theta'_{W_{t-1}^4} \mathbf{dh}_{t-1}^2 \boldsymbol{\chi}_{t-1}, \quad (95)$$

tth Feature Updates:

$$\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1 = (\mathbf{W}_t^1 - \mathbf{W}_{t-1}^1)\boldsymbol{\xi} \quad (96)$$

$$= -\theta_{\delta W_t^1} \eta d\mathbf{h}_{t-1}^1 \boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi} \quad (97)$$

$$(\delta \mathbf{h}_t^1)_\alpha = \phi \left(\{(\mathbf{d}\mathbf{h}_{t-1}^1)_{\alpha j}\}_{j=1}^{d_{\text{out}}}; -\eta, \{(\boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi})_j\}_{j=1}^{d_{\text{out}}}\right) \quad (\text{Nonlin}) \quad (98)$$

$$\phi(\dots) := -\eta \sum_{j=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^1)_{\alpha j} (\boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi})_j \quad (99)$$

$$\theta_{\delta h_t^1} := \theta_{\delta W_t^1}. \quad (100)$$

$$\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2 = \mathbf{W}_{t-1}^2 (\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1) + (\mathbf{W}_t^2 - \mathbf{W}_{t-1}^2) \mathbf{h}_t^1 + (\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2) \quad (101)$$

$$= \theta'_{W_{t-1}^2} \theta_{\delta h_t^1} \frac{\overline{\mathbf{W}_{t-1}^2} \delta \mathbf{h}_t^1}{n} - \theta'_{\delta W_t^2} \eta d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1} \frac{(\mathbf{h}_{t-1}^1)^T \mathbf{h}_t^1}{n} - \eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1}, \quad (102)$$

$$\theta_{\delta h_t^2} (\delta \mathbf{h}_t^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}_{t-1}^2})_{k\alpha}, (\delta \mathbf{h}_t^1)_\alpha, (\mathbf{h}_{t-1}^1)_\alpha, (\mathbf{h}_t^1)_\alpha; \right. \\ \left. \theta'_{W_{t-1}^4}, \theta'_{W_{t-1}^2}, \theta_{\delta h_t^1}, \theta'_{\delta W_t^2}, \eta, \{(\mathbf{d}h_{t-1,k}^2)_\beta, (\boldsymbol{\chi}_{t-1})_\beta\}_{\beta=1}^{d_{\text{out}}}\right), \quad (\text{Moment}) \quad (103)$$

$$\phi(\dots) := \theta'_{W_{t-1}^2} \theta_{\delta h_t^1} (\overline{\mathbf{W}_{t-1}^2})_{k\alpha} (\delta \mathbf{h}_t^1)_\alpha - \theta'_{\delta W_t^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}h_{t-1,k}^2)_\beta (\boldsymbol{\chi}_{t-1})_\beta (\mathbf{h}_{t-1}^1)_\alpha (\mathbf{h}_t^1)_\alpha \\ - \eta \theta'_{W_{t-1}^4} \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta} (\boldsymbol{\chi}_{t-1})_\beta. \quad (104)$$

$$\theta_{\delta h_t^2} := \max(\theta'_{W_{t-1}^2} \theta_{\delta h_t^1}, \theta'_{\delta W_t^2}, \theta'_{W_{t-1}^4}). \quad (105)$$

$$\tilde{\mathbf{h}}_t^2 - \tilde{\mathbf{h}}_{t-1}^2 = \mathbf{G}(\mathbf{h}_t^2) - \mathbf{G}(\mathbf{h}_{t-1}^2) \quad (106)$$

$$= \mathbf{G} \left(\mathbf{h}_{t-1}^2 + \theta_{\delta h_t^2} \delta \mathbf{h}_t^2 \right) - \mathbf{G}(\mathbf{h}_{t-1}^2) \quad (107)$$

$$\theta_{\delta \tilde{h}_t^2} (\delta \tilde{\mathbf{h}}_t^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi_k \left(-; \{(\mathbf{h}_{t-1}^2)_k, (\delta \mathbf{h}_t^2)_k\}_{k=1}^m, \theta_{\delta h_t^2} \right), \quad (\text{Moment}), \quad (108)$$

$$\theta_{\delta \tilde{h}_t^2} := \theta_{\delta h_t^2}, \quad (109)$$

where \mathbf{G} is the gating function (either softmax or sigmoid) (see Appendix G) (Yang and Hu, 2021).

$$\mathbf{h}_t^{3,k} - \mathbf{h}_{t-1}^{3,k} = \mathbf{W}_{t-1}^{3,k}(\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1) + (\mathbf{W}_t^{3,k} - \mathbf{W}_{t-1}^{3,k})\mathbf{h}_t^1, \quad (110)$$

$$= \theta_{\delta h_t^1} \mathbf{W}_0^{3,k} \delta \mathbf{h}_t^1 - \sum_{s=0}^{t-2} \theta'_{\delta W_{s+1}^3} \theta_{\delta h_t^1} \eta d\mathbf{h}_s^{3,k} \boldsymbol{\chi}_s \frac{(\mathbf{h}_s^1)^T (\delta \mathbf{h}_t^1)}{n} - \theta'_{\delta W_t^3} \eta d\mathbf{h}_{t-1}^{3,k} \boldsymbol{\chi}_{t-1} \frac{(\mathbf{h}_{t-1}^1)^T (\mathbf{h}_t^1)}{n}, \quad (111)$$

$$(\delta \mathbf{h}_t^{3,k})_\alpha = \phi \left(\Upsilon_\alpha, \{\Lambda_{s,\alpha}\}_{s=0}^{t-1}, \mathbf{V}_\alpha; \theta_{\delta h_t^1 / \delta h_t^{3,k}}, \{\theta'_{\delta W_s^3 \delta h_t^1 / \delta h_t^{3,k}}\}_{s=0}^{t-1}, \theta'_{\delta W_t^3 / \delta h_t^{3,k}} \right) \quad (\text{Nonlin}) \quad (112)$$

$$\phi(\dots) := \theta_{\delta h_t^1 / \delta h_t^{3,k}} \Upsilon_\alpha - \sum_{s=0}^{t-2} \theta'_{\delta W_{s+1}^3 \delta h_t^1 / \delta h_t^{3,k}} \Lambda_{s,\alpha} - \theta'_{\delta W_t^3 / \delta h_t^{3,k}} \mathbf{V}_\alpha, \quad (113)$$

$$\Upsilon := \mathbf{W}_0^{3,k} \delta \mathbf{h}_t^1, \quad (\text{MatMul}) \quad (114)$$

$$\Lambda_{s,\alpha} := \psi \left(\{(\mathbf{d}\mathbf{h}_s^{3,k})_{\alpha l}\}_{l=1}^{d_{\text{out}}}; \eta, \{(\boldsymbol{\chi}_s)_l\}_{l=1}^{d_{\text{out}}}, d_s \right) = \eta d_s \sum_{l=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_s^{3,k})_{\alpha l} (\boldsymbol{\chi}_s)_l, \quad (\text{Nonlin}) \quad (115)$$

$$d_s := \frac{1}{n} \sum_{\alpha=1}^n \nu \left((\mathbf{h}_s^1)_\alpha, (\delta \mathbf{h}_t^1)_\alpha; - \right), \quad (\text{Moment}) \quad (116)$$

$$\mathbf{V}_\alpha := \psi \left(\{(\mathbf{d}\mathbf{h}_{t-1}^{3,k})_{\alpha l}\}_{l=1}^{d_{\text{out}}}; \{(\boldsymbol{\chi}_{t-1})_l\}_{l=1}^{d_{\text{out}}}, \eta, c_{t-1} \right) = \eta c_{t-1} \sum_{l=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^{3,k})_{\alpha l} (\boldsymbol{\chi}_{t-1})_l, \quad (\text{Nonlin}) \quad (117)$$

$$c_{t-1} := \frac{1}{n} \sum_{\alpha=1}^n \psi \left((\mathbf{h}_{t-1}^1)_\alpha, (\mathbf{h}_t^1)_\alpha; - \right), \quad (\text{Moment}) \quad (118)$$

$$\theta_{\delta h_t^{3,k}} = \max \left(\theta_{\delta h_t^1}, \max_{0 \leq s \leq t-2} \left(\theta'_{\delta W_{s+1}^3} \theta_{\delta h_t^1} \right), \theta'_{\delta W_t^3} \right), \quad (119)$$

$$\theta_{\delta h_t^1 / \delta h_t^{3,k}} := \theta_{\delta h_t^1} / \theta_{\delta h_t^{3,k}}, \quad (120)$$

$$\theta'_{\delta W_s^3 \delta h_t^1 / \delta h_t^{3,k}} := \theta'_{\delta W_s^3} \theta_{\delta h_t^1} / \theta_{\delta h_t^{3,k}}, \quad (121)$$

$$\theta'_{\delta W_t^3 / \delta h_t^{3,k}} := \theta'_{\delta W_t^3} / \theta_{\delta h_t^{3,k}}. \quad (122)$$

We repeat the above process for each expert k .

$$\mathbf{h}_t^3 - \mathbf{h}_{t-1}^3 = \sum_{k=1}^m \left(\tilde{h}_{t-1,k}^2 (\mathbf{h}_t^{3,k} - \mathbf{h}_{t-1}^{3,k}) + (\tilde{h}_{t,k}^2 - \tilde{h}_{t-1,k}^2) \mathbf{h}_t^{3,k} \right), \quad (123)$$

$$= \sum_{k=1}^m \theta_{\delta h_t^{3,k}} \tilde{h}_{t-1,k}^2 \delta \mathbf{h}_t^{3,k} + \theta_{\delta \tilde{h}_t^2} \delta \tilde{h}_{t,k}^2 \mathbf{h}_t^{3,k}, \quad (124)$$

$$(\delta \mathbf{h}_t^3)_\alpha = \phi \left(\{(\delta \mathbf{h}_t^{3,k})_\alpha\}_{k=1}^m, \{(\mathbf{h}_t^{3,k})_\alpha\}_{k=1}^m; \theta_{\delta h_t^{3,k}/\delta h_t^3}, \theta_{\delta \tilde{h}_t^2/\delta h_t^3}, \{\tilde{h}_{t-1,k}^2\}_{k=1}^m, \{\delta \tilde{h}_{t,k}^2\}_{k=1}^m \right), \quad (\text{Nonlin}) \quad (125)$$

$$\theta_{\delta h_t^3} := \max(\theta_{\delta h_t^{3,k}}, \theta_{\delta \tilde{h}_t^2}), \quad (126)$$

$$\theta_{\delta h_t^{3,k}/\delta h_t^3} := \theta_{\delta h_t^{3,k}} / \theta_{\delta h_t^3}, \quad (127)$$

$$\theta_{\delta \tilde{h}_t^2/\delta h_t^3} := \theta_{\delta \tilde{h}_t^2} / \theta_{\delta h_t^3}. \quad (128)$$

Note that in the above, all experts have the same scale by symmetry so we do not write a maximum over experts.

$$\mathbf{f}_t - \mathbf{f}_{t-1} = \mathbf{W}_{t-1}^4 (\mathbf{h}_t^3 - \mathbf{h}_{t-1}^3) + (\mathbf{W}_t^4 - \mathbf{W}_{t-1}^4) \mathbf{h}_t^3, \quad (129)$$

$$= \theta'_{W_{t-1}^4} \theta_{\delta h_t^3} \frac{\overline{\mathbf{W}_{t-1}^4} \delta \mathbf{h}_t^3}{n} - \theta'_{\delta W_t^4} \eta \chi_{t-1} \frac{(\mathbf{h}_{t-1}^3)^T (\mathbf{h}_t^3)}{n}, \quad (130)$$

$$\theta_{\delta f_t} (\delta \mathbf{f}_t)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}_{t-1}^4})_{\beta\alpha}, (\delta \mathbf{h}_t^3)_\alpha, (\mathbf{h}_{t-1}^3)_\alpha, (\mathbf{h}_t^3)_\alpha; \theta'_{W_{t-1}^4}, \theta_{\delta h_t^3}, \theta'_{\delta W_t^4}, \eta, (\chi_{t-1})_\beta \right) \quad (\text{Moment}) \quad (131)$$

$$\phi(\dots) := \theta'_{W_{t-1}^4} \theta_{\delta h_t^3} (\overline{\mathbf{W}_{t-1}^4})_{\beta\alpha} (\delta \mathbf{h}_t^3)_\alpha - \theta'_{\delta W_t^4} \eta (\chi_{t-1})_\beta (\mathbf{h}_{t-1}^3)_\alpha (\mathbf{h}_t^3)_\alpha \quad (132)$$

$$\theta_{\delta f_t} := \max(\theta'_{W_{t-1}^4} \theta_{\delta h_t^3}, \theta'_{\delta W_t^4}). \quad (133)$$

C.2. Infinite Width Limit for Soft MoE with SGD

Under stable parametrisations (Theorem 9), all the θ parameters in the program will converge to either zero or one.

Z random variables: According to the Master Theorem (Yang and Hu, 2021), the distribution of the coordinates of each vector z in the program will follow an iid distribution which we refer to as Z^z . Each scalar in the program \bullet converges to a deterministic limit that we label $\dot{\bullet}$. We recursively define each Z distribution based on the operation in the program used to create it.

Z Init: Let \mathcal{V}_I be the set of initial vectors. For each $z \in \mathcal{V}_I$, we set $\hat{Z}^z := Z^z$ and $\dot{Z}^z := 0$ (defined below).

Z Nonlin: Given $\phi : \mathbb{R}^k \times \mathbb{R}^l \rightarrow \mathbb{R}$, scalars $\theta_1, \dots, \theta_l \in \mathbb{R}$, and vectors $x^1, \dots, x^k \in \mathbb{R}^n$, we have

$$Z \phi(x^1, \dots, x^k; \theta_1, \dots, \theta_l) := \phi(Z^{z^1}, \dots, Z^{z^k}; \dot{\theta}_1, \dots, \dot{\theta}_l). \quad (134)$$

Z Moment: Given a scalar $\theta = \frac{1}{n} \sum_{\alpha=1}^n \phi(x_\alpha^1, \dots, x_\alpha^k; \theta_1, \dots, \theta_l)$, we have

$$\dot{\theta} := \mathbb{E}[\phi(Z^{z^1}, \dots, Z^{z^k}; \dot{\theta}_1, \dots, \dot{\theta}_l)], \quad (135)$$

where the expectation is over Z^{z^1}, \dots, Z^{z^k} .

Z MatMul: $Z^{Wx} := \hat{Z}^{Wx} + \dot{Z}^{Wx}$ for every matrix $W \sim \mathcal{N}(0, \sigma_W^2/n)$ and vector x .

- \hat{Z}^{Wx} is Gaussian with zero mean. Let $\mathcal{V}_W := \{Wy \text{ for some vector } y\}$. Then $\{\hat{Z}^{Wy} : Wy \in \mathcal{V}_W\}$ is jointly Gaussian with zero mean and covariance,

$$\text{Cov}(\hat{Z}^{Wx}, \hat{Z}^{Wy}) := \sigma_W^2 \mathbb{E}[Z^x Z^y] \text{ for any } Wx, Wy \in \mathcal{V}_W. \quad (136)$$

- Z^x has been computed by some operation that we unwind $Z^x = \phi(\{\hat{Z}^{W^T y^i}\}_{i=1}^k, \{\hat{Z}^{z^i}\}_{i=1}^j; \{\hat{\theta}_i\}_{i=1}^l)$ with $z^i \notin \mathcal{V}_{W^T}$ ¹. Define $\partial Z^x / \partial \hat{Z}^{W^T y^i} := \partial_i \phi(\dots)$. Then

$$\dot{Z}^{Wx} := \sigma_W^2 \sum_{i=1}^k Z^{y^i} \mathbb{E} \left[\frac{\partial Z^x}{\partial \hat{Z}^{W^T y^i}} \right]. \quad (137)$$

Initial vectors: All our initial vectors have coordinates distributed according to $\mathcal{N}(0, 1)$ by construction. The set of initial vectors are $\mathcal{V} := \{\overline{\mathbf{W}}^1, \overline{\mathbf{W}}^2, \overline{\mathbf{W}}^4\}$. Thus, $Z^W = \mathcal{N}(0, 1)$ for $W \in \mathcal{V}$.

First Forward Pass:

$$Z^{h_0^1} = \sum_{j=1}^{d_{\text{in}}} \xi_j Z^{W_{0,j}^1}, \quad (138)$$

$$\mathring{\mathbf{h}}_0^2 = \mathbf{b}_0^2, \quad (139)$$

$$\mathring{\mathbf{h}}_0^2 = G(\mathbf{b}_0^2), \quad (140)$$

$$Z^{h_0^{3,k}} = Z^{W_0^{3,k} h_0^1} \text{ for every expert } k, \quad (141)$$

$$Z^{h_0^3} = \sum_{k=1}^m \mathring{h}_{0,k}^2 Z^{h_0^{3,k}} \quad (142)$$

$$\mathring{\mathbf{f}}_0 = \mathbf{0}. \quad (143)$$

Here, $Z^{W_0^1}$ is standard Gaussian by construction. $\mathring{\mathbf{h}}_0^2 = \mathbf{b}_0^2$ as the non bias router logits tend to zero $\mathbb{E}[Z^{W_0^2}] \mathbb{E}[Z^{h_0^1}] = 0$. $Z^{W_0^{3,k} h_0^1} = \hat{Z}^{W_0^{3,k} h_0^1}$ from the Z MatMul rule. Finally $\mathring{\mathbf{f}}_0 = \mathbf{0}$ as $\mathring{f}_0 = \mathbb{E}[Z^{W_0^4}] \mathbb{E}[Z^{h_0^3}] = 0$.

First Backward Pass: For $d_{\text{out}} > 1$, we index the distributions by β ,

$$Z^{(dh_0^3)_\beta} = Z^{(\overline{W_0^4})_\beta}, \quad (144)$$

$$d\tilde{h}_{0,k}^2 = 0 \text{ for every expert } k, \quad (145)$$

$$dh_{0,k}^2 = 0 \text{ for every expert } k, \quad (146)$$

$$Z^{(dh_0^{3,k})_\beta} = \mathring{h}_{0,k}^2 Z^{(dh_0^3)_\beta} \text{ for every expert } k, \quad (147)$$

$$Z^{(dh_0^1)_\beta} = \sum_{k=1}^m Z^{(W_0^{3,k})^T (dh_0^{3,k})_\beta} \quad (148)$$

$$(149)$$

1. z^i is not a function of W^T .

$Z^{\overline{W_0^4}}$ is a standard Gaussian by construction. $d\tilde{h}_{0,k}^2 = \mathbb{E}[Z^{h_0^{3,k}}] \mathbb{E}[\overline{W_0^4}] = 0$. From this, $dh_{0,k}^2 = 0$ follows. $Z^{(W_0^{3,k})^T dh_0^3} = \hat{Z}^{(W_0^{3,k})^T dh_0^3}$ from the Z MatMul rule.

First Feature Update:

$$Z^{\delta h_1^1} = -\eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_0^1)_\beta} (\dot{\chi}_0)_\beta \boldsymbol{\xi}^T \boldsymbol{\xi}, \quad (150)$$

$$\delta \dot{h}_{1,k}^2 = \dot{\theta}_{\delta h_1^1} \mathbb{E}[Z^{(\overline{W_0^2})_{k,:}} Z^{\delta h_1^1}] - \dot{\theta}'_{\delta W_1^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (dh_{0,k}^2)_\beta (\dot{\chi}_0)_\beta \mathbb{E}[Z^{h_0^1} Z^{h_1^1}] - \eta d\mathbf{h}_0^2 \dot{\chi}_0, \quad (151)$$

$$\delta \dot{\mathbf{h}}_1^2 = \begin{cases} G(\dot{\mathbf{h}}_1^2) - G(\dot{\mathbf{h}}_0^2) & \text{if } \dot{\theta}_{\delta h_1^2} = 1 \\ G'(\dot{\mathbf{h}}_0^2)(\delta \dot{\mathbf{h}}_1^2) & \text{if } \dot{\theta}_{\delta h_1^2} = 0 \end{cases}, \quad (152)$$

$$Z^{\delta h_1^{3,k}} = \dot{\theta}_{\delta h_1^1 / \delta h_1^{3,k}} Z^{W_0^{3,k} \delta h_1^1} - \dot{\theta}'_{\delta W_1^3 / \delta h_1^{3,k}} \eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_0^{3,k})_\beta} (\dot{\chi}_0)_\beta \mathbb{E}[Z^{h_0^1} Z^{h_1^1}] \text{ for every expert } k, \quad (153)$$

$$Z^{\delta h_1^3} = \sum_{k=1}^m \dot{\theta}_{\delta h_1^{3,k} / \delta h_1^3} \delta \dot{h}_{0,k}^2 Z^{\delta h_1^{3,k}} + \dot{\theta}_{\delta \tilde{h}_1^2} \dot{\theta}_{h_1^{3,k} / \delta h_3} \delta \tilde{h}_{1,k}^2 Z^{h_1^{3,k}}, \quad (154)$$

$$\delta f_1^3 = \dot{\theta}_{\delta h_1^3} \mathbb{E}[Z^{\overline{W_0^4}} Z^{\delta h_1^3}] - \dot{\theta}'_{\delta W_1^4} \eta \dot{\chi}_0 \mathbb{E}[Z^{h_0^3} Z^{h_1^3}], \quad (155)$$

where G is either softmax or sigmoid and G' is the derivative. Further, $Z^{W_0^{3,k} \delta h_1^1}$ is given by the Z MatMul rule,

$$Z^{W_0^{3,k} \delta h_1^1} = \hat{Z}^{W_0^{3,k} \delta h_1^1} + \dot{Z}^{W_0^{3,k} \delta h_1^1} \quad (156)$$

$$= \hat{Z}^{W_0^{3,k} \delta h_1^1} + \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_0^{3,k})_\beta} \mathbb{E} \left[\frac{\partial Z^{\delta h_1^1}}{\partial \hat{Z}^{(W_0^{3,k})^T (dh_0^{3,k})_\beta}} \right]. \quad (157)$$

We write the limit of each vector in the program as

$$Z^{h_t} = Z^{h_0} + \dot{\theta}_{\delta h_1} Z^{\delta h_1} + \dots + \dot{\theta}_{\delta h_t} Z^{\delta h_t}. \quad (158)$$

Weight updates: The distribution of certain weight updates are required in the backward pass and thus the forward passes.

$$Z^{\delta \overline{W}_t^4} = -\eta \dot{\chi}_{t-1} Z^{h_{t-1}^3}, \quad (159)$$

$$Z^{(\delta \overline{W}_t^2)_{k,:}} = -\eta \sum_{\beta=1}^{d_{\text{out}}} (dh_{t-1,k}^2)_\beta (\dot{\chi}_{t-1})_\beta Z^{h_{t-1}^1}, \quad (160)$$

$$\delta \mathbf{b}_t^2 = -\eta \dot{\theta}'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \dot{\chi}_{t-1}. \quad (161)$$

***t*th Backward Pass:**

$$Z(\text{dh}_t^3)_\beta = Z(\overline{W_t^4})_\beta^T, \quad (162)$$

$$(\mathring{\text{d}}\hat{h}_{t,k}^2)_\beta = \mathbb{E}[Z^{h_t^{3,k}} Z(\text{dh}_t^3)_\beta], \quad (163)$$

$$(\mathring{\text{d}}\hat{h}_{t,k}^2)_\beta = \sum_{j=1}^m J(\hat{h}_t^2)_{jk} (\mathring{\text{d}}\hat{h}_{t,j}^2)_\beta, \quad (164)$$

$$Z(\text{dh}_t^{3,k})_\beta = \mathring{\hat{h}}_{t,k}^2 Z(\text{dh}_t^3)_\beta, \quad (165)$$

$$Z(\text{dh}_t^1)_\beta = \sum_{k=1}^m \left(Z(\overline{W_t^2})_{k,:} (\mathring{\text{d}}\hat{h}_{t,k}^2)_\beta + Z^{(W_0^{3,k})^T} (\text{dh}_t^{3,k})_\beta - \sum_{s=1}^{t-1} Z^{h_s^1} \sum_{\gamma=1}^{d_{\text{out}}} (\dot{\chi}_s)_\gamma \mathbb{E}[Z^{(\text{dh}_s^{3,k})_\gamma} Z(\text{dh}_t^{3,k})_\beta] \right), \quad (166)$$

where $Z^{(W_0^{3,k})^T} \text{dh}_t^{3,k}$ is given by the Z MatMul rule,

$$Z^{(W_0^{3,k})^T} \text{dh}_t^{3,k} = \hat{Z}^{(W_0^{3,k})^T} \text{dh}_t^{3,k} + \dot{Z}^{(W_0^{3,k})^T} \text{dh}_t^{3,k} \quad (167)$$

$$= \hat{Z}^{(W_0^{3,k})^T} \text{dh}_t^{3,k} + \sum_{v \in \mathcal{V}: W_0^{3,k} v \in \mathcal{V}} Z^v \mathbb{E} \left[\frac{\partial Z^{\text{dh}_t^{3,k}}}{\partial \hat{Z}^{W_0^{3,k} v}} \right], \quad (168)$$

where $\mathcal{V} := \{v \text{ for some vector in the program } v\}$ is the set of all vectors in the program.

***t*th Forward Pass:**

$$Z^{\delta h_t^1} = -\eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(\text{dh}_{t-1}^1)_\beta} (\dot{\chi}_{t-1})_\beta \xi^T \xi, \quad (169)$$

$$\delta \hat{h}_{t,k}^2 = \theta'_{W_{t-1}^2} \theta_{\delta h_t^1} \mathbb{E}[Z^{(\overline{W_{t-1}^2})_{k,:}} Z^{\delta h_t^1}] - \theta'_{\delta W_t^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathring{\text{d}}\hat{h}_{t-1,k}^2)_\beta (\dot{\chi}_{t-1})_\beta \mathbb{E}[Z^{h_{t-1}^1} Z^{h_t^1}] + \delta \mathbf{b}_t^2, \quad (170)$$

$$\delta \hat{\mathbf{h}}_t^2 = \begin{cases} \mathbf{G}(\hat{\mathbf{h}}_t^2) - \mathbf{G}(\hat{\mathbf{h}}_{t-1}^2) & \text{if } \theta_{\delta h_t^2} = 1 \\ \mathbf{G}'(\hat{\mathbf{h}}_{t-1}^2)(\delta \hat{\mathbf{h}}_t^2) & \text{if } \theta_{\delta h_t^2} = 0 \end{cases}, \quad (171)$$

$$Z^{\delta h_t^{3,k}} = \theta_{\delta h_t^1 / \delta h_t^{3,k}} Z^{W_0^{3,k} \delta h_t^1} - \sum_{s=0}^{t-2} \theta'_{\delta W_{s+1}^{3,k} \delta h_t^1 / \delta h_t^{3,k}} \eta \sum_{\beta}^{d_{\text{out}}} Z^{(\text{dh}_s^{3,k})_\beta} (\dot{\chi}_s)_\beta \mathbb{E}[Z^{h_s^1} Z^{\delta h_t^1}] \quad (172)$$

$$- \theta'_{\delta W_t^{3,k} / \delta h_t^{3,k}} \eta \sum_{\beta}^{d_{\text{out}}} Z^{(\text{dh}_{t-1}^{3,k})_\beta} (\dot{\chi}_{t-1})_\beta \mathbb{E}[Z^{h_{t-1}^1} Z^{h_t^1}], \quad (173)$$

$$Z^{\delta h_t^3} = \sum_{k=1}^m \theta_{\delta h_t^{3,k} / \delta h_t^3} \mathring{\hat{h}}_{k,t-1}^2 Z^{\delta h_t^{3,k}} + \theta_{\delta \hat{h}_t^2 / \delta h_t^3} \delta \hat{h}_{t,k}^2 Z^{h_t^{3,k}}, \quad (174)$$

$$(\delta \hat{f}_t)_\beta = \theta'_{W_{t-1}^4} \theta_{\delta h_t^3} \mathbb{E}[Z^{(\overline{W_{t-1}^4})_\beta} Z^{\delta h_t^3}] - \theta'_{\delta W_t^4} \eta (\dot{\chi}_{t-1})_\beta \mathbb{E}[Z^{h_{t-1}^3} Z^{h_t^3}], \quad (175)$$

where G is either softmax or sigmoid and G' is the derivative, and

$$Z^{W_0^{3,k}} \delta h_t^1 = \hat{Z}^{W_0^{3,k}} \delta h_t^1 + \dot{Z}^{W_0^{3,k}} \delta h_t^1 \quad (176)$$

$$= \hat{Z}^{W_0^{3,k}} \delta h_t^1 + \sum_{s=0}^{t-1} \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_s^{3,k})_\beta} \mathbb{E} \left[\frac{\partial Z^{\delta h_t^1}}{\partial \hat{Z}^{(W_0^{3,k})^T (dh_s^{3,k})_\beta}} \right] \quad (177)$$

C.3. Tensor Program for Soft MoE with Adam

We consider Adam optimiser updates here.

Tensor Program: We use $\text{NE} \otimes \text{ORT}$ for the Adam limit (Yang and Littwin, 2023). This has the same operations as NETSORT^+ (Appendix C.1), but with a generalisation of the `Nonlin`.

- `OuterNonlin`: Given $r \geq 0$, vectors $\mathbf{x}^0, \dots, \mathbf{x}^r$, and function $\psi : \mathbb{R}^{(r+1)+\ell} \rightarrow \mathbb{R}$, generate a new vector (\mathbb{R}^n) as

$$y_\alpha = \frac{1}{n^r} \sum_{\beta_1, \dots, \beta_r=1}^n \psi(x_\alpha^0; x_{\beta_1}^1; \dots; x_{\beta_r}^r; \theta_1, \dots, \theta_l)$$

Definition 13 (Adam update) Given gradients g_0, \dots, g_t for a parameter \mathbf{W} , the Adam update at step t with learning rate η is,

$$\begin{aligned} \hat{m} &\leftarrow \frac{1}{1 - \beta_1^{t+1}} \sum_{s=0}^t (1 - \beta_1) \beta_1^{t-s} g_s, \\ \hat{v} &\leftarrow \sqrt{\frac{1}{1 - \beta_2^{t+1}} \sum_{s=0}^t (1 - \beta_2) \beta_2^{t-s} g_s^2 + \epsilon^2}, \\ \mathbf{W} &\leftarrow \mathbf{W} - \eta \left(\frac{\hat{m}}{\sqrt{\hat{v} + \epsilon^2}} \right), \end{aligned}$$

where (β_1, β_2) are momentum parameters and $\epsilon > 0$ is a constant.

bcd parametrisation for Adam updates: The parameters b and c play the same role as in Appendix C.1. We parametrise the weights for layer l as $\mathbf{W}^l = n^{-b_l} \overline{\mathbf{W}}^l$. The learning rate for the parameter \mathbf{W}^l is ηn^{-c_l} , where η is some width independent constant. The update function for layer l at time t is Q_t^l , which takes in scaled gradients $Q_t^l(n^{d_l} g_0, \dots, n^{d_l} g_t)$. For Adam the update function is

$$Q_t^l(n^{d_l} g_0, \dots, n^{d_l} g_t) = \frac{\frac{1}{1 - \beta_1^{t+1}} \sum_{s=0}^t (1 - \beta_1) \beta_1^{t-s} n^{d_l} g_s}{\sqrt{\frac{1}{1 - \beta_2^{t+1}} \sum_{s=0}^t (1 - \beta_2) \beta_2^{t-s} n^{2d_l} g_s^2 + \epsilon^2}}. \quad (178)$$

Note that instead of scaling the gradients g by n^{d_l} , we can equivalently scale $\epsilon \rightarrow n^{-d_l} \epsilon$.

Faithfulness (Theorem 7): To ensure that the inputs to the update function Q are neither too large (causing diverging behavior) or too small (trivialising the update function), we require that the inputs to the update function scale $\Theta(1)$. Equivalently, for our Adam update function, we require that the scale of $n^{-d_i}\epsilon$ is the same as the scale of the gradients g_0, \dots, g_t .

We assume that the parametrisation is stable (Theorem 9), and parameter values are chosen such that the network is stable at initialisation (Theorem 11).

First Backward Pass: Same as Appendix C.1.

First Weight Updates: Recall that we define the gradient of the loss with respect to the function as χ_t . For the vectors $\mathbf{W} \in \{\mathbf{W}^1, \mathbf{W}^2, \mathbf{W}^4\}$ we write $\delta\mathbf{W}_t := \theta_{\delta W_t}^{-1}(\mathbf{W}_t - \mathbf{W}_{t-1})$, and thus the weights at step t as $\mathbf{W}_t = \theta_{W_0} \overline{\mathbf{W}}_0 + \dots + \theta_{\delta W_t} \delta W_t$.

$$\mathbf{W}_1^4 - \mathbf{W}_0^4 = -n^{-c_4} \eta Q_1^4 \left(n^{d_4} \chi_0(\mathbf{h}_0^3)^T \right) \quad (179)$$

$$\theta_{\delta W_1^4} := n^{-c_4}. \quad (180)$$

$$\mathbf{W}_1^{3,k} - \mathbf{W}_0^{3,k} = -n^{-c_3} \eta Q_1^3 \left(n^{(d_3-1)} d\mathbf{h}_0^{3,k} \chi_0(\mathbf{h}_0^1)^T \right), \quad (181)$$

$$\theta_{\delta W_1^3} := n^{-c_3}. \quad (182)$$

$$\mathbf{W}_1^2 - \mathbf{W}_0^2 = -n^{-c_2} \eta Q_1^2 \left(n^{d_2} n \theta_{W_0^2} d\mathbf{h}_0^2 \chi_0(\mathbf{h}_0^1)^T \right) \quad (183)$$

$$= -n^{-c_2} \eta Q_1^2 \left(n^{d_2} d\mathbf{h}_0^2 \chi_0(\mathbf{h}_0^1)^T \right) \quad (184)$$

$$\theta_{\delta W_1^2} := n^{-c_2}. \quad (185)$$

$$\mathbf{W}_1^1 - \mathbf{W}_0^1 = -n^{-c_1} \eta Q_1^1 \left(n^{d_1} n \theta_{W_0^2} \theta_{W_0^4} d\mathbf{h}_0^1 \chi_0(\boldsymbol{\xi})^T \right) \quad (186)$$

$$= -n^{-c_1} \eta Q_1^1 \left(n^{(d_1-1)} d\mathbf{h}_0^1 \chi_0(\boldsymbol{\xi})^T \right) \quad (187)$$

$$\theta_{\delta W_1^1} := n^{-c_1}. \quad (188)$$

The router bias in the learnable case updates as

$$\mathbf{b}_1^2 - \mathbf{b}_0^2 = -\eta \theta'_{W_0^4} d\mathbf{h}_0^2 \chi_0, \quad (189)$$

where due to stability at initialisation (Theorem 11) $\theta'_{W_0^4} = 1$.

First Feature Updates: The feature updates are affected by the change in the previous layers features (*propagating updates*) and the update in the weights at the current layer (*effective updates*). We write the feature updates for the vectors $\mathbf{h} \in \{\mathbf{h}^1, \{\mathbf{h}^{3,k}\}_{k=1}^m, \mathbf{h}^3\}$ as $\delta\mathbf{h}_t := \theta_{\delta h_t}^{-1}(\mathbf{h}_t - \mathbf{h}_{t-1})$. This gives each features for the vectors at step t as $\mathbf{h}_t = \mathbf{h}_0 + \theta_{\delta h_1} \delta\mathbf{h}_1 + \dots + \theta_{\delta h_t} \delta\mathbf{h}_t$. For the scalars $\mathbf{h} \in \{\mathbf{h}^2, \tilde{\mathbf{h}}^2, \mathbf{f}\}$ we simply write $\delta\mathbf{h}_t := (\mathbf{h}_t - \mathbf{h}_{t-1})$ and $\mathbf{h}_t = \mathbf{h}_0 + \dots + \delta\mathbf{h}_t$.

$$\mathbf{h}_1^1 - \mathbf{h}_0^1 = (\mathbf{W}_1^1 - \mathbf{W}_0^1)\boldsymbol{\xi} \quad (190)$$

$$= -\theta_{\delta W_1^1} \eta Q_1^1 \left(n^{d_1} n^{-1} d\mathbf{h}_0^1 \boldsymbol{\chi}_0 \boldsymbol{\xi}^T \right) \boldsymbol{\xi} \quad (191)$$

$$(\delta \mathbf{h}_1^1)_\alpha = \phi \left(\{ (d\mathbf{h}_0^1)_{\alpha\beta} \}_{\beta=1}^{d_{\text{out}}}; -\eta, n^{(d_1-1)}, \{ (\boldsymbol{\chi}_0)_\beta \}_{\beta=1}^{d_{\text{out}}}, \boldsymbol{\xi} \right) \quad (\text{Nonlin}) \quad (192)$$

$$\phi(\dots) := -\eta Q_1^1 \left(n^{(d_1-1)} \sum_{\beta=1}^{d_{\text{out}}} (d\mathbf{h}_0^1)_{\alpha\beta} (\boldsymbol{\chi}_0)_\beta \boldsymbol{\xi}^T \right) \boldsymbol{\xi}, \quad (193)$$

$$\theta_{\delta h_1^1} := \theta_{\delta W_1^1}, \quad (194)$$

$$\mathbf{h}_1^2 - \mathbf{h}_0^2 = \mathbf{W}_0^2 (\mathbf{h}_1^1 - \mathbf{h}_0^1) + (\mathbf{W}_1^2 - \mathbf{W}_0^2) \mathbf{h}_1^1 + (\mathbf{b}_1^2 - \mathbf{b}_0^2) \quad (195)$$

$$= \theta_{\delta h_1^1} \theta'_{W_0^2} \frac{\overline{\mathbf{W}_0^2} \delta \mathbf{h}_1^1}{n} - \theta'_{\delta W_1^2} \eta \frac{1}{n} Q_1^2 \left(n^{d_2} d\mathbf{h}_0^2 \boldsymbol{\chi}_0 (\mathbf{h}_0^1)^T \right) \mathbf{h}_1^1 - \eta d\mathbf{h}_0^2 \boldsymbol{\chi}_0, \quad (196)$$

$$\theta_{\delta h_1^2} (\delta \mathbf{h}_1^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_0^2}_{k\alpha}, (\delta \mathbf{h}_1^1)_\alpha, (\mathbf{h}_0^1)_\alpha, (\mathbf{h}_1^1)_\alpha; \theta_{\delta h_1^1}, \theta'_{W_0^2}, \theta'_{\delta W_1^2}, n^{d_2}, \eta, \{ (d\mathbf{h}_0^2)_{k\beta}, (\boldsymbol{\chi}_0)_\beta \}_{\beta=1}^{d_{\text{out}}} \right) \quad (\text{Moment}) \quad (197)$$

$$\phi(\dots) := \theta_{\delta h_1^1} \theta'_{W_0^2} \overline{\mathbf{W}_0^2}_{k\alpha} (\delta \mathbf{h}_1^1)_\alpha - \theta'_{\delta W_1^2} \eta \frac{1}{n} Q_1^2 \left(n^{d_2} \sum_{\beta=1}^{d_{\text{out}}} (d\mathbf{h}_0^2)_{k\beta} (\boldsymbol{\chi}_0)_\beta (\mathbf{h}_0^1)_\alpha \right) (\mathbf{h}_1^1)_\alpha \quad (198)$$

$$- \eta \theta'_{W_0^2} d\mathbf{h}_0^2 \boldsymbol{\chi}_0, \quad (199)$$

$$\theta_{\delta h_1^2} := \max(\theta_{\delta h_1^1} \theta'_{W_0^2}, \theta'_{\delta W_1^2}), \quad (200)$$

$$\tilde{\mathbf{h}}_1^2 - \tilde{\mathbf{h}}_0^2 = \mathbf{G}(\mathbf{h}_1^2) - \mathbf{G}(\mathbf{h}_0^2) \quad (201)$$

$$= \mathbf{G} \left(\mathbf{h}_0^2 + \theta_{\delta h_1^2} \delta \mathbf{h}_1^2 \right) - \mathbf{G}(\mathbf{h}_0^2) \quad (202)$$

$$\theta_{\delta \tilde{h}_1^2} (\delta \tilde{\mathbf{h}}_1^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi_k \left(-; \{ (\mathbf{h}_0^2)_k, (\delta \mathbf{h}_1^2)_k \}_{k=1}^m \right), \quad (\text{Moment}), \quad (203)$$

$$\theta_{\delta \tilde{h}_1^2} := \theta_{\delta h_1^2}. \quad (204)$$

where G is either softmax or sigmoid and can be represented as a Moment (see Appendix G) (Yang and Hu, 2021).

$$\mathbf{h}_1^{3,k} - \mathbf{h}_0^{3,k} = \mathbf{W}_0^{3,k}(\mathbf{h}_1^1 - \mathbf{h}_0^1) + (\mathbf{W}_1^{3,k} - \mathbf{W}_0^{3,k})\mathbf{h}_1^1, \quad (205)$$

$$= \theta_{\delta h_1^1} \mathbf{W}_0^{3,k} \delta \mathbf{h}_1^1 - \theta'_{\delta W_1^3} \eta \frac{1}{n} Q_1^3 \left(n^{(d_3-1)} d\mathbf{h}_0^{3,k} \boldsymbol{\chi}_0(\mathbf{h}_0^1)^T \right) (\mathbf{h}_1^1), \quad (206)$$

$$\begin{aligned} (\delta \mathbf{h}_1^{3,k})_\alpha &= \frac{1}{n} \sum_{\gamma=1}^n \phi \left(\Upsilon_\alpha, \{(\delta \mathbf{h}_0^{3,k})_{\alpha l}\}_{l=1}^{d_{\text{out}}}; (\mathbf{h}_0^1)_\gamma, (\mathbf{h}_1^1)_\gamma \right. \\ &\quad \left. ; \theta_{\delta h_1^1 / \delta h_1^{3,k}}, \theta'_{\delta W_1^3 / \delta h_1^{3,k}}, n^{(d_3-1)}, \{(\boldsymbol{\chi}_0)_l\}_{l=1}^{d_{\text{out}}}, \eta \right), \end{aligned} \quad (\text{OuterNonlin}) \quad (207)$$

$$\Upsilon := \mathbf{W}_0^{3,k} \delta \mathbf{h}_1^1, \quad (\text{MatMul}) \quad (208)$$

$$\phi(\dots) := \theta_{\delta h_1^1 / \delta h_1^{3,k}} \Upsilon_\alpha - \theta'_{\delta W_1^3 / \delta h_1^{3,k}} \eta Q_1^3 \left(n^{(d_3-1)} \sum_{l=1}^{d_{\text{out}}} (\delta \mathbf{h}_0^{3,k})_{\alpha l} (\boldsymbol{\chi}_0)_l (\mathbf{h}_0^1)_\gamma \right) (\mathbf{h}_1^1)_\gamma, \quad (209)$$

$$\theta_{\delta h_1^{3,k}} = \max(\theta_{\delta h_1^1}, \theta'_{\delta W_1^3}), \quad (210)$$

$$\theta_{\delta h_1^1 / \delta h_1^{3,k}} := \theta_{\delta h_1^1} / \theta_{\delta h_1^{3,k}}, \quad (211)$$

$$\theta'_{\delta W_1^3 / \delta h_1^{3,k}} := \theta'_{\delta W_1^3} / \theta_{\delta h_1^{3,k}}. \quad (212)$$

We repeat the above process for each expert k to give $\{\delta \mathbf{h}_1^{3,k}\}_{k=1}^m$.

$$\mathbf{h}_1^3 - \mathbf{h}_0^3 = \sum_{k=1}^m \left(\tilde{h}_{0,k}^2 (\mathbf{h}_1^{3,k} - \mathbf{h}_0^{3,k}) + (\tilde{h}_{1,k}^2 - \tilde{h}_{0,k}^2) \mathbf{h}_1^{3,k} \right) \quad (213)$$

$$= \sum_{k=1}^m \left(\theta_{\delta h_1^{3,k}} \tilde{h}_{0,k}^2 \delta \mathbf{h}_1^{3,k} + \theta_{\delta \tilde{h}_1^2} \delta \tilde{h}_{1,k}^2 \mathbf{h}_1^{3,k} \right), \quad (214)$$

$$(\delta \mathbf{h}_1^3)_\alpha = \phi \left(\{(\delta \mathbf{h}_1^{3,k})_\alpha\}_{k=1}^m, \{(\mathbf{h}_1^{3,k})_\alpha\}_{k=1}^m; \theta_{\delta h_1^{3,k} / \delta h_1^3}, \theta_{\delta \tilde{h}_1^2 / \delta h_1^3}, \{\tilde{h}_{0,k}^2\}_{k=1}^m, \{\delta \tilde{h}_{1,k}^2\}_{k=1}^m \right), \quad (\text{Nonlin}) \quad (215)$$

$$\phi(\dots) := \sum_{k=1}^m \left(\theta_{\delta h_1^{3,k} / \delta h_1^3} \tilde{h}_{0,k}^2 (\delta \mathbf{h}_1^{3,k})_\alpha + \theta_{\delta \tilde{h}_1^2 / \delta h_1^3} \delta \tilde{h}_{1,k}^2 (\mathbf{h}_1^{3,k})_\alpha \right), \quad (216)$$

$$\theta_{\delta h_1^3} := \max(\theta_{\delta h_1^{3,k}}, \theta_{\delta \tilde{h}_1^2}), \quad (217)$$

$$\theta_{\delta h_1^{3,k} / \delta h_1^3} := \theta_{\delta h_1^{3,k}} / \theta_{\delta h_1^3}, \quad (218)$$

$$\theta_{\delta \tilde{h}_1^2 / \delta h_1^3} := \theta_{\delta \tilde{h}_1^2} / \theta_{\delta h_1^3}. \quad (219)$$

$$\mathbf{f}_1 - \mathbf{f}_0 = \mathbf{W}_0^4(\mathbf{h}_1^3 - \mathbf{h}_0^3) + (\mathbf{W}_1^4 - \mathbf{W}_0^4)\mathbf{h}_1^3 \quad (220)$$

$$= \theta'_{W_0^4} \theta_{\delta h_1^3} \frac{\overline{\mathbf{W}_0^4} \delta \mathbf{h}_1^3}{n} - \theta'_{\delta W_1^4} \eta \frac{1}{n} Q_1^4 \left(n^{d_4} \boldsymbol{\chi}_0 (\mathbf{h}_0^3)^T \right) \mathbf{h}_1^3, \quad (221)$$

$$\theta_{\delta f_1} (\delta \mathbf{f}_0)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_0^4}_{\beta\alpha}, (\delta \mathbf{h}_1^3)_\alpha, (\mathbf{h}_0^3)_\alpha, (\mathbf{h}_1^3)_\alpha; \theta'_{W_0^4}, \theta_{\delta h_1^3}, \theta'_{\delta W_1^4}, n^{d_4}, \eta, (\boldsymbol{\chi}_0)_\beta \right), \quad (\text{Moment}) \quad (222)$$

$$\phi(\dots) := \theta'_{W_0^4} \theta_{\delta h_1^3} (\overline{\mathbf{W}_0^4}_{\beta\alpha}) (\delta \mathbf{h}_1^3)_\alpha - \theta'_{\delta W_1^4} \eta Q_1^4 \left(n^{d_4} (\boldsymbol{\chi}_0)_\beta (\mathbf{h}_0^3)_\alpha \right) (\mathbf{h}_1^3)_\alpha, \quad (223)$$

$$\theta_{\delta f_1} := \max(\theta'_{W_0^4} \theta_{\delta h_1^3}, \theta'_{\delta W_1^4}). \quad (224)$$

***t*th Backward Pass:** Same as Appendix C.1 except for the gradient $d\mathbf{h}_t^1$.

$$\begin{aligned} d\mathbf{h}_t^1 &= \sum_{k=1}^m (\overline{\mathbf{W}_t^2})_{:,k}^T dh_{t,k}^2 + \sum_{k=1}^m (\mathbf{W}_0^{3,k})^T d\mathbf{h}_t^{3,k} \\ &\quad - \sum_{k=1}^m \sum_{s=0}^{t-1} \eta \frac{1}{n} Q_s^1 \left(n^{(d_3-1)} \mathbf{h}_s^1 (\boldsymbol{\chi}_s)^T (d\mathbf{h}_s^{3,k})^T \right) d\mathbf{h}_t^{3,k}, \quad \in \mathbb{R}^{n \times d_{\text{out}}}, \end{aligned} \quad (225)$$

$$\begin{aligned} (d\mathbf{h}_t^1)_{\alpha\beta} &= \frac{1}{n} \sum_{\gamma=1}^n \phi \left(\{(\overline{\mathbf{W}_t^2})_{\alpha k}^T, (\mathbf{V}_t^k)_{\alpha\beta}\}_{k=1}^m, \{(\mathbf{h}_s^1)_\alpha\}_{s=1}^{t-1}; \right. \\ &\quad \left. \{(\mathbf{h}_s^{3,k})_{\gamma\beta}\}_{s=1}^{t-1}, (d\mathbf{h}_t^{3,k})_{\gamma\beta}; n^{(d_3-1)}, \{(dh_{t,k}^2)_\beta\}_{k=1}^m, \{(\boldsymbol{\chi}_s)_\beta\}_{s=1}^{t-1} \right), \end{aligned} \quad (\text{Outer}) \quad (226)$$

$$(\mathbf{V}_t^k)_{:, \beta} := (\mathbf{W}_0^{3,k})^T (d\mathbf{h}_t^{3,k})_{:, \beta}, \quad (\text{MatMu}) \quad (227)$$

$$\phi(\dots) := \sum_{k=1}^m (\overline{\mathbf{W}_t^2})_{\alpha k}^T (dh_{t,k}^2)_\beta + \sum_{k=1}^m (\mathbf{V}_t^k)_{\alpha\beta} - \sum_{k=1}^m \sum_{s=0}^{t-1} \eta \frac{1}{n} Q_s^1 \left(n^{(d_3-1)} (\mathbf{h}_s^1)_\alpha (\boldsymbol{\chi}_s)_\beta (d\mathbf{h}_s^{3,k})_{\gamma\beta} \right) (d\mathbf{h}_t^{3,k})_{\gamma\beta}, \quad (228)$$

***t*th Weight Updates:**

$$\mathbf{W}_t^4 - \mathbf{W}_{t-1}^4 = -n^{-c_4} \eta Q_t^4 \left(n^{d_4} \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^3)^T \right) \quad (229)$$

$$\theta_{\delta W_t^4} := n^{-c_4}. \quad (230)$$

$$\mathbf{W}_t^{3,k} - \mathbf{W}_{t-1}^{3,k} = -n^{-c_3} \eta Q_t^3 \left(n^{d_3} \theta_{W_{t-1}^4} d\mathbf{h}_{t-1}^{3,k} \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^1)^T \right) \quad (231)$$

$$\theta_{\delta W_t^3} := n^{-c_3}. \quad (232)$$

$$\mathbf{W}_t^2 - \mathbf{W}_{t-1}^2 = -n^{-c_2} \eta Q_t^2 \left(n^{d_2} n \theta_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^1)^T \right) \quad (233)$$

$$\theta_{\delta W_t^2} := n^{-c_2}. \quad (234)$$

$$\mathbf{W}_t^1 - \mathbf{W}_{t-1}^1 = -n^{-c_1} \eta Q_t^1 \left(n^{d_1} n \theta_{W_{t-1}^4} \theta_{W_{t-1}^2} d\mathbf{h}_{t-1}^1 \boldsymbol{\chi}_{t-1} (\boldsymbol{\xi})^T \right) \quad (235)$$

$$\theta_{\delta W_t^1} := n^{-c_1}. \quad (236)$$

$$\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2 = -\eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1}, \quad (237)$$

*t*th Feature Updates:

$$\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1 = (\mathbf{W}_t^1 - \mathbf{W}_{t-1}^1) \boldsymbol{\xi} \quad (238)$$

$$= -\theta_{\delta W_t^1} \eta Q_t^1 \left(n^{d_1} d\mathbf{h}_{t-1}^1 \boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi} \right) \quad (239)$$

$$(\delta \mathbf{h}_t^1)_\alpha = \phi \left(\{(\mathbf{d}\mathbf{h}_{t-1}^1)_{\alpha\beta}\}_{\beta=1}^{d_{\text{out}}}; n^{d_1}, -\eta, \{(\boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi})_\beta\}_{\beta=1}^{d_{\text{out}}}\right) \quad (\text{Nonlin}) \quad (240)$$

$$\phi(\dots) := -\eta Q_t^1 \left(n^{d_1} \theta_{d\mathbf{h}_{t-1}^1} \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^1)_{\alpha\beta} (\boldsymbol{\chi}_{t-1} \boldsymbol{\xi}^T \boldsymbol{\xi})_\beta \right), \quad (241)$$

$$\theta_{\delta h_t^1} := \theta_{\delta W_t^1}, \quad (242)$$

$$\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2 = \mathbf{W}_{t-1}^2 (\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1) + (\mathbf{W}_t^2 - \mathbf{W}_{t-1}^2) \mathbf{h}_t^1 + (\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2) \quad (243)$$

$$= \theta'_{W_t^2} \theta_{\delta h_t^1} \frac{\overline{\mathbf{W}_{t-1}^2} \delta \mathbf{h}_t^1}{n} - \theta'_{\delta W_t^2} \eta \frac{1}{n} Q_t^2 \left(n^{d_2} \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^1)^T \right) \mathbf{h}_t^1 \quad (244)$$

$$- \eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1} \quad (245)$$

$$\theta_{h_t^2} (\delta \mathbf{h}_t^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_{t-1}^2}{}_{k\alpha}, (\delta \mathbf{h}_t^1)_\alpha, (\mathbf{h}_{t-1}^1)_\alpha, (\mathbf{h}_{t-1}^1)_\alpha; \right. \\ \left. \theta'_{W_t^2}, \theta_{\delta h_t^1}, \theta'_{\delta W_t^2}, n^{d_2}, \theta'_{W_{t-1}^4}, \eta, \{(\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta}, (\boldsymbol{\chi}_{t-1})_\beta\}_{\beta=1}^{d_{\text{out}}}\right) \quad (\text{Moment}) \quad (246)$$

$$\phi(\dots) := \theta'_{W_{t-1}^2} \theta_{\delta h_t^1} \overline{\mathbf{W}_{t-1}^2}{}_{k\alpha} (\delta \mathbf{h}_t^1)_\alpha \quad (247)$$

$$- \theta'_{\delta W_t^2} \eta Q_t^2 \left(n^{d_2} \theta'_{W_{t-1}^4} \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta} (\boldsymbol{\chi}_{t-1})_\beta (\mathbf{h}_{t-1}^1)_\alpha \right) (\mathbf{h}_{t-1}^1)_\alpha \quad (248)$$

$$- \eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \boldsymbol{\chi}_{t-1} \quad (249)$$

$$\theta_{\delta h_t^2} := \max(\theta'_{W_t^2} \theta_{\delta h_t^1}, \theta'_{\delta W_t^2}, \theta'_{W_{t-1}^4}), \quad (250)$$

$$\tilde{\mathbf{h}}_t^2 - \tilde{\mathbf{h}}_{t-1}^2 = \mathbf{G}(\mathbf{h}_t^2) - \mathbf{G}(\mathbf{h}_{t-1}^2) \quad (251)$$

$$= \mathbf{G} \left(\mathbf{h}_{t-1}^2 + \theta_{h_t^2} \delta \mathbf{h}_t^2 \right) - \text{Softmax}(\mathbf{h}_{t-1}^2) \quad (252)$$

$$\theta_{\tilde{h}_t^2} (\delta \tilde{\mathbf{h}}_t^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi_k \left(-; \{(\mathbf{h}_{t-1}^2)_k, (\delta \mathbf{h}_t^2)_k\}_{k=1}^m \right), \quad (\text{Moment}), \quad (253)$$

$$\theta_{\delta \tilde{h}_t^2} := \theta_{\delta h_t^2}. \quad (254)$$

$$\mathbf{h}_t^{3,k} - \mathbf{h}_{t-1}^{3,k} = \mathbf{W}_{t-1}^{3,k}(\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1) + (\mathbf{W}_t^{3,k} - \mathbf{W}_{t-1}^{3,k})\mathbf{h}_t^1, \quad (255)$$

$$= \theta_{\delta h_t^1} \mathbf{W}_0^{3,k} \overline{\delta \mathbf{h}_t^1} - \sum_{s=0}^{t-2} \theta'_{\delta W_s^3} \theta_{\delta h_t^1} \frac{1}{n} Q_t^3 \left(n^{d_3} \theta_{W_s^4} \eta d\mathbf{h}_s^{3,k} \boldsymbol{\chi}_s(\mathbf{h}_s^1)^T \right) \delta \mathbf{h}_t^1 \quad (256)$$

$$- \theta'_{\delta W_t^3} \frac{1}{n} Q_t^3 \left(n^{d_3} \theta_{W_{t-1}^4} \eta d\mathbf{h}_{t-1}^{3,k} \boldsymbol{\chi}_{t-1}(\mathbf{h}_{t-1}^1)^T \right) (\mathbf{h}_t^1), \quad (257)$$

$$\begin{aligned} (\delta \mathbf{h}_t^{3,k})_\alpha = & \frac{1}{n} \sum_{\gamma=1}^n \phi \left(\Upsilon_\alpha, \{ \{ (d\mathbf{h}_s^{3,k})_{\alpha l} \}_{l=1}^{d_{\text{out}}} \}_{s=0}^{t-2}, \{ (d\mathbf{h}_{t-1}^{3,k})_{\alpha l} \}_{l=1}^{d_{\text{out}}}; \right. \\ & \{ (\mathbf{h}_s^1)_{\gamma} \}_{s=0}^{t-2}, (\delta \mathbf{h}_t^1)_{\gamma}, (\mathbf{h}_{t-1}^1)_{\gamma}, (\mathbf{h}_t^1)_{\gamma}; \\ & \left. \theta_{\delta h_t^1 / \delta h_t^{3,k}}, \{ \theta'_{\delta W_s^3 \delta h_t^1 / \delta h_t^{3,k}} \}_{s=0}^{t-2}, \theta'_{\delta W_t^3 / \delta h_t^{3,k}}, n^{d_3}, \{ \theta_{W_s^4} \}_{s=0}^{t-1} \right) \end{aligned} \quad (\text{OuterNonlin}) \quad (258)$$

$$\Upsilon := \mathbf{W}_0^{3,k} \delta \mathbf{h}_1^1, \quad (\text{MatMul}) \quad (259)$$

$$\theta_{\delta h_t^{3,k}} = \max \left(\theta_{\delta h_t^1}, \max_s \left(\theta'_{\delta W_s^3} \theta_{\delta h_t^1} \right), \theta'_{\delta W_t^3} \right), \quad (260)$$

$$\theta_{\delta h_t^1 / \delta h_t^{3,k}} := \theta_{\delta h_t^1} / \theta_{\delta h_t^{3,k}}, \quad (261)$$

$$\theta'_{\delta W_s^3 \delta h_t^1 / \delta h_t^{3,k}} := \theta'_{\delta W_s^3} \theta_{\delta h_t^1} / \theta_{\delta h_t^{3,k}}, \quad (262)$$

$$\theta'_{\delta W_t^3 / \delta h_t^{3,k}} := \theta'_{\delta W_t^3} / \theta_{\delta h_t^{3,k}}. \quad (263)$$

We repeat the above process for each expert k .

$$\mathbf{h}_t^3 - \mathbf{h}_{t-1}^3 = \sum_{k=1}^m \left(\tilde{h}_{t-1,k}^2 (\mathbf{h}_t^{3,k} - \mathbf{h}_{t-1}^{3,k}) + (\tilde{h}_{t,k}^2 - \tilde{h}_{t-1,k}^2) \mathbf{h}_t^{3,k} \right) \quad (264)$$

$$= \sum_{k=1}^m \left(\theta_{\delta h_t^{3,k}} \tilde{h}_{t-1,k}^2 \delta \mathbf{h}_t^{3,k} + \theta_{\delta \tilde{h}_t^2} \delta \tilde{h}_{t,k}^2 \mathbf{h}_t^{3,k} \right), \quad (265)$$

$$(\delta \mathbf{h}_t^3)_\alpha = \phi \left(\{ (\delta \mathbf{h}_t^{3,k})_\alpha \}_{k=1}^m, \{ (\mathbf{h}_t^{3,k})_\alpha \}_{k=1}^m; \theta_{\delta h_t^{3,k} / \delta h_t^3}, \theta_{\delta \tilde{h}_t^2 / \delta h_t^3}, \{ \tilde{h}_{t-1,k}^2 \}_{k=1}^m, \{ \delta \tilde{h}_{t,k}^2 \}_{k=1}^m \right), \quad (\text{Nonlin}) \quad (266)$$

$$\phi(\dots) := \sum_{k=1}^m \left(\theta_{\delta h_t^{3,k}} \tilde{h}_{t-1,k}^2 (\delta \mathbf{h}_t^{3,k})_\alpha + \theta_{\delta \tilde{h}_t^2} \delta \tilde{h}_{t,k}^2 (\mathbf{h}_t^{3,k})_\alpha \right), \quad (267)$$

$$\theta_{\delta h_t^3} := \max(\theta_{\delta h_t^{3,k}}, \theta_{\delta \tilde{h}_t^2}), \quad (268)$$

$$\theta_{\delta h_t^{3,k} / \delta h_t^3} := \theta_{\delta h_t^{3,k}} / \theta_{\delta h_t^3}, \quad (269)$$

$$\theta_{\delta \tilde{h}_t^2 / \delta h_t^3} := \theta_{\delta \tilde{h}_t^2} / \theta_{\delta h_t^3}. \quad (270)$$

$$\mathbf{f}_t - \mathbf{f}_{t-1} = \mathbf{W}_{t-1}^4 (\mathbf{h}_t^3 - \mathbf{h}_{t-1}^3) + (\mathbf{W}_t^4 - \mathbf{W}_{t-1}^4) \mathbf{h}_t^3 \quad (271)$$

$$= \theta'_{W_{t-1}^4} \theta_{\delta h_t^3} \frac{\overline{\mathbf{W}_{t-1}^4} \delta \mathbf{h}_t^3}{n} - \theta'_{\delta W_t^4} \eta \frac{1}{n} Q_t^4 \left(n^{d_4} \boldsymbol{\chi}_{t-1} (\mathbf{h}_{t-1}^3)^T \right) \mathbf{h}_t^3, \quad (272)$$

$$\theta_{f_t} (\delta \mathbf{f}_{t-1})_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_{t-1}^4} \beta_\alpha, (\delta \mathbf{h}_t^3)_\alpha, (\mathbf{h}_{t-1}^3)_\alpha, (\mathbf{h}_t^3)_\alpha; \theta'_{W_{t-1}^4}, \theta_{\delta h_t^3}, \theta'_{\delta W_t^4}, n^{d_4}, \eta, (\boldsymbol{\chi}_{t-1})_\beta \right), \quad (\text{Moment}) \quad (273)$$

$$\phi(\dots) := \theta'_{W_{t-1}^4} \theta_{\delta h_t^3} (\overline{\mathbf{W}_{t-1}^4})_{\beta\alpha} (\delta \mathbf{h}_t^3)_\alpha - \theta'_{\delta W_t^4} \eta Q_t^4 \left(n^{d_4} (\boldsymbol{\chi}_{t-1})_\beta (\mathbf{h}_{t-1}^3)_\alpha \right) (\mathbf{h}_t^3)_\alpha, \quad (274)$$

$$\theta_{\delta f_t} := \max(\theta'_{W_{t-1}^4} \theta_{\delta h_t^3}, \theta'_{\delta W_t^4}). \quad (275)$$

C.4. Infinite Width Limit for Soft MoE with Adam

See Appendix C.2 for an explanation of the limit notation.

Adam requires defining the limit of an extra operation, `OuterNonlin` (Yang and Littwin, 2023):

Z OuterNonlin: Given $r \geq 0$ and function $\psi : \mathbb{R}^{(r+1)+\ell} \rightarrow \mathbb{R}$ such that $y_\alpha = \frac{1}{n^r} \sum_{\beta_1, \dots, \beta_r=1}^n \psi(x_\alpha^0, x_{\beta_1}^1; \dots; x_{\beta_r}^r; \theta_1, \dots, \theta_\ell)$ then

$$Z^y := \mathbb{E}_{Z_1^{x^1}, \dots, Z_1^{x^r}} \left[\psi \left(Z^{x^0}; Z_1^{x^1}; \dots; Z_1^{x^r}; \theta_1, \dots, \theta_\ell \right) \right], \quad (276)$$

where $(Z_1^{x^1}, \dots, Z_1^{x^r}) \stackrel{d}{=} (Z^{x^1}, \dots, Z^{x^r})$ is an iid copy of $(Z^{x^1}, \dots, Z^{x^r})$ and is independent of Z^{x^0} .

Assumption 14 (Pseudo Lipschitz Q, Assumption 2.3.2 (Yang and Littwin, 2023))

Assume that Q_t^ℓ, ϵ are pseudo-Lipschitz for all ℓ, t .

Assumption 15 (Faithfulness at initialisation) If we want faithfulness in the first weight updates, it already restricts the values of d to $d_1 = 1, d_2 = 0, d_3 = 1, d_4 = 0$.

First Forward Pass:

$$Z^{h_0^1} = \sum_{j=1}^{d_{\text{in}}} \xi_j Z^{W_{0,j}^1}, \quad (277)$$

$$\mathring{\mathbf{h}}_0^2 = \mathbf{b}_0^2, \quad (278)$$

$$\mathring{\mathbf{h}}_0^2 = G(\mathbf{b}_0^2) \quad (279)$$

$$Z^{h_0^{3,k}} = Z^{W_0^{3,k} h_0^1} \text{ for every expert } k, \quad (280)$$

$$Z^{h_0^3} = \sum_{k=1}^m \mathring{h}_{0,k}^2 Z^{h_0^{3,k}} \quad (281)$$

$$\mathring{\mathbf{f}}_0 = \mathbf{0}. \quad (282)$$

Here, $Z^{W_0^1}$ is standard Gaussian by construction. $\mathring{\mathbf{h}}_0^2 = \mathbf{0}$ as $\mathring{h}_0^2 = \mathbb{E}[Z^{W_0^2}] \mathbb{E}[Z^{h_0^1}] = 0$. $\mathring{\mathbf{h}}_0^2$ then follows by the definition of `Softmax(0)`. $Z^{W_0^{3,k} h_0^1} = \hat{Z}^{W_0^{3,k} h_0^1}$ from the Z MatMul rule. Finally $\mathring{\mathbf{f}}_0 = \mathbf{0}$ as $\mathring{f}_0 = \mathbb{E}[Z^{W_0^4}] \mathbb{E}[Z^{h_0^3}] = 0$.

First Backward Pass:

$$Z^{\text{dh}_0^3} = Z^{\overline{W_0^4}}, \quad (283)$$

$$\mathring{d}\tilde{h}_{0,k}^2 = 0 \text{ for every expert } k, \quad (284)$$

$$\mathring{d}h_{0,k}^2 = 0 \text{ for every expert } k, \quad (285)$$

$$Z^{\text{dh}_0^{3,k}} = \mathring{\mathbf{h}}_{0,k}^2 Z^{\text{dh}_0^3} \text{ for every expert } k, \quad (286)$$

$$Z^{\text{dh}_0^1} = \sum_{k=1}^m Z^{(W_0^{3,k})^T \text{dh}_0^{3,k}} \quad (287)$$

$$(288)$$

$Z^{\overline{W_0^4}}$ is a standard Gaussian by construction. $\mathring{d}\tilde{h}_{0,k}^2 = \mathbb{E}[Z^{h_0^{3,k}}] \mathbb{E}[\overline{W_0^4}] = 0$. From this, $\mathring{d}h_{0,k}^2 = 0$ follows. $Z^{(W_0^{3,k})^T \text{dh}_0^3} = \hat{Z}^{(W_0^{3,k})^T \text{dh}_0^3}$ from the Z MatMul rule.

First Feature Update:

$$Z^{\delta h_1^1} = -\eta Q_1^1 \left(\sum_{\beta=1}^{d_{\text{out}}} Z^{(\text{dh}_0^1)_\beta} (\mathring{\chi}_0)_\beta \boldsymbol{\xi}^T \right) \boldsymbol{\xi}, \quad (289)$$

$$\delta \mathring{h}_{1,k}^2 = \mathring{\theta}_{\delta h_1^1} \mathbb{E}[Z^{\overline{W_0^2}}]_{k,:} Z^{\delta h_1^1} - \mathring{\theta}'_{\delta W_1^2} \eta \mathbb{E} \left[Q_1^2 \left(\sum_{\beta=1}^{d_{\text{out}}} (\text{dh}_{0,k}^2)_\beta (\mathring{\chi}_0)_\beta Z^{h_0^1} \right) Z^{h_1^1} \right], \quad (290)$$

$$\delta \mathring{\mathbf{h}}_1^2 = \begin{cases} \mathbf{G}(\mathring{\mathbf{h}}_1^2) - \mathbf{G}(\mathring{\mathbf{h}}_0^2) & \text{if } \mathring{\theta}_{\delta h_1^2} = 1 \\ \mathbf{G}'(\mathring{\mathbf{h}}_0^2)(\delta \mathring{\mathbf{h}}_0^2) & \text{if } \mathring{\theta}_{\delta h_1^2} = 0 \end{cases}, \quad (291)$$

$$\begin{aligned} Z^{\delta h_1^{3,k}} &= \mathring{\theta}_{\delta h_1^1 / \delta h_1^{3,k}} Z^{W_0^{3,k} \delta h_1^1} \\ &\quad - \mathring{\theta}'_{\delta W_1^3 / \delta h_1^{3,k}} \eta \mathbb{E} \left(Z_1^{h_0^1}, Z_1^{h_1^1} \right) \left[Q_1^3 \left(\sum_{\beta=1}^{d_{\text{out}}} Z^{(\text{dh}^{3,k})_\beta} (\mathring{\chi}_0)_\beta Z_1^{h_0^1} \right) Z_1^{h_1^1} \right], \text{ for every expert } k, \end{aligned} \quad (292)$$

$$Z^{\delta h_1^3} = \sum_{k=1}^m \mathring{\theta}_{\delta h_1^{3,k} / \delta h_1^3} \mathring{h}_{0,k}^2 Z^{\delta h_1^{3,k}} - \mathring{\theta}_{\delta \tilde{h}_1^2} \mathring{\theta}_{h_1^{3,k} / \delta h^3} \delta \mathring{h}_{1,k}^2 Z^{h_1^{3,k}}, \quad (293)$$

$$\delta \mathring{f}_1 = \mathring{\theta}_{\delta h_1^3} \mathbb{E}[Z^{\overline{W_0^4}} Z^{\delta h_1^3}] - \mathring{\theta}'_{\delta W_1^4} \eta \mathbb{E} \left[Q_1^4 \left(\mathring{\chi}_0 Z^{h_0^3} \right) \right] Z^{h_1^3}, \quad (294)$$

where $Z^{W_0^{3,k} \delta h_1^1}$ is given by the Z MatMul rule,

$$Z^{W_0^{3,k} \delta h_1^1} = \hat{Z}^{W_0^{3,k} \delta h_1^1} + \dot{Z}^{W_0^{3,k} \delta h_1^1} \quad (295)$$

$$= \hat{Z}^{W_0^{3,k} \delta h_1^1} + Z^{\text{dh}_0^{3,k}} \mathbb{E} \left[\frac{\partial Z^{\delta h_1^1}}{\partial \hat{Z}^{(W_0^{3,k})^T \text{dh}_0^{3,k}}} \right]. \quad (296)$$

Furthermore, the random variables $(Z_1^{h_0^1}, Z_1^{h_1^1}) \stackrel{d}{=} (Z^{h_0^1}, Z^{h_1^1})$ is an iid copy and is independent from $Z^{\text{dh}^{3,k}}$.

We write the limit of each vector in the program as

$$Z^{h_t} = Z^{h_0} + \overset{\circ}{\theta}_{\delta h_1} Z^{\delta h_1} + \dots + \overset{\circ}{\theta}_{\delta h_t} Z^{\delta h_t}. \quad (297)$$

Weight updates: The distribution of certain weight updates are required in the backward pass and thus the forward passes.

$$Z^{\overline{\delta W_t^4}} = -Q_t^4 \left(\eta \overset{\circ}{\chi}_{t-1} Z^{h_{t-1}^3} \right), \quad (298)$$

$$Z^{(\overline{\delta W_t^2})_{k,:}} = -Q_t^2 \left(\eta \sum_{\beta=1}^{d_{\text{out}}} (\overset{\circ}{d}h_{t,k}^2)_{\beta} (\overset{\circ}{\chi}_{t-1})_{\beta} Z^{h_{t-1}^1} \right). \quad (299)$$

t th Backward Pass:

$$Z^{\text{d}h_t^3} = Z^{(\overline{W_t^4})^T}, \quad (300)$$

$$\overset{\circ}{d}h_{t,k}^2 = \mathbb{E}[Z^{h_t^{3,k}} Z^{\text{d}h_t^3}], \quad (301)$$

$$\overset{\circ}{d}h_{t,k}^2 = \sum_{j=1}^m J(\overset{\circ}{h}_t^2)_{kj} \overset{\circ}{d}h_{t,k}^2, \quad (302)$$

$$Z^{\text{d}h_t^{3,k}} = \overset{\circ}{h}_{t,k}^2 Z^{\text{d}h_t^3}, \quad (303)$$

$$Z^{\text{d}h_t^1} = \sum_{k=1}^m Z^{(\overline{W_t^2})_{k,:}} \overset{\circ}{d}h_{t,k}^2 + Z^{(W_0^{3,k})^T} \text{d}h_t^{3,k} - \eta \sum_{s=1}^{t-1} \mathbb{E} \left(Z_1^{\text{d}h_s^{3,k}}, Z_1^{\text{d}h_t^{3,k}} \right) \left[Q_s^1 \left(Z^{h_s^1} \overset{\circ}{\chi}_s Z_1^{\text{d}h_s^{3,k}} \right) Z_1^{\text{d}h_t^{3,k}} \right], \quad (304)$$

where $Z^{(W_0^{3,k})^T} \text{d}h_t^{3,k}$ is given by the Z MatMul rule,

$$Z^{(W_0^{3,k})^T} \text{d}h_t^{3,k} = \hat{Z}^{(W_0^{3,k})^T} \text{d}h_t^{3,k} + \dot{Z}^{(W_0^{3,k})^T} \text{d}h_t^{3,k} \quad (305)$$

$$= \hat{Z}^{(W_0^{3,k})^T} \text{d}h_t^{3,k} + \sum_{v \in \mathcal{V}: W_0^{3,k} v \in \mathcal{V}} Z^v \mathbb{E} \left[\frac{\partial Z^{\text{d}h_t^{3,k}}}{\partial \hat{Z}^{W_0^{3,k} v}} \right], \quad (306)$$

where $\mathcal{V} := \{v \text{ for some vector in the program } v\}$ is the set of all vectors in the program. Further the random variables $(Z_1^{\text{d}h_s^{3,k}}, Z_1^{\text{d}h_t^{3,k}}) \stackrel{d}{=} (Z^{\text{d}h_s^{3,k}}, Z^{\text{d}h_t^{3,k}})$ is an iid copy and is independent of $Z^{h_s^1}$.

t th Forward Pass:

$$Z^{\delta h_t^1} = -\eta Q_t^1 \left(\sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_{t-1}^1)_\beta} (\hat{\chi}_{t-1})_\beta \xi^T \right) \xi, \quad (307)$$

$$\delta \hat{h}_{t,k}^2 = \hat{\theta}'_{W_{t-1}^2} \hat{\theta}_{\delta h_t^1} \mathbb{E}[Z^{(\overline{W_{t-1}^2})_k} Z^{\delta h_t^1}] - \hat{\theta}'_{\delta W_t^2} \eta \mathbb{E} \left[Q_t^2 \left(\sum_{\beta=1}^{d_{\text{out}}} (d\hat{h}_{t-1,k}^2)_\beta (\hat{\chi}_{t-1})_\beta Z^{h_{t-1}^1} \right) Z^{h_t^1} \right], \quad (308)$$

$$\delta \hat{\mathbf{h}}_t^2 = \begin{cases} \mathbf{G}(\hat{\mathbf{h}}_t^2) - \mathbf{G}(\hat{\mathbf{h}}_{t-1}^2) & \text{if } \hat{\theta}_{\delta h_t^2} = 1 \\ \mathbf{G}'(\hat{\mathbf{h}}_{t-1}^2)(\delta \hat{\mathbf{h}}_{t-1}^2) & \text{if } \hat{\theta}_{\delta h_t^2} = 0 \end{cases}, \quad (309)$$

$$Z^{\delta h_t^{3,k}} = \hat{\theta}_{\delta h_t^1 / \delta h_t^{3,k}} Z^{W_0^{3,k} \delta h_t^1} \quad (310)$$

$$- \sum_{s=0}^{t-2} \hat{\theta}'_{\delta W_s^{3,k} \delta h_t^1 / \delta h_t^{3,k}} \eta \mathbb{E} \left(Z_1^{h_s^1}, Z_1^{\delta h_t^1} \right) \left[Q_t^3 \left(\sum_{\beta}^{d_{\text{out}}} Z^{(dh_s^{3,k})_\beta} (\hat{\chi}_s)_\beta Z_1^{h_s^1} \right) Z_1^{\delta h_t^1} \right] \quad (311)$$

$$- \hat{\theta}'_{\delta W_t^{3,k} / \delta h_t^{3,k}} \eta \mathbb{E} \left(Z_1^{h_{t-1}^1}, Z_1^{\delta h_t^1} \right) \left[Q_t^3 \left(\sum_{\beta}^{d_{\text{out}}} Z^{(dh_{t-1}^{3,k})_\beta} (\hat{\chi}_{t-1})_\beta Z_1^{h_{t-1}^1} \right) Z_1^{\delta h_t^1} \right], \quad (312)$$

$$Z^{\delta h_t^3} = \sum_{k=1}^m \hat{\theta}_{\delta h_t^{3,k} / \delta h_t^3} \hat{h}_{k,t-1}^2 Z^{\delta h_t^{3,k}} - \hat{\theta}_{\delta \hat{h}_t^2 / \delta h_t^3} \delta \hat{h}_{t,k}^2 Z^{h_t^{3,k}}, \quad (313)$$

$$\delta \hat{f}_t = \hat{\theta}'_{W_{t-1}^4} \hat{\theta}_{\delta h_t^3} \mathbb{E}[Z^{\overline{W_{t-1}^4}} Z^{\delta h_t^3}] - \hat{\theta}'_{\delta W_t^4} \eta \mathbb{E} \left[Q_t^4 \left(\hat{\chi}_{t-1} Z^{h_{t-1}^3} \right) Z^{h_t^3} \right], \quad (314)$$

where

$$Z^{W_0^{3,k} \delta h_t^1} = \hat{Z}^{W_0^{3,k} \delta h_t^1} + \dot{Z}^{W_0^{3,k} \delta h_t^1} \quad (315)$$

$$= \hat{Z}^{W_0^{3,k} \delta h_t^1} + \sum_{s=0}^{t-1} Z^{dh_s^1} \mathbb{E} \left[\frac{\partial Z^{dh_s^{3,k}}}{\partial \hat{Z}^{(W_0^{3,k})^T dh_s^{3,k}}} \right] \quad (316)$$

Appendix D. Derivation of Soft MoE Results

We derive the results from Section 3 here. The results rely on the Tensor Program derivation in Appendix C.

D.1. SP: Router Collapse

In this section, we derive the results for the soft MoE network with SP parametrisation. We use the Tensor Program derivations in Appendix C.1 and Appendix C.3 to do this.

Recall that the SP parametrisation refers to initialisation parameters $b_1 = 0$ and $b_\ell = 0.5$ for all layers $\ell \neq 1$. Furthermore, the learning rate and epsilon scaling parameters are the same for each layer $c = c_\ell$ and $d = d_\ell$ for all layers ℓ . The initial values of b_ℓ set the scalars

for the network in Theorem 1

$$\theta_{W_0^1} = 1, \quad (317)$$

$$\theta_{W_0^2} = n^{-0.5}, \quad (318)$$

$$\theta_{W_0^{3,k}} = 1, \quad (319)$$

$$\theta_{W_0^4} = n^{-0.5}. \quad (320)$$

$$(321)$$

We write the Tensor Program for the first forward pass with SP as follows (the rest of the tensor program is similar to Appendix C.1). We write this without the bias as it does not affect the result.

First Forward Pass:

$$\mathbf{h}_0^1 = \overline{\mathbf{W}}_0^1 \boldsymbol{\xi}, \quad (\mathbf{h}_0^1)_\alpha = \phi\left(\left(\overline{\mathbf{W}}_0^1\right)_{\alpha 1}, \dots, \left(\overline{\mathbf{W}}_0^1\right)_{\alpha d_{\text{in}}}; \xi_1, \dots, \xi_{d_{\text{in}}}\right) = \sum_{j=1}^{d_{\text{in}}} \left(\overline{\mathbf{W}}_0^1\right)_{\alpha j} \xi_j, \quad (\text{Nonlin})$$

$$\mathbf{h}_0^2 = \sqrt{n} \overline{\mathbf{W}}_0^2 \mathbf{h}_0^1, \quad (\mathbf{h}_0^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi\left(\left(\overline{\mathbf{W}}_0^2\right)_{k\alpha}, (\mathbf{h}_0^1)_\alpha; \sqrt{n}\right) = \frac{1}{n} \sqrt{n} \sum_{\alpha=1}^n \left(\overline{\mathbf{W}}_0^2\right)_{k\alpha} (\mathbf{h}_0^1)_\alpha, \quad (\text{Moment})$$

$$\tilde{\mathbf{h}}_0^2 = G(\mathbf{h}_0^2), \quad (\text{Moment, See G})$$

$$\hat{\mathbf{h}}_0^{3,k} = \overline{\mathbf{W}}_0^{3,k} \mathbf{h}_0^1, \quad (\text{MatMul})$$

$$\mathbf{h}_0^{3,k} = \hat{\mathbf{h}}_0^{3,k}, \quad (\mathbf{h}_0^{3,k})_\alpha = \phi\left(\left(\hat{\mathbf{h}}_0^{3,k}\right)_\alpha; -\right) = \left(\hat{\mathbf{h}}_0^3\right)_\alpha \quad (\text{Nonlin})$$

$$\mathbf{h}_0^3 = \sum_{k=1}^m \mathbf{h}_0^{3,k} \tilde{h}_{0,k}^2, \quad (\mathbf{h}_0^3)_\alpha = \phi\left(\left(\mathbf{h}_0^{3,1}\right)_\alpha, \dots, \left(\mathbf{h}_0^{3,m}\right)_\alpha; \tilde{h}_{0,1}^2, \dots, \tilde{h}_{0,m}^2\right) = \sum_{k=1}^m \left(\mathbf{h}_0^{3,k}\right)_\alpha \tilde{h}_{0,k}^2 \quad (\text{Nonlin})$$

$$\mathbf{f}_0 = \sqrt{n} \frac{1}{n} \overline{\mathbf{W}}_0^4 \mathbf{h}_0^3, \quad (\mathbf{f}_0)_\beta = \frac{1}{n} \sum_{\alpha=1}^n \phi\left(\left(\overline{\mathbf{W}}_0^4\right)_{\beta\alpha}, (\mathbf{h}_0^3)_\alpha; \sqrt{n}\right) = \frac{1}{n} \sqrt{n} \sum_{\alpha=1}^n \left(\overline{\mathbf{W}}_0^4\right)_{\beta\alpha} (\mathbf{h}_0^3)_\alpha, \quad (\text{Moment})$$

where $G(\cdot)$ is the gating function, either Softmax or Sigmoid..

To find the scale of the output-like layers \mathbf{h}^2 and \mathbf{f} we use Master's theorem on the variance of the coordinates

$$\mathbb{E}\left[(\mathbf{h}_0^2)_k\right]^2 = \sum_{\alpha=1}^n \frac{1}{n} \mathbb{E}\left[\left(\overline{\mathbf{W}}_0^2\right)_{k\alpha}^2\right] \mathbb{E}\left[(\mathbf{h}_0^1)_\alpha\right]^2 = \mathbb{E}\left[\frac{1}{n} \sum_{\alpha=1}^n (\mathbf{h}_0^1)_\alpha^2\right], \quad (\text{Moment}) \quad (322)$$

$$\mathbb{E}\left[(\mathbf{f}_0)_\beta\right]^2 = \sum_{\alpha=1}^n \frac{1}{n} \mathbb{E}\left[\left(\overline{\mathbf{W}}_0^4\right)_{\beta\alpha}^2\right] \mathbb{E}\left[(\mathbf{h}_0^3)_\alpha\right]^2 = \mathbb{E}\left[\frac{1}{n} \sum_{\alpha=1}^n (\mathbf{h}_0^3)_\alpha^2\right]. \quad (\text{Moment}) \quad (323)$$

We consider SGD updates here for ease of exposition. The result trivially holds for Adam updates as well (see Theorem 18).

Using the Master’s theorem (Yang and Hu, 2021), the coordinates of each layer will have scale $\Theta(1)$,

$$\theta_{h_0^1} = 1, \quad (324)$$

$$\theta_{h_0^2} = 1, \quad (325)$$

$$\theta_{\tilde{h}_0^2} = 1, \quad (326)$$

$$\theta_{h_0^{3,k}} = 1, \quad (327)$$

$$\theta_{h_0^3} = 1, \quad (328)$$

$$\theta_{f_0} = 1. \quad (329)$$

Without loss of generality, we assume that values of c is chosen such that $\delta \mathbf{h}_1^1$ and $\mathbf{h}_1^1 = \mathbf{h}_0^1 + \delta \mathbf{h}_1^1$ are $\Theta(1)$. That is $\theta_{\delta h_1^1} = \theta_{h_1^1} = 1$. It is safe to make this assumption as $\delta \mathbf{h}_1^1$ exploding or vanishing with width are both undesirable. The update to \mathbf{h}^2 then has the form (Appendix C.1)

$$\delta \mathbf{h}_1^2 = \theta'_{W_0^2} \frac{\overline{\mathbf{W}_0^2} \delta \mathbf{h}_1^1}{n} - \theta'_{\delta W_1^2} \eta d\mathbf{h}_0^2 \boldsymbol{\chi}_0 \frac{(\mathbf{h}_0^1)^T \mathbf{h}_1^1}{n}, \quad (330)$$

where $\overline{\mathbf{W}_0^2}, \delta \mathbf{h}_1^1$ are the unscaled ($\Theta(1)$) terms. The above operation is written as a Moment in the Tensor Program. By the TP Master Theorem, $\frac{\overline{\mathbf{W}_0^2} \delta \mathbf{h}_1^1}{n} \rightarrow \mathbb{E}[Z^{\overline{W_0^2}} Z^{\delta h_1^1}]$ (Appendix C.1). As the gradients required to calculate δh_1^1 depend on W_0^2 , the limiting random variables $Z^{\overline{W_0^2}}, Z^{\delta h_1^1}$ are correlated, thus the expectation is non-zero and the term scales $\Theta(1)$. Thus, the first term has the scale $\theta'_{W_0^2} = n^{0.5}$, and the second scales $\theta'_{\delta W_1^2} = n^{1-c} = n^{1-c}$. Hence, the first term scales $\Theta(n^{0.5})$ irrespective of the scaling of the second term, which means that the entire term scales $\Omega(n^{0.5})$. As $\mathbf{h}_1^2 = \delta \mathbf{h}_1^2 + \mathbf{h}_0^2$, where the first term scales $\Omega(n^{0.5})$, and the second term scales $\Theta(1)$, hence \mathbf{h}_1^2 scales $\Omega(n^{0.5})$.

Saturation of Softmax $\tilde{\mathbf{h}}^2$: In the softmax network, $\tilde{\mathbf{h}}_1^2 = \text{Softmax}(\mathbf{h}_1^2)$ where \mathbf{h}_1^2 scales $\Omega(n^{0.5})$. Using the well known property that softmax of diverging inputs results in a one-hot distribution, we get,

$$\lim_{n \rightarrow \infty} \frac{e^{(h_i^2)_k}}{\sum_{j=1}^m e^{(h_i^2)_j}} \rightarrow 1,$$

for some expert index $k \in [m]$. The gradient of the softmax layer is

$$\frac{\partial (\tilde{\mathbf{h}}_1^2)_i}{\partial (\mathbf{h}_1^2)_j} = (\tilde{h}_1^2)_i (\delta_{ij} - (\tilde{h}_1^2)_j),$$

where $(\tilde{\mathbf{h}}_1^2)_i$ is the i^{th} component of the vector $\tilde{\mathbf{h}}_1^2$, and δ_{ij} is the Kronecker delta function. If $\tilde{\mathbf{h}}_1^2$ is one-hot, then the gradient $\frac{\partial \tilde{\mathbf{h}}_1^2}{\partial \mathbf{h}_1^2}$ is a zero matrix of size $m \times m$. As $\frac{\partial \tilde{\mathbf{h}}_1^2}{\partial \mathbf{h}_1^2}$ is required to calculate the gradient of the router $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^2}$,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_1^2} = \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{h}}_1^2} \frac{\partial \tilde{\mathbf{h}}_1^2}{\partial \mathbf{h}_1^2} \frac{\partial \mathbf{h}_1^2}{\partial \mathbf{W}_1^2}$$

this implies that the router gradient will approach a matrix of zeroes as the width is increased. Hence $\mathbf{W}_t^2 = \mathbf{W}_1^2$ for all subsequent steps of training.

Saturation of Sigmoid $\tilde{\mathbf{h}}^2$: In the sigmoid network, $\tilde{\mathbf{h}}_1^2 = \sigma(\mathbf{h}_1^2)$ where \mathbf{h}_1^2 scales $\Omega(n^{0.5})$. With n increasing we get

$$\lim_{n \rightarrow \infty} \sigma((h_t^2)_k) \rightarrow \begin{cases} 1 & \text{if } (h_t^2)_k > 0 \\ 0 & \text{if } (h_t^2)_k < 0 \end{cases},$$

for every expert $k \in [m]$. The gradient of the sigmoid layer is

$$\frac{\partial(\tilde{\mathbf{h}}_1^2)_i}{\partial(\mathbf{h}_1^2)_i} = \sigma((\mathbf{h}_1^2)_i) (1 - \sigma((\mathbf{h}_1^2)_i)),$$

which approaches a zero matrix as n increases. Similarly to the softmax case, this results in a vanishing gradient for the router.

Remark 16 (Zero initialisation of router with SGD) *Zero intitialisation of the router with SGD ensures that the interaction between \mathbf{W}_0^2 and $\delta\mathbf{h}_1^1$ does not occur. However, here the results from Yang and Hu (2021) hold, primarily that it is not possible for the “output-like” layers (router and function) to feature learn. This can easily be seen by the fact that $\theta_{\delta\mathbf{h}_1^1} = n^{-(c+1)}$, which means that $c = -1$ is required for \mathbf{h}_1^1 to feature learn. This however causes the second term (effective updates) in $\delta\mathbf{h}_1^2$ to diverge as $\theta'_{\delta\mathbf{W}_1^2} = n \cdot n^{-c} \cdot n \cdot \theta_{\mathbf{W}_0^2} = n^3$.*

Remark 17 (Results hold for Top-K) *The router logits diverge due to the interaction between \mathbf{W}_0^2 and $\delta\mathbf{h}_1^1$, and thus holds irrespective of the gating function, or sparse execution. If Softmax or Sigmoid activations are used with Top-K, the router probabilities will still saturate with SP.*

Remark 18 (Router diverges with Adam) *The analysis above relies on the propagating updates $\mathbf{W}_0^2\delta\mathbf{h}_1^1$ of the router blowing up. Adam updates only affect the effective updates of the router $\delta\mathbf{W}_1^2\mathbf{h}_1^2$. Hence, the same analysis holds for Adam.*

D.2. μ P-MoE: Stability but with Router Stagnation without Router Bias

Table 1: *bcd* parametrisation vlaue for μ P-MoE SGD

Layer (ℓ)	\mathbf{b}_ℓ	\mathbf{c}_ℓ
1	0	-1
2	1	1
3	0.5	0
4	1	1

Stability and Feature Learning for SGD: We show that the values in Table 1 for the network in Theorem 1 satisfy our desired properties of stability and feature learning for SGD. The values in Theorem 11 already satisfy stability. To show feature learning, we show that

the limiting values of the scalars $\hat{\theta}_\bullet$ in the feature updates converge to one. Stability ($\mathcal{O}(1)$ features) combined with $\Theta(1)$ feature updates then imply that all the features scale $\Theta(1)$.

$$\hat{\theta}_{\delta h_1^1} = 1, \quad \hat{\theta}'_{\delta W_1^2} = 1, \quad (331)$$

$$\hat{\theta}_{\delta h_1^2} = 1, \quad (332)$$

$$\hat{\theta}_{\delta h_1^1 / \delta h_1^{3,k}} = 1, \quad \hat{\theta}'_{\delta W_1^3 / \delta h_1^{3,k}} = 1, \quad (333)$$

$$\hat{\theta}_{\delta h_1^{3,k} / \delta h_1^3} = 1, \quad \hat{\theta}_{\delta \tilde{h}_1^2 / \delta h_1^3} = 1, \quad (334)$$

$$\hat{\theta}_{\delta h_1^3} = 1, \quad \hat{\theta}'_{\delta W_1^4} = 1. \quad (335)$$

For the next forward passes we have $\theta_{W_t^2} = \theta_{W_0^2}$ and $\theta_{W_t^4} = \theta_{W_0^4}$, thus we have

$$\hat{\theta}'_{W_{t-1}^2} \hat{\theta}_{\delta h_t^1} = 1, \quad \hat{\theta}'_{\delta W_t^2} = 1, \quad (336)$$

$$\hat{\theta}_{\delta h_t^2} = 1, \quad (337)$$

$$\hat{\theta}_{\delta h_t^1 / \delta h_t^{3,k}} = 1, \quad \hat{\theta}'_{\delta W_t^3 \delta h_t^1 / \delta h_t^{3,k}} = 1, \quad \hat{\theta}'_{\delta W_t^3 / \delta h_t^{3,k}} = 1, \quad (338)$$

$$\hat{\theta}_{\delta h_t^{3,k} / \delta h_t^3} = 1, \quad \hat{\theta}_{\delta \tilde{h}_t^2 / \delta h_t^3} = 1, \quad (339)$$

$$\hat{\theta}'_{W_{t-1}^4} \hat{\theta}_{\delta h_t^3} = 1, \quad \hat{\theta}'_{\delta W_t^4} = 1. \quad (340)$$

For a formal argument to show feature learning (such that the limits are not degenerate), we refer to (Yang and Hu, 2021, H.7).

Table 2: *bcd* parametrisation vlaue for μ P-MoE Adam

Layer (ℓ)	\mathbf{b}_ℓ	\mathbf{c}_ℓ	\mathbf{d}_ℓ
1	0	0	1
2	1	1	0
3	0.5	1	1
4	1	1	0

Stability, Feature Learning, and Faithfulness for Adam: The faithfulness follows from the values at the first weight updates Theorem 15. Similar to above, it is easy to see that all $\hat{\theta}_\bullet$ converge to one.

Lack of specialisation: We now show that for μ P-MoE without router bias, specialisation of the experts fails. We first show that, in the infinite-width limit, the gradient received by each router logit converges to the same deterministic value. Next, we show that this implies that the router logits, and hence the probabilities evolve synchronously, keeping the same value for each expert. In Softmax-Soft MoEs, due to the form of the softmax gradient, the above results imply that the gradient received by the router is zero. We finish with remarks that these theorems hold for Adam, and we find the finite-width scaling of the router gradient in Softmax-Soft MoEs.

Proposition 19 Consider a soft MoE network without a router bias with m experts with either Softmax or Sigmoid gating parametrised such that stability holds throughout training (Theorem 9), and such that the router logits converge to zero in the infinite width limit, $\mathbf{h}_0^2 = 0$ (Theorem 11, Table 1). For any time step $t \geq 0$, each router activation $\tilde{\mathbf{h}}^2$ component receives the same deterministic unscaled ($\Theta(1)$) gradient. That is,

$$d\tilde{h}_{t,k}^2 := \mathbb{E} \left[Z^{h_t^{3,k}} Z^{dh_t^3} \right], \quad (341)$$

is identical for all experts $k \in \{1, \dots, m\}$.

Proof First, observe that $Z^{dh_t^3} = Z^{\overline{W}_t^4}$ is a global variable shared by the entire network; it depends only on the output layer weights, making it independent of the specific expert index k . Therefore, it suffices to prove that the marginal distribution of the expert features $Z^{h_t^{3,k}}$ is identical for all k . We proceed by induction.

Base Case ($t = 0$): At initialization, the expert features are defined by the random variable:

$$Z^{h_0^{3,k}} = Z^{W_0^{3,k} h_0^1}. \quad (342)$$

The input features \mathbf{h}_0^1 are shared across all experts. The weights $\mathbf{W}_0^{3,k}$ are initialized i.i.d. such that $(\hat{W}_k^3)_{ij} \sim \mathcal{N}(0, n^{-1})$. Thus $Z^{W_0^{3,k} h_0^1}$ (derived by ZMatMul rule) is the same for every expert. Consequently, the random variables $Z^{h_0^{3,k}}$ are identically distributed for all k :

$$Z^{h_0^{3,k}} \stackrel{d}{=} Z^{h_0^{3,j}} \quad \forall k, j. \quad (343)$$

Thus, the expectation $\mathbb{E}[Z^{h_0^{3,k}} Z^{dh_0^3}]$ is constant with respect to k , and thus $d\tilde{h}_{0,k}^2$ is the same for each expert.

Note that for softmax or sigmoid gating, the router probabilities are the same for each expert at initialisation due to $\tilde{h}_{0,k}^2 = 0$ and without a router bias,

$$\tilde{h}_{0,k}^2 = \begin{cases} \frac{1}{m} \mathbf{1} & \text{if G = Softmax} \\ \frac{1}{2} \mathbf{1} & \text{if G = Sigmoid} \end{cases}.$$

Hence $\tilde{h}_{0,i}^2 = \tilde{h}_{0,j}^2$ for every $i, j \in [m]$.

Inductive Step: Assume that for all time steps $s < t$, the expert features are identically distributed across experts:

$$Z^{h_s^{3,k}} \stackrel{d}{=} Z^{h_s^{3,j}} \implies d\tilde{h}_{s,k}^2 = \tilde{h}_{s,j}^2 \quad \forall k, j. \quad (344)$$

The router logits at a given time step $s < t$ are

$$\mathring{h}_{s,k}^2 = \tilde{h}_{0,k}^2 + \sum_{p=1}^s \delta \mathring{h}_{p,k}^2,$$

where $\mathring{h}_{0,k}^2 = \mathring{h}_{0,j}^2$ for any experts k, j due to the initialisation assumption. We now show that the updates $\delta \mathring{h}_{p,k}^2$ are the same for each expert index k . The infinite width limit router weight updates and corresponding logit updates for each expert is

$$Z_{k,:}^{(\delta W_p^2)} = -\eta \sum_{\beta=1}^{d_{\text{out}}} (\mathring{h}_{p,k}^2)_{\beta} (\mathring{\chi}_{p-1})_{\beta} Z^{h_{p-1}^1}, \quad (345)$$

$$\delta \mathring{h}_{p,k}^2 = \mathring{\theta}_{W_{p-1}^2} \mathring{\theta}_{\delta h_p^1} \mathbb{E}[Z^{(\overline{W_{p-1}^2})_{k,:}} Z^{\delta h_p^1}] - \mathring{\theta}'_{\delta W_p^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathring{h}_{p-1,k}^2)_{\beta} (\mathring{\chi}_{p-1})_{\beta} \mathbb{E}[Z^{h_{p-1}^1} Z^{h_p^1}]. \quad (346)$$

As the inductive hypothesis implies $\mathring{h}_{p,k}^2 = \mathring{h}_{p,j}^2, \forall k, j$, we see that updates are the same for each expert a every time step $p < s < t$ and hence that

$$\mathring{h}_{s,k}^2 = \mathring{h}_{s,j}^2 \quad \forall k, j.$$

As the gradient for each expert k at time step s is

$$Z^{\text{d}h_s^{3,k}} = \mathring{h}_{s,k}^2 Z^{\text{d}h_s^3},$$

where $Z^{\text{d}h_s^3}$ is shared for all the experts. Using the fact that $\mathring{h}_{s,k}^2 = \mathring{h}_{s,j}^2$ for all experts k, j , we get

$$Z^{\text{d}h_s^{3,k}} \stackrel{d}{=} Z^{\text{d}h_s^{3,j}} \quad \forall k, j.$$

Consider the update at step t , defined in the infinite width limit as:

$$\mathbf{h}_t^{3,k} = \mathbf{h}_{t-1}^{3,k} + \delta \mathbf{h}_t^{3,k}, \quad (347)$$

$$\begin{aligned} Z^{\delta h_t^{3,k}} &= \mathring{\theta}_{\delta h_t^1 / \delta h_t^{3,k}} Z^{W_0^{3,k} \delta h_t^1} - \sum_{s=0}^{t-2} \mathring{\theta}'_{\delta W_s^{3,k} \delta h_t^1 / \delta h_t^{3,k}} \eta \sum_{\beta=1}^{d_{\text{out}}} Z_{\beta}^{(\text{d}h_s^{3,k})} (\mathring{\chi}_s)_{\beta} \mathbb{E}[Z^{h_s^1} Z^{\delta h_t^1}] \\ &\quad - \mathring{\theta}'_{\delta W_t^{3,k} / \delta h_t^{3,k}} \eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(\text{d}h_{t-1}^{3,k})_{\beta}} (\mathring{\chi}_{t-1})_{\beta} \mathbb{E}[Z^{h_{t-1}^1} Z^{h_t^1}], \end{aligned} \quad (348)$$

Note that $\mathring{\theta}'_{\delta W_t^{3,k}}$ is the same for each expert. As we have shown, inductive hypothesis implies that $Z^{\text{d}h_s^{3,k}} \stackrel{d}{=} Z^{\text{d}h_s^{3,j}}$, for all experts k, j at every time step $s < t$. Furthermore, $Z^{W_0^{3,k} \delta h_t^1}$ (from the ZMatMul rule) is the same for each expert as $Z^{W_0^{3,k}}$ is the same for each expert and $Z^{\delta h_t^1}$ is shared by all the experts. Thus, using the inductive hypothesis that $Z^{h_{t-1}^{3,k}} \stackrel{d}{=} Z^{h_{t-1}^{3,j}}$ we have,

$$Z^{\delta h_t^{3,k}} \stackrel{d}{=} Z^{\delta h_t^{3,j}} \implies Z^{h_t^{3,k}} \stackrel{d}{=} Z^{h_t^{3,j}} \quad \forall k, j. \quad (349)$$

This in turn implies that

$$\mathring{d} \mathring{h}_{t,k}^2 = \mathring{d} \mathring{h}_{t,j}^2 \quad \forall k, j. \quad (350)$$

■

Corollary 20 Consider a soft MoE network without a router bias with m experts with either softmax or sigmoid gating parametrised such that stability holds throughout training (Theorem 9), and such that the router logits converge to zero in the infinite width limit, $\mathring{\mathbf{h}}_0^2 = 0$ (Theorem 11, Table 1). For any time step $t \geq 0$, the router probability is the same for every expert in the infinite width limit. That is,

$$\mathring{\tilde{h}}_{t,i}^2 = \mathring{\tilde{h}}_{t,j}^2 \quad (351)$$

for every expert $i, j \in [m]$.

Proof Result follows from the inductive proof in Theorem 19. For the last step of the induction, the result of Theorem 19, which implies that the gradient of the router logits are the same for each expert index $d\mathring{h}_{t,k}^2 = d\mathring{h}_{t,j}^2$. This, combined with the initialisation $\mathring{\mathbf{h}}_0^2 = 0$ implies that $\mathring{h}_{t,i}^2 = \mathring{h}_{t,j}^2$ for every expert $i, j \in [m]$ and time step t . The result for $\mathring{\tilde{h}}_{t,i}^2 = \mathring{\tilde{h}}_{t,j}^2$ for every expert $i, j \in [m]$ follows. ■

Corollary 21 Consider a Softmax-Soft MoE network without a router bias. For any time step $t \geq 0$, in the infinite width limit, the gradient with respect to the router vanishes,

$$d\mathring{h}_{t,k}^2 = 0, \quad \text{for every expert } k. \quad (352)$$

Consequently, the router weights do not evolve from initialisation,

$$Z^{(\delta W_t^2)_{k,:}} = 0, \quad (353)$$

$$\implies \mathring{\tilde{h}}_{t,k}^2 = \mathbb{E}[Z^{(\overline{W}_0^2)_{k,:}} Z^{\delta h_t^1}], \quad (354)$$

for every k . That is, in the limit, the router is stuck at initialisation.

Proof The limiting value of the gradient with respect to the router is

$$d\mathring{h}_{t,k}^2 = \sum_{j=1}^m J(\mathring{h}_t^2)_{kj} d\mathring{h}_{t,j}^2, \quad (355)$$

where $J(h^2)_{kj} = \tilde{h}_k^2(\delta_{kj} - \tilde{h}_j^2)$ is the Jacobian of the Softmax and $\tilde{h}^2 = \text{Softmax}(h^2)$. From Theorem 19, $d\mathring{h}_{t,j}^2 = d\mathring{h}_{t,k}^2$ for all experts k, j . Let $C := d\mathring{h}_{t,j}^2$. For every k , we have

$$d\mathring{h}_{t,k}^2 = \sum_{j=1}^m J(\mathring{h}_t^2)_{kj} C \quad (356)$$

$$= C \mathring{\tilde{h}}_{t,k}^2 \sum_{j=1}^m (\delta_{kj} - \mathring{\tilde{h}}_{t,j}^2) \quad (357)$$

$$= 0, \quad (358)$$

where we use the fact that $\sum_{j=1}^m \mathring{\tilde{h}}_{t,j}^2 = 1$. The rest follows from looking at the infinite width limits of $\delta \mathbf{W}_t^2$ and $\delta \mathbf{h}_{t,k}^2$. ■

Remark 22 (Router bias prevents the stagnation in Theorem 20) *The reason for the router stagnation is the symmetric initialisation due to μP -MoE and the resulting identical deterministic gradients for each expert. A router bias (or noise) ensures that at initialisation the router is not symmetric, hence leading to varied deterministic gradient for each expert. The router thus does not stagnate.*

Remark 23 (Results hold for any stable parametrisation) *For the above results, we only assume stability and the results are shown using unscaled quantities (quantities that scale $\Theta(1)$). As such, any values of bcd do not change our results. We also assume $\mathring{\mathbf{h}}_0^2 = 0$ — the only stable parametrisation where this is not true is with $b_2 = 0.5$. We refer to (Yang and Hu, 2021, H.4.1) for discussion that. Furthermore, $b_2 = 1$ is required for feature learning considerations (Theorem 11).*

Remark 24 (Extension to Adam) *The above theorems hold trivially for Adam. The only change required is in the proof for Theorem 19. Here, the infinite width limit update to $\mathbf{h}_t^{3,k}$ will be (Appendix C.4)*

$$Z^{\delta h_t^{3,k}} = \mathring{\theta}_{\delta h_t^1 / \delta h_t^{3,k}} Z^{W_0^{3,k} \delta h_t^1} \quad (359)$$

$$- \sum_{s=0}^{t-2} \mathring{\theta}'_{\delta W_s^{3,k} / \delta h_t^1 / \delta h_t^{3,k}} \eta \mathbb{E} \left(Z_1^{h_s^1}, Z_1^{\delta h_t^1} \right) \left[Q_t^3 \left(\sum_{\beta}^{d_{out}} Z^{(dh_s^{3,k})_{\beta}} (\mathring{\chi}_s)_{\beta} Z_1^{h_s^1} \right) Z_1^{\delta h_t^1} \right] \quad (360)$$

$$- \mathring{\theta}'_{\delta W_t^{3,k} / \delta h_t^{3,k}} \eta \mathbb{E} \left(Z_1^{h_{t-1}^1}, Z_1^{h_t^1} \right) \left[Q_t^3 \left(\sum_{\beta}^{d_{out}} Z^{(dh_{t-1}^{3,k})_{\beta}} (\mathring{\chi}_{t-1})_{\beta} Z_1^{h_{t-1}^1} \right) Z_1^{h_t^1} \right], \quad (361)$$

where Q is the Adam update function (Theorem 13). Here it is easy to see that $Z^{dh_s^{3,k}} \stackrel{d}{=} Z^{dh_s^{3,j}}$, for all experts k, j at every time step $s < t$ implies the same result as Theorem 19. The result from Theorem 21 follows.

Remark 25 (Scale of $dh_{t,k}^2$ for finite n) *In the tensor program, for each output coordinate β we have*

$$(d\tilde{h}_{t,k}^2)_{\beta} = \frac{1}{n} \sum_{\alpha=1}^n \theta'_{W_t^4} \theta_{h_t^{3,k}} (h_t^{3,k})_{\alpha} (dh_t^3)_{\alpha\beta} = \theta'_{W_t^4} \theta_{h_t^{3,k}} \frac{1}{n} \sum_{\alpha=1}^n (Y_{t,k})_{\alpha\beta}, \quad (\text{Moment}) \quad (362)$$

where we defined

$$(Y_{t,k})_{\alpha\beta} := (h_t^{3,k})_{\alpha} (dh_t^3)_{\alpha\beta}. \quad (363)$$

The gradient with respect to the router logits is

$$(dh_{t,k}^2)_{\beta} = \sum_{j=1}^m J(\tilde{h}_t^2)_{kj} (d\tilde{h}_{t,j}^2)_{\beta} = \theta'_{W_t^4} \theta_{h_t^{3,k}} \frac{1}{n} \sum_{\alpha=1}^n X_{t,k,\alpha\beta}, \quad (\text{Moment}) \quad (364)$$

where

$$X_{t,k,\alpha\beta} := \sum_{j=1}^m J(\tilde{h}_t^2)_{kj} (Y_{t,j})_{\alpha\beta}. \quad (365)$$

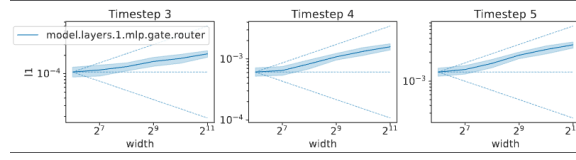


Figure 2: Average norm of router features as width is increased with router $d = 0.5$. Shows that for high ϵ , this causes the router logits to diverge.

In the infinite-width limit, we showed above that $\mathring{d}h_{t,k}^2 = 0$, so in particular $\mathbb{E}[X_{t,k,\alpha\beta}] = 0$. To capture the scale, consider the second-moment ‘‘Moment’’ variable

$$\theta'_{W_t^4} \theta_{h_t^{3,k}} \frac{1}{n} \sum_{\alpha=1}^n (X_{t,k,\alpha\beta})^2 = \theta'_{W_t^4} \theta_{h_t^{3,k}} \sum_{j=1}^m \sum_{i=1}^m J(\tilde{h}_t^2)_{kj} J(\tilde{h}_t^2)_{ki} \frac{1}{n} \sum_{\alpha=1}^n (Y_{t,j})_{\alpha\beta} (Y_{t,i})_{\alpha\beta}, \quad (366)$$

whose infinite-width limit, by the Master theorem, is

$$\hat{\theta}'_{W_t^4} \hat{\theta}_{h_t^{3,k}} \sum_{j=1}^m \sum_{i=1}^m J(\mathring{h}_t^2)_{kj} J(\mathring{h}_t^2)_{ki} \mathbb{E}[Z^{Y_{t,j}} Z^{Y_{t,i}}]. \quad (367)$$

Thus the variance of $(dh_{t,k}^2)_{\beta} = n^{-1} \sum_{\alpha=1}^n X_{t,k,\alpha\beta}$ scales as

$$\text{Var}[(dh_{t,k}^2)_{\beta}] = \frac{1}{n} (\theta'_{W_t^4} \theta_{h_t^{3,k}})^2 \text{Var}(X_{t,k,\alpha\beta}), \quad (368)$$

where $\text{Var}(X_{t,k,\alpha\beta})$ is some non-zero scalar. Consequently, we have that $\mathring{d}h_{t,k}^2$ scales $\Theta(n^{-\frac{1}{2}})$.

Remark 26 (Scale of $\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2$ for Softmax-Soft MoE networks under $\mu\text{P-Moe}$ parametrisation)

Here, we find the scale of $\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2$ for Softmax-Soft MoE networks under the $\mu\text{P-MoE}$ parametrisation. We find the scale under zero-initialisation of \mathbf{W}_0^2 , so that the initialisation of the network matches the infinite-width limit, as recommended in Yang et al. (2021). Here, only the effective updates decide the scale of the feature update. The scale of the router weight update is $\theta_{\delta W_t^2} = n^{-1} \cdot n^{-0.5} = n^{-1.5}$, where the extra $n^{-0.5}$ comes from the finite width scale of the router gradient, due to the softmax Jacobian (Theorem 25). Thus, the scale of $\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2$ is $\theta'_{\delta W_t^2} = n \cdot n^{-1.5} = n^{-0.5}$.

Remark 27 (Benign Scaling for Adam) In Softmax-Soft MoEs, the scale of the gradient has an extra $n^{-0.5}$ factor. This implies setting $d = 0.5$ for the router can be used to rescale the variance such that it is $\Theta(1)$. However, as this results in a random scalar in the program, this limit is not expressible in Tensor Programs (Yang and Hu, 2021, Theorem 7.4). In practice we find this can lead to diverging logits Figure 2.

Appendix E. Proof of Theorem 4

From Appendix C.2, we see that at initialisation the router logits converge deterministically to $C1$ (with $C = 0$) under $\mu\text{P-MoE}$ scaling without a bias. A simple CLT argument, (formalised

in Yang (2019)) shows that $\sqrt{n}(\mathbf{h}^2(n) - C\mathbf{1})$ converges in distribution to a Gaussian with variance proportional to $\mathbb{E}[(Z^{h^1})^2]$. Thus, we can write

$$\mathbf{h}^2(n) = C\mathbf{1} + \sigma_n \epsilon_n,$$

with $\epsilon_n = \sqrt{n}(\mathbf{h}^2(n) - C\mathbf{1})$ and $\sigma_n = \frac{1}{\sqrt{n}}$.

Theorem 28 (Limit of the Top- K masked vector) *Let*

$$S_n := \text{TopK}(\epsilon_n) = \text{TopK}(\{h_i^2(n)\}_{i=1}^m),$$

and define

$$\tilde{h}_i^2(n) := h_i^2(n) \mathbf{1}\{i \in S_n\}, \quad \tilde{\mathbf{h}}^2(n) := (\tilde{h}_1^2(n), \dots, \tilde{h}_M^2(n)).$$

Then,

$$\tilde{\mathbf{h}}^2(n) \xrightarrow{d} C \mathbf{1}_S \quad \text{in } \mathbb{R}^M, \quad \text{where } S := \text{TopK}(\epsilon)$$

where $\mathbf{1}_S$ is the indicator function indicating whether index $j \in S$.

Proof Define the Top- K indicator map

$$T : \mathbb{R}^M \rightarrow \{0, 1\}^M, \quad T(z) := \mathbf{1}_{\text{TopK}(z)}.$$

This map is discontinuous only at points where ties occur among the order statistics. As the probability of ties when sampling from ϵ is zero (as Gaussian as a continuous joint density),

$$\mathbb{P}(\epsilon \in D_T) = 0,$$

where D_T denotes the discontinuity set of T . Hence, by the extended continuous mapping theorem (Billingsley, 2013),

$$T(\epsilon_n) \xrightarrow{d} T(\epsilon), \quad \text{i.e. } \mathbf{1}_{S_n} \xrightarrow{d} \mathbf{1}_S.$$

Now write

$$\tilde{\mathbf{h}}^2(n) = (C\mathbf{1} + \sigma_n \epsilon_n) \odot T(\epsilon_n) = C T(\epsilon_n) + \sigma_n (\epsilon_n \odot T(\epsilon_n)).$$

Since $\epsilon_n \xrightarrow{d} \epsilon$, the sequence $\{\epsilon_n\}$ is tight, and therefore

$$\|\epsilon_n\|_\infty = O_p(1).$$

As $\sigma_n \rightarrow 0$, it follows that

$$\left\| \sigma_n (\epsilon_n \odot T(\epsilon_n)) \right\|_\infty \leq \sigma_n \|\epsilon_n\|_\infty \xrightarrow{p} 0.$$

Combining this with $\mathbf{1}_{S_n} \xrightarrow{d} \mathbf{1}_S$ and applying Slutsky's theorem yields

$$\tilde{\mathbf{h}}^2(n) \xrightarrow{d} C \mathbf{1}_S.$$

■

Thus, in the first forward pass (at initialisation), where the router is symmetric, the Top- K index set is not deterministic as $S := \text{TopK}(\epsilon)$.

Appendix F. Tensor Program for Top- K MoE

In this section we derive the Tensor program and the corresponding infinite width limit for Top- K MoEs (??). Note we only derive this for SGD, corresponding limit to Adam is analogous to Appendix C.3.

F.1. Tensor Program for Top- K MoE with SGD

We use the same notation as in Appendix C.1. Throughout we make the assumption that the parametrisation is stable (Theorem 9).

Network:

$$\mathbf{h}^1 = \mathbf{W}^1 \boldsymbol{\xi}, \quad \in \mathbb{R}^{n \times 1} \quad (369)$$

$$\mathbf{h}^2 = \mathbf{W}^2 \mathbf{h}^1 + \mathbf{b}^2, \quad \in \mathbb{R}^{m \times 1} \quad (370)$$

$$\mathbf{s} = \mathsf{T}(\mathbf{h}^2) \quad \text{where } \mathsf{T}_k(\mathbf{h}^2) := \mathbf{1}\{k \in \text{TopK}(\mathbf{h}^2)\} \quad (371)$$

$$\tilde{\mathbf{h}}^2 = \mathsf{G}(\mathbf{h}^2, \mathbf{s}), \quad \in [0, 1]^{m \times 1} \quad (372)$$

$$\mathbf{h}^{3,k} = \mathbf{W}^{3,k} \mathbf{h}^1, \quad \in \mathbb{R}^{n \times 1} \quad (373)$$

$$\mathbf{h}^3 = \sum_{k=1}^m \mathbf{h}^{3,k} \tilde{h}_k^2, \quad \in \mathbb{R}^{n \times 1} \quad (374)$$

$$\mathbf{f} = \mathbf{W}^4 \mathbf{h}^3, \quad \in \mathbb{R}^{d_{\text{out}} \times 1} \quad (375)$$

where $\mathbf{b}^2 \in \mathbb{R}^m$ is a bias parameter that is represented as an initial scalar in the program, \mathbf{s} is the Top- K mask, and $\mathsf{G}(\cdot, \mathbf{s})$ is the activations applied to the index set \mathbf{s} . For Softmax

$$\mathsf{G}_k(u, \mathbf{s}) := \frac{s_k \exp(u_k)}{\sum_{j=1}^m s_j \exp(u_j)}, \quad (376)$$

and for Sigmoid

$$\mathsf{G}_k(u, \mathbf{s}) := s_k \sigma(u_k).$$

Conditioned on noise/bias: We derive the Tensor program conditioned on the bias term noises $\{\mathbf{b}_t^2\}_{t=0}^T$ or the learnable bias \mathbf{b}^2 . We assume that they have been sampled from a distribution such that $\mathbb{P}(b_i^2 = b_j^2) = 0$. Practically, we sample $b_i \sim \mathcal{N}(0, 1)$.

Remark 29 *Top- K MoEs are not admissible in the Tensor program framework without a router bias. This is because the router history is not deterministic, whereas the Tensor program Master Theorem requires scalars (such as router mask) to converge to deterministic constants.*

First Forward Pass:

$$\begin{aligned}
\mathbf{h}_0^1 &= \theta_{W_0^1} \overline{\mathbf{W}}_0^1 \boldsymbol{\xi}, & (\mathbf{h}_0^1)_\alpha &= \phi \left(\{(\overline{\mathbf{W}}_0^1)_{\alpha\beta}\}_{\beta=1}^{d_{\text{in}}}; \xi_1, \dots, \xi_{d_{\text{in}}}, \theta_{W_0^1} \right) = \sum_{j=1}^{d_{\text{in}}} \theta_{W_0^1} (\overline{\mathbf{W}}_0^1)_{\alpha j} \xi_j, \\
\mathbf{h}_0^2 &= \theta'_{W_0^2} \frac{1}{n} \overline{\mathbf{W}}_0^2 \mathbf{h}_0^1 + \mathbf{b}_0^2, & (\mathbf{h}_0^2)_k &= \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}}_0^2)_{k\alpha}, (\mathbf{h}_0^1)_\alpha; \theta'_{W_0^2}, \{(b_0^2)_k\}_{k=1}^m \right) = \frac{1}{n} \theta'_{W_0^2} \sum_{\alpha=1}^n (\overline{\mathbf{W}}_0^2)_{k\alpha} (\mathbf{h}_0^1)_\alpha + (b_0^2)_k, \\
\mathbf{s} &= \mathbf{T}(\mathbf{h}^2), & & \text{(Moment)} \\
\tilde{\mathbf{h}}_0^2 &= \mathbf{G}(\hat{\mathbf{h}}_0^2, \mathbf{s}), & & \text{(Moment, See G)} \\
\hat{\mathbf{h}}_0^{3,k} &= \overline{\mathbf{W}}_0^{3,k} \mathbf{h}_0^1, & & \text{(MatMul)} \\
\mathbf{h}_0^{3,k} &= \theta_{W_0^3} \hat{\mathbf{h}}_0^{3,k}, & (\mathbf{h}_0^{3,k})_\alpha &= \phi \left((\hat{\mathbf{h}}_0^{3,k})_\alpha; \theta_{W_0^3} \right) = \theta_{W_0^3} (\hat{\mathbf{h}}_0^3)_\alpha \\
\mathbf{h}_0^3 &= \sum_{k=1}^m \mathbf{h}_0^{3,k} \tilde{h}_{0,k}^2, & (\mathbf{h}_0^3)_\alpha &= \phi \left((\mathbf{h}_0^{3,1})_\alpha, \dots, (\mathbf{h}_0^{3,m})_\alpha; \tilde{h}_{0,1}^2, \dots, \tilde{h}_{0,m}^2 \right) = \sum_{k=1}^m (\mathbf{h}_0^{3,k})_\alpha \tilde{h}_{0,k}^2 \\
\mathbf{f}_0 &= \theta'_{W_0^4} \frac{1}{n} \overline{\mathbf{W}}_0^4 \mathbf{h}_0^3, & (\mathbf{f}_0)_\beta &= \frac{1}{n} \sum_{\alpha=1}^n \phi \left((\overline{\mathbf{W}}_0^4)_{\beta\alpha}, (\mathbf{h}_0^3)_\alpha; \theta'_{W_0^4} \right) = \frac{1}{n} \theta'_{W_0^4} \sum_{\alpha=1}^n (\overline{\mathbf{W}}_0^4)_{\beta\alpha} (\mathbf{h}_0^3)_\alpha.
\end{aligned}$$

We assume stability at initialisation and thus the initialisation values in Theorem 11. For these values, the scale of $\mathbf{W}_0^2 \mathbf{h}_0^1$ is $n^{-0.5}$.

Remark 30 (Initial mask is deterministic conditioned on the bias) *With the initialisation values in Theorem 11, as the width increases $\mathbf{W}_0^2 \mathbf{h}_0^1 \rightarrow \mathbf{0}$, with vanishing finite-width fluctuation (Yang, 2019). Consequently, if \mathbf{h}_0^2 denotes the limit of \mathbf{h}_0^2*

$$\mathbf{T}(\mathring{\mathbf{h}}_0^2) = \mathbf{T}(\mathbf{b}_0^2).$$

Hence the Top-K operator is determined by the initial bias and is deterministic in the limit. This can be shown with the result in Appendix E.

Noise and Learnable Bias: We consider two cases here:

1. Case 1: A different noise is added at each time step $\{\mathbf{b}_t^2\}_{t=1}^T$.
2. Case 2: The term \mathbf{b}_0^2 is learned. Note that the updates below are written only for learnable biases.

First Backward Pass: Same as Appendix C.1, except the gradient $dh_{0,k}^2$ requires gradient of the Top-K mask.

$$dh_{0,k}^2 = \mathbf{s}_0 \sum_{j \in \mathcal{S}} J(\tilde{\mathbf{h}}_0^2)_{jk} d\tilde{h}_{0,j}^2, \quad \in \mathbb{R}^{1 \times d_{\text{out}}} \quad \text{(Moment)} \quad (377)$$

where $J(\tilde{\mathbf{h}}_0^2)_{jk} = \frac{\partial (\tilde{\mathbf{h}}_0^2)_j}{\partial (\tilde{\mathbf{h}}_0^2)_k}$ is the Jacobian of the activation function on the sparse logits.

First Weight Updates: We need to consider the update to the learnable \mathbf{b}^2

$$\mathbf{b}_1^2 - \mathbf{b}_0^2 = -\eta \theta'_{W_0^4} d\mathbf{h}_0^2 \boldsymbol{\chi}_0. \quad (378)$$

The rest is the same as Appendix C.1.

First Feature Updates: The only feature update that will change from Appendix C.1 is \mathbf{h}^2 and $\tilde{\mathbf{h}}^2$:

$$\mathbf{h}_1^2 - \mathbf{h}_0^2 = \mathbf{W}_0^2(\mathbf{h}_1^1 - \mathbf{h}_0^1) + (\mathbf{W}_1^2 - \mathbf{W}_0^2)\mathbf{h}_1^1 + \mathbf{b}_1^2 - \mathbf{b}_0^2 \quad (379)$$

$$= \theta_{\delta h_1^1} \theta'_{W_0^2} \frac{\overline{\mathbf{W}_0^2} \delta \mathbf{h}_1^1}{n} - \theta'_{\delta W_1^2} \eta d\mathbf{h}_0^2 \chi_0 \frac{(\mathbf{h}_0^1)^T \mathbf{h}_1^1}{n} - \eta \theta'_{W_0^4} d\mathbf{h}_0^2 \chi_0 \quad (380)$$

$$\theta_{h_1^2}(\delta \mathbf{h}_1^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_{0k\alpha}^2}, (\delta \mathbf{h}_1^1)_\alpha, (\mathbf{h}_0^1)_\alpha, (\mathbf{h}_1^1)_\alpha; \theta_{\delta h_1^1}, \theta'_{W_0^2}, \theta'_{\delta W_1^2}, \eta, \{(\mathbf{d}\mathbf{h}_0^2)_{k\beta}, (\chi_0)_\beta\}_\beta^{d_{\text{out}}} \right) \quad (\text{Moment}) \quad (381)$$

$$\phi(\dots) := \theta_{\delta h_1^1} \theta'_{W_0^2} \overline{\mathbf{W}_{0k\alpha}^2} (\delta \mathbf{h}_1^1)_\alpha - \theta'_{\delta W_1^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_0^2)_{k\beta} (\chi_0)_\beta (\mathbf{h}_0^1)_\alpha (\mathbf{h}_1^1)_\alpha - \eta \theta'_{W_0^4} \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_0^2)_{k\beta} (\chi_0)_\beta, \quad (382)$$

$$\theta_{\delta h_1^2} := \max(\theta_{\delta h_1^1} \theta'_{W_0^2}, \theta'_{\delta W_1^2}, \theta'_{W_0^4}). \quad (383)$$

$$\tilde{\mathbf{h}}_1^2 - \tilde{\mathbf{h}}_0^2 = \mathbf{G}(\mathbf{h}_1^2, \mathbf{s}_1) - \mathbf{G}(\mathbf{h}_0^2, \mathbf{s}_0), \quad (\text{Moment}) \quad (384)$$

The rest of the feature updates are the same as Appendix C.1.

t th Backward Pass: Same as Appendix C.1, except for the mask

$$d\mathbf{h}_{t,k}^2 = \mathbf{s}_t \sum_{j \in \mathcal{S}} J(\tilde{\mathbf{h}}_t^2)_{jk} d\tilde{h}_{t,j}^2, \quad \in \mathbb{R}^{1 \times d_{\text{out}}} \quad (\text{Moment}) \quad (385)$$

where $J(\tilde{\mathbf{h}}_t^2)_{jk} = \frac{\partial (\tilde{\mathbf{h}}_t^2)_j}{\partial (\tilde{\mathbf{h}}_t^2)_k}$ is the Jacobian of the activation function on the sparse logits.

t th Weight Updates: We need to consider the update to \mathbf{b}^2

$$\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2 = -\eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \chi_{t-1}. \quad (386)$$

The rest is the same as Appendix C.1.

t th Feature Updates: The only feature update that will change from Appendix C.1 is \mathbf{h}^2 and $\hat{\mathbf{h}}^2$:

$$\mathbf{h}_t^2 - \mathbf{h}_{t-1}^2 = \mathbf{W}_{t-1}^2(\mathbf{h}_t^1 - \mathbf{h}_{t-1}^1) + (\mathbf{W}_t^2 - \mathbf{W}_{t-1}^2)\mathbf{h}_t^1 + \mathbf{b}_t^2 - \mathbf{b}_{t-1}^2 \quad (387)$$

$$= \theta_{\delta h_t^1} \theta'_{W_{t-1}^2} \frac{\overline{\mathbf{W}_{t-1}^2} \delta \mathbf{h}_t^1}{n} - \theta'_{\delta W_t^2} \eta d\mathbf{h}_{t-1}^2 \chi_{t-1} \frac{(\mathbf{h}_{t-1}^1)^T \mathbf{h}_t^1}{n} - \eta \theta'_{W_{t-1}^4} d\mathbf{h}_{t-1}^2 \chi_{t-1} \quad (388)$$

$$\theta_{h_t^2}(\delta \mathbf{h}_t^2)_k = \frac{1}{n} \sum_{\alpha=1}^n \phi \left(\overline{\mathbf{W}_{t-1k\alpha}^2}, (\delta \mathbf{h}_t^1)_\alpha, (\mathbf{h}_{t-1}^1)_\alpha, (\mathbf{h}_t^1)_\alpha; \theta_{\delta h_t^1}, \theta'_{W_{t-1}^2}, \theta'_{\delta W_t^2}, \eta, \{(\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta}, (\chi_{t-1})_\beta\}_\beta^{d_{\text{out}}} \right) \quad (389)$$

$$\phi(\dots) := \theta_{\delta h_t^1} \theta'_{W_{t-1}^2} \overline{\mathbf{W}_{t-1k\alpha}^2} (\delta \mathbf{h}_t^1)_\alpha - \theta'_{\delta W_t^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta} (\chi_{t-1})_\beta (\mathbf{h}_{t-1}^1)_\alpha (\mathbf{h}_t^1)_\alpha - \eta \theta'_{W_{t-1}^4} \sum_{\beta=1}^{d_{\text{out}}} (\mathbf{d}\mathbf{h}_{t-1}^2)_{k\beta} (\chi_{t-1})_\beta \quad (390)$$

$$\theta_{\delta h_t^2} := \max(\theta_{\delta h_t^1} \theta'_{W_{t-1}^2}, \theta'_{\delta W_t^2}, \theta'_{W_{t-1}^4}). \quad (391)$$

$$\tilde{\mathbf{h}}_t^2 - \tilde{\mathbf{h}}_{t-1}^2 = \mathbf{G}(\mathbf{h}_t^2) - \mathbf{G}(\mathbf{h}_{t-1}^2), \quad (\text{Moment}) \quad (392)$$

The rest will be the same as Appendix C.1.

F.2. Infinite Width Limit for TopK MoE with SGD

We use the same notation as in Appendix C.2.

Remark 31 (Deterministic Limiting Mask Conditioned on Noise/Bias) *To find the limit, we need to show that the Top-K mask is deterministic in both cases: 1) where noise is added at each step to the router logits, 2) where a learnable bias is added to the router logits.*

1. *Note that the noise \mathbf{b}_t^2 is drawn independently of $\mathbf{W}_t \mathbf{h}_t^1$ at each step. This implies that the probability of ties $(\mathbf{h}^2)_i = (\mathbf{h}^2)_j$ is 0. $\mathbf{W}_t \mathbf{h}_t^1$ converges deterministically according to the Master Theorem: If $\mathbf{W}_t \mathbf{h}_t^1$ converges to the same value for each dimension such that the setting of Appendix E holds for $\mathbf{W}_t \mathbf{h}_t^1$, then the mask is $\mathbf{T}(\mathbf{b}_t^2)$ in the limit and thus deterministic.*
2. *With a learnable bias, the mask is deterministic at initialisation (Theorem 30). However, for subsequent steps, we must make the additional assumption that for widths $n > N$ for some $N > 0$, in the resulting logits $h_{t,(K)}^2 - h_{t,(K+1)}^2 > \epsilon$ where ϵ is width independent, for all t . As TP only allows for pseudo-Lipschitz non-linearities, scalars in the program converge to a deterministic value, and the logits are bounded under μP -MoE, this can be proven using induction. However, we leave this for future work.*

First Forward Pass:

$$Z^{h_0^1} = \xi Z^{W_0^1}, \quad (393)$$

$$\mathring{\mathbf{h}}_0^2 = \mathbf{b}_0^2, \quad (394)$$

$$\mathring{\mathbf{s}}_0 = \mathbf{T}(\mathring{\mathbf{h}}_0^2) \quad (395)$$

$$\mathring{\mathbf{h}}_0^2 = \mathbf{G}(\mathring{\mathbf{h}}_0^2, \mathring{\mathbf{s}}_0) \quad (396)$$

$$Z^{h_0^{3,k}} = Z^{W_0^{3,k} h_0^1} \text{ for every expert } k, \quad (397)$$

$$Z^{h_0^3} = \sum_{k=1}^m \mathring{h}_{0,k}^2 Z^{h_0^{3,k}} \quad (398)$$

$$\mathring{\mathbf{f}}_0 = \mathbf{0}. \quad (399)$$

Here, $Z^{W_0^1}$ is standard Gaussian by construction. The gating function \mathbf{G} is either a Softmax or Sigmoid. $Z^{W_0^{3,k} h_0^1} = \hat{Z}^{W_0^{3,k} h_0^1}$ from the Z MatMul rule. Finally $\mathring{\mathbf{f}}_0 = \mathbf{0}$ as $\mathring{f}_0 = \mathbb{E}[Z^{W_0^4}] \mathbb{E}[Z^{h_0^3}] = 0$.

First Backward Pass:

$$Z^{dh_0^3} = Z^{\overline{W_0^4}}, \quad (400)$$

$$\overset{\circ}{d}\tilde{h}_{0,k}^2 = 0 \text{ for every expert } k, \quad (401)$$

$$d\overset{\circ}{h}_{0,k}^2 = 0 \text{ for every expert } k, \quad (402)$$

$$Z^{dh_0^{3,k}} = \tilde{h}_{0,k}^2 Z^{dh_0^3} \text{ for every expert } k, \quad (403)$$

$$Z^{dh_0^1} = \sum_{k=1}^m Z^{(W_0^{3,k})^T dh_0^{3,k}} \quad (404)$$

$$(405)$$

$Z^{\overline{W_0^4}}$ is a standard Gaussian by construction. $d\tilde{h}_{0,k}^2 = \mathbb{E}[Z^{h_0^{3,k}}] \mathbb{E}[\overline{W_0^4}] = 0$. From this, $d\overset{\circ}{h}_{0,k}^2 = 0$ follows. $Z^{(W_0^{3,k})^T dh_0^3} = \hat{Z}^{(W_0^{3,k})^T dh_0^3}$ from the Z MatMul rule.

First Feature Update:

$$Z^{\delta h_1^1} = -\eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_0^1)_\beta} (\hat{\chi}_0)_\beta \boldsymbol{\xi}^T \boldsymbol{\xi}, \quad (406)$$

$$\delta \overset{\circ}{h}_{1,k}^2 = \overset{\circ}{\theta}_{\delta h_1^1} \mathbb{E}[Z^{(\overline{W_0^2})_{k,:}} Z^{\delta h_1^1}] - \overset{\circ}{\theta}'_{\delta W_1^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (d\overset{\circ}{h}_{0,k}^2)_\beta (\hat{\chi}_0)_\beta \mathbb{E}[Z^{h_0^1} Z^{h_1^1}] - \eta \overset{\circ}{\theta}'_{W_0^4} \sum_{\beta=1}^{d_{\text{out}}} (d\overset{\circ}{h}_{0,k}^2)_\beta (\hat{\chi}_0)_\beta, \quad (407)$$

$$\delta \overset{\circ}{s}_1^2 = \text{T}(\overset{\circ}{h}_1^2) - \text{T}(\overset{\circ}{h}_0^2) \text{ if } \overset{\circ}{\theta}_{\delta h_1^2} = 1, \quad (408)$$

$$\delta \overset{\circ}{h}_1^2 = \text{G}(\overset{\circ}{s}_1^2, \overset{\circ}{s}_1) - \text{G}(\overset{\circ}{h}_0^2, \overset{\circ}{s}_0) \text{ if } \overset{\circ}{\theta}_{\delta h_1^2} = 1, \quad (409)$$

$$Z^{\delta h_1^{3,k}} = \overset{\circ}{\theta}_{\delta h_1^1 / \delta h_1^{3,k}} Z^{W_0^{3,k} \delta h_1^1} - \overset{\circ}{\theta}'_{\delta W_1^3 / \delta h_1^{3,k} \eta} \sum_{\beta=1}^{d_{\text{out}}} Z^{(dh_0^{3,k})_\beta} (\hat{\chi}_0)_\beta \mathbb{E}[Z^{h_0^1} Z^{h_1^1}] \text{ for every expert } k, \quad (410)$$

$$Z^{\delta h_1^3} = \sum_{k=1}^m \overset{\circ}{\theta}_{\delta h_1^{3,k} / \delta h_1^3} \overset{\circ}{h}_{0,k}^2 Z^{\delta h_1^{3,k}} - \overset{\circ}{\theta}_{\delta \tilde{h}_1^2} \overset{\circ}{\theta}_{h_1^{3,k} / \delta h_3} \delta \tilde{h}_{1,k}^2 Z^{h_1^{3,k}}, \quad (411)$$

$$\delta \overset{\circ}{f}_1 = \overset{\circ}{\theta}_{\delta h_1^3} \mathbb{E}[Z^{\overline{W_0^4}} Z^{\delta h_1^3}] - \overset{\circ}{\theta}'_{\delta W_1^4} \eta \hat{\chi}_0 \mathbb{E}[Z^{h_0^3} Z^{h_1^3}], \quad (412)$$

where G is either softmax or sigmoid and G' is the derivative. Further, $Z^{W_0^{3,k} \delta h_1^1}$ is given by the Z MatMul rule,

$$Z^{W_0^{3,k} \delta h_1^1} = \hat{Z}^{W_0^{3,k} \delta h_1^1} + \dot{Z}^{W_0^{3,k} \delta h_1^1} \quad (413)$$

$$= \hat{Z}^{W_0^{3,k} \delta h_1^1} + Z^{dh_0^{3,k}} \mathbb{E} \left[\frac{\partial Z^{\delta h_1^1}}{\partial \hat{Z}^{(W_0^{3,k})^T dh_0^{3,k}}} \right]. \quad (414)$$

We write the limit of each vector in the program as

$$Z^{h_t} = Z^{h_0} + \overset{\circ}{\theta}_{\delta h_1} Z^{\delta h_1} + \dots + \overset{\circ}{\theta}_{\delta h_t} Z^{\delta h_t}. \quad (415)$$

Weight updates: The distribution of certain weight updates are required in the backward pass and thus the forward passes. Note that \mathbf{b} is a scalar in the program ($\in \mathbb{R}^m$) and hence does not have a limit.

$$Z^{\overline{\delta W_t^4}} = -\eta \dot{\chi}_{t-1} Z^{h_{t-1}^3}, \quad (416)$$

$$Z^{(\overline{\delta W_t^2})_{k,:}} = -\eta \sum_{\beta=1}^{d_{\text{out}}} (\dot{h}_{t,k}^2)_{\beta} (\dot{\chi}_{t-1})_{\beta} Z^{h_{t-1}^1}, \quad (417)$$

$$\mathbf{b}_t^2 - \mathbf{b}_{t-1}^2 = -\eta \dot{\theta}'_{W_{t-1}^4} \sum_{\beta=1}^{d_{\text{out}}} (\dot{h}_{t-1,k}^2)_{\beta} (\dot{\chi}_{t-1})_{\beta}. \quad (418)$$

t th Backward Pass: Let $\mathcal{S} = \text{TopK}(\mathbf{h}_t^2)$ be the set of active experts,

$$Z dh_t^3 = Z^{(\overline{W_t^4})^T}, \quad (419)$$

$$\dot{h}_{t,k}^2 = \mathbb{E}[Z h_t^{3,k} Z dh_t^3], \quad (420)$$

$$\dot{h}_{t,k}^2 = \dot{s}_{t,k} \sum_{j=1}^m J(\dot{h}_t^2)_{kj} \dot{h}_{t,k}^2, \quad (421)$$

$$Z dh_t^{3,k} = \dot{h}_{t,k}^2 Z dh_t^3, \quad (422)$$

$$Z dh_t^1 = \sum_{k=1}^m Z^{(\overline{W_t^2})_{k,:}} \dot{h}_{t,k}^2 + Z^{(W_0^{3,k})^T} dh_t^{3,k} - \sum_{s=1}^{t-1} Z^{h_s^1} \dot{\chi}_s \mathbb{E}[Z dh_s^{3,k} Z dh_t^{3,k}], \quad (423)$$

where $Z^{(W_0^{3,k})^T} dh_t^{3,k}$ is given by the Z MatMul rule,

$$Z^{(W_0^{3,k})^T} dh_t^{3,k} = \hat{Z}^{(W_0^{3,k})^T} dh_t^{3,k} + \dot{Z}^{(W_0^{3,k})^T} dh_t^{3,k} \quad (424)$$

$$= \hat{Z}^{(W_0^{3,k})^T} dh_t^{3,k} + \sum_{v \in \mathcal{V}: W_0^{3,k} v \in \mathcal{V}} Z^v \mathbb{E} \left[\frac{\partial Z dh_t^{3,k}}{\partial \hat{Z}^{W_0^{3,k} v}} \right], \quad (425)$$

where $\mathcal{V} := \{v \text{ for some vector in the program } v\}$ is the set of all vectors in the program.

tth Forward Pass:

$$Z^{\delta h_t^1} = -\eta \sum_{\beta=1}^{d_{\text{out}}} Z^{(\text{dh}_{t-1}^1)_\beta} (\dot{\chi}_{t-1})_\beta \xi^T \xi, \quad (426)$$

$$\delta \dot{h}_{t,k}^2 = \hat{\theta}'_{W_{t-1}^2} \hat{\theta}_{\delta h_t^1} \mathbb{E}[Z^{(\overline{W_{t-1}^2})_{k,\cdot}} Z^{\delta h_t^1}] - \hat{\theta}'_{\delta W_t^2} \eta \sum_{\beta=1}^{d_{\text{out}}} (\text{dh}_{t-1,k}^2)_\beta (\dot{\chi}_{t-1})_\beta \mathbb{E}[Z^{h_{t-1}^1} Z^{h_t^1}] \quad (427)$$

$$- \eta \hat{\theta}'_{W_{t-1}^4} \sum_{\beta=1}^{d_{\text{out}}} (\text{dh}_{t-1,k}^2)_\beta (\dot{\chi}_{t-1})_\beta, \quad (428)$$

$$\delta \dot{s}_t^2 = \text{T}(\mathbf{h}_t^2) - \text{T}(\mathbf{h}_{t-1}^2) \text{ if } \hat{\theta}_{\delta h_t^2} = 1, \quad (429)$$

$$\delta \dot{\mathbf{h}}_t^2 = \text{G}(\mathbf{h}_t^2, \mathbf{s}_t) - \text{G}(\mathbf{h}_{t-1}^2, \mathbf{s}_{t-1}) \text{ if } \hat{\theta}_{\delta h_t^2} = 1, \quad (430)$$

$$Z^{\delta h_t^{3,k}} = \hat{\theta}_{\delta h_t^1 / \delta h_t^{3,k}} Z^{W_0^{3,k} \delta h_t^1} - \sum_{s=0}^{t-2} \hat{\theta}'_{\delta W_s^{3,k} \delta h_t^1 / \delta h_t^{3,k}} \eta \sum_{\beta}^{d_{\text{out}}} Z^{(\text{dh}_s^{3,k})_\beta} (\dot{\chi}_s)_\beta \mathbb{E}[Z^{h_s^1} Z^{\delta h_t^1}] \quad (431)$$

$$- \hat{\theta}'_{\delta W_t^{3,k} / \delta h_t^{3,k}} \eta \sum_{\beta}^{d_{\text{out}}} Z^{(\text{dh}_{t-1}^{3,k})_\beta} (\dot{\chi}_{t-1})_\beta \mathbb{E}[Z^{h_{t-1}^1} Z^{h_t^1}], \quad (432)$$

$$Z^{\delta h_t^3} = \sum_{k=1}^m \hat{\theta}_{\delta h_t^{3,k} / \delta h_t^3} \hat{h}_{k,t-1}^2 Z^{\delta h_t^{3,k}} - \hat{\theta}_{\delta \bar{h}_t^2 / \delta h_t^3} \delta \dot{h}_{t,k}^2 Z^{h_t^{3,k}}, \quad (433)$$

$$\delta \dot{f}_t = \hat{\theta}'_{W_{t-1}^4} \hat{\theta}_{\delta h_t^3} \mathbb{E}[Z^{\overline{W_{t-1}^4}} Z^{\delta h_t^3}] - \hat{\theta}'_{\delta W_t^4} \eta \dot{\chi}_{t-1} \mathbb{E}[Z^{h_{t-1}^3} Z^{h_t^3}], \quad (434)$$

where G is either softmax or sigmoid and G' is the derivative, and

$$Z^{W_0^{3,k} \delta h_t^1} = \hat{Z}^{W_0^{3,k} \delta h_t^1} + \dot{Z}^{W_0^{3,k} \delta h_t^1} \quad (435)$$

$$= \hat{Z}^{W_0^{3,k} \delta h_t^1} + \sum_{s=0}^{t-1} Z^{\text{dh}_s^1} \mathbb{E} \left[\frac{\partial Z^{\text{dh}_s^{3,k}}}{\partial \hat{Z}^{(W_0^{3,k})^T \text{dh}_s^{3,k}}} \right] \quad (436)$$

Remark 32 ($\mu\text{P-MoE}$ gives feature learning and stability) *As the infinite-width limit does not change much for the Top-K case. It is easy to see that the $\mu\text{P-MoE}$ values in Table 1 for SGD and Table 2 for Adam hold for Top-K MoEs.*

Appendix G. Expressing softmax and sigmoid over a multi-scalar as a moment operation in TP

We can express the sum of l scalars $\theta_1, \dots, \theta_k \in \mathbb{R}$ using the **Moment** operations with no vectors ($k = 0$):

$$\phi : \mathbb{R}^0 \times \mathbb{R}^l \rightarrow \mathbb{R}, \quad \phi(-; \theta_1, \dots, \theta_k) = \sum_{i=1}^k \theta_i.$$

Then

$$\frac{1}{n} \sum_{\alpha=1}^n \phi(-; \theta_1, \dots, \theta_l) = \frac{1}{n} \sum_{\alpha=1}^n \left(\sum_{i=1}^l \theta_i \right) = \sum_{i=1}^l \theta_i,$$

since the expression inside the sum is constant with respect to α .

Softmax: We can do the same with a softmax. Simply define the function ϕ as the softmax function over the scalars.

$$\phi : \mathbb{R}^0 \times \mathbb{R}^l \rightarrow \mathbb{R}, \quad \phi(-; \theta_1, \dots, \theta_k) = \frac{\exp \theta_i}{\sum_{i=1}^k \exp \theta_i}.$$

Then

$$\frac{1}{n} \sum_{\alpha=1}^n \phi(-; \theta_1, \dots, \theta_l) = \frac{1}{n} \sum_{\alpha=1}^n \frac{\exp \theta_i}{\sum_{i=1}^k \exp \theta_i} = \frac{\exp \theta_i}{\sum_{i=1}^k \exp \theta_i},$$

so each of these operations generates a valid scalar. **Sanity check.** Note that if a sequence of random vectors (multi-scalars in TP language)

$$\theta^{(n)} = (\theta_1^{(n)}, \theta_2^{(n)}, \dots, \theta_l^{(n)}) \longrightarrow \theta^{(*)} = (\theta_1^{(*)}, \theta_2^{(*)}, \dots, \theta_l^{(*)})$$

where the convergence may be in probability, distribution, or almost sure, then the continuous mapping theorem states that if $\phi : \mathbb{R}^l \rightarrow \mathbb{R}$ is continuous at $\theta^{(*)}$, then

$$\phi(\theta^{(n)}) = \phi(\theta_1^{(n)}, \theta_2^{(n)}, \dots, \theta_l^{(n)}) \longrightarrow \phi(\theta^{(*)}) = \phi(\theta_1^{(*)}, \theta_2^{(*)}, \dots, \theta_l^{(*)})$$

in the same mode of convergence (probability, distribution, a.s). So continuity at the limit alone suffices for a.s convergence to hold.

TP requires that the non-linearities have polynomially bounded derivatives which are of course much stronger and clearly satisfied by the softmax.

Sigmoid: We can simply express the sigmoid of a scalar with the function ϕ defined as

$$\phi : \mathbb{R}^0 \times \mathbb{R} \rightarrow \mathbb{R}, \quad \phi(-; \theta_1) = \frac{1}{1 + \exp(-\theta_1)}.$$

Then

$$\frac{1}{n} \sum_{\alpha=1}^n \phi(-; \theta_1) = \frac{1}{1 + \exp(-\theta_1)}.$$

Appendix H. Additional Experiments and Details

H.1. Details for Figures

Figures in the main paper: All figures in the main paper are created by training for $t = 20$ steps on Fashion-MNIST. The network used is the same as Theorem 1 but with non-linear experts with ReLU activation. The models contained 4 experts, and $K = 2$ were executed in the Top- K case. Hyperparameters of the model such as multiplier (for μ P-MoE) and learning rate were tuned so that the figures showed effective learning. The router and function output are zero-initialised to show limiting behaviour (Yang et al., 2021).

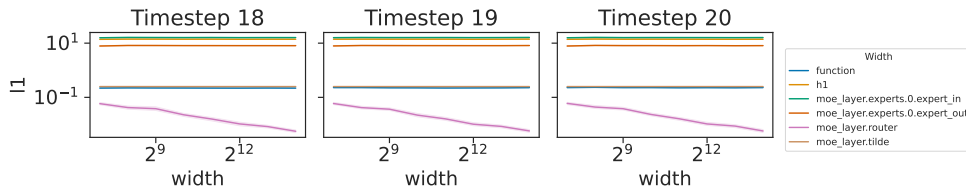


Figure 3: Coordinate check for Softmax-Soft MoE. Shows that the router feature scales $n^{-0.5}$ as predicted in Theorem 25. Tilde signifies the router probability.

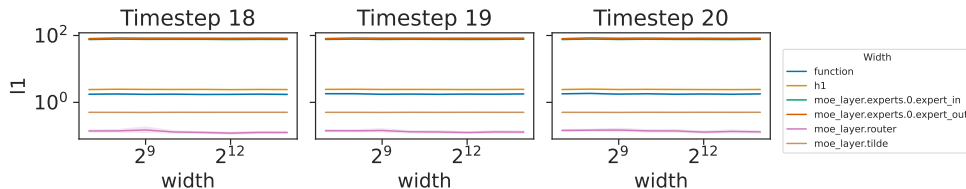


Figure 4: Coordinate check for Sigmoid-Soft MoE. Shows that the router feature learns in comparison with Softmax gating. Tilde signifies the router probability.

HP Transfer in Appendix H.4: HP transfer of learning rate in Appendix H.4 was carried out using a Mixtral style transformer language model (Jiang et al., 2024) trained on Wikitext-103-v1 for 1 epoch. The model was made smaller due to compute constraints. The model contained three transformer layers with feedforward width of 4 times the width dimension. A default load balancing loss coefficient of 0.01 was fixed. Hyperparameters such as the multipliers (for output, router, inputs) and learning rate were randomly sampled 100 times and were tuned based on the validation loss. The model consisted of 8 experts and 2 experts were executed per token. QK Layernorm was implemented as it has been shown to be more stable for HP transfer (Haas et al., 2025).

Compute Resources: The experiments for this paper (including the Mixtral run) were run on a single H100 GPU.

H.2. Coordinate Checks for Soft Routing

We show additional coordinate checks in for sigmoid soft in Figure 4 and softmax soft in Figure 3.

H.3. Coordinate Checks for Top- K Routing

We show additional coordinate checks for sigmoid Top- K in Figure 6 and Figure 5. We additionally verify that μ P-MoE scaling (Figure 8) allows for feature learning in Mixtral style transformer MoEs compare to SP-MoE (Figure 7).

H.4. Additional HP Transfer Experiments

We show in Algorithm 1 how to transfer hyperparameters conditional on noise terms/bias.

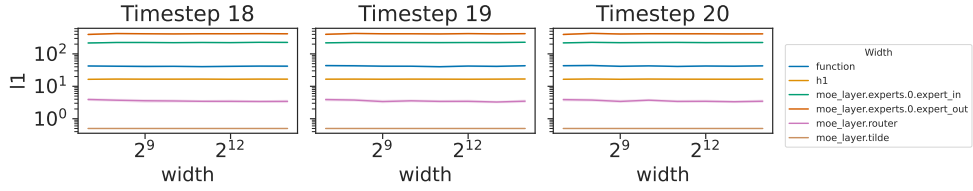


Figure 5: Coordinate check for Softmax-Top- K MoE with 2/4 active experts per forward pass. Shows that in comparison to the soft routing case, the Top- K operation can allow for expert specialisation and hence route feature learning. Tilde signifies the router probability.

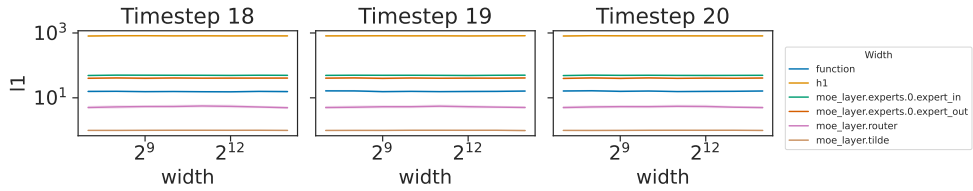


Figure 6: Coordinate check for Sigmoid-Top- K MoE with 2/4 active experts per forward pass. Shows that like in the soft routing case, the router exhibits feature learning. Tilde signifies the router probability.

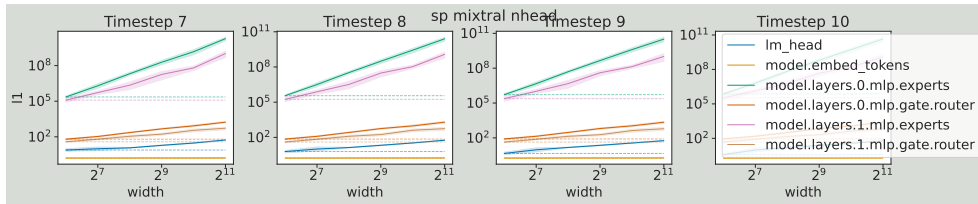


Figure 7: Coordinate check for Top- K Mixtral MoE with SP. Shows that the network diverges with increasing width.

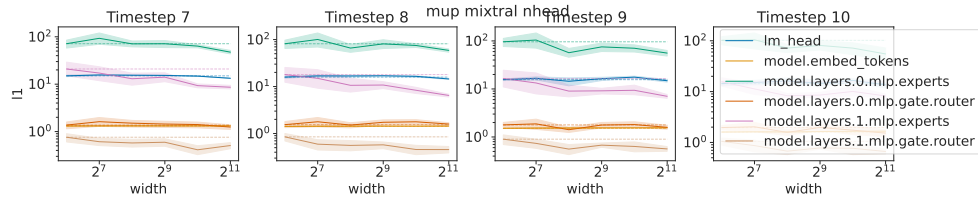


Figure 8: Coordinate check for Top- K Mixtral MoE with μ P-MoE. Shows that the network feature learns as with increases.

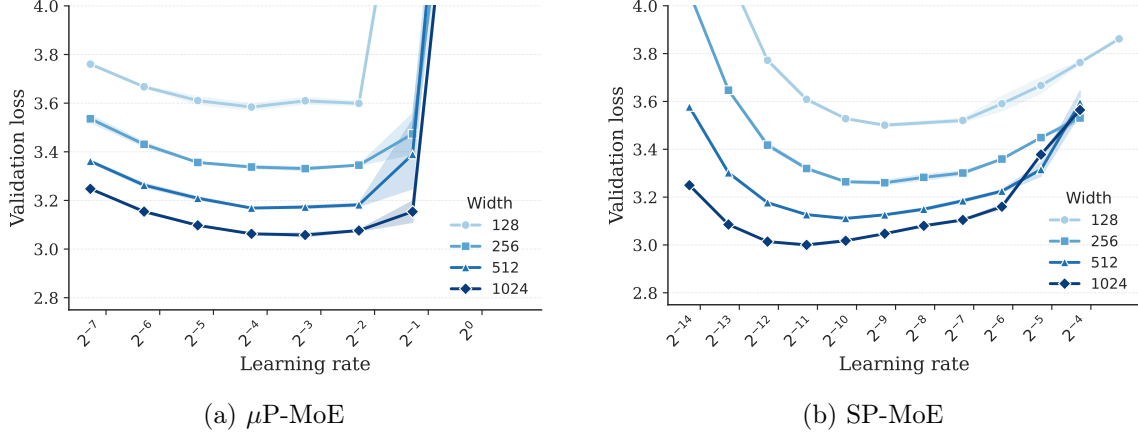


Figure 9: Learning rate transfer for Top- K MoE models parametrised according to (a) μ P-MoE and (b) SP-MoE. Hyperparameters are tuned on the smallest model.

Algorithm 1: μ P HP-transfer for Top- K MoE with a shared router-bias schedule

Input: base width n_0 , target width n , steps T , experts m , Top- K , bias scale σ_{b^2}

Bias schedule (fix once, reuse for all widths)

| Sample $\{b_t^2\}_{t=0}^T$ with $b_t^2 \sim \mathcal{N}(0, \sigma_{b^2} I_m)$ using a fixed seed.

end

Training loop at width $w \in \{n_0, n\}$, given HPs $\mathcal{H}(\cdot)$

| **for** $t = 0$ to T **do**

| | $h_t^2 \leftarrow W_t^2 h_t^1 + b_t^2$ $I_t \leftarrow \text{TopK}(h_t^2)$ Route to experts in I_t and update parameters.

| **end**

end

Tune at width n_0

| Run the training loop with μ P-heuristic scaling \Rightarrow tuned HPs $\mathcal{H}(n_0)$.

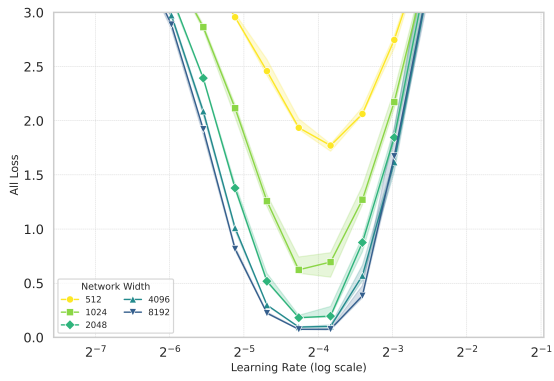
end

Transfer to width n

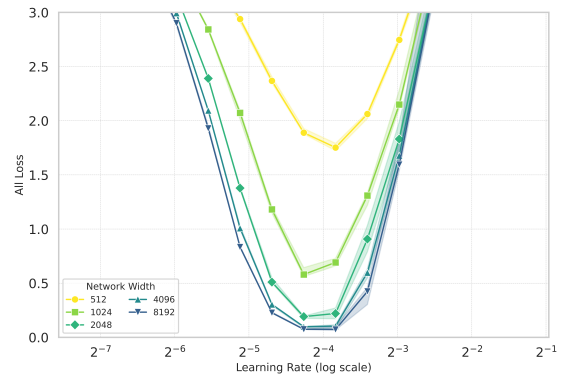
| Initialise width- n model with μ P-heuristic and $\mathcal{H}(n)$; run the same training loop.

end

We carry out additional HP transfer with a similar network to Theorem 1, but with non-linear experts and Top-2 out of 4 expert activation. In Figure 10 we show HP transfer with and without the bias term on tinyimagenet. In Figure 11 we show HP transfer with and without the bias term on emnist.

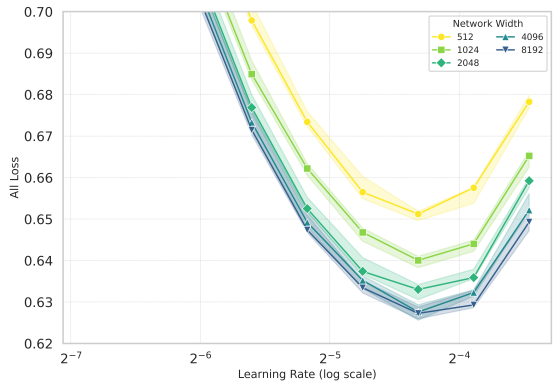


(a) With Bias

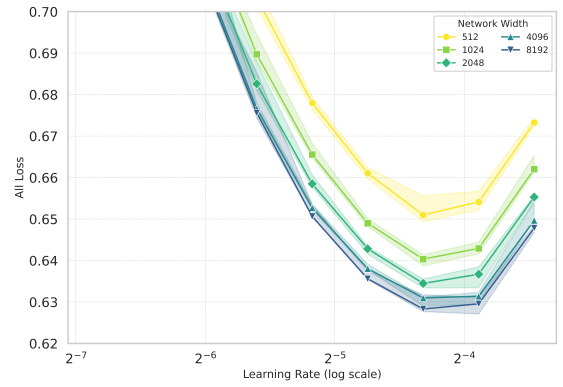


(b) Without Bias

Figure 10



(a) With Bias



(b) Without Bias

Figure 11