

# How to Scale Mixture-of-Experts: From $\mu$ P to the Maximally Scale-Stable Parameterization

author names withheld

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

## Abstract

Recent frontier large language models predominantly rely on Mixture-of-Experts (MoE) architectures. Despite empirical progress, there is still no principled understanding of how hyperparameters should scale with network width  $N$ , expert width  $N_e$ , number of experts  $M$ , sparsity  $K$ , and depth  $L$  to ensure both stability and optimal performance at scale. We take a principled step toward resolving this gap by analyzing three different scaling regimes: (I) co-scaling  $N \asymp N_e$ , (II) co-scaling  $N \asymp M \asymp K$ , and (III) full proportional scaling of  $N$ ,  $N_e$ ,  $M$ , and  $K$ . For each regime, we develop a novel Dynamical Mean Field Theory (DMFT) description of the limiting training dynamics of MoEs that provides a formal foundation for our analysis. Within this framework, we derive the unique parameterization for SGD and Adam satisfying all maximal-update ( $\mu$ ) desiderata. We then show that the resulting  $\mu$ P prescription does not reliably induce monotonic improvement with scale or robust learning-rate transfer. We trace these pathologies to scale-dependent observables in the aggregation dynamics, which motivates a refined set of desiderata that we term *maximal scale stability*. Guided by this principle, we derive a *Maximally Scale-Stable Parameterization* (MSSP) for both SGD and Adam in all three scaling regimes, and characterize the corresponding limiting dynamics - qualitatively distinct from the  $\mu$ P limit - through a separate DMFT analysis. Experiments verify that MSSP robustly recovers learning rate transfer and monotonic improvement with scale across regimes. Combined with existing depth-scaling theory, these results provide a complete scaling prescription for MoE architectures as a function of width, depth, expert width, and number of experts.

## 1. Introduction

While scale-invariant parametrizations for dense models, such as the Maximal Update Parameterization ( $\mu$ P), are well-developed [13, 15, 58], modern frontier large language models increasingly employ Mixture-of-Experts (MoE) layers [19, 27], which decouple parameter count from computational cost via sparse routing [47]. In doing so, they introduce additional scaling axes that must be coordinated jointly with width  $N$  and depth  $L$ : *the number of experts  $M$ , the expert width  $N_e$ , and the routing sparsity  $K/M$* . Recent engineering advances render scaling along all of these dimensions practically tractable [43], and a growing body of evidence indicates that increasingly fine-grained experts outperform configurations comprising a few large experts [11, 12, 30]. Yet, little is known about how to optimally scale fine-grained, sparse MoE architectures, which motivates the question:

*How should MoE architectures and training hyperparameters be scaled to yield scale-invariant, non-degenerate feature and prediction dynamics in various co-scaling regimes of  $M$ ,  $N$ ,  $N_e$ , and  $K$ ?*

As we will see in this work, developing scaling theories for MoEs is substantially more involved than for dense networks. One cannot assume commutativity across axes ( $M, N, N_e, K$ ); instead, joint limits governed by their relative rates must be analyzed. Routing and aggregation in MoEs couple

the dynamics of the router, the experts, and update statistics [16, 47]. This leads to a combinatorial space of co-scaling regimes, each, as we demonstrate, exhibiting qualitatively distinct behaviour.

**Main Contributions.** We study three co-scaling regimes for MoEs: (Regime I)  $N \asymp N_e \rightarrow \infty$  with  $M, K$  fixed; (Regime II)  $N \asymp M \asymp K \rightarrow \infty$  with  $N_e$  fixed — the fine-grained / bottleneck regime increasingly favored in modern MoE designs; and (Regime III) joint scaling  $N \asymp N_e \asymp M \asymp K \rightarrow \infty$ . Across these regimes, we make the following contributions:

- **$\mu$ P for MoEs across regimes and optimizers.** In each scaling regime, we derive the parameterizations that satisfy the  $\mu$ P principles for SGD and Adam via signal propagation analyses which are formally justified by a novel Dynamical Mean Field Theory (DMFT) for each scaling regime.
- **Scale-dependence of  $\mu$ P in MoEs.** Despite formally satisfying the  $\mu$ -desiderata, we find that  $\mu$ P does not reliably deliver learning-rate transfer or monotonic improvement with scale in MoEs. We trace this to scale-dependence in the training dynamics.
- **Maximally Scale-Stable Parameterization (MSSP).** We propose *maximal scale stability* as a refined principle generalizing  $\mu$ P, and derive the Maximally Scale-Stable Parameterization (MSSP) for SGD and Adam in each regime. The required corrections to  $\mu$ P are structurally distinct across regimes: zero router initialization in Regime I, amplified expert-output initialization variance ( $1/N_e \rightarrow M/N_e$ ) in Regime II, and shared expert weights at initialization in Regime III.
- **DMFT under MSSP.** We derive self-consistent DMFT equations for the limiting training dynamics in each regime under MSSP, including a four-level conditional mean-field hierarchy in Regime III induced by the shared expert initialization, qualitatively distinct from the DMFT under  $\mu$ P.
- **Empirical validation.** We verify on MLP and Transformer MoEs that MSSP robustly outperforms  $\mu$ P, recovers learning-rate transfer, and restores monotonic improvement with scale across all three regimes. We provide a complete scaling prescription for modern Transformer MoE architectures along width  $N$ , depth  $L$ , number of experts  $M$ , expert width  $N_e$  and number of active experts  $K$ .

Overall, this paper takes important steps towards stable, predictable and optimal MoE scaling that preserves reliable hyperparameter transfer and monotonic improvement with scale.

**Independent and concurrent work.** Jiang et al. [28] concurrently derive a  $\mu$ P for Sign SGD in Regime III. We address SGD and Adam across all three regimes, identify that  $\mu$ P alone does not yield learning rate transfer or monotonic improvement with scale in Regimes II and III, and propose MSSP as a resolution. A detailed comparison is given in App. H. Limitations and avenues for future work are discussed in App F.

## 2. Setting: Architecture and Scaling Regimes

**MoE architecture.** Given input  $x_t \in \mathbb{R}^D$ , the residual stream  $h_t^0 = W_t^{\text{in}} x_t$  is transformed as:

$$h_t^l = h_t^{l-1} + K^{-\alpha_{\text{agg}}} \sum_{i \in \text{top-}K} \phi_{t,i}^l \cdot W_t^{l,\text{out},i} \varphi(W_t^{l,\text{in},i} h_t^{l-1}), \quad \phi_{t,i}^l = \sigma(\beta (Q_t^l h_t^{l-1})_i),$$

for  $l \in [L]$ , and is finally transformed into output logits  $f_t = (W_t^{\text{out}})^\top h_t^L$ . Here  $\varphi$  is a coordinate-wise nonlinearity,  $\beta$  a tunable inverse temperature, and  $\sigma$  either sigmoid or softmax. Our scaling prescriptions are *local* to the MoE block and apply unchanged when stacked with attention, residual connections, and normalization (with non-MoE components parameterized via  $\mu$ P/mean-field), as validated by our Transformer MoE experiments in Appendices Q.1 and Q.2.

**Coordinate-wise update rules.** We consider optimizers of the form  $W_t = W_{t-1} - \eta_t \Psi_t(g_0, \dots, g_t)$  with  $g_t = \nabla_W \mathcal{L}_t$  and  $\Psi_t$  acting entrywise [57]; this covers SGD and Adam.

**Scaling axes, regimes, and parameterizations** We study joint scaling in width  $N$ , expert width  $N_e$ , number of experts  $M$ , and active experts  $K$  ( $1 \leq K \leq M$ ). All  $\Theta, \mathcal{O}, \Omega$  are taken as  $n \rightarrow \infty$  along a trajectory  $\mathcal{S}(n) = (N, N_e, M, K)(n)$ . We organize results by three regimes: (I) *fixed number of experts*:  $M, K = \Theta(1)$  and  $N \asymp N_e \rightarrow \infty$ ; (II) *infinite number of experts with fixed expert width*:  $N_e = \Theta(1)$  and  $N \asymp M \asymp K \rightarrow \infty$ ; (III) *joint proportional scaling*:  $N \asymp N_e \asymp M \asymp K \rightarrow \infty$ . Here we focus on Regimes II and III; the analysis and empirics for Regime I are collected in App. C. Depth requires separate analytical treatment; we defer depth scaling to the end of Section 4. Unless stated otherwise, we treat all other quantities as fixed  $\Theta(1)$  quantities.

**$bc\alpha$ -parametrization.** A  $bc\alpha$ -parametrization fixes exponents  $(b_W, c_W, d_W)$  for each trainable tensor  $W$  and an aggregation exponent  $\alpha_{\text{agg}} \in [0, 1]$  such that  $W_0 \sim \mathcal{N}(0, n^{-2b_W})$ ,

$$W_t = W_{t-1} - \eta n^{-c_W} \Psi_t(n^{d_W} g_0(W), \dots, n^{d_W} g_t(W)), \quad h_t^l = n^{-\alpha_{\text{agg}}} \sum_{i=1}^M \phi_{t,i}^l h_{t,i}^{l,\text{out}},$$

with  $\Psi_t$  applied entrywise. Our goal is to identify a set of exponents  $\{\alpha_{\text{agg}}\} \cup \{b_W, c_W, d_W\}_W$  that yields stable, performant dynamics at large model scale.

**Experiment Setup.** We train MLP MoEs on TinyImageNet [32] and GPT MoEs on Dolma3 [39]. All details can be found in Appendix P and additional experiments in Appendix Q.

### 3. Maximal Update Desiderata for MoEs and their Shortcomings

The maximal-update ( $\mu$ ) desiderata have served as the primary principled scaling rules for dense networks [7, 13, 53, 56]. Hence they serve as a natural starting point for MoE scaling. Here, we will derive  $\mu$ P for each MoE scaling regime, and uncover its pathologies when applied to MoE scaling.

**Norm choice.** For any vector, matrix or tensor  $T$ , we measure average entry size with the RMS norm  $\|T\|_{\text{RMS}} := \left(\frac{1}{|T|} \sum_{i \in \text{entries}} T_i^2\right)^{1/2}$ , and measure spectral properties of matrices  $W$  using the operator norm when equipping input and output space with  $\|\cdot\|_{\text{RMS}}$ :  $\|W\|_{\text{op}} := \sup_{x \neq 0} \frac{\|Wx\|_{\text{RMS}}}{\|x\|_{\text{RMS}}}$ .

**Recapitulating the Maximal Update Desiderata.** Intuitively, the  $\mu$ -desiderata require  $\Theta(1)$  feature updates that propagate forward through all linear layers without vanishing or diverging. For Adam, faithfulness ensures  $\epsilon$  scales as the layerwise gradient RMS norm.

**Desideratum  $\mu$ -1: Stability and feature learning.** For every layer (including router and expert sub-layers), activations and their updates neither vanish nor diverge with model scale:  $\|h_t\|_{\text{RMS}} = \mathcal{O}(1)$  for all  $t \geq 0$  and  $\|\Delta h_{t+1}\|_{\text{RMS}} = \Theta(1)$  for some  $t \geq 0$ .

**Desideratum  $\mu$ -2: Maximal effective and propagating updates.** For every trainable weight  $W_t$  and input activations  $x_t$ , write  $h_t := W_t x_t$ . The update  $\Delta h_t$  decomposes into an *effective update* (from updates to the current layer’s parameters) and a *propagating update* (from updates to upstream layers):  $\Delta h_t = \Delta W_t x_t + W_0 \Delta x_t$ . Both contributions are  $\Theta(1)$  at some  $t \in \mathcal{O}(1)$ .

**Desideratum  $\mu$ -3: Faithfulness.** The input to the update function  $\Psi_t$  is  $\Theta(1)$  at some  $t > 0$ .

A forward and backward signal propagation analysis yields  $\mu$ P in each regime (App. J). For SGD the scaling is formally justified by a corresponding DMFT analysis (App. K–App. O).

**Result 1 ( $\mu$ -parameterization for MoEs).** For each scaling regime (I–III) and each optimizer (SGD, Adam), there exists a unique  $bc\alpha$ -parameterization under which the MoE training dynamics satisfy Desiderata  $\mu$ -(1-3). The exponents  $(b_W, c_W, d_W, \alpha_{\text{agg}})$  are summarized in App. B.1.

**Empirical degeneracies of  $\mu$ P in MoEs.** We find that  $\mu$ P *can fail* empirically in Regimes II and III: (i) Effective and propagating updates do not remain  $\Theta(1)$  at initialization, and can take many steps to

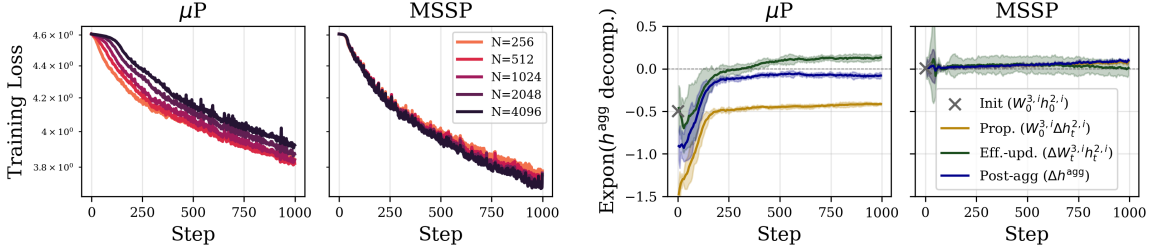


Figure 1: **Delayed learning and scale dependent dynamics of  $\mu P$  resolved by MSSP (SGD, Regime II).** Training loss (left) is worse at large scale in  $\mu P$  (the darker, the wider), but monotonically improves in MSSP. Scaling exponents of sub-terms of the aggregated MoE activations  $h_t^l$  (right) are approximately 0 in all time steps in MSSP, signaling scale-independent training dynamics. In  $\mu P$ , initially vanishing sub-terms cause cascading time-dependent scale dependence of all terms. Even after 1000 updates, post-aggregation propagating updates still vanish in  $\mu P$ , as we predict.

drift toward their predicted asymptotic values (Figure Q.26). (ii) performance often degrades with scale (Fig. 1 and 2, App. Q.4.2 and Q.4.3) and (iii) learning rates that are optimal at small scales do not remain optimal at larger scales (App. Q.3.2 and Q.3.3). *In short,  $\mu P$  does not reliably deliver the scaling benefits in MoEs that motivate it in dense networks.*

**Diagnosing the mechanism for failure in Regimes II and III.** The post-aggregation activations  $h_t^{l+1} = h_t^l + h_t^{agg}$  admit the decomposition

$$h_t^{agg} = \frac{1}{M} \sum_{i=1}^M \phi_i^l \cdot \left( \underbrace{W_0^{l,out,i} h_{0,i}^{l,in}}_{\text{init.}} + \underbrace{\Delta W_t^{l,out,i} h_{t,i}^{l,in}}_{\text{effective}} + \underbrace{W_0^{l,out,i} \Delta h_{t,i}^{l,in}}_{\text{propagating}} \right). \quad (\text{MoE})$$

A single principle organizes the scaling of each term (explained in Appendix D): the cross-expert average is  $\Theta(1)$  (LLN-like) when the per-expert summands share a coherent direction, and  $\Theta(1/\sqrt{M})$  (CLT-like) when their directions are independent across experts. App. J shows that the training-induced update  $\Delta h_t^{agg}$ , contains unbalanced contributions in both regimes, most pronounced in Regime II. If, at finite scale,  $\Delta h_t^{agg}$  is dominated by the CLT-suppressed contributions, feature learning vanishes early, recovers slowly, and performance may degrade with scale (Fig. 1 and 2).

#### 4. A More Fundamental Desideratum Beyond Maximal Updates: Scaling MoEs Requires Maximal Scale Stability

In Section 3, we traced the failure modes of  $\mu P$  to scale-dependent contributions arising in the training dynamics. This motivates a more general and fundamental desideratum, which we call *maximal scale stability*: *every primitive interaction obtained by decomposing each weight as  $W = W_0 + \Delta W$  and propagating this split through the network should have a  $\Theta(1)$  impact on the forward and (appropriately normalized) backward dynamics.* For most settings of interest, this reduces to a compact set of operational conditions; we instantiate the desideratum for MoEs as follows.

##### Maximal Scale Stability Desiderata for MoEs.

1. *Forward.* The maximal update Desiderata  $\mu$ -1,  $\mu$ -2,  $\mu$ -3 hold.
2. *Backward.* Along any linear map  $h_t^\ell = W_t^\ell x_t^{\ell-1}$ , analogous to activations, gradients admit the transpose recursion  $\bar{\delta}_t^{\ell-1} = (W_t^\ell)^\top \delta_t^\ell$  and consequently the analogous effective/propagating decomposition  $\bar{\delta}_t^{\ell-1} = (W_0^\ell)^\top \delta_t^\ell + (\Delta W_t^\ell)^\top \delta_t^\ell$ , where  $\bar{\delta}_t^\ell := \nabla_{x_t^\ell} \mathcal{L}$  and  $\delta_t^\ell := \nabla_{h_t^\ell} \mathcal{L}$ . Both contributions are balanced and remain  $\Theta(1)$  under appropriate normalization.

3. *Forward and Backward Aggregation.* Each of the init, propagating, and effective contributions of the cross-expert aggregation (MoE) is  $\Theta(1)$ . The analogous requirement applies to the expert-aggregated gradient at the shared input  $h^l$  (see Eq. (I.6)).

We refer to a parameterization that satisfies these conditions as the *Maximally Scale-Stable Parameterization* (MSSP). For linear layers with  $\text{fan-in} \asymp \text{fan-out} \rightarrow \infty$ , forward and backward desiderata coincide by the symmetry of the Gaussian initialization. MSSP desiderata reduce to  $\mu$ -desiderata in dense architectures such as MLPs and Transformers. But in general, forward scale stability does not imply backward scale stability as we will see below.

**MSSP for MoEs.** As for  $\mu$ P, signal propagation in MSSP is analyzed in App. J, with a corresponding DMFT in App. L and N, and qualitative insights from these DMFT limits in App. E.

**Result 2 (MSSP for MoEs).** The parameterization summarized in App. B.1 satisfies all MSSP desiderata in Regimes I and III. In Regime II, it satisfies Desiderata 1–3 except that the propagating update of  $h^{3,i}$  is of order  $\Theta(\sqrt{M})$ ; however, its impact on the aggregate dynamics remains  $\Theta(1)$ .

*Regime II.* MSSP rectifies the unbalanced  $\mu$ P dynamics with a single intervention: **the expert output initialization variance is amplified from  $1/N_e$  to  $M/N_e$** . Although failure modes in  $\mu$ P arise through structurally distinct mechanisms, this variance amplification simultaneously rebalances every forward and backward cross-expert sub-aggregate, every previously sub-leading contribution to per-expert gradients as well as to the gradient flowing through the router (cf. tables in App. J.3.2 for  $\mu$ P versus App. J.3.4 for MSSP). The only scale-dependent quantity under MSSP is the per-expert propagating update  $W_0^{l,\text{out},i} h_t^{l,\text{in},i} = \Theta(\sqrt{M})$ , which is impossible to avoid. We argue that this divergence is benign since its impact on the training dynamics is  $\Theta(1)$ , and no training observable inherits this divergence. The dynamics admits a well-defined DMFT limit under MSSP (App. L). We empirically verify that all relevant terms remain scale-stable over time (Figure Q.26, App. Q.4.2). Consequently, MSSP-Regime-II restores monotonic performance improvement with scale and learning rate transfer from small to large scale for both SGD and Adam (Fig. 2, App. Q.1 and Q.3.2).

*Regime III.* The viable fix in Regime III is structural: **share the expert initialization,  $W_0^{1,\text{out},i} = W_0^{1,\text{out}}$  and  $W_0^{1,\text{in},i} = W_0^{1,\text{in}}$  for all  $i$**  and let the routing mechanism to diversify experts over time. This restores complete scale stability (see theory App. J.4.5 and experiment App. Q.4.3). Performance of SGD and Adam on TinyImageNet reliably improves with scale only under MSSP but not  $\mu$ P (Fig. 2), all relevant quantities remain  $\Theta(1)$  throughout training (see App. Q.4.3), and LR reliably transfers across scales for both MLP MoEs and GPT MoEs (App. Q.1).

Hence, both regimes call for structurally distinct interventions: increased expert output initialization in Regime II, and shared initialization in Regime III. Neither solution transfers to the other regime. A more detailed, self-contained but intuitive mechanism-by-mechanism account of how the MSSP prescription rectifies each  $\mu$ P failure is provided in App. I.

**Scaling depth  $L$ .** The analysis of Bordelon et al. [6] extends with minimal modification to MoEs with residual scaling  $\Theta(1/L)$ . Our complete MoE scaling prescription is summarized in Table B.2.

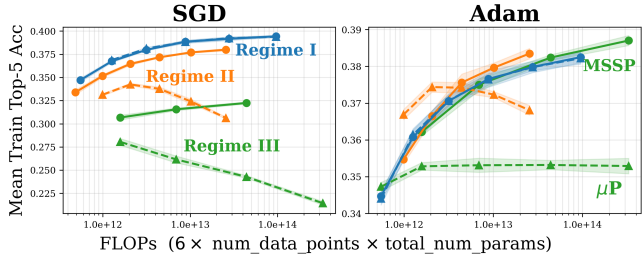


Figure 2: **MSSP outperforms  $\mu$ P.** Across optimizers and regimes, MSSP (solid lines) outperforms  $\mu$ P (dashed lines) at large scale. All details in App. P.3.

## References

- [1] Greg W. Anderson, Alice Guionnet, and Ofer Zeitouni. *An Introduction to Random Matrices*, volume 118 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, 2010.
- [2] Zhidong Bai and Jack W. Silverstein. *Spectral Analysis of Large Dimensional Random Matrices*. Springer Series in Statistics. Springer, 2 edition, 2010.
- [3] Charlie Blake, Constantin Eichenberg, Josef Dean, Lukas Balles, Luke Yuri Prince, Björn Deiseroth, Andres Felipe Cruz-Salinas, Carlo Luschi, Samuel Weinbach, and Douglas Orr.  $u\text{-}\mu$ p: The unit-scaled maximal update parametrization. In *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025. URL <https://openreview.net/forum?id=P7KRiiLM8T>.
- [4] Enric Boix-Adsera and Philippe Rigollet. The power of fine-grained experts: Granularity boosts expressivity in mixture of experts. *arXiv preprint arXiv:2505.06839*, 2025.
- [5] Blake Bordelon and Cengiz Pehlevan. Self-consistent dynamical field theory of kernel evolution in wide neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 35: 32240–32256, 2022.
- [6] Blake Bordelon, Hamza Tahir Chaudhry, and Cengiz Pehlevan. Infinite limits of multi-head transformer dynamics. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024. URL <https://openreview.net/forum?id=p0BBKhD5aI>.
- [7] Blake Bordelon, Lorenzo Noci, Mufan Bill Li, Boris Hanin, and Cengiz Pehlevan. Depth-wise hyperparameter transfer in residual networks: Dynamics and scaling limit. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=KZJehvRKGD>.
- [8] Louis-Pierre Chaintron, Lénaïc Chizat, and Javier Maas. Resnets of all shapes and sizes: Convergence of training dynamics in the large-scale limit. *arXiv preprint arXiv:2603.18168*, 2026.
- [9] Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding mixture of experts in deep learning. *arXiv preprint arXiv:2208.02813*, 2022.
- [10] Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- [11] Aidan Clark, Diego De Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann, Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George Bm Van Den Driessche, Eliza Rutherford, Tom Hennigan, Matthew J Johnson, Albin Cassirer, Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol Vinyals, Marc’Aurelio Ranzato, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen Simonyan. Unified scaling laws for routed language models. In *Proceedings of the 39th International*

*Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*. PMLR, 2022.

- [12] Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Yu Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- [13] 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:2505.01618*, 2025.
- [14] David Eigen, Marc’Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013.
- [15] Katie Everett, Lechao Xiao, Mitchell Wortsman, Alexander A Alemi, Roman Novak, Peter J Liu, Izzeddin Gur, Jascha Sohl-Dickstein, Leslie Pack Kaelbling, Jaehoon Lee, et al. Scaling exponents across parameterizations and optimizers. *arXiv:2407.05872*, 2024.
- [16] 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):1–39, 2022.
- [17] Nikhil Ghosh, Denny Wu, and Alberto Bietti. Understanding the mechanisms of fast hyperparameter transfer. *arXiv preprint arXiv:2512.22768*, 2025.
- [18] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010.
- [19] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [20] 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. URL <https://openreview.net/forum?id=hTxnm6H93P>.
- [21] Soufiane Hayou and Greg Yang. Width and depth limits commute in residual networks. In *International Conference on Machine Learning (ICML)*, pages 12700–12723. PMLR, 2023.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *IEEE international conference on computer vision (ICCV)*, pages 1026–1034, 2015.
- [23] Xu Owen He. Mixture of a million experts. *arXiv preprint arXiv:2407.04153*, 2024.
- [24] J. Hubbard. Calculation of partition functions. *Physical Review Letters*, 3(2):77–78, 1959.
- [25] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.

- [26] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural Tangent Kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8571–8580, 2018.
- [27] 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.
- [28] Tianze Jiang, Blake Bordelon, Cengiz Pehlevan, and Boris Hanin. Hyperparameter transfer with mixture-of-expert layers. *arXiv preprint arXiv:2601.20205*, 2026.
- [29] Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 6(2):181–214, 1994.
- [30] Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michał Krutul, Szymon Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdź, Piotr Sankowski, et al. Scaling laws for fine-grained mixture of experts. *arXiv preprint arXiv:2402.07871*, 2024.
- [31] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [32] Yann Le, Xuan Yang, et al. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [34] Jan Małaśnicki, Kamil Ciebiera, Mateusz Boruń, Maciej Pióro, Jan Ludziejewski, Maciej Stefaniak, Michał Krutul, Sebastian Jaszczur, Marek Cygan, Kamil Adamczewski, et al. mu-parametrization for mixture of experts. *arXiv preprint arXiv:2508.09752*, 2025.
- [35] V. A. Marchenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mat. Sbornik (N.S.)*, 72(114)(4):507–536, 1967.
- [36] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [37] Mansour Zoubeirou Mayaki. Generalization and scaling laws for mixture-of-experts transformers, 2026. Manuscript.
- [38] 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):E7665–E7671, 2018.
- [39] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. Olmoe: Open mixture-of-experts language models. *arXiv preprint arXiv:2409.02060*, 2024.
- [40] Radford M. Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages 29–53. Springer New York, 1996.

- [41] Lorenzo Noci, Chuning Li, Mufan Li, Bobby He, Thomas Hofmann, Chris J Maddison, and Dan Roy. The shaped transformer: Attention models in the infinite depth-and-width limit. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024.
- [42] Lorenzo Noci, Alexandru Meterez, Thomas Hofmann, and Antonio Orvieto. Super consistency of neural network landscapes and learning rate transfer. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- [43] NVIDIA. Scalable training of mixture-of-experts models with megatron core. *arXiv preprint arXiv:2603.07685*, 2026.
- [44] Team Olmo, Allyson Ettinger, Amanda Bertsch, Bailey Kuehl, David Graham, David Heineman, Dirk Groeneveld, Faeze Brahman, Finbarr Timbers, Hamish Ivison, Jacob Morrison, Jake Poznanski, Kyle Lo, Luca Soldaini, Matt Jordan, Mayee Chen, Michael Noukhovitch, Nathan Lambert, Pete Walsh, Pradeep Dasigi, Robert Berry, Saumya Malik, Saurabh Shah, Scott Geng, Shane Arora, Shashank Gupta, Taira Anderson, Teng Xiao, Tyler Murray, Tyler Romero, Victoria Graf, Akari Asai, Akshita Bhagia, Alexander Wettig, Alisa Liu, Aman Rangapur, Chloe Anastasiades, Costa Huang, Dustin Schwenk, Harsh Trivedi, Ian Magnusson, Jaron Lochner, Jiacheng Liu, Lester James V. Miranda, Maarten Sap, Malia Morgan, Michael Schmitz, Michal Guerquin, Michael Wilson, Regan Huff, Ronan Le Bras, Rui Xin, Rulin Shao, Sam Skjonsberg, Shannon Zejiang Shen, Shuyue Stella Li, Tucker Wilde, Valentina Pyatkin, Will Merrill, Yapei Chang, Yuling Gu, Zhiyuan Zeng, Ashish Sabharwal, Luke Zettlemoyer, Pang Wei Koh, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. Olmo 3, 2025. URL <https://arxiv.org/abs/2512.13961>.
- [45] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [46] Babak Shahbaba and Radford Neal. Nonlinear models using dirichlet process mixtures. *Journal of Machine Learning Research*, 10(8), 2009.
- [47] 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 (ICLR)*, 2017.
- [48] Charles M. Stein. Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, 9(6):1135–1151, 1981.
- [49] Wenfang Sun, Xinyuan Song, Pengxiang Li, Lu Yin, Yefeng Zheng, and Shiwei Liu. The curse of depth in large language models. *arXiv preprint arXiv:2502.05795*, 2025.
- [50] Volker Tresp. Mixtures of gaussian processes. *Advances in neural information processing systems*, 13, 2000.

- [51] 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 (NeurIPS)*, 2024. URL <https://openreview.net/forum?id=aQv5AbN1wF>.
- [52] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- [53] Nikhil Vyas, Alexander Atanasov, Blake Bordelon, Depen Morwani, Sabarish Sainathan, and Cengiz Pehlevan. Feature-learning networks are consistent across widths at realistic scales. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024.
- [54] Mitchell Wortsman, Peter J Liu, Lechao Xiao, Katie E Everett, Alexander A Alemi, Ben Adlam, John D Co-Reyes, Izzeddin Gur, Abhishek Kumar, Roman Novak, Jeffrey Pennington, Jascha Sohl-Dickstein, Kelvin Xu, Jaehoon Lee, Justin Gilmer, and Simon Kornblith. Small-scale proxies for large-scale transformer training instabilities. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=d8w0pmvXbZ>.
- [55] Greg Yang. Wide feedforward or recurrent neural networks of any architecture are gaussian processes. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- [56] Greg Yang and Edward J. Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning (ICML)*, 2021.
- [57] Greg Yang and Etai Littwin. Tensor programs ivb: Adaptive optimization in the infinite-width limit. *arXiv:2308.01814*, 2023.
- [58] 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. *arXiv:2203.03466*, 2022.
- [59] Greg Yang, James B Simon, and Jeremy Bernstein. A spectral condition for feature learning. *arXiv preprint arXiv:2310.17813*, 2023.
- [60] 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.

# Appendices

**Appendix Contents.**

<b>I Quick Reference</b>	<b>15</b>
<b>A Appendix Roadmap and Suggested Reading Paths</b>	<b>15</b>
<b>B Complete Maximally Scale-Stable Parameterization for Mixtral-Style MoEs</b>	<b>17</b>
<b>C Regime I: Empirics and MSSP Correction</b>	<b>19</b>
<b>D Diagnosing the Mechanism for Failure in Regimes II and III</b>	<b>19</b>
<b>E Summary of our DMFT Analyses for MoE training dynamics</b>	<b>21</b>
<b>F Discussion and Future Work</b>	<b>22</b>
<b>II Background and Extended Related Work</b>	<b>23</b>
<b>G Background on Width-scaling Parameterizations</b>	<b>23</b>
<b>H Detailed Related Work</b>	<b>24</b>
<b>III Theory</b>	<b>26</b>
<b>I Intuitive Explanations of the Shortcomings of <math>\mu</math>P for MoEs and How MSSP Fixes Them</b>	<b>26</b>
I.1 Overview . . . . .	26
I.2 Per-layer alignment: scaling of linear (forward) maps $W \cdot x$ . . . . .	27
I.3 MSSP Desiderata . . . . .	30
I.4 Backward scaling . . . . .	31
I.5 Aggregation layers . . . . .	31
I.6 Width dependence in $\mu$ P and how MSSP fixes it . . . . .	31
I.7 Cross-expert aggregation: a taxonomy of mechanisms . . . . .	32
I.8 Summary: regime imbalances under $\mu$ P and how MSSP resolves them . . . . .	37
I.9 Justification of MSSP in Regime II . . . . .	38
I.10 Composability with depth: stacked MoE blocks under MSSP-Regime-II . . . . .	39
I.11 Approaches We Tried That Do Not Restore Scale Stability . . . . .	41
<b>J Signal Propagation Analysis</b>	<b>43</b>
J.1 Preliminaries . . . . .	43
J.2 Standing assumptions and notation . . . . .	44
J.3 Scaling derivation for Regime II . . . . .	44
J.4 Scaling derivation for Regime III . . . . .	70

<b>K</b>	<b>DMFT Analysis for Regime I</b>	<b>89</b>
K.1	Setup . . . . .	89
K.2	Dynamics . . . . .	91
K.3	Mean field theory . . . . .	94
K.4	Averages of dual variables . . . . .	103
K.5	The coupling kernels A and B . . . . .	106
K.6	Final DMFT . . . . .	106
<b>L</b>	<b>DMFT Analysis for Regime II (MSSP)</b>	<b>109</b>
L.1	Architectural definitions . . . . .	109
L.2	Learning rates and initialization . . . . .	109
L.3	Gradient definitions . . . . .	109
L.4	DMFT kernels and order parameters . . . . .	110
L.5	Learning Dynamics and the Neural Tangent Kernel . . . . .	111
L.6	Decomposition of the MoE NTK . . . . .	111
L.7	Evolution of weights, preactivations, and pregradients . . . . .	112
L.8	Stochastic initial fields . . . . .	114
L.9	Deriving the DMFT Action . . . . .	116
L.10	Softmax . . . . .	119
L.11	Partition function . . . . .	119
L.12	Saddle point approximation . . . . .	121
L.13	Saddle-Point Equations . . . . .	122
L.14	Expert-local order parameters . . . . .	123
L.15	Hubbard-Stratonovich transformation . . . . .	124
L.16	DMFT Dynamics . . . . .	124
<b>M</b>	<b>DMFT Analysis for Regime II (<math>\mu</math>P)</b>	<b>126</b>
<b>N</b>	<b>DMFT Analysis for Regime III (MSSP)</b>	<b>128</b>
N.1	Architectural Definitions . . . . .	128
N.2	Learning Rates and Initialization . . . . .	128
N.3	Gradient Definitions . . . . .	128
N.4	DMFT Kernels and Order Parameters . . . . .	129
N.5	Learning Dynamics and the Neural Tangent Kernel . . . . .	130
N.6	Decomposition of the MoE NTK . . . . .	130
N.7	Evolution of weights, preactivations, and pregradients . . . . .	131
N.8	Stochastic Initial Fields . . . . .	133
N.9	Deriving the DMFT Action . . . . .	135
N.10	Softmax . . . . .	137
N.11	Partition function . . . . .	137
N.12	Saddle Point Approximation . . . . .	141
N.13	Saddle-Point Equations . . . . .	142
N.14	Expert-Local Order Parameters . . . . .	143
N.15	Hubbard-Stratonovich transformation . . . . .	144
N.16	DMFT Dynamics . . . . .	144

<b>O</b>	<b>DMFT Analysis for Regime III (<math>\mu</math>P)</b>	<b>146</b>
<b>IV</b>	<b>Experiments</b>	<b>149</b>
<b>P</b>	<b>Experimental Setup</b>	<b>149</b>
P.1	MLP MoE experiments on TinyImagenet . . . . .	149
P.2	Transformer MoE experiments . . . . .	150
P.3	Figure details . . . . .	152
<b>Q</b>	<b>Additional Experiments</b>	<b>153</b>
Q.1	Learning rate sweeps for Transformer MoEs . . . . .	153
Q.2	Refined coordinate checks for Transformer MoEs . . . . .	153
Q.3	Learning rate sweeps for MLP MoEs . . . . .	156
Q.4	Fine-grained scaling evaluations in MLP MoEs . . . . .	160
Q.5	Soft softmax routing collapses to uniform for $\mu$ P in Regime I . . . . .	176
Q.6	MoEs require layerwise learning rate tuning . . . . .	178
Q.7	Random search and 2D multiplier tuning do not suffice . . . . .	179
Q.8	Global Adam $\varepsilon$ induces width dependence at sufficient scale . . . . .	180

## Part I

# Quick Reference

This part serves as a concise entrypoint to the appendix that enables readers to orient themselves quickly and to navigate the subsequent parts in a structured way. Appendix A provides a roadmap and suggested reading paths depending on the reader’s background and interests.

For practitioners, Appendix B concisely provides our *complete MSSP scaling prescription for modern MoE architectures*.

Then we summarize our empirical and theoretical findings in Regime I, concisely explain the failure mechanism of  $\mu$ P in Regimes II and III, and summarize our DMFT analyses.

### A. Appendix Roadmap and Suggested Reading Paths

The remaining parts of this appendix are self-contained and can be summarized follows:

- **Part II (Background and Extended Related Work)** provides background on the established infinite width literature, in particular the most common *abc*-parameterizations for standard dense networks, necessary to understand our corrected parameterizations for MoEs (Appendix G) and a detailed account of related work (Appendix H).
- **Part III (Theory)** provides three complementary levels of formalism: an intuitive scaling explanation (Appendix I), followed by a full forward and backward signal propagation analysis for deriving the  $\mu$ P and MSSP hyperparameter scaling rules in each scaling regime (Appendix J) followed by a rigorous self-consistent DMFT for  $\mu$ P and MSSP in each of the three scaling regimes (Appendices K to O). Note that Regimes II and III provide differing challenges that require their own resolutions, both in terms of scaling and analysis technique. Whenever the DMFT equations for a differing parameterization our routing mechanism follows from an analogous derivation, we omit the duplication and just provide the final set of self-consistent DMFT equations.
- **Part IV (Experiments)** contains the experimental setup for our MLP MoE and Transformer MoE experiments (Appendix P), followed by extensive empirical evidence for our claims about learning rate transfer and scaling properties of  $\mu$ P and MSSP for both SGD and Adam in all 3 scaling regimes (Appendix Q). We close with further empirical scaling insights that practitioners should be aware of in Sections Q.5 to Q.8.

**If you only have 45 minutes:** read (i) the remaining subsections of this part for concise summaries of our empirical and theoretical insights in each regime, and (ii) Appendix I.8 for a more detailed regime-by-regime summary of how  $\mu$ P breaks and how MSSP fixes it.

**For conceptual understanding:** We recommend familiarizing oneself with the background on width scaling (Appendix G), then reading Appendix I for a comprehensive introduction into the scaling arguments necessary for understanding the signal propagation pathologies of MoEs and how to fix them. These arguments will provide the reader with tools that are more broadly applicable to similar scaling derivations for other non-standard architectures and optimizers. We then recommend reading the complete MSSP recipe for modern MoEs (Appendix B), and ending with a quick pass over the empirical results in Appendix Q. In particular, we recommend skipping Appendices J to P during this first read.

**Deep dive:** For the ambitious reader, we recommend starting with the above path for generating sufficient conceptual understanding, before diving into the detailed signal propagation analyses in Appendix J, and the DMFT analyses in Appendices K to O. The signal propagation analysis decomposes all terms in the training dynamics and determines their scaling with arguments from random matrix theory. This makes the order of contributions of all observables explicit. The DMFT provides complementary value by clarifying how different variables interact via evolution of their kernels.

**Practitioners:** For readers that are primarily interested in the practical takeaways, we recommend starting with our complete Mixtral-style MSSP recipe in Appendix B. Our Transformer setup is introduced in Appendix P.2 and our empirical findings are provided in Appendix Q.

**Notation mapping.** The notation in the appendix differs slightly from the main paper. For readability, we write out the theory for a single block, but, as is common in the related DMFT and  $\mu$ P literature [28, 58], every derivation follows in the same way for the architecture that stacks multiple blocks that is considered in the main paper. A detailed explanation is provided in Appendix I.10.

Recall the following main paper notation.

**MoE architecture.** Given input  $x_t \in \mathbb{R}^D$ , the residual stream  $h_t^0 = W_t^{\text{in}} x_t$  is transformed as:

$$h_t^l = h_t^{l-1} + K^{-\alpha_{\text{agg}}} \sum_{i \in \text{top-K}} \phi_{t,i}^l \cdot W_t^{l,\text{out},i} \varphi(W_t^{l,\text{in},i} h_t^{l-1}), \quad \phi_{t,i}^l = \sigma(\beta (Q_t^l h_t^{l-1})_i),$$

for  $l \in [L]$ , and is finally transformed into output logits  $f_t = (W_t^{\text{out}})^\top h_t^L$ . Here  $\varphi$  is a coordinatewise nonlinearity,  $\beta$  a tunable inverse temperature, and  $\sigma$  either sigmoid or softmax.

For theoretical purposes, it suffices to analyze the signal propagation through the following single MoE block without residual connection.

**Minimal MoE architecture.** For readability, we work with the minimal MoE architecture

$$h^1 = W^1 x, \tag{A.1}$$

$$\psi = Q h^1, \tag{A.2}$$

$$\phi = \sigma(\psi), \tag{A.3}$$

$$h^{2,i} = W^{2,i} h^1, \quad i = 1, \dots, M, \tag{A.4}$$

$$h^{3,i} = W^{3,i} h^{2,i}, \tag{A.5}$$

$$h^3 = \frac{1}{M} \sum_{i=1}^M \phi_i h^{3,i}, \tag{A.6}$$

$$f = W^4 h^3, \tag{A.7}$$

where  $x \in \mathbb{R}^D$  is the input,  $f \in \mathbb{R}$  the scalar output, and the trainable parameters are

$$\theta = (W^1, Q, \{W^{2,i}, W^{3,i}\}_{i=1}^M, W^4)$$

with shapes  $W^1 \in \mathbb{R}^{N \times D}$ ,  $Q \in \mathbb{R}^{M \times N}$ ,  $W^{2,i} \in \mathbb{R}^{N_e \times N}$ ,  $W^{3,i} \in \mathbb{R}^{N \times N_e}$ ,  $W^4 \in \mathbb{R}^{1 \times N}$ . The same scaling exponents hold under sigmoid gating with aggregation factor  $\alpha_{\text{agg}} = 1$  and under softmax gating without an explicit aggregation multiplier.

In particular, here we write  $h^1$  instead of  $h^l$ , and  $h^3 = \frac{1}{M} \sum_{i=1}^M \phi_i h^{3,i}$  instead of  $h^{l+1} = h^l + h^{\text{agg}}$ .

## B. Complete Maximally Scale-Stable Parameterization for Mixtral-Style MoEs

Here, we provide our  $\mu$ P and MSSP scaling rules for scaling MoEs trained with SGD and Adam in Table B.1, and then provide our complete MSSP scaling prescription for scaling Adam training of GPT MoEs with width  $N$ , depth  $L$ , expert width  $N_e$ , number of experts  $M$  and number of active experts  $K$ .

Group	Qty.	Regime I ( $N, N_e \asymp n \rightarrow \infty$ )		Regime II ( $N, M, K \asymp n \rightarrow \infty$ )		Regime III ( $N, N_e, M, K \asymp n \rightarrow \infty$ )	
		SGD	Adam	SGD	Adam	SGD	Adam
Router $Q$	Init $\sigma$	$N^{-1}   0$	$N^{-1}   0$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$
	LR $\eta$	$N^{-1}$	$N^{-1}$	$MN^{-1}$	$N^{-1}$	1	$N^{-1}$
	Adam $\epsilon$	-	1	-	$M^{-1}$	-	$M^{-1}$
Expert in	Init $\sigma$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}   (N^{-1/2})^{\text{tied}}$	$N^{-1/2}   (N^{-1/2})^{\text{tied}}$
	LR $\eta$	1	$N^{-1}$	$MN^{-1}$	$N^{-1}$	$M$	$N^{-1}$
	Adam $\epsilon$	-	$N^{-1}$	-	$M^{-1}$	-	$M^{-1}N^{-1}$
Expert out	Init $\sigma$	$N_e^{-1/2}$	$N_e^{-1/2}$	$N_e^{-1/2}   M^{1/2}N_e^{-1/2}$	$N_e^{-1/2}   M^{1/2}N_e^{-1/2}$	$N_e^{-1/2}   (N_e^{-1/2})^{\text{tied}}$	$N_e^{-1/2}   (N_e^{-1/2})^{\text{tied}}$
	LR $\eta$	1	$N_e^{-1}$	$MN$	$N_e^{-1}$	$M$	$N_e^{-1}$
	Adam $\epsilon$	-	$N^{-1}$	-	$M^{-1}N^{-1}$	-	$M^{-1}N^{-1}$

Table B.1:  $\mu$ P and MSSP for SGD and Adam across scaling regimes. Scaling prescriptions as a function of  $(N, N_e, M, K)$  for initialization standard deviation  $\sigma$ , learning rate  $\eta$ , and Adam  $\epsilon$  for each MoE weight matrix, using expert aggregation factor  $\alpha^{\text{agg}} = 1$  for sigmoid routers. Entries where  $\mu$ P and MSSP differ are highlighted as  $\mu$ P | MSSP. All other weight matrices should be scaled as in  $\mu$ P. In addition, decoupled weight decay [33] and Adam’s  $\beta_1, \beta_2$  should be scale-independent. *tied*: weights should be shared across experts *at initialization*.

Table B.2: **MoE AdamW hyperparameter scaling.** Columns correspond to three MoE scaling regimes. **Width** and **depth** control terms are highlighted. For example in the bottleneck Regime II, width  $N$ , number of experts  $M$  and top- $K$  are all scaling proportionally, thus could be replaced by the width multiplier  $m_N = N/N_0$  in relation to some base width  $N_0$ . In AdamW as proposed by Loshchilov and Hutter [33], weight decay should stay fixed across model scales. In the PyTorch implementation of AdamW, the weight decay multiplier should always be the inverse of the learning rate multiplier so that the effective weight decay  $\eta \cdot \text{wd}$  stays scale-independent. The minimal fine-graining of multiplier tuning at small model size recommended includes a base residual multiplier  $\alpha_{\text{base}}$ , a global initialization multiplier, and learning rate multipliers based on layer type (embedding, readout, gain, router, expert in, expert out, other).

\*: Router zero initialization requires initial randomness in the routing mechanism.

*tied*: **In Regime III, expert weights should be shared at initialization.**

SCALING MOES: FROM  $\mu$ P TO MSSP

Parameterization	Regime I	Regime II	Regime III
	$N, N_e \rightarrow \infty$ $M, K$ fixed	$N, M, K \rightarrow \infty$ $N_e$ fixed	$N, M, K, N_e \rightarrow \infty$
<b>Emb. Init. Std.</b>	$d_{\text{in}}^{-1/2}$	$d_{\text{in}}^{-1/2}$	$d_{\text{in}}^{-1/2}$
<b>Emb. Adam LR</b>	$d_{\text{in}}^{-1}$	$d_{\text{in}}^{-1}$	$d_{\text{in}}^{-1}$
<b>Emb. Adam <math>\epsilon</math></b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Pre-LN Init. Std.</b>	1	1	1
<b>Pre-LN Adam LR</b>	1	1	1
<b>Pre-LN Adam <math>\epsilon</math></b>	$N^{-1}L^{-1}$	$N^{-1}L^{-1}$	$N^{-1}L^{-1}$
<b>Hidden Init. Std.</b>	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$
<b>Hidden Adam LR</b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Hidden Bias Adam LR</b>	1	1	1
<b>Hidden Adam <math>\epsilon</math></b>	$N^{-1}L^{-1}$	$N^{-1}L^{-1}$	$N^{-1}L^{-1}$
<b>MoE routing &amp; experts</b>			
<b>Router (gating) Init. Std.</b>	0*	$N^{-1/2}$	$N^{-1/2}$
<b>Router (gating) Adam LR</b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Router Adam <math>\epsilon</math></b>	$L^{-1}$	$M^{-1}L^{-1}$	$M^{-1}L^{-1}$
<b>Expert Layer 1 Init. Std.</b>	$N^{-1/2}$	$N^{-1/2}$	$(N^{-1/2})^{\text{tied}}$
<b>Expert Layer 1 Adam LR</b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Expert Layer 1 Adam <math>\epsilon</math></b>	$N^{-1}L^{-1}$	$M^{-1}L^{-1}$	$N^{-1}M^{-1}L^{-1}$
<b>Expert Layer 2 Init. Std.</b>	$N_e^{-1/2}$	$M^{1/2}N_e^{-1/2}$	$(N_e^{-1/2})^{\text{tied}}$
<b>Expert Layer 2 Adam LR</b>	$N_e^{-1}$	$N_e^{-1}$	$N_e^{-1}$
<b>Expert Layer 2 Adam <math>\epsilon</math></b>	$N^{-1}L^{-1}$	$N^{-1}M^{-1}L^{-1}$	$N^{-1}M^{-1}L^{-1}$
<b>Aggregation multiplier</b>	$K^{-1}$	$K^{-1}$	$K^{-1}$
<b>Aux load-balancing loss multiplier</b>	1	1	1
<b>Router z-loss multiplier</b>	1	1	1
<b>MHA Residual</b>	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MHA(LN(X^l))$	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MHA(LN(X^l))$	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MHA(LN(X^l))$
<b>MoE FFN Residual</b>	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MoE(LN(X^l))$	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MoE(LN(X^l))$	$X^l + \alpha_{\text{base}} \cdot L^{-1} \cdot$ $MoE(LN(X^l))$
<b>Final-LN Init. Std.</b>	-	-	-
<b>Final-LN Adam LR</b>	1	1	1
<b>Final-LN Adam <math>\epsilon</math></b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Unemb. Init. Std.</b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Unemb. Adam LR</b>	$N^{-1}$	$N^{-1}$	$N^{-1}$
<b>Unemb. AdamW <math>\epsilon</math></b>	1	1	1
<b>Unemb. Fwd.</b>	1	1	1

### C. Regime I: Empirics and MSSP Correction

This appendix collects the empirical observations and the MSSP correction for Regime I ( $N, N_e \rightarrow \infty$  with  $M, K = \Theta(1)$ ), which were deferred from the main text because  $\mu\text{P}$  already broadly delivers the expected scaling behaviour in this regime, and the MSSP correction is marginal in practice.

**Empirics of  $\mu\text{P}$  in Regime I.** Under  $\mu\text{P}$ , training in Regime I broadly behaves as expected: loss improves monotonically with scale, and learning-rate transfer holds across scales (Section Q.3.1). One empirical anomaly stands out: propagating updates of the router typically vanish with scale (Fig. Q.22), even though our derivation predicts them to be  $\Theta(1)$ . We address this discrepancy with the MSSP correction below.

**MSSP correction in Regime I.**  $\mu\text{P}$  satisfies our operational MSSP desiderata for MoEs in Regime I. However, splitting  $\Delta h^{l-1}$  into the contributions of the router-gradient and the expert-gradient pathways reveals that the former is correlated with  $Q_0$  and contributes coherently at  $\Theta(1)$ , while the latter is weakly correlated with  $Q_0$  and is CLT-suppressed at  $\Theta(1/\sqrt{N})$ . To recover residual balance in the router’s propagating update  $Q_0^l \Delta h^{l-1}$ , **MSSP initializes the router to zero**. The empirical impacts of this adjustment are marginal (Fig. 2, App. Q.4.1).

**Relation to the shape-based  $\mu\text{P}$ -heuristic.** For dense networks, it is common to scale the layerwise initialization variance and learning rate of a weight matrix purely based on its shape [58, 59], distinguishing input-like (fixed  $\rightarrow \infty$ ), hidden-like ( $\infty \rightarrow \infty$ ), and output-like ( $\infty \rightarrow$  fixed) layers. Małaśnicki et al. [34] apply this heuristic to MoEs, treating the router output-like under  $M, K \in \Theta(1)$ . Table B.1 shows that the  $\mu\text{P}$ -heuristic coincides with  $\mu\text{P}$  in Regime I (whereas it fails to induce all  $\mu$ -desiderata in Regimes II and III, see Section 3 of the main text).

**Pointers.** DMFT for Regime I under  $\mu\text{P}$ : App. K; learning-rate transfer plots: App. Q.3.1; coordinate checks: App. Q.4.1.

### D. Diagnosing the Mechanism for Failure in Regimes II and III

This appendix provides a mechanism-by-mechanism walkthrough of the cross-expert aggregation imbalances summarized in §3 of the main text. To understand these empirical failures under  $\mu\text{P}$ , examine the aggregation over experts unique to MoEs. The post-aggregation activations  $h_t^{l+1}$  can be written as  $h_t^l + h_t^{\text{agg}}$ , where  $h_t^{\text{agg}}$  admits the decomposition (MoE), reproduced here for convenience:

$$h_t^{\text{agg}} = \frac{1}{M} \sum_{i=1}^M \phi_i^l \cdot \left( \underbrace{W_0^{l,\text{out},i} h_{0,i}^{l,\text{in}}}_{\text{init.}} + \underbrace{\Delta W_t^{l,\text{out},i} h_{t,i}^{l,\text{in}}}_{\text{effective}} + \underbrace{W_0^{l,\text{out},i} \Delta h_{t,i}^{l,\text{in}}}_{\text{propagating}} \right).$$

A single principle organizes the scaling properties of each of these quantities. The cross-expert average exhibits LLN-like behaviour and contributes at  $\Theta(1)$  whenever the per-expert summands share a *coherent direction across experts*; when the per-expert summands have independent random directions, the corresponding aggregate exhibits CLT-like behaviour which suppresses a factor of  $\sqrt{M}$ . The underlying mechanism that produces or destroys this coherence varies across these aggregate terms and across regimes, so they can occupy distinct orders in  $N$ , and the order of  $h_t^{l+1}$  at any finite scale is set by whichever sub-term dominates. We now provide intuition for how each of these aggregates behaves. We recommend non-technical readers to skip to ‘*Empirical consequences*’.

*Init aggregate.* Each per-expert summand  $\phi_i^l \cdot W_0^{l,\text{out},i} h_{0,i}^{l,\text{in}}$  is a chain  $W_0^{l,\text{out},i} W_0^{l,\text{in},i}$  of two independent Gaussian matrices acting on the shared input  $x_0^l$ . By rotation invariance, the cross-expert dependence between summands is mediated by its scaled norm,  $\|x_0^l\|^2/M$ , which as  $M \rightarrow \infty$  converges to a deterministic scalar. Since summands are asymptotically i.i.d., the init aggregate exhibits CLT-like behaviour and is  $\Theta(1/\sqrt{M})$  in both Regimes II and III.

*Propagating aggregate.* The propagating aggregate admits a further decomposition. At leading order,  $\Delta h_t^{l,\text{in},i} \approx W_0^{l,\text{in},i} \Delta x_t^l + \Delta W_t^{l,\text{in},i} x_t^l$ , splitting the propagating aggregate into two contributions:  $\frac{1}{M} \sum_i \phi_i^l W_0^{l,\text{out},i} W_0^{l,\text{in},i} \Delta x_t^l$  and  $\frac{1}{M} \sum_i \phi_i^l W_0^{l,\text{out},i} \Delta W_t^{l,\text{in},i} x_t^l$ . The former is governed by the same CLT mechanism as the init aggregate, with one caveat:  $\Delta x_t^l$  depends on each  $W_0^{l,\text{in},i}$  through backpropagation. However, a single vector in  $\mathbb{R}^N$  cannot align nontrivially with  $M \rightarrow \infty$  independent random Gaussian chains simultaneously. The per-expert summands therefore remain approximately independent across experts, and the same CLT argument applies. This contribution is therefore also of order  $\Theta(1/\sqrt{M})$  in both Regimes II and III. The second contribution behaves differently in different scaling regimes. Substituting the rank-one gradient  $\Delta W_t^{l,\text{in},i} \propto (W_0^{l,\text{out},i})^\top \delta_t^{l+1} (x_t^l)^\top$  (where  $\delta_t^{l+1} = \nabla_{h_t^{l+1}} \mathcal{L}$  represents the backpropagated error signal from the subsequent layer) yields the form  $\frac{1}{M} \sum_{i=1}^M \phi_i^l G_i u$ , where  $G_i := W_0^{l,\text{out},i} (W_0^{l,\text{out},i})^\top$  and  $u \propto \delta_t^{l+1}$  is shared across experts.

- (a) *Regime II.* When  $N_e$  is held fixed,  $G_i$  has rank  $N_e = \Theta(1)$  and after rescaling, it is the orthogonal projector onto a  $O(1)$  dimensional subspace of  $\mathbb{R}^N$  that is Haar-distributed on the Grassmannian  $\text{Gr}(N_e, N)$  (i.e., uniformly random among  $N_e$ -dimensional subspaces), and these subspaces are independent across experts. Therefore, the per-expert vectors  $\{G_i u\}_{i=1}^M$  are approximately i.i.d. across experts with no coherent direction, and the cross-expert average exhibits (CLT-like behaviour). Both contributions to the propagating aggregate are therefore  $\Theta(1/\sqrt{M})$  under  $\mu$ P.
- (b) *Regime III.* When  $N_e \asymp N$ ,  $W_0^{l,\text{out},i}$  is generically full rank and the empirical spectrum of  $G_i = W_0^{l,\text{out},i} (W_0^{l,\text{out},i})^\top$  concentrates around its mean [35]. Consequently  $G_i \approx c I_N$  for a deterministic scalar  $c = \Theta(1)$ , so  $G_i u \approx c u$  for every  $i$ . The per-expert summands therefore share a coherent direction  $u$ , and this contribution exhibits LLN-like behaviour at  $\Theta(1)$  scale. Together with the  $\Theta(1/\sqrt{M})$  first contribution, the propagating aggregate contains a strictly subleading contribution and exhibits a milder version of the same CLT/LLN imbalance. Moreover, the init term in both regimes exhibits CLT behaviour so remains subleading.

*Effective aggregate.* Substituting the rank-one gradient  $\Delta W_t^{l,\text{out},i} \propto \delta_t^{l+1} (h_{t-1}^{l,\text{in},i})^\top$  reveals that the summands of the effective aggregate share a coherent direction  $\delta_t^{l+1}$ . The effective aggregate therefore exhibits LLN-like behaviour at  $\Theta(1)$  along  $\delta_t^{l+1}$  in both Regimes II and III.

*Empirical consequences.* In summary,  $h_t^{\text{agg}}$  is composed of unbalanced contributions: the init term is of order  $\Theta(1/\sqrt{M})$  in both regimes; the effective term is of order  $\Theta(1)$  in both; the propagating term decomposes into a first contribution of  $\Theta(1/\sqrt{M})$  in both regimes, and a second contribution of  $\Theta(1/\sqrt{M})$  in Regime II but  $\Theta(1)$  in Regime III. The training-induced update  $\Delta h_t^{\text{agg}}$ , formed from the effective and propagating terms, is therefore unbalanced in both regimes, with the imbalance more pronounced in Regime II. This is confirmed empirically in Regime II under  $\mu$ P (Fig. 1): at finite scale  $\Delta h_t^{\text{agg}}$  is dominated by the scale-suppressed init and propagating terms. Feature learning is therefore vanishing early in training and takes many optimization steps to become  $\Theta(1)$  and performance consequently degrades with increasing scale (Fig. 2).

## E. Summary of our DMFT Analyses for MoE training dynamics

For each regime, we derive self-consistent equations describing the MoE training dynamics. Each admits a well-defined, scale-independent DMFT limit under both  $\mu$ P and MSSP. In each regime, the limiting dynamics under MSSP enjoy several analytical properties described below. The full set of DMFT equations is provided in App.K-N. We defer the comparison to  $\mu$ P limits to App. O.

**Regime I** ( $M, K \in \Theta(1)$ ). The DMFT has a finite-expert mean-field factorization. Residual-stream neurons and their backward signals are i.i.d. draws from a global single-site distribution; within each expert, hidden-layer neurons are i.i.d. draws from an expert-local single-site distribution. Both distributions are characterized self-consistently by a finite set of macroscopic order parameters. Because the number of experts is finite, the expert index is not replaced by a population mean field: the experts and router variables remain explicitly represented. Averages over the global distribution define the global feature and gradient kernels; Together with router variables, these order parameters close the limiting dynamics. Under soft routing, the stochastic router field vanishes in this limit, inducing deterministic routing. Without an additional symmetry-breaking mechanism (Top- $K$ , routing noise, or random router biases), the router remains uniform and experts stay identical throughout training.

**Regime II.** The simultaneous  $N, M, K \asymp n \rightarrow \infty$  limit with  $N_e \in \Theta(1)$  reveals a *two-level, conditional, mean-field hierarchy* in MSSP. Residual-stream neurons and their backward signals are i.i.d. draws from an effective global single-site distribution; experts together with their router variables are i.i.d. draws from an effective expert/router distribution. Both are self-consistently characterized by a finite set of order parameters. Since per-expert hidden width is fixed, the hidden layer of a typical expert is a finite-dimensional component of the expert/router process. Averages over the global distribution define the global feature and gradient kernels; averages over the expert/router distribution define the routed expert-level kernels. These order parameters close the limiting dynamics.

**Regime III.** The simultaneous  $N, N_e, M, K \asymp n \rightarrow \infty$  limit reveals a *four-level, conditional, mean-field hierarchy* in MSSP. Schematically,

$$\begin{aligned} \mathcal{X}_{\text{glob}} &\sim P_{\text{glob}}(\cdot; \mathcal{K}), & \mathcal{F}_{\text{sh}} &\sim P_{\text{sh}}(\cdot; \mathcal{K}), \\ \mathcal{E} \mid \mathcal{F}_{\text{sh}} &\sim P_{\text{ex/r}}(\cdot \mid \mathcal{F}_{\text{sh}}; \mathcal{K}), & \mathcal{U} \mid (\mathcal{F}_{\text{sh}}, \mathcal{E}) &\sim P_{\text{loc}}(\cdot \mid \mathcal{F}_{\text{sh}}, \mathcal{E}; \mathcal{K}), \end{aligned}$$

where  $\mathcal{X}_{\text{glob}}$  denotes the global single-site process,  $\mathcal{F}_{\text{sh}}$  the shared expert-hidden single-site process,  $\mathcal{E}$  the expert/router single-site process, and  $\mathcal{U}$  the within-expert hidden-neuron single-site process.

*Global single-site process.* Residual-stream coordinates and their backward signals are i.i.d. draws from an effective distribution characterized self-consistently by a finite set of order parameters; averages over this process define the global feature and gradient kernels.

*Shared expert-hidden single-site process.* Coordinates of *expert-averaged* hidden-layer fields, arising from expert-averaged hidden activations and gradients, are i.i.d. draws from an effective shared distribution. Averages define the shared expert-hidden kernels, which summarize the correlation structure of these fields induced by the shared initial expert weights.

*Expert/router single-site process.* Conditional on the shared expert-hidden fields, expert/router sites are i.i.d. draws from an effective expert/router distribution; averages define expert-level kernels.

*Within-expert hidden-neuron single-site process.* Conditional on the shared expert-hidden fields and the corresponding router state, neurons in the expert hidden layers are i.i.d. draws from an effective distribution; averages define the expert-local forward and backward kernels.

The four processes form a closed self-consistent system: kernels defined at each level enter as parameters of the single-site distributions at the others, and the limiting dynamics are determined by solving all four levels simultaneously.

## F. Discussion and Future Work

Scaling laws and hyperparameter transfer have become the backbone of modern model development. Yet, for MoE architectures there remained a critical gap of how to scale the optimization hyperparameters jointly with MoE model dimensions  $M, N, N_e, K, L$ . The Maximally Scale-Stable Parameterization (MSSP) introduced here closes this gap. We believe generalizing the  $\mu$ -desiderata to maximally scale-stable desiderata provides a framework for scale-invariant dynamics across a broader set of architectures, optimizers, and scaling dimensions, with practical benefits including predictable performance gains and hyperparameter transfer.

MSSP enables the first principled compute-optimal comparison of the three co-scaling regimes. Earlier work on optimal expert granularity [11, 12, 30] was conducted under suboptimal model scaling and merits reevaluation. More fundamentally, why hyperparameters transfer under scale-invariant parameterizations, in MoEs, dense networks, or beyond, remains an open theoretical problem.

The DMFT limits are of independent interest, opening routes to studying finite-size corrections, expert specialization, and the role of routing. They also provide a framework for analyzing auxiliary load-balancing and  $z$ -losses, which we verify empirically do not alter the scaling exponents but may shed light on router collapse through their interaction with the limiting dynamics.

**Limitations.** In Regime II, MSSP causes the per-expert pre-aggregation  $W_0^{l,\text{out},i} h_t^{l,\text{in},i}$  to scale as  $\Theta(\sqrt{M})$ . Though its dynamical impact remains  $\Theta(1)$ , at very large  $M$  this may raise numerical-precision concerns, which a fused kernel computing the cross-expert aggregate without materializing the per-expert quantity could address. Our DMFT derivations are at the level of physical rigor; fully rigorous proofs remain open. Of practical interest would also be studies that capture increasing sparsity  $M/K \rightarrow \infty$ .

## Part II

## Background and Extended Related Work

This part collects background on existing width-scaling parameterizations and a detailed account of related-work.

## G. Background on Width-scaling Parameterizations

Table G.1 gives an overview over the most common  $abc$ -parameterizations for training dense networks with SGD. Yang and Hu [56] use width-dependent weight multipliers  $n^{-a} \cdot W$  to scale gradients and avoid layerwise learning rates. This extra degree of freedom introduces an equivalence class of valid  $\mu\text{P}$   $abc$ -parameterizations. Mean field parameterization [5, 10, 38] is equivalent to  $\mu\text{P}$ , using the multipliers  $\text{fan-in}^{-1/2}$  except  $N^{-1}$  in the output layer. For practical readability, we remove this extra degree of freedom in the main paper by using naive weight multipliers  $a = 0$  for all trainable weights.

		Weight-multiplier version			Weight-multiplier-free version		
		Input-like	Hidden-like	Output-like	Input-like	Hidden-like	Output-like
SP	$\alpha_l \cdot W^l$ , $\alpha_l \propto$				1	1	1
	$\mathcal{N}(0, \sigma_l^2)$ , $\sigma_l \propto$		-		1	$N^{-1/2}$	$N^{-1/2}$
	$\eta_l \cdot \nabla_{W^l} \mathcal{L}$ , $\eta_l \propto$				$N^{-c}$	$N^{-c}$	$N^{-c}$
NTP	$\alpha_l \cdot W^l$ , $\alpha_l \propto$	1	$N^{-1/2}$	$N^{-1/2}$	1	1	1
	$\mathcal{N}(0, \sigma_l^2)$ , $\sigma_l \propto$	1	1	1	1	$N^{-1/2}$	$N^{-1/2}$
	$\eta_l \cdot \nabla_{W^l} \mathcal{L}$ , $\eta_l \propto$	1	1	1	1	$N^{-1}$	$N^{-1}$
$\mu\text{P}$	$\alpha_l \cdot W^l$ , $\alpha_l \propto$	$N^{1/2}$	1	$N^{-1/2}$	1	1	1
	$\mathcal{N}(0, \sigma_l^2)$ , $\sigma_l \propto$	$N^{-1/2}$	$N^{-1/2}$	$N^{-1/2}$	1	$N^{-1/2}$	$N^{-1}$
	$\eta_l \cdot \nabla_{W^l} \mathcal{L}$ , $\eta_l \propto$	1	1	1	$n$	1	$N^{-1}$

Table G.1: **(Common  $abc$ -parameterizations)** Standard parameterization (SP), neural tangent parameterization (NTP) and the maximal update parameterization ( $\mu\text{P}$ ) for MLPs trained with SGD in their multiplier version which purely adapts the architecture and allows width-independent global learning rates (*left*) and in their weight multiplier-free version (*right*). Parameterizations differ in their layerwise choice of width-dependent weight multipliers  $\alpha_l$ , initialization variances  $\sigma_l$  and learning rates  $\eta_l$ . Weight multiplier-free representatives of an  $abc$ -equivalence class purely adapt the optimization algorithm highlighting the fact that parameterizations effectively only induce layerwise learning rates. Knowing that  $\mu\text{P}$  correctly scales the updates in all layers, observe that the input- and hidden-layer learning rates in NTP induce vanishing updates. The same holds in SP when choosing  $c \geq 1$  as is necessary for avoiding logit blowup in the infinite-width limit. SP with large learning rates  $c = 1/2$  recovers stable hidden-layer feature learning despite logit divergence [20].

**$\mu\text{P}$  for Adam in standard architectures.** Adam requires the same layerwise initialization variances as SGD. For clarity, consider the weight-multiplier-free version of  $\mu\text{P}$ . Since Adam normalizes its update, but the update direction remains correlated with the incoming activations, Adam’s layerwise learning rate scaling simplifies to  $\eta(W) = 1/\text{fan-in}(W)$  for any linear weight matrix  $W : \mathbb{R}^{\text{fan-in}} \rightarrow \mathbb{R}^{\text{fan-out}}$ . For faithfulness, Adam  $\varepsilon(W)$  should always scale as the gradient RMS norm of  $W$ , that is  $\Theta(1)$  for the readout layer and  $\Theta(N^{-1})$  otherwise.

**Baseline Parameterizations: SP and the  $\mu$ P heuristic.** Standard parametrization (SP) uses Xavier/Glorot-type variance scaling [18, 22], i.e. typical dense layers take  $b_W = \frac{1}{2}$  (with common exceptions for input embeddings), and uses the same optimizer hyperparameters such as learning rate and Adam  $\varepsilon$  for all layers.

In contrast,  $\mu$ P is designed so that maximal-update behavior (and, for adaptive methods, faithfulness) hold in the infinite-width limit for dense networks [56, 57]. Practitioners often use the following  $\mu$ P heuristic: classify a weight tensor by whether it maps fixed  $\rightarrow n$  (*input-like*),  $n \rightarrow n$  (*matrix-like*), or  $n \rightarrow$ fixed (*output-like*), and assign  $(b, c, d)$  accordingly as summarized in Table G.2. For MoEs in Regime I, this heuristic treats expert MLP weights matrix-like, while router projections are output-like due to the fixed number of experts  $M$ . In Regime I,  $\mu$ P heuristic indeed satisfies all  $\mu$  desiderata, but it does not suffice in Regimes II and III, when co-scaling the number of experts  $M \rightarrow \infty$ .

Table G.2:  $\mu$ P-heuristic  $(b, c_{\text{SGD}}, c_{\text{Adam}}, d_{\text{Adam}})$  by layer type. Both SGD and Adam share the same initialization  $\mathcal{N}(0, N^{-b})$ . SGD requires learning rates  $\eta_N = N^{-c_{\text{SGD}}}$ , and Adam requires learning rates  $\eta_N = N^{-c_{\text{Adam}}}$  and gradient scaling  $N^{-d_{\text{Adam}}} \cdot \nabla_W \mathcal{L}$ .

	$b$	$c_{\text{SGD}}$	$c_{\text{Adam}}$	$d_{\text{Adam}}$
Input-like	0	-1	0	1
Matrix-like	0.5	0	1	1
Output-like	1	1	1	0

## H. Detailed Related Work

**Mixture of Experts models.** Mixture of experts models have a long history in machine learning [25], with instantiations such as hierarchical models [29], Gaussian processes [50] or Dirichlet processes [46]. After the first application in deep networks by Eigen et al. [14], Shazeer et al. [47] sparked the use of MoEs in modern architectures, using the top- $k$  operation for sparsification and trainable noise for load balancing. Fedus et al. [16] lowered communication and computational cost for improved scalability, using Kaiming initialization and simplifying the load balancing auxiliary loss. Most MoE architectures still use a width-independent initialization of 0.006 [12]. The trend in recent frontier models goes toward increased fine-graining, from 2 out of 8 in Mixtral 8x7B and 8x22B [27] to 8 out of 256 in Deepseek-R1 [19]. Zoph et al. [60] introduce the router z-loss for router logit regularization; Haas et al. [20] suggest that miss-scaling causes logit divergence, hence the additional regularization might be unnecessary in MSSP. Insightful empirical investigations can be found in the fully open OIMoE report [39]. They identify removing the need for the load balancing loss as an important future direction, since it significantly constrains the model flexibility and may prevent experts from sufficiently specializing over the course of training. To date, the theory of MoEs is very limited. Chen et al. [9] show that MoEs are able to learn data with cluster structure using two-layer nonlinear convolutional neural networks as experts.

**Scaling theory.** Infinite-width theory dates back to Jacot et al. [26], Matthews et al. [36], Neal [40]. The Tensor Program series [55–57] has served as a crucial tool for developing flexible and general infinite-width theory allowing to study non-vanishing feature learning limits [15, 20, 51]. This theory gave rise to the maximal update parameterization ( $\mu$ P) [58] under which the optimal learning rate and further dynamical properties [42] have been observed to transfer to larger model width. The mean-field parameterization is equivalent to  $\mu$ P. DMFT has proven useful as it allows to simulate from the limit [5, 6, 10, 38]. Similar techniques allow infinite depth limits [7, 8, 21, 41]. These works apply depth-dependent scaling factors like  $L^{-1/2}$  or  $L^{-1}$  [13] to the residual stream; it remains open

whether these interventions suffice. Dey et al. [13] show that  $L^{-1}$  outperforms  $L^{-1/2}$  in Transformer training, but Sun et al. [49] argue that scaling should depend on the layer index.

**Scaling mixture of experts models.** Clark et al. [11], Krajewski et al. [30] and Dai et al. [12, Figure 3] all suggest that increasingly fine-grained experts outperform fewer larger ones in terms of learning performance, with saturating gains. Boix-Adsera and Rigollet [4] prove an exponentially improved expressivity due to fine graining. Mayaki [37] provides theory showing that performance can improve either by scaling active capacity or by increasing the number of experts, depending on the dominant bottleneck. In practice, eventually the increased communication cost of routing dominates compute-optimality considerations, but increasingly efficient implementations are being developed. NVIDIA [43] provides a highly optimized distributed training framework across scaling regimes, that also allows efficient training with many fine-grained experts. He [23] proposes PEER layers that take fine graining to the extreme. However, all existing works use suboptimal model scaling, hence it will be crucial to evaluate the impact of granularity on scaling in MSSP.

**Independent and concurrent work.** Jiang et al. [28] derive a parameterization for *Sign SGD* in Regime III via a DMFT analysis, and present Transformer MoE experiments showing approximate learning rate transfer when scaling one axis at a time with all others fixed. We derive  $\mu\text{P}$  for *SGD and Adam across all three Regimes I, II, and III* using a signal propagation analysis complemented by a DMFT for each regime, and our experiments scale jointly along the axes each regime prescribes. Our  $\mu\text{P}$  in Regime III coincides with theirs up to a larger router initialization. More importantly, we identify that in Regimes II and III,  $\mu\text{P}$  principles do not suffice to robustly yield learning rate transfer or monotonic improvement with scale in MoEs, due to scale-dependent dynamics, and we propose MSSP as a resolution. In Regime III, our DMFT for MSSP is qualitatively distinct: it exhibits a four-layer mean-field hierarchy (versus a three-layer hierarchy under  $\mu\text{P}$ ), and we additionally derive novel DMFT limits for Regimes I and II under MSSP. We provide a detailed account of further related work in Appendix H.

**Under Regimes II and III ( $M, K \rightarrow \infty$ ),  $\mu\text{P}$  differs from shape-based  $\mu\text{P}$ -heuristic.** For dense networks, it is common to scale the layerwise initialization variance and learning rate of a weight matrix purely based on its shape [58, 59], distinguishing input-like (fixed  $\rightarrow \infty$ ), hidden-like ( $\infty \rightarrow \infty$ ) and output-like ( $\infty \rightarrow \text{fixed}$ ) layers. Małaśnicki et al. [34] apply this heuristic to MoEs, treating the router output-like in Regime I ( $M, K \in \Theta(1)$ ). Table B.1 shows that  $\mu\text{P}$ -heuristic coincides with  $\mu\text{P}$  in Regime I, but fails to induce all  $\mu$ -desiderata in Regimes II and III.

## Part III

# Theory

This part of the appendix presents our scaling analysis at three levels of formalism, each progressively more precise. We start with an intuitive explanation in Appendix I, then provide the detailed signal propagation analysis in Appendix J, and finally provide a full DMFT analysis for each scaling regime in Appendices K to O.

### I. Intuitive Explanations of the Shortcomings of $\mu$ P for MoEs and How MSSP Fixes Them

To set the stage for our formal results, we provide intuition for how to derive scaling rules for MoE blocks that satisfy Desiderata 1 and 2. In each scaling regime, these desiderata determine a unique  $(b, c, d, \alpha)$ -parameterization (up to the rescaling invariance of homogeneous optimizers discussed in Section 2). Table B.1 reports a canonical parameterization of this equivalence class. The SGD scaling in each regime is formally justified by the corresponding DMFT derivation. In addition, Adam scaling is heuristically derived but is directly guided by the same DMFT analysis. The conditions under which these parameterizations satisfy Desideratum 3 are discussed in Section J. Here and throughout, asymptotic notation (e.g.,  $\Theta(\cdot)$ ) is understood with respect to the joint scaling trajectory  $\mathcal{S}(n)$  as  $n \rightarrow \infty$ , specialized to the scaling regime (I–III) under consideration. For the intuitive explanation, we focus mainly on Regimes II and III. Regime I is more straightforward and is omitted for brevity.

#### I.1. Overview

Let’s start by recalling the Maximal update desiderata<sup>1</sup>.

##### I.1.1. MAXIMAL UPDATE DESIDERATA

The Maximal Update Desiderata require that in each layer, the effective and propagating updates remain  $\Theta(1)$  in RMS at some  $t > 0$ :

$$\|\mathcal{T}_{\text{eff}}^\ell(t)\|_{\text{RMS}} = \|\Delta W_t^\ell x_t^{\ell-1}\|_{\text{RMS}} = \Theta(1), \quad \|\mathcal{T}_{\text{prop}}^\ell(t)\|_{\text{RMS}} = \|W_0^\ell \Delta x_t^{\ell-1}\|_{\text{RMS}} = \Theta(1). \quad (\text{I.1})$$

For the MoE aggregation layer, we additionally require that the updates to the aggregated representation remains non-vanishing and non-diverging:

$$\|\Delta h_t^{\text{agg}}(t)\|_{\text{RMS}} = \Theta(1). \quad (\text{I.2})$$

##### I.1.2. EFFECTIVE AND PROPAGATING UPDATES

**Alignment exponents for effective and propagating updates.** In all regimes, the asymptotic scales of forward and backward effective and propagating terms are governed by the training-induced correlation structure between the two factors appearing in each product. Following Everett et al.

1. We omit an explicit discussion of the optimizer faithfulness desideratum here for brevity but we will briefly discuss  $\epsilon$  scaling.

[15], Yang et al. [59], we summarise this effect via *alignment exponents*  $p_\ell$  and  $q_\ell$  ( $p_\ell^\nabla$  and  $q_\ell^\nabla$  analogously for backward), defined by

$$\begin{aligned} \|\Delta W_t^\ell x_t^{\ell-1}\|_{\text{RMS}} &= \Theta\left((n_{\text{in}}^\ell)^{p_\ell} \|\Delta W_t^\ell\|_{\text{RMS}} \|x_t^{\ell-1}\|_{\text{RMS}}\right), \\ \|W_0^\ell \Delta x_t^{\ell-1}\|_{\text{RMS}} &= \Theta\left((n_{\text{in}}^\ell)^{q_\ell} \|W_0^\ell\|_{\text{RMS}} \|\Delta x_t^{\ell-1}\|_{\text{RMS}}\right), \end{aligned} \quad (\text{I.3})$$

where  $n_{\text{in}}^\ell := \dim(x^{\ell-1})$  denotes the fan-in of layer  $\ell$ . Informally,  $p_\ell$  quantifies the degree of alignment between  $\Delta W_t^\ell$  and the current features  $x_t^{\ell-1}$ , while  $q_\ell$  quantifies the degree of alignment between the initial weights  $W_0^\ell$  and the upstream feature change  $\Delta x_t^{\ell-1}$ .

**From  $(p_\ell, q_\ell)$  to scaling rules for optimizers.** Fix a layer  $\ell$  and suppose the alignment exponents  $(p_\ell, q_\ell)$  are known. Assume inductively that at some time  $t > 0$  the previous-layer quantities satisfy  $\|x_t^{\ell-1}\|_{\text{RMS}} = \Theta(1)$  and  $\|\Delta x_t^{\ell-1}\|_{\text{RMS}} = \Theta(1)$ . Enforcing Desideratum 1 at layer  $\ell$  then reduces, via (I.3), to choosing scales for  $\|W_0^\ell\|_{\text{RMS}}$  and  $\|\Delta W_t^\ell\|_{\text{RMS}}$  so that the effective and propagating updates remain  $\Theta(1)$ . These marginal scales are directly controlled by optimizer hyperparameters: initialization variances set the scaling of  $\|W_0^\ell\|_{\text{RMS}}$ , while learning-rate scaling sets the scaling of  $\|\Delta W_t^\ell\|_{\text{RMS}}$  once the scale of the layerwise gradients is known. For Adam, optimizer faithfulness further requires  $\epsilon$  to remain  $\Theta(1)$  relative to the typical gradient scale; accordingly, we choose  $\epsilon$  to scale with the layerwise gradient RMS norm. The same reduction applies to the backward pass, with backward alignment exponents  $(p_\ell^\nabla, q_\ell^\nabla)$  playing the role of  $(p_\ell, q_\ell)$  in setting the leading-order scales of  $(W_0^\ell)^\top \delta_t^\ell$  and  $(\Delta W_t^\ell)^\top \delta_t^\ell$ , where  $\delta_t^\ell := \nabla_{h_t^\ell} \mathcal{L}$  denotes the pre-activation gradient passed backward.

### I.1.3. DMFT PREDICTIONS.

Our DMFT analysis provides the regime-specific asymptotic inputs needed to instantiate the above program. It predicts the forward alignment behaviour  $(p_\ell, q_\ell)$ , the backward alignment behaviour  $(p_\ell^\nabla, q_\ell^\nabla)$ , and the leading-order scaling of expert-indexed aggregates that enter both  $\Delta h^{\text{agg}}(t)$  and the gradients at shared representations when  $M$  and  $K$  scale. Together, these predictions determine the regime-appropriate scaling of initialization variances, learning rates, and (for Adam)  $\epsilon$ , yielding the canonical parameterizations reported in Table B.1.

The remainder of this section follows the programme described above in two stages. Section I.2 treats single linear maps  $W \cdot x$  and the alignment exponents  $(p_\ell, q_\ell)$  that govern them. Section I.7 treats the MoE-specific cross-expert aggregations, organized as a small taxonomy of structural mechanisms that determine whether each aggregation term survives at leading scale or is suppressed by the cross-expert CLT effect. The same taxonomy makes the imbalances under  $\mu\text{P}$  in Regimes II and III transparent, and identifies, for each mechanism, the structural change that MSSP introduces to restore balance.

## I.2. Per-layer alignment: scaling of linear (forward) maps $W \cdot x$

Each of the products entering the alignment-exponent decomposition (I.3) is a single linear map  $W x$ , and the prefactors in  $\|W x\|_{\text{RMS}}$  are governed by random-matrix and rank-1-alignment arguments below. We focus mainly on the forward pass; the backward arguments are analogous.

I.2.1. RANDOM-MATRIX PRELIMINARIES

Let  $A \in \mathbb{R}^{n \times n}$  have i.i.d. entries  $A_{ij} \sim \mathcal{N}(0, \sigma^2)$ . Then  $A$  acts as a near-isometry rescaled by  $\sigma\sqrt{n}$  along almost every direction in  $\mathbb{R}^n$ : its singular values are  $\Theta(\sigma\sqrt{n})$  throughout the Marchenko–Pastur bulk [2, 35], and the corresponding singular vectors are uniformly distributed on  $S^{n-1}$  [1, Ch. 2]. The directions along which  $A$  fails to act as such an isometry, namely those aligned with its smallest singular values, span only a vanishing fraction of  $\mathbb{R}^n$  and form an exponentially thin set on the sphere [e.g. 52, §3.1]. Therefore, for generic  $x \in \mathbb{R}^n$  (even when  $x$  depends on  $A$ , provided it is not adversarially chosen to lie in this subspace),

$$\|Ax\|_2 \asymp \sigma\sqrt{n} \|x\|_2.$$

Equivalently, in RMS norm,

$$\|Ax\|_{\text{RMS}} \asymp \sqrt{n} \|A\|_{\text{RMS}} \|x\|_{\text{RMS}}.$$

We refer to this as *CLT-like scaling*: each entry of  $Ax$  is a sum of  $n$  uncorrelated mean-zero terms, so its typical magnitude grows only as  $\sqrt{n}$  relative to the per-term scale, in analogy with the Central Limit Theorem. Later, we will contrast this with *LLN-like scaling*, where the prefactor becomes  $n$  rather than  $\sqrt{n}$ .

I.2.2. PROPAGATING UPDATES

**Standard MLP hidden layer: CLT-like scaling.** Applying the arguments above to a standard MLP hidden layer, with initial weights ( $W_0$ ) an i.i.d. Gaussian matrix whose fan-in and fan-out both diverge at rate  $n$ , a propagating perturbation  $\Delta x$  inherits this CLT-like scaling:

$$\|W_0 \Delta x\|_{\text{RMS}} \asymp \sqrt{n} \|W_0\|_{\text{RMS}} \|\Delta x\|_{\text{RMS}}.$$

Non-trivial feature learning requires the previous layer’s updates to satisfy  $\|\Delta x\|_{\text{RMS}} = \Theta(1)$ , and for the propagating update to remain  $\Theta(1)$  after passing through the weight matrix we therefore need

$$\|W_0\|_{\text{RMS}} \asymp \frac{1}{\sqrt{n}}, \quad \text{equivalently,} \quad \sigma \asymp \frac{1}{\sqrt{n}}.$$

The CLT-like scaling above hinges on  $A$  mapping between two diverging dimensions: for an  $n \times n$  i.i.d. Gaussian, the operator norm is  $\Theta(\sigma\sqrt{n})$ , which caps the prefactor at  $\sqrt{n}$  regardless of how  $x$  depends on  $A$ . No alignment, however structured, can lift the scaling to  $n$  in the  $\infty \rightarrow \infty$  case. When at least one dimension of  $A$  is finite or  $A$  has an extremely skewed aspect ratio, this argument no longer applies and alignment-induced LLN scaling becomes possible, although whether it actually occurs depends on the surrounding network dynamics; we shall give some examples below to illustrate this. In the  $\infty \rightarrow \infty$  case,  $x$  has too few degrees of freedom to align coherently with all  $n$  rows of  $A$  at once: maximal alignment with any single row boosts only the corresponding entry of  $Ax$  to LLN order, while the remaining  $n - 1$  entries stay at CLT order and dominate the norm. There is no way for one  $x$  to align with  $n$  independent rows simultaneously.

**Remark.** In all three regimes, propagating updates (both forward and backward) of all the layers that have diverging input and output dimensions admit CLT-like scaling. For instance, in Regime III all the layers (including the experts and the router), except the input and the output layer, fall into this category.

**Readout layer: the canonical LLN-like scaling.** For the readout layer of a standard network,  $W_0 \in \mathbb{R}^{k \times n}$  maps the diverging width  $n$  to a fixed output dimension  $k$ . The perturbation  $\Delta x \in \mathbb{R}^n$  entering the layer has been shaped by backpropagation through  $W_0^\top$ , so its entries acquire systematic alignment with the columns of  $W_0$ . The inner products  $\langle a^{(j)}, \Delta x \rangle$  no longer cancel, and the LLN-like prefactor is realized in the diverging dimension being summed over:

$$\|W_0 \Delta x\|_{\text{RMS}} \asymp n \|W_0\|_{\text{RMS}} \|\Delta x\|_{\text{RMS}}.$$

**MoE first expert layer (Regime II): alignment broken by expert independence.** The MoE bottleneck regime (Regime II) is an instructive intermediate case. The first expert layer  $W_0^{(i)} \in \mathbb{R}^{N_e \times n}$  maps the diverging width  $n$  to the finite expert hidden width  $N_e$ , geometrically the same situation as a readout layer, so alignment-induced LLN scaling is in principle permitted. Yet the alignment mechanism is broken here: the  $M$  expert first-layers  $\{W_0^{(i)}\}_{i=1}^M$  are independent at initialization, and a single shared input perturbation  $\Delta x$  cannot simultaneously align with the column structure of every  $W_0^{(i)}$ . CLT-like scaling therefore survives,

$$\|W_0^{(i)} \Delta x\|_{\text{RMS}} \asymp \sqrt{n} \|W_0^{(i)}\|_{\text{RMS}} \|\Delta x\|_{\text{RMS}}.$$

**MoE second expert layer (Regime II): alignment is restored at finite dimension.** The second expert layer  $W_0^{(i,2)} \in \mathbb{R}^{n \times N_e}$  maps the finite hidden width  $N_e$  back to the diverging width  $n$ . Here the relevant sum runs over only  $N_e$  terms, so the LLN prefactor cannot exceed  $N_e$  regardless of alignment. The perturbation entering this layer is private to expert  $i$  and is propagated by training through its own  $W_0^{(i,2)\top}$ , so the alignment is realized:

$$\|W_0^{(i,2)} \Delta x\|_{\text{RMS}} \asymp N_e \|W_0^{(i,2)}\|_{\text{RMS}} \|\Delta x\|_{\text{RMS}} \asymp \|W_0^{(i,2)}\|_{\text{RMS}} \|\Delta x\|_{\text{RMS}} \quad \text{since } N_e \in \Theta(1).$$

LLN-like in the fixed dimension being summed over.

### I.2.3. EFFECTIVE UPDATES: RANK-1 ALIGNMENT

The SGD update  $\Delta W = -\eta (\nabla_x L) x^\top$  is rank-1, so applying it to  $x$  saturates Cauchy-Schwarz:  $\|\Delta W \cdot x\|_2 = \|\Delta W\|_F \|x\|_2$ , equivalently

$$\|\Delta W \cdot x\|_{\text{RMS}} \asymp n \|\Delta W\|_{\text{RMS}} \|x\|_{\text{RMS}},$$

where  $n$  is the layer’s fan-in. The prefactor  $n$  (vs. the  $\sqrt{n}$  of CLT) reflects the rank-1 alignment between  $\Delta W$  and  $x$  due to gradient descent, which precludes any mean-zero cancellation. When the fan-in is fixed (as for the second expert layer in Regime II, with  $N_e = \Theta(1)$ ), the prefactor is  $\Theta(1)$  rather than diverging. By transpose, the backward effective piece  $(\Delta W)^\top \delta$  satisfies the same alignment relation with fan-out replacing fan-in. Together with the gradient scale at the layer, the rank-1 relation places a constraint on the SGD learning rate at every layer in any regime; the gradient scale is the only missing input to instantiate it.

In contrast, for Adam, entrywise normalization sets  $\Delta W$  entries to  $\Theta(\eta_{\text{Adam}})$  regardless of the gradient scale, so the alignment argument alone constrains  $\eta_{\text{Adam}}$ : the input-side LLN factor gives  $\eta_{\text{Adam}} = 1/n$  when the input dimension diverges and  $\eta_{\text{Adam}} = 1$  when it is fixed. Setting Adam’s  $\epsilon$ , however, still requires the gradient scale:  $\epsilon$  must remain  $\Theta$  of the typical per-coordinate gradient

magnitude (else either  $\varepsilon$  dominates and the update becomes trivial, or the gradients dominate and the normalization becomes vacuous), so it scales with the gradient at the layer.

Together with the propagating analysis above, this yields per-layer constraints relating  $\sigma$ ,  $\eta$  (and  $\varepsilon$  for Adam) to the layerwise gradient scale, which is in turn determined by the cross-expert aggregations of Section I.7. Our goal is a parameterization that simultaneously satisfies these per-layer constraints together with the cross-expert balance constraints below; trade-offs between competing constraints can in principle require expanding the hyperparameter space (e.g., layer-specific scalings) to admit a feasible solution.

Crucially, the  $\mu$ P desideratum (that the layerwise effective and propagating terms in the *forward* pass remain  $\Theta(1)$ ) already pins down a unique parameterization: the per-layer initialization variances, learning rates, and (for Adam)  $\varepsilon$ . We refer to this parameterization as the  $\mu$ P *baseline*. The alignment-exponent framework (I.3) furnishes the analytical machinery for solving the corresponding  $\Theta(1)$  conditions. The  $\mu$ P desideratum, however, imposes no explicit constraints on the backward pass or on the cross-expert aggregations (I.4), (I.6); as discussed in §I.1.1, balance across the terms of these decompositions constitutes the additional desiderata required of an MoE parameterization. In what follows, we proceed case by case through the  $\mu$ P baseline, identifying where it produces imbalances under our extended desiderata and showing how MSSP resolves each.

### I.3. MSSP Desiderata

For MoE blocks, the analogous requirement applies to each term in the effective/propagating decomposition of the aggregated representation. Writing  $W^{3,i}(t) = W^{3,i}(0) + \Delta W^{3,i}(t)$  and  $h^{2,i}(t) = h^{2,i}(0) + \Delta h^{2,i}(t)$ , the aggregation  $h^{\text{agg}}(t) = \sum_{i=1}^M \phi_i(t) W^{3,i}(t) h^{2,i}(t)$  decomposes as

$$h^{\text{agg}}(t) = \underbrace{\sum_{i=1}^M \phi_i W^{3,i}(0) h^{2,i}(0)}_{\text{A1: init.}} + \underbrace{\sum_{i=1}^M \phi_i W^{3,i}(0) \Delta h^{2,i}(t)}_{\text{A2: propagating}} + \underbrace{\sum_{i=1}^M \phi_i \Delta W^{3,i}(t) h^{2,i}(t)}_{\text{A3: effective}}. \quad (\text{I.4})$$

The MoE-specific extension of Desideratum 1 requires each of these three contributions (init., propagating, effective) to remain  $\Theta(1)$  in RMS at some  $t > 0$ . Expanding the propagating term one step further (substituting  $\Delta h^{2,i}(t) \approx W^{2,i}(0) \Delta h^1 + \Delta W^{2,i}(t) h^1(0)$ ) yields the four-term decomposition  $h^{\text{agg}}(t) = A_1 + A_{2,1} + A_{2,2} + A_3$  which will be used in the case studies below.

**Backward desiderata.** Recall that the same structure carries over to the backward pass. Along any linear map  $h_t^\ell = W_t^\ell x_t^{\ell-1}$ , writing  $\bar{\delta}_t^\ell := \nabla_{x_t^\ell} \mathcal{L}$  and  $\delta_t^\ell := \nabla_{h_t^\ell} \mathcal{L}$ , the transpose recursion  $\bar{\delta}_t^{\ell-1} = (W_t^\ell)^\top \delta_t^\ell$  admits the analogous effective/propagating split

$$\bar{\delta}_t^{\ell-1} = \underbrace{(W_0^\ell)^\top \delta_t^\ell}_{\text{propagating}} + \underbrace{(\Delta W_t^\ell)^\top \delta_t^\ell}_{\text{effective}}. \quad (\text{I.5})$$

When backpropagation crosses a representation that fans out into  $M$  expert branches, the gradient at the shared node aggregates expert contributions in the same form as the forward aggregation. Writing  $g^{2,i}(t) := \nabla_{h^{2,i}(t)} \mathcal{L}$  for the per-expert gradient at the bottleneck representation, the expert-pathway contribution to the gradient at the input shared across experts is  $\bar{\delta}^{1,\text{exp}}(t) = \sum_{i=1}^M \phi_i(t) (W^{2,i}(t))^\top g^{2,i}(t)$ , and expanding  $W^{2,i}(t) = W^{2,i}(0) + \Delta W^{2,i}(t)$  and

$g^{2,i}(t) = g^{2,i}(0) + \Delta g^{2,i}(t)$  gives

$$\bar{\delta}^{1,\text{exp}}(t) = \underbrace{\sum_{i=1}^M (W^{2,i}(0))^\top g^{2,i}(0)}_{\text{A4: init.}} + \underbrace{\sum_{i=1}^M (W^{2,i}(0))^\top \Delta g^{2,i}(t)}_{\text{A5: propagating}} + \underbrace{\sum_{i=1}^M (\Delta W^{2,i}(t))^\top g^{2,i}(t)}_{\text{A6: effective}} \quad (\text{I.6})$$

Desideratum 1 extends naturally to the backward pass: each piece of the backward decomposition (the layerwise effective/propagating split (I.5) and each term in the expert-aggregation decomposition (I.6)) is required to be *balanced*, i.e. to share a common scale within each decomposition. (Under an appropriate per-layer normalization that places backward signals on a scale comparable to forward signals, the common scale is  $\Theta(1)$ ; since this normalization depends on architecture, algorithm, and layer type, we work with the normalization-free *balance* statement throughout, and omit the explicit normalization.)

#### I.4. Backward scaling

The forward arguments above transpose directly to the backward pass at every layer with diverging input and output dimensions:  $(W_0^\ell)^\top \delta$  inherits the same CLT-like scaling as  $W_0^\ell \Delta x$ , and  $(\Delta W^\ell)^\top \delta$  inherits the same rank-1 alignment as  $\Delta W^\ell \cdot x$  (with fan-out replacing fan-in). We do not repeat these case studies.

#### I.5. Aggregation layers

In MoEs, aggregations across experts arise in both the forward pass (Eq. (I.4)) and the backward pass (Eq. (I.6)). The aggregation operator introduces sums over *expert-indexed* contributions (weighted by routing), and the leading-order scale of these sums depends on the evolving cross-expert correlation structure. In particular, expert sums can exhibit CLT-type scaling (weak cross-expert correlations), LLN-type scaling (strongly correlated contributions), or mixed behaviour across different contributions; our analysis makes these distinctions regime-by-regime.

#### I.6. Width dependence in $\mu\text{P}$ and how MSSP fixes it

Below, we provide some case studies on the key imbalances in MoE dynamics under  $\mu\text{P}$  and follow this by how MSSP fixes them.

##### I.6.1. EXPERT HIDDEN LAYER GRADIENT $\partial f / \partial h^{2,i}$

Consider the four-piece decomposition

$$\partial f / \partial h^{2,i} \propto (W_0^{3,i} + \Delta W^{3,i})^\top (W_0^4 + \Delta W^4)^\top.$$

In Regime III,  $(W_0^{3,i})^\top (\Delta W^4)^\top$  admits CLT-like scaling, since  $W_0^{3,i}$  has diverging input and output dimensions. In Regime II,  $W_0^{3,i}$  is a low-rank matrix and  $\Delta W^4$  is shaped by  $W_0^{3,i}$  through backpropagation, so LLN-like scaling is in principle possible. However,  $\Delta W^4$  is built from sums over all  $M \rightarrow \infty$  experts'  $W_0^{3,i}$ , and hence cannot simultaneously align with any single  $W_0^{3,i}$ ; the term therefore also admits CLT-like scaling.

$$\|(W_0^{3,i})^\top (\Delta W^4)^\top\|_{\text{RMS}} = \sqrt{N} \|(W_0^{3,i})^\top\|_{\text{RMS}} \|(\Delta W^4)^\top\|_{\text{RMS}} \asymp \sqrt{N} \cdot 1 / \sqrt{N_e} \cdot 1/N \quad \text{under } \mu\text{P}.$$

Now consider the term  $(\Delta W^{3,i})^\top (W_0^4)^\top$ . As discussed in the effective-updates section, since the updates to the second expert-layer weights are rank-1 in structure, the rank-1 alignment argument applies, yielding LLN-like scaling:

$$\|(\Delta W^{3,i})^\top (W_0^4)^\top\|_{\text{RMS}} = N \|(\Delta W^{3,i})^\top\|_{\text{RMS}} \|(W_0^4)^\top\|_{\text{RMS}} \asymp N \cdot 1/N_e \cdot 1/N \quad \text{under } \mu\text{P}.$$

The remaining two contractions in the four-piece expansion follow the same scaling laws:  $(W_0^{3,i})^\top (W_0^4)^\top$  behaves identically to the CLT-like term analysed above, and  $(\Delta W^{3,i})^\top (\Delta W^4)^\top$  behaves identically to the rank-1 term.

In Regime III we have  $N_e \asymp N$ , so both terms are of order  $1/N$  and are therefore balanced (and, under the appropriate normalization, both terms are of order 1). In Regime II,  $N_e \in \Theta(1)$ , and the two terms have scales  $1/\sqrt{N}$  and 1 respectively.

In summary, in Regime III,  $|W_0^{3,i}| = \Theta(1/\sqrt{N})$  and  $|\Delta W^{3,i}| = \Theta(1/N)$ , leading to an entry-size gap of  $\sqrt{N}$ . This  $\sqrt{N}$  exactly compensates the CLT-vs.-rank-1 alignment-type gap, so both contractions land at  $\Theta(1/N)$ : *balanced*. In Regime II,  $|W_0^{3,i}| = \Theta(1)$  and  $|\Delta W^{3,i}| = \Theta(1)$ , so the entry sizes are equal (since  $\sigma_3^2 = 1/N_e = \Theta(1)$ ). The CLT-vs.-rank-1 alignment-type gap is no longer cancelled, so the two pieces split by  $\sqrt{N}$ : *imbalanced*.

**MSSP correction in Regime II.** The MSSP-Regime-II boost  $\sigma_3^2 = M/N_e$  increases the scale of  $|W_0^{3,i}|$  from  $\Theta(1)$  to  $\Theta(\sqrt{M}) = \Theta(\sqrt{N})$ , while the update  $|\Delta W^{3,i}| = \Theta(1)$  is unchanged (the SGD learning rate  $\eta_3 = \eta MN$  is identical to its  $\mu\text{P}$ -Regime-II value). The entry-size ratio  $|\Delta W^{3,i}|/|W_0^{3,i}|$  thus moves from  $\Theta(1)$  under  $\mu\text{P}$  to  $\Theta(1/\sqrt{N})$  under MSSP, restoring the buffer to the same  $\sqrt{N}$  value as in Regime III. The CLT-vs.-rank-1 alignment-type gap is then cancelled exactly, and the four pieces of  $\partial f/\partial h^{2,i}$  all sit at the leading  $\Theta(1/N)$  scale.

## I.7. Cross-expert aggregation: a taxonomy of mechanisms

We now turn to the MoE-specific cross-expert aggregations. Two such aggregations arise: the forward aggregation  $h^{\text{agg}}(t)$  of (I.4), with expanded four-term form  $h^{\text{agg}}(t) = A_1 + A_{2,1} + A'_{2,2} + A_3$ ; and the backward aggregation  $\bar{\delta}^{1,\text{exp}}(t)$  of (I.6), which collects the expert-pathway contributions to the gradient with respect to the input  $h^1$  shared across experts, with an analogous decomposition. Each summand depends on the per-expert weights  $W_0^{2,i}$ ,  $W_0^{3,i}$  and their updates; the leading-order scale of the aggregate is determined by the cross- $i$  correlation structure of these summands. The scaling of these aggregation terms is determined by the crude CLT–LLN dichotomy: summands sharing a coherent direction aggregate by LLN (no  $1/\sqrt{M}$  suppression), whereas summands with i.i.d. random directions aggregate by CLT ( $1/\sqrt{M}$  suppression). The structurally interesting question, however, is: *What generates a coherent direction in the first place, or prevents one from emerging?* The four mechanisms identified below (A–D) classify the case studies that follow according to the algebraic feature that controls this question, and we examine each in turn under the  $\mu\text{P}$  baseline and under MSSP. For each mechanism we (i) state the algebraic feature and case-study term it governs, (ii) state its scaling under the  $\mu\text{P}$  baseline, and (iii) close with a paragraph describing the corresponding behaviour under MSSP.

### I.7.1. MECHANISM A: CROSS-EXPERT CLT UNDER INDEPENDENT OR WEAKLY CORRELATED SUMMANDS

Mechanism A governs aggregates  $A := \frac{1}{M} \sum_i \phi_i W_0^{3,i} W_0^{2,i} v$  with  $v$  shared across experts. Terms  $A_1$  and  $A_{2,1}$  in the forward aggregation are of this form, with  $v = h^1$  and  $v = \Delta h^1$  respectively. In

both cases a cross-expert central limit theorem applies, suppressing the aggregate by  $1/\sqrt{M}$  relative to the per-summand scale.

**Independent case ( $A_1$ ).** Take  $v \in \mathbb{R}^N$  independent of all expert weights. With  $W_0^{2,i} \in \mathbb{R}^{N_e \times N}$  Gaussian (entries  $\mathcal{N}(0, \sigma_2^2)$ ),

$$W_0^{2,i} v \sim \mathcal{N}(0, \sigma_2^2 \|v\|^2 I_{N_e}),$$

isotropic on  $\mathbb{R}^{N_e}$  and depending on  $v$  only through  $\|v\|^2$ . Conditional on  $W_0^{2,i} v$ ,

$$u_i := W_0^{3,i} W_0^{2,i} v \mid W_0^{2,i} v \sim \mathcal{N}(0, \sigma_3^2 \|W_0^{2,i} v\|^2 I_N).$$

Marginalising over  $W_0^{2,i}$  preserves rotation invariance: the distribution of  $u_i$  depends on  $v$  only through  $\|v\|^2$ . Across experts,  $(W_0^{2,i}, W_0^{3,i})_{i=1}^M$  are i.i.d., so the only source of cross-expert dependence among  $\{u_i\}$  is the shared  $\|v\|^2$ .

The scalar  $\|v\|^2 = \Theta(N)$  does not itself converge. Combined with the  $1/M$  aggregation prefactor and  $M \asymp N$ , it appears in the aggregate variance  $\mathbb{E}\|A\|^2 \asymp M^{-1} \mathbb{E}\|u_i\|^2$  as  $\|v\|^2/M \asymp \|v\|^2/N = \|v\|_{\text{RMS}}^2$ , which converges deterministically under standard initialisation. The summands are asymptotically independent at the aggregate level, and a multivariate CLT applies.

**Weakly dependent case ( $A_{2,1}$ ).**  $\Delta h^1$  is shaped by gradient backflow through every expert and depends on each  $W_0^{2,i}$ , so strict independence fails. A single vector in  $\mathbb{R}^N$  cannot, however, align nontrivially with  $M \rightarrow \infty$  independent Gaussian matrices simultaneously: the dependence of  $\Delta h^1$  on any one  $W_0^{2,i}$  enters through a single summand of the cross-expert loss gradient and is of relative order  $1/M$ . Decompose  $\Delta h^1 = \Delta h^{1,(\neq i)} + \Delta h^{1,(i)}$  with  $\Delta h^{1,(\neq i)}$  independent of  $W_0^{2,i}$  and  $\|\Delta h^{1,(i)}\|/\|\Delta h^1\| = O(1/M)$ ; the dominant term  $W_0^{3,i} W_0^{2,i} \Delta h^{1,(\neq i)}$  is governed by the previous paragraph, the residual by a  $1/M$ -suppressed correction.

**Per-summand and aggregate magnitudes.**

$$\mathbb{E}\|u_i\|^2 = N\sigma_3^2 \cdot N_e\sigma_2^2 \|v\|^2 = \sigma_2^2 \sigma_3^2 N N_e \|v\|^2,$$

$$\|u_i\|_{\text{RMS}} \asymp (\sigma_2 \sqrt{N})(\sigma_3 \sqrt{N_e}) \|v\|_{\text{RMS}}, \quad \|A\|_{\text{RMS}} \asymp \frac{\sigma_2 \sigma_3 \sqrt{N N_e}}{\sqrt{M}} \|v\|_{\text{RMS}}.$$

This holds in both Regime II ( $N_e$  fixed) and Regime III ( $N_e \asymp N$ ). Under  $\mu\text{P}$ ,  $\sigma_2 = 1/\sqrt{N}$  and  $\sigma_3 = 1/\sqrt{N_e}$ , so  $\|A_1\|_{\text{RMS}}, \|A_{2,1}\|_{\text{RMS}}$  are of order  $1/\sqrt{M}$ , strictly subleading to  $A'_{2,2}$  and  $A_3$  at  $\Theta(1)$ .

**MSSP correction.** The two MSSP fixes resolve this in regime-specific ways. In Regime III, MSSP shares expert weights at initialization ( $W_0^{2,i} = W_0^2$  and  $W_0^{3,i} = W_0^3$  for all  $i$ ); the per-expert summands then become strongly correlated across  $i$  rather than independent, the cross-expert CLT collapses to an LLN, and both terms are lifted to  $\Theta(1)$  along the shared direction  $W_0^3 W_0^2 h^1$  (or its  $\Delta h^1$  analogue). In Regime II, MSSP boosts  $\sigma_3$  from  $1/\sqrt{N_e}$  to  $\sqrt{M/N_e}$ ; the per-summand RMS scale is amplified by  $\sqrt{M}$ , exactly compensating the  $1/\sqrt{M}$  cross-expert CLT suppression and lifting both terms to  $\Theta(1)$  in expert-specific directions.

## I.7.2. MECHANISM B: GRAM OPERATORS ON A SHARED VECTOR

We will use the quantity  $A'_{2,2} = \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,i} \Delta_t W^{2,i} h_0^1$  to illustrate this mechanism. Each summand has the form  $G_i u$ , where  $G_i = W_0^{3,i} (W_0^{3,i})^\top$  acts on a vector  $u$  that is shared across experts (no  $i$  index). The cross-expert direction of  $G_i u$  is then governed by whether  $G_i$  is full-rank (so  $G_i \approx N_e \sigma_3^2 \cdot I$  and  $G_i u$  inherits the shared direction  $u$ ) or rank-deficient (so  $G_i$  is, after rescaling, a projection onto a random  $N_e$ -dimensional subspace of  $\mathbb{R}^N$ , distinct across experts). The same rank threshold separates Regime III ( $N_e \asymp N$ , isotropic, LLN) from Regime II ( $N_e = \Theta(1)$ , anisotropic, CLT).

**Intuitive scaling of  $A'_{2,2}$  across regimes.** Substituting the rank-1 update  $\Delta_t W^{2,i} \propto (W_0^{3,i})^\top (W_0^4)^\top (h_0^1)^\top$  and using  $\Delta_t W^{2,i} h_0^1 \propto \|h_0^1\|^2 (W_0^{3,i})^\top (W_0^4)^\top$ , each summand of  $A'_{2,2}$  reduces to

$$W_0^{3,i} \Delta_t W^{2,i} h_0^1 \propto \|h_0^1\|^2 \underbrace{W_0^{3,i} (W_0^{3,i})^\top}_{=: G_i} (W_0^4)^\top.$$

Writing  $u := (W_0^4)^\top \in \mathbb{R}^N$ , the analysis splits into a within-expert step (the action of  $G_i$  on the fixed vector  $u$ ) and a cross-expert step (the aggregation  $(1/M) \sum_i$ ).

*Within-expert structure (fix  $i$ ).* The matrix  $W_0^{3,i} \in \mathbb{R}^{N \times N_e}$  has rank  $\min(N, N_e)$  generically, and so does  $G_i$ . The image of  $G_i$  is the column space of  $W_0^{3,i}$  in  $\mathbb{R}^N$ ; let  $\Pi_i$  denote the orthogonal projection onto it. By rotational invariance of the entries of  $W_0^{3,i}$ , the random subspace  $\text{col}(W_0^{3,i})$  is Haar-distributed on the Grassmannian  $\text{Gr}(\min(N, N_e), N)$ .

The non-zero eigenvalues of  $G_i$  are the squared singular values of  $W_0^{3,i}$ , which concentrate around  $N\sigma_3^2$  by Marchenko–Pastur [35]. Consequently  $G_i$  admits the operator-level approximation

$$G_i \approx (N\sigma_3^2) \Pi_i,$$

giving  $G_i u \approx (N\sigma_3^2) \Pi_i u$ . Since  $u$  is independent of  $W_0^{3,i}$  and  $\text{col}(W_0^{3,i})$  is Haar-distributed, each projection retains

$$\|\Pi_i u\|_2^2 \asymp \frac{\min(N, N_e)}{N} \|u\|_2^2.$$

*Cross-expert structure.* The matrices  $\{W_0^{3,i}\}_{i=1}^M$  are i.i.d., so the projections  $\{\Pi_i u\}_{i=1}^M$  are independent across experts, each marginally Haar-distributed in direction within  $\mathbb{R}^N$ .

*Regime III ( $N_e \asymp N$ ): isotropic Gram matrix preserves LLN scaling.* For each fixed  $i$ , since  $W_0^{3,i}$  is generically full rank at leading order, so  $\Pi_i \approx I_N$  and

$$G_i u \approx (N_e \sigma_3^2) u = \Theta(1) \cdot u$$

(taking  $\sigma_3^2 = 1/N_e$ ). The Gram acts as a near-scalar multiple of the identity, and the result lies along  $u = (W_0^4)^\top$  (shared direction across all experts). Cross-expert aggregation admits a law of large numbers like scaling:

$$\|A'_{2,2}\|_{\text{RMS}} \asymp \|h_0^1\|^2 \cdot N_e \sigma_3^2 \cdot \|(W_0^4)^\top\|_{\text{RMS}} = \Theta(1) \text{ along } (W_0^4)^\top \quad \text{under } \mu\text{P}$$

---

2. The same argument holds under weak correlations

*Regime II* ( $N_e$  fixed,  $N \rightarrow \infty$ ): *anisotropic Gram yields CLT scaling.* For each fixed  $i$ , the column span  $\text{col}(W_0^{3,i})$  is a fixed-dimensional random subspace within a diverging ambient space, on which  $G_i$  has eigenvalue  $N\sigma_3^2 = N/N_e$  under  $\mu$ P. By rotational invariance, the Haar distribution of this subspace implies  $\|\Pi_i u\|_2^2 \asymp (N_e/N) \|u\|_2^2$ : each orthonormal basis vector of  $\text{col}(W_0^{3,i})$  is uniform on  $S^{N-1}$  by rotation invariance, contributing  $\|u\|^2/N$  in expectation, and there are  $N_e$  such basis vectors. Therefore,

$$\|G_i u\|_2 \asymp (N\sigma_3^2) \sqrt{N_e/N} \|u\|_2 = \sigma_3^2 \sqrt{NN_e} \|u\|_2,$$

with  $G_i u$  lying entirely in the expert-specific subspace  $\text{col}(W_0^{3,i})$ . Across experts, these subspaces are independent random elements of the Grassmannian, so the directions  $\{G_i u / \|G_i u\|_2\}_{i=1}^M$  are independent and there is no shared coherent direction along which  $\{G_i u\}_i$  can accumulate. Cross-expert aggregation therefore proceeds by the central limit theorem rather than the law of large numbers:  $A'_{2,2}$  is suppressed by an additional  $1/\sqrt{M}$  relative to the per-summand scale, and is not aligned with  $(W_0^4)^\top$ .

**MSSP correction.** The Regime II imbalance is removed under MSSP-Regime-II by the amplification  $\sigma_3^2 = M/N_e$ , which raises the eigenvalue of  $G_i$  on its rank- $N_e$  column span from  $N\sigma_3^2 = N/N_e$  under  $\mu$ P to  $N\sigma_3^2 = NM/N_e$  under MSSP. The per-summand scale  $\|G_i u\|_2 \asymp \sigma_3^2 \sqrt{NN_e} \|u\|_2$  is thereby amplified by exactly  $\sqrt{M}$ , and the cross-expert  $1/\sqrt{M}$  CLT cancellation is exactly absorbed;  $A'_{2,2}$  then sits at  $\Theta(1)$ , in expert-specific directions rather than along  $(W_0^4)^\top$ . In Regime III,  $A'_{2,2}$  is already  $\Theta(1)$  under  $\mu$ P (the Gram is isotropic and aligned along  $u$ ); MSSP-Regime-III's shared experts collapse the marginal randomness of  $\{\Pi_i\}$  entirely, but do not change the  $\Theta$ -class.

### I.7.3. MECHANISM C: EXPERT SUMS OF EFFECTIVE CONTRIBUTIONS

Let us use the term  $A_3$  of the forward aggregation to illustrate. Each summand has the form  $\Delta W^{3,i} x_i$  with  $\Delta W^{3,i}$  rank-1 of the form  $u_{\text{shared}} v_i^\top$ , where  $u_{\text{shared}}$  is a backpropagated factor shared across experts and  $v_i$  is expert-specific. The product  $\Delta W^{3,i} x_i = \langle v_i, x_i \rangle u_{\text{shared}}$  then writes the summand as a (per-expert) scalar times a shared direction. Cross-expert aggregation reduces to a scalar mean, and the LLN survives whenever that scalar mean is dominated by a strictly positive contribution. Unlike Mechanism B, the shared direction is forced by the rank-1 structure of the gradient itself rather than recovered from a Gram acting on a fixed vector, so this mechanism is regime-insensitive.

**Scaling  $A_3$ .** Recall that

$$A_3 = \frac{1}{M} \sum_i \phi_{i,t} \Delta_t W^{3,i} h_t^{2,i},$$

where

$$\Delta_t W^{3,i} h_t^{2,i} = -\frac{\eta_3 \chi \phi_{i,t-1}}{M} \langle h_{t-1}^{2,i}, h_t^{2,i} \rangle (W_{t-1}^4)^\top.$$

Therefore, the per-summand is a single inner product times the shared readout direction. Decomposing the inner product,

$$\langle h_{t-1}^{2,i}, h_t^{2,i} \rangle = \|h_{t-1}^{2,i}\|_2^2 + \langle h_{t-1}^{2,i}, \Delta_t h_t^{2,i} \rangle.$$

The first term,  $\|h_{t-1}^{2,i}\|_2^2$ , is strictly positive for every expert (non-zero mean across experts). The second term is a smaller fluctuation around it: at  $t = 1$  it is random across  $i$ , while at  $t = 2$  it

becomes coherent via the alignment chain. In either case the sum is dominated by the strictly positive  $\|h_{t-1}^{2,i}\|_2^2$ , so the cross-expert scalar average is LLN-preserved:

$$\frac{1}{M} \sum_i \phi_{i,t} \phi_{i,t-1} \langle h_{t-1}^{2,i}, h_t^{2,i} \rangle \asymp \|h^{2,i}\|_2^2.$$

Substituting back,

$$A_3 \asymp -\frac{\eta_3 \chi}{M} \|h^{2,i}\|_2^2 (W_{t-1}^4)^\top$$

along the shared readout direction and in  $\Theta(1)$  in Regime II under both  $\mu\text{P}$  and MSSP. In Regime III, the same argument shows that these terms are  $\Theta(1)$  under  $\mu\text{P}$ . Under MSSP (which uses shared experts), similar arguments show that these terms are  $\Theta(1)$  as well.

#### I.7.4. MECHANISM D: SCALAR SELF-PAIRING STRUCTURAL INVARIANCE

Certain quantities in the aggregation are built from summands that comprise scalar self-pairings  $u^\top G_i u$  rather than a vector  $G_i u$ . The expectation  $\mathbb{E}[u^\top G_i u] = \|u\|_2^2$  depends only on the diagonal mean  $\mathbb{E}[(G_i)_{aa}] = N_e \sigma_3^2$ , which is  $\Theta(1)$  in both Regimes II and III under  $\mu\text{P}$  (and is therefore regime-symmetric). The per-layer parameterization, however, is calibrated against the (regime-asymmetric) vector form, so when the structure happens to be scalar the calibration finds nothing to act on and passes through unattenuated, leading to width dependence.

**Mechanism: The structural mechanism underlying the imbalance of  $A_6$ .** The term  $A_6$  arises in the expansion of the backward expert-pathway gradient  $\bar{\delta}^{1,\text{exp}}(t)$  of (I.6). Specifically,  $A_6 = (1/M) \sum_i \phi_{i,t} (\Delta W^{2,i})^\top (W_0^{3,i})^\top (W_t^4)^\top$ , the contribution in which  $\Delta W^{2,i}$  is the perturbed expert factor while  $W^{3,i}$  remains at initialization. Substituting  $\Delta W^{2,i}$  into  $A_6$  produces the scalar self-pairing

$$S_i := u^\top G_i u = \|(W_0^{3,i})^\top u\|_2^2, \quad u := (W_0^4)^\top, \quad G_i := W_0^{3,i} (W_0^{3,i})^\top.$$

Since  $(G_i)_{ab} = \sum_{k=1}^{N_e} (W_0^{3,i})_{ak} (W_0^{3,i})_{bk}$  with entries i.i.d.  $\mathcal{N}(0, \sigma_3^2)$ ,

$$\mathbb{E}[(G_i)_{ab}] = N_e \sigma_3^2 \delta_{ab} = \delta_{ab} \quad (\sigma_3^2 = 1/N_e),$$

hence

$$\mathbb{E}[S_i] = \|u\|_2^2 = \Theta(1/N) \quad \text{in both regimes.}$$

The remaining terms of the decomposition of the aggregated gradient depend on  $W_0^{3,i}$  only through linear, vector-form quantities, which are regime-asymmetric by a factor of  $\sqrt{N}$  in per-coordinate magnitude and which are accounted for by the parameterization. Due to the regime symmetry of  $A_6$ , the parameterization therefore passes through it unattenuated, yielding *imbalanced* scaling in Regime II.

$$A_6 \in \Theta(1/N^2) \text{ in Regime II,} \quad A_6 \in \Theta(1/N) \text{ in Regime III.}$$

**MSSP correction.** The MSSP-Regime-II boost  $\sigma_3^2 = M/N_e$  promotes the diagonal mean to  $\mathbb{E}[(G_i)_{aa}] = N_e \sigma_3^2 = M = \Theta(N)$ , so the scalar self-pairing satisfies  $\mathbb{E}[S_i] = M \|u\|_2^2 = \Theta(N) \cdot \Theta(1/N) = \Theta(1)$ . The  $\Theta(N)$  amplification of the scalar exactly fills the  $\sqrt{N} \cdot \sqrt{N}$  deficit that  $\mu\text{P}$ -Regime-II’s vector-form-calibrated parameterization had been unable to address, and  $A_6$  moves from  $\Theta(1/N^2)$  to  $\Theta(1/N)$ , the leading scale in Regime II. The structural symmetry of  $\mathbb{E}[G_i] \propto I$  that made  $A_6$  inert under  $\mu\text{P}$ -Regime-II is exploited in reverse by MSSP: the boost is precisely tuned so that the regime-symmetric part of  $G_i$  becomes regime-asymmetric in the way the parameterization expects. In Regime III,  $A_6$  is already at the leading scale under  $\mu\text{P}$ , and MSSP-Regime-III’s shared experts do not change its  $\Theta$ -class.

### I.8. Summary: regime imbalances under $\mu\text{P}$ and how MSSP resolves them

**Regime II.** In Regime II, MSSP introduces a single change to the  $\mu\text{P}$  baseline: the initialization variance of the second expert layer is increased to  $\sigma_3^2 = M/N_e$ . Although the change operates through a single parameter, it restores width independence in all four sources of imbalance present under  $\mu\text{P}$  in Regime II. The cross-expert CLT under independent or weakly correlated summands (Mechanism A) is corrected because the increased variance amplifies each per-summand by  $\sqrt{M}$ , exactly offsetting the  $1/\sqrt{M}$  suppression of the cross-expert sum. The Gram operators on a shared vector (Mechanism B) and the scalar self-pairing structural invariance (Mechanism D) are corrected by the corresponding amplification of the Gram’s eigenvalues and diagonal entries. The gradient at the expert hidden layer (§I.6.1) is corrected because the increased  $\sigma_3$  raises  $\|W_0^{3,i}\|$  to a scale that allows the four contractions of  $\partial f / \partial h^{2,i}$  to balance. The same change to the parameterization therefore corrects all four structurally distinct sources of width dependence in Regime II simultaneously.

**Regime III.** In Regime III, MSSP introduces a different single change: expert weights are shared at initialization, so that all experts begin from a common pair  $(W_0^2, W_0^3)$ . Under  $\mu\text{P}$ , the only mechanism that produces a width dependence in Regime III is the cross-expert CLT under independent or weakly correlated summands (Mechanism A). With shared experts, the per-expert summands are identical across  $i$  and accumulate by an LLN rather than cancelling by a CLT, and the affected terms are lifted to the leading scale. The remaining mechanisms (the Gram operators on a shared vector, the expert sums of effective contributions, the scalar self-pairing structural invariance, and the gradient at the expert hidden layer) are already balanced under  $\mu\text{P}$  in Regime III, and shared experts leave their leading behaviour unchanged.

**The two fixes cannot be interchanged.** The two fixes are matched to the two distinct sources of width dependence in their respective regimes, and cannot be interchanged. Applying the Regime-II fix in Regime III would amplify the second expert layer’s initial weights too aggressively, causing the backward gradient at the expert hidden layer to scale as  $\sqrt{N}$  above its target; this is the very imbalance that the Regime-II fix was introduced to repair in Regime II. Conversely, sharing experts in Regime II would correct only Mechanism A: the Gram operators on a shared vector, the scalar self-pairing structural invariance, and the gradient at the expert hidden layer would all remain imbalanced, since shared experts do not affect the  $\sigma_3$ -dependent scales through which those mechanisms operate. Regime II requires the boosted variance; Regime III requires shared experts; each is the appropriate fix for its regime and would cause width dependence in the other.

### I.9. Justification of MSSP in Regime II

We close the Regime II analysis by addressing a structural feature of the MSSP parameterization (Definition 5) that distinguishes it from conventional feature-learning parameterizations: the per-expert hidden state  $h^{3,i}$  (specifically the initial and propagating contributions  $W_0^{3,i} h_0^{2,i}$  and  $W_0^{3,i} \Delta h^{2,i}$ ) has coordinate scale  $\Theta(\sqrt{N})$  vs the standard  $\Theta(1)$  scale. We argue in this subsection that this is admissible because  $h^{3,i}$  never functions as a primitive variable in the dynamics, and that MSSP-Regime-II is, up to structurally equivalent reformulations, the unique parameterization that achieves the correct  $\Theta(1)$  scale for the decomposition of the aggregated activation  $h^3$  in Regime II while confining the resulting amplification to a quantity that the dynamics avoid as a primitive.

**The dynamics never realize  $h^{3,i}$  as a primitive variable.** We claim that, throughout both forward and backward passes, every place at which  $h^{3,i}$  would arise in a natural decomposition of the dynamics admits a contraction-order rewriting that bypasses it. The forward aggregate is computable directly:

$$h^3 = \frac{1}{M} \sum_{i=1}^M \phi_i W^{3,i} h^{2,i}, \quad (\text{I.7})$$

without forming any  $h^{3,i}$  explicitly. The backward gradients factor through small-dimensional contractions:

$$\frac{\partial f}{\partial h^{2,i}} = \frac{\phi_i}{M} (W^{3,i})^\top (W^4)^\top, \quad (\text{I.8})$$

$$\frac{\partial L}{\partial W^{3,i}} = \frac{\chi \phi_i}{M} (W^4)^\top (h^{2,i})^\top, \quad (\text{I.9})$$

$$\frac{\partial f}{\partial \phi_i} = \frac{1}{M} (W^4 W^{3,i}) h^{2,i}. \quad (\text{I.10})$$

The gating gradient (I.10) is the only quantity in which  $h^{3,i}$  would appear under the natural decomposition  $\partial f / \partial \phi_i = (1/M) (h^{3,i})^\top (W^4)^\top$ . The contraction reordering on the right-hand side replaces the formation of  $h^{3,i} \in \mathbb{R}^N$  with the inner-product chain  $W^4 W^{3,i} \in \mathbb{R}^{1 \times N_e}$ , an  $N_e$ -dimensional row vector. Because  $N_e = \Theta(1)$ , this object is small-dimensional and may be formed at numerical magnitude  $\Theta(\sqrt{N})$  without storage or precision concerns.

The full forward and backward dynamics therefore close on the variable set

$$\mathcal{V} = \{h^1, \psi, \phi, \{h^{2,i}\}_{i=1}^M, h^3, f, \{W^{2,i}\}_{i=1}^M, \{W^{3,i}\}_{i=1}^M, W^4, Q\}, \quad (\text{I.11})$$

none of whose elements simultaneously carry both large coordinate magnitude ( $\Theta(\sqrt{N})$ ) and large dimension ( $\Theta(N)$ ). The amplified per-expert hidden state  $h^{3,i}$  — the unique object in the architecture that combines both — is structurally absent from  $\mathcal{V}$  as a primitive variable.

**Mean-field interpretation.** The structure above is the signature of a mean-field theory: per-particle quantities are amplified, the aggregate is the natural order parameter, and the closed dynamics involve only aggregated moments of per-particle quantities together with low-dimensional inputs. The DMFT effective fields associated with MSSP-Regime-II can be expressed entirely in terms of aggregated weight-product moments such as

$$\bar{G}_t = \frac{1}{M} \sum_{i=1}^M \phi_{i,t} W_t^{3,i} W_t^{2,i}, \quad \bar{K}_t = \frac{1}{M} \sum_{i=1}^M \phi_{i,t}^2 (W_t^{2,i})^\top W_t^{2,i},$$

together with the per-expert layer-1 activations  $\{h^{2,i}\}$ , all of which sit at the standard  $\Theta(1)$  coordinate scale. The per-expert hidden state  $h^{3,i}$  does not appear in the effective theory.

### I.10. Composability with depth: stacked MoE blocks under MSSP-Regime-II

The boost-and-cancel mechanism that defines MSSP-Regime-II is structurally local to a single MoE block: it operates within that block's cross- $i$  CLT on the aggregated activation, and its scaling justification (Definition 5) does not refer to any architectural element outside the block. We formalize this observation by showing that stacking  $K$  MoE blocks under MSSP-Regime-II preserves the leading scales of every primitive variable in the forward pass and every backward gradient at every block boundary. The parameterization composes with depth without per-block modification.

**Stacked architecture.** Let  $K \geq 2$ . We consider the composition

$$x \xrightarrow{W^1} h^1 \xrightarrow{\text{block 1}} h^{3,(1)} \xrightarrow{\text{block 2}} h^{3,(2)} \rightarrow \dots \xrightarrow{\text{block } K} h^{3,(K)} \xrightarrow{W^4} f,$$

where each block  $k \in [K]$  has its own gating  $Q^{(k)}$ , expert weights  $\{W^{2,i,(k)}, W^{3,i,(k)}\}_{i \in [M]}$ , and aggregation

$$\begin{aligned} \psi^{(k)} &= Q^{(k)} h^{3,(k-1)}, & \phi_i^{(k)} &= \sigma(\psi_i^{(k)}), \\ h^{2,i,(k)} &= W^{2,i,(k)} h^{3,(k-1)}, \\ h^{3,i,(k)} &= W^{3,i,(k)} h^{2,i,(k)}, \\ h^{3,(k)} &= \frac{1}{M} \sum_{i=1}^M \phi_i^{(k)} h^{3,i,(k)}, \end{aligned}$$

with the convention  $h^{3,(0)} := h^1$ . All weights are drawn independently across blocks at initialization. Each block adopts the MSSP-Regime-II parameterization (Definition 5): per-block init variances  $\sigma_Q^2 = 1/N$ ,  $\sigma_2^2 = 1/N$ ,  $\sigma_3^2 = M/N_e$ ; output  $W_0^4 = 0$ ; and SGD learning rates  $\eta_Q = \eta_2 = \eta M/N$ ,  $\eta_3 = \eta MN$ ,  $\eta_4 = \eta/N$  within each block.

#### Forward composability.

**Proposition 1 (Forward depth composition)** *Under the stacked MSSP-Regime-II architecture above, for every  $k \in [K]$  at initialization:*

- (i)  $h_0^{2,i,(k)} \in \Theta(1)$  entry-wise.
- (ii)  $h_0^{3,i,(k)} \in \Theta(\sqrt{N})$  entry-wise.
- (iii) The aggregate  $h_0^{3,(k)} \in \Theta(1)$  entry-wise via cross- $i$  CLT.

**Proof [Sketch]** Induction on  $k$ . The base case  $k = 1$  is the single-block analysis (§J.3.3) with input  $h^1 \in \Theta(1)$ . For the inductive step, suppose  $h_0^{3,(k-1)} \in \Theta(1)$  entry-wise. Then  $h_0^{2,i,(k)} = W_0^{2,i,(k)} h_0^{3,(k-1)}$  has variance per entry  $\sigma_2^2 \|h_0^{3,(k-1)}\|^2 = (1/N) \cdot \Theta(N) = \Theta(1)$ , giving (i). Next,  $h_0^{3,i,(k)} = W_0^{3,i,(k)} h_0^{2,i,(k)}$  has variance per entry  $N_e \sigma_3^2 \Theta(1) = N_e \cdot (M/N_e) = M = \Theta(N)$ , giving (ii). Finally, the aggregate variance is  $(1/M^2) \sum_i \phi_{i,0}^2 N_e \sigma_3^2 \Theta(1) = \Theta(\sigma_3^2/M) = \Theta(1)$  via cross- $i$  CLT, giving (iii). The inductive hypothesis enters only through the  $\Theta(1)$  entry scale of the input  $h_0^{3,(k-1)}$  — exactly the property required to start the per-block boost-and-cancel.  $\blacksquare$

**Backward composability.**

**Proposition 2 (Backward depth composition)** *At  $t = 1$  — when only  $W^4$  has updated, with  $W_1^4 = -\eta_4 \chi_0 (h_0^{3,(K)})^\top$  having entries  $\Theta(1/N)$  — the propagated gradient at every inter-block aggregate satisfies*

$$\partial f_1 / \partial h_0^{3,(k)} \in \Theta(1/N) \quad \text{entry-wise,} \quad k \in \{0, 1, \dots, K-1\},$$

where  $h_0^{3,(0)} := h^1$ .

**Proof [Sketch]** Reverse induction on  $k$ . Base case  $k = K - 1$ : the expert pathway through block  $K$  is

$$\partial f_1 / \partial h_0^{3,(K-1)} \Big|_{\text{exp}} = \frac{1}{M} \sum_{i=1}^M \phi_{i,0}^{(K)} (W_0^{2,i,(K)})^\top (W_0^{3,i,(K)})^\top (W_1^4)^\top.$$

The intermediate  $(W_0^{3,i,(K)})^\top (W_1^4)^\top \in \mathbb{R}^{N_e}$  has entries  $\Theta(1)$  — the boost in  $\sigma_3^2 = M/N_e$  produces  $(W_0^{3,i,(K)})^\top h_0^{3,(K)} \in \Theta(N)$ , multiplied by  $-\eta_4 \chi_0 = -1/N$  to give  $\Theta(1)$ . Then  $(W_0^{2,i,(K)})^\top$  acting on this  $\Theta(1)$  vector via the rank-deficient initialization Gram (§J.3.1) produces per-summand entries  $\Theta(1/\sqrt{N})$  in mostly-random directions; cross- $i$  CLT over  $M$  approximately-independent summands cancels this to  $\Theta(1/\sqrt{MN}) = \Theta(1/N)$ . The router pathway through block  $K$  contributes  $\Theta(1/N)$  similarly.

Inductive step: assume  $\partial f_1 / \partial h_0^{3,(k+1)} \in \Theta(1/N)$  entries. The expert pathway through block  $k + 1$  is

$$\partial f_1 / \partial h_0^{3,(k)} \Big|_{\text{exp}} = \frac{1}{M} \sum_{i=1}^M \phi_{i,0}^{(k+1)} (W_0^{2,i,(k+1)})^\top (W_0^{3,i,(k+1)})^\top \partial f_1 / \partial h_0^{3,(k+1)}.$$

The intermediate  $(W_0^{3,i,(k+1)})^\top \partial f_1 / \partial h_0^{3,(k+1)} \in \mathbb{R}^{N_e}$  has variance per entry  $N \cdot \sigma_3^2 \cdot \Theta(1/N^2) = N \cdot N \cdot (1/N^2) = 1$ , hence coordinate scale  $\Theta(1)$  — the boost in  $\sigma_3^2$  exactly compensates the inductive  $\Theta(1/N)$  scale of the incoming gradient. The remainder of the calculation reduces to the base case:  $(W_0^{2,i,(k+1)})^\top$  acting on a  $\Theta(1)$  vector produces per-summand  $\Theta(1/\sqrt{N})$ , cross- $i$  CLT yields aggregate  $\Theta(1/N)$ . Router pathway analogous.  $\blacksquare$

**Theorem 3 (Depth composability of MSSP-Regime-II)** *Under the stacked MSSP-Regime-II architecture above, for every  $K \geq 1$  and every block  $k \in [K]$ :*

- (i) *Every primitive variable at initialization in block  $k$  has the same coordinate scale as in the single-block ( $K = 1$ ) MSSP-Regime-II analysis (§J.3.3).*
- (ii) *Every backward gradient at every primitive variable at  $t = 1$  in block  $k$  has the same coordinate scale as in the single-block analysis (§J.3.3), and consequently the SGD updates  $\Delta W^{2,i,(k)}$ ,  $\Delta W^{3,i,(k)}$ ,  $\Delta Q^{(k)}$  within block  $k$  retain their single-block scales.*
- (iii) *Every per-expert hidden state  $h^{3,i,(k)}$  remains non-primitive: the contraction-reordering identity*

$$(h^{3,i,(k)})^\top v = (v^\top W^{3,i,(k)}) h^{2,i,(k)}$$

*bypasses formation of  $h^{3,i,(k)}$  for any vector  $v \in \mathbb{R}^N$  at coordinate scale  $\Theta(1/N)$ , with the small-dimensional intermediate  $v^\top W^{3,i,(k)} \in \mathbb{R}^{N_e}$  at coordinate scale  $\Theta(1)$ .*

Consequently, *MSSP-Regime-II* composes with arbitrary depth  $K$  without any per-block parameter modification.

**Proof** [Sketch] (i) is Proposition 1. (ii) follows from Proposition 2 together with the single-block backward analysis: once the gradient at block  $k$ 's output aggregate  $\partial f_1 / \partial h_0^{3,(k)} \in \Theta(1/N)$  is established, all backward computations within block  $k$  proceed exactly as in §J.3.3, since they depend only on the input gradient's coordinate scale and the in-block weight structure. (iii) is the algebraic identity together with the entry-scale calculation  $v^\top W^{3,i,(k)} \in \Theta(1)$  for  $v \in \Theta(1/N)$  entries, which is precisely the intermediate computation in the proof of Proposition 2. ■

### I.11. Approaches We Tried That Do Not Restore Scale Stability

Before settling on the parameterization-level prescription that defines MSSP, we explored a number of architectural and parameterization-level interventions intended to remove the delayed-learning phenomenon and the underlying width-dependent subterms in MoE training under  $\mu$ P. None of the interventions enumerated below, applied on its own, restored clean refined-coordinate-check exponents while avoiding delayed learning at scale. We document them here because each is a natural intervention that practitioners may consider, and the reasons they fail clarify why the underlying problem is structural and why a parameterization-level fix is necessary.

#### I.11.1. RMSNORM (OR LAYERNORM) AFTER THE AGGREGATION

A natural first response to the observation that  $\|h^{\text{agg}}\|_{\text{RMS}}$  vanishes with width under  $\mu$ P-Regime-II is to insert a normalization layer after the aggregation, rescaling  $h^{\text{agg}}$  back to  $\Theta(1)$ . We tested this intervention with both RMSNorm and post-aggregation sigmoid and observed that neither restores width-independent dynamics. The reasons fall into four distinct points.

**Norm layers preserve ratios, not balance.** The post-aggregation activation  $h^{\text{agg}}$  decomposes into multiple subterms with distinct width scalings under  $\mu$ P-Regime-II — most notably the LLN-aligned effective contribution  $A_3 \in \Theta(1)$  and the CLT-suppressed propagating contribution  $A_{2,1} \in \Theta(1/\sqrt{N})$  (§J.3.2). RMSNorm divides the entire aggregate by its own RMS norm, rescaling every subterm by the same factor. The relative magnitudes of the subterms are preserved: after normalization, the post-norm activation is dominated by the effective contribution at the leading scale, while the propagating contribution sits at  $\Theta(1/\sqrt{N})$  of the signal — a vanishing fraction. The propagating channel that should carry information from upstream-layer updates into the current block has been effectively erased.

**The norm's  $\epsilon$  has to be width-scaled, and the right scaling is exactly the imbalanced-subterm structure we are trying to avoid having to track.** Practical RMSNorm computes  $\tilde{h} = h / \sqrt{\|h\|^2 + \epsilon^2}$ , where  $\epsilon$  prevents over- and under-flow. With  $\|h_0^{\text{agg}}\|^2 \in \Theta(1/N)$  in  $\mu$ P-Regime-II, there is a width-dependent crossover between two regimes: when  $\|h\|^2 \ll \epsilon^2$  the norm acts as a constant rescaling  $h/\epsilon$  (no normalization in the limit), while when  $\|h\|^2 \gg \epsilon^2$  the norm acts as a full RMS normalization. A fixed  $\epsilon$  therefore behaves qualitatively differently at small versus large widths, and the threshold between the two behaviors moves with  $N$ . To recover width-independent behavior,  $\epsilon$  would need to be scaled in proportion to  $\|h^{\text{agg}}\|$ ; but  $\|h^{\text{agg}}\|$  in  $\mu$ P-Regime-II is itself a sum of subterms with mixed scaling exponents, and the dominant subterm changes during training.

There is no single power-of- $N$  schedule for  $\epsilon$  that works across widths without first having solved the underlying scaling problem. The norm therefore moves the scaling burden into a different hyperparameter rather than eliminating it.

**The backward pass through the norm introduces its own width-dependent rescaling.** The Jacobian of RMSNorm scales as  $1/\|h^{\text{agg}}\|$ , so gradients flowing backward through the norm layer are inflated by that factor. At initialization in  $\mu$ P-Regime-II this factor is  $\sqrt{N}$ , and gradients upstream of the norm are amplified by a width-dependent multiplier whose magnitude also evolves over training as  $\|h^{\text{agg}}\|$  climbs. The effective layerwise learning rate on every parameter upstream of the norm therefore becomes a function of both width and step number — the time-varying width-dependence that scale stability is designed to remove. So even granting the loss of the propagating channel discussed above and the hyperparameter burden of scaling  $\epsilon$ , the backward pass through a normalization layer is itself a fresh source of instability rather than a fix for the existing one.

**A post-aggregation norm cannot reach imbalances that live elsewhere in the dynamics.** The imbalanced subterms in  $\mu$ P-Regime-II are not localized to  $h^{\text{agg}}$ . The four-piece expert-hidden-gradient buffer  $\partial f/\partial h^{2,i}$  (§I.6.1) and the scalar self-pairing imbalance in  $A_6$  (Mechanism D, §I.7.4) sit on the backward pathway and never pass through any post- $h^{\text{agg}}$  operation. Even granting that a post-aggregation norm somehow fixed the forward aggregation, three of the four mechanisms identified in §I.7 would remain broken. A single normalization layer placed at one point in the architecture cannot address imbalances that arise across multiple structurally distinct points in the dynamics; this is one of the empirical reasons we converged on a parameterization-level fix rather than an architectural patch.

## J. Signal Propagation Analysis

Here we provide a self-consistent forward and backward signal propagation analysis through an MoE block for the two more interesting Regimes II and III.

This analysis complements the rigorous DMFT treatment in the subsequent sections. We use the minimal MoE architecture defined below, which places the MoE block between an input and an output layer so that its incoming and outgoing signals carry the same asymptotic scaling as in a standard dense architecture. The arguments therefore extend to MoE blocks embedded in larger networks, such as a Transformer MoEs, whenever additional modules preserve this scaling. This is precisely what  $\mu$ P’s layer-wise parametrization guarantees [56, 58]. This compositionality is verified by our Transformer MoE coordinate checks in Appendix Q.2. For brevity, we omit the analysis for Regime I, which follows the known arguments for standard architectures.

### J.1. Preliminaries

**Minimal MoE architecture.** We work with the minimal MoE architecture

$$h^1 = W^1 x, \tag{J.1}$$

$$\psi = Q h^1, \tag{J.2}$$

$$\phi = \sigma(\psi), \tag{J.3}$$

$$h^{2,i} = W^{2,i} h^1, \quad i = 1, \dots, M, \tag{J.4}$$

$$h^{3,i} = W^{3,i} h^{2,i}, \tag{J.5}$$

$$h^3 = \frac{1}{M} \sum_{i=1}^M \phi_i h^{3,i}, \tag{J.6}$$

$$f = W^4 h^3, \tag{J.7}$$

where  $x \in \mathbb{R}^D$  is the input,  $f \in \mathbb{R}$  the scalar output, and the trainable parameters are

$$\theta = (W^1, Q, \{W^{2,i}, W^{3,i}\}_{i=1}^M, W^4)$$

with shapes  $W^1 \in \mathbb{R}^{N \times D}$ ,  $Q \in \mathbb{R}^{M \times N}$ ,  $W^{2,i} \in \mathbb{R}^{N_e \times N}$ ,  $W^{3,i} \in \mathbb{R}^{N \times N_e}$ ,  $W^4 \in \mathbb{R}^{1 \times N}$ . Here  $N$  is the hidden width,  $N_e$  the expert hidden width,  $M$  the number of experts, and  $D$  the input dimension. The derivations below are written for *sigmoid* gating,  $\sigma(\psi)_i = (1 + e^{-\psi_i})^{-1}$ , paired with the explicit  $1/M$  aggregation factor in (J.1); the same scaling exponents hold under softmax gating without an explicit aggregation multiplier.

**Scaling regimes.** As described in the main paper, we analyze three asymptotic regimes, distinguished by which of the width parameters  $N, N_e, M$  diverge:

- **Regime I:**  $(N, N_e \asymp n \rightarrow \infty$  with  $M, K = \Theta(1))$ ,
- **Regime II:**  $(N, M, K \asymp n \rightarrow \infty$  with  $N_e = \Theta(1))$ ,
- **Regime III:**  $(N, N_e, M, K \asymp n \rightarrow \infty)$ .

The heuristic derivation for Regimes II and III is provided in Sections J.3, and J.4 respectively. In each of these regimes, we want to understand the scale of all quantities that arise in the forward and backward computations in the network, both at initialization and during training. While the

initial random parameters are independent of the data distribution, correlations develop over time since updates carry information about the data; this is the case for both SGD and Adam [15, 57, 58]. To derive the correct scaling it is crucial to consider the correlations that arise during training and how these correlations propagate forward and backward through the network. In general, after one step of SGD the full set of post-update correlations is already present [5, 15, 56, 58], so analysing the first two forward and backward passes suffices to determine the asymptotic width scaling of every quantity arising in the computation<sup>3</sup>.

## J.2. Standing assumptions and notation

We track leading-order coordinate-scale exponents only, ignoring  $\Theta(1)$  prefactors.

### Standing assumptions.

**(B1)**  $\chi_t = \partial L / \partial f_t \in \Theta(1)$  at every step  $t$ ; absorbed into  $\Theta(1)$  scalar prefactors throughout.

**(B2)** Nonlinearities preserve leading scaling exponents. For sigmoid gating with  $1/M$  aggregation,  $\phi_{i,t}, \dot{\phi}_{i,t} \in \Theta(1)$ .

**Notation.**  $X \in \Theta(\alpha)$  denotes entry-wise (coordinate) scale. Cumulative SGD updates:  $W_t^\ell = W_0^\ell + \Delta_t W^\ell$ ,  $\Delta_t X = X_t - X_{t-1}$ , with init values subscripted by 0.

**Standard tools.** We invoke the following standard probabilistic facts as needed, without per-instance citation:

- Central limit theorem (CLT) for sums of approximately independent terms;
- Law of large numbers (LLN) for empirical averages;
- Operator-norm concentration for iid (sub-)Gaussian matrices: for  $W \in \mathbb{R}^{n \times m}$  with iid  $\mathcal{N}(0, \sigma^2)$  entries,  $\|W\|_{op} = \sigma(\sqrt{n} + \sqrt{m})(1 + o(1))$  with high probability [52].

**Feature update decomposition.** For  $h^\ell = W^\ell h^{\ell-1}$ ,

$$\Delta h^\ell = \underbrace{W_0^\ell \Delta h^{\ell-1}}_{\text{propagating}} + \underbrace{\Delta W^\ell h_0^{\ell-1}}_{\text{effective}} + \underbrace{\Delta W^\ell \Delta h^{\ell-1}}_{\text{cross term}}.$$

## J.3. Scaling derivation for Regime II

In this section, we provide the heuristic scaling derivation for Regime II.

### Definition 4 ( $\mu$ P baseline, Regime II)

- **Initialization.** All parameters are drawn independently according to:

$$W_0^1 \sim \mathcal{N}(0, D^{-1}), Q_0 \sim \mathcal{N}(0, N^{-1}), W_0^{2,i} \sim \mathcal{N}(0, N^{-1}),$$

$$W_0^{3,i} \sim \mathcal{N}(0, N_e^{-1}), W_0^4 \sim \mathcal{N}(0, N^{-2}), \text{ for } i \in [M].$$

3. When  $W_{L+1}^{(0)} = 0$ , only the readout updates in the first SGD step, since the gradient into every non-readout parameter carries a factor of  $W_{L+1}^\top$  that vanishes at initialization. A third forward and backward pass at  $\theta^{(2)}$  is then needed, since the non-readout parameters first move at  $t = 2$ .

- **SGD learning rates.**  $\eta_1 = \eta N$ ,  $\eta_Q = \eta(M/N)$ ,  $\eta_2 = (M/N)\eta$ ,  $\eta_3 = (MN)\eta$ ,  $\eta_4 = \eta N^{-1}$ .
- **Adam learning rates.**  $\eta_1 = \eta$ ,  $\eta_Q = \eta_2 = \eta_4 = \eta N^{-1}$ ,  $\eta_3 = 1/N_e\eta$ .
- **Adam epsilon.**  $\epsilon_1 = \epsilon N^{-1}$ ,  $\epsilon_Q = \epsilon M^{-1}$ ,  $\epsilon_2 = \epsilon M^{-1}$ ,  $\epsilon_3 = \epsilon N^{-1}M^{-1}$ ,  $\epsilon_4 = \epsilon$ .

**Definition 5 (SSP, Regime II)**

- **Initialization.** All parameters are drawn independently according to:

$$W_0^1 \sim \mathcal{N}(0, D^{-1}), \quad Q_0 \sim \mathcal{N}(0, N^{-1}), \quad W_0^{2,i} \sim \mathcal{N}(0, N^{-1}),$$

$$W_0^{3,i} \sim \mathcal{N}(0, MN_e^{-1}), \quad W_0^4 \sim \mathcal{N}(0, N^{-2}), \quad \text{for } i \in [M].$$

- **SGD learning rates.**  $\eta_1 = \eta N$ ,  $\eta_Q = \eta(M/N)$ ,  $\eta_2 = (M/N)\eta$ ,  $\eta_3 = (MN)\eta$ ,  $\eta_4 = \eta N^{-1}$ .
- **Adam learning rates.**  $\eta_1 = \eta$ ,  $\eta_Q = \eta_2 = \eta_4 = \eta N^{-1}$ ,  $\eta_3 = 1/N_e\eta$ .
- **Adam epsilon.**  $\epsilon_1 = \epsilon N^{-1}$ ,  $\epsilon_Q = \epsilon M^{-1}$ ,  $\epsilon_2 = \epsilon M^{-1}$ ,  $\epsilon_3 = \epsilon N^{-1}M^{-1}$ ,  $\epsilon_4 = \epsilon$ .

 J.3.1. DERIVING  $\mu\text{P}$  IN REGIME II

This subsection presents a self-contained derivation for the  $\mu\text{P}$  baseline (Definition 4) in Regime II.

**Setup and standing assumptions**

**Architecture.** As in §J.1: minimal MoE with sigmoid gating and explicit  $1/M$  aggregation  $h^3 = (1/M) \sum_i \phi_i h^{3,i}$ .

**Initialization ( $\mu\text{P}$  baseline).** Per Definition 4, with all parameters drawn independently:

$$(W_0^1)_{ab} \sim \mathcal{N}(0, 1/D), \quad (Q_0)_{ia} \sim \mathcal{N}(0, 1/N),$$

$$(W_0^{2,i})_{ab} \sim \mathcal{N}(0, 1/N), \quad (W_0^{3,i})_{ab} \sim \mathcal{N}(0, 1/N_e),$$

$$(W_0^4)_a \sim \mathcal{N}(0, 1/N^2).$$

Resulting entry scales:  $W_0^4$  entries  $\Theta(1/N)$ ,  $W_0^{2,i}$  entries  $\Theta(1/\sqrt{N})$ ,  $W_0^{3,i}$  entries  $\Theta(1)$  (since  $\sigma_3^2 = 1/N_e = \Theta(1)$ ).

**Learning rates (SGD).**  $\eta_1 = \eta N$ ,  $\eta_Q = \eta M/N$ ,  $\eta_2 = \eta M/N$ ,  $\eta_3 = \eta MN$ ,  $\eta_4 = \eta/N$ ,  $\eta \in \Theta(1)$ .

**Gating.** Sigmoid:  $\phi_{i,t}, \dot{\phi}_{i,t} \in \Theta(1)$ .

**Auxiliary scaling lemmas** We reuse Lemma 10 from §J.4.4.

**Specialization of Lemma 10 to the rank-deficient initialization Gram matrices in Regime II.**

For the random initialization  $W_0^{3,i} \in \mathbb{R}^{N \times N_e}$  with  $\sigma_W^2 = 1/N_e$  and  $N_e = \Theta(1)$ , the Gram matrix  $W_0^{3,i}(W_0^{3,i})^\top \in \mathbb{R}^{N \times N}$  has rank at most  $N_e = \Theta(1)$ : its diagonal mean is 1 and its off-diagonal entries have coordinate scale  $\Theta(1)$  (variance  $1/N_e = \Theta(1)$ ). Consequently, acting on a vector  $v = (W_0^4)^\top$  with entries of coordinate scale  $\Theta(1/N)$ , the product  $W_0^{3,i}(W_0^{3,i})^\top v$  has entries of coordinate scale  $\Theta(1/\sqrt{N})$  in mostly-random directions, not aligned along  $v$ .

Analogously, for the random initialization  $W_0^{2,i} \in \mathbb{R}^{N_e \times N}$  with  $\sigma_W^2 = 1/N$ , the Gram matrix  $(W_0^{2,i})^\top W_0^{2,i} \in \mathbb{R}^{N \times N}$  has rank at most  $N_e = \Theta(1)$  and entries of coordinate scale  $\Theta(1/N)$ . Acting on  $h_0^1$ , it produces entries of coordinate scale  $\Theta(1/\sqrt{N})$  in mostly-random directions, not aligned along  $h_0^1$ .

We refer to  $W_0^{3,i}(W_0^{3,i})^\top$  and  $(W_0^{2,i})^\top W_0^{2,i}$  collectively as the *rank-deficient initialization Gram matrices* of Regime II.

### First forward pass

$$h_0^1 = W_0^1 x \in \Theta(1). \quad (\text{CLT}; \sigma_1^2 = 1/D, \|x\|^2 \in \Theta(D)) \quad (\text{R2}\mu\text{-F1.1})$$

$$\psi_0 = Q_0 h_0^1 \in \Theta(1). \quad (\text{CLT}; \sigma_Q^2 = 1/N, \|h_0^1\|^2 \in \Theta(N)) \quad (\text{R2}\mu\text{-F1.2})$$

$$\phi_{i,0} = \sigma(\psi_{i,0}) \in \Theta(1). \quad (\text{sigmoid bounded; gating assumption}) \quad (\text{R2}\mu\text{-F1.3})$$

$$\begin{aligned} h_0^{2,i} &= W_0^{2,i} h_0^1 \in \Theta(1), \quad \text{independent across } i \\ &(\text{CLT}; \sigma_2^2 = 1/N, \|h_0^1\|^2 \in \Theta(N); \\ &\|h_0^{2,i}\|^2 = N_e \cdot \Theta(1) = \Theta(1) \text{ since } N_e = \Theta(1), \text{ not } \Theta(N_e) \rightarrow \infty). \end{aligned} \quad (\text{R2}\mu\text{-F1.4})$$

$$\begin{aligned} h_0^{3,i} &= W_0^{3,i} h_0^{2,i} \in \Theta(1) \\ &(\text{variance per entry } \sigma_3^2 \|h_0^{2,i}\|^2 = (1/N_e) \cdot \Theta(1) = \Theta(1); \\ &\text{direct variance, not asymptotic CLT, since } N_e = O(1)). \end{aligned} \quad (\text{R2}\mu\text{-F1.5})$$

$$\begin{aligned} h_0^3 &= (1/M) \sum_i \phi_{i,0} h_0^{3,i} \in \Theta(1/\sqrt{M}) = \Theta(1/\sqrt{N}) \\ &(\text{cross-}i \text{ CLT on entries of coordinate scale } \Theta(1); \\ &\{h_0^{3,i}\} \text{ independent across } i). \end{aligned} \quad (\text{R2}\mu\text{-F1.6})$$

$$\begin{aligned} f_0 &= W_0^4 h_0^3 = \langle W_0^4, h_0^3 \rangle \in \Theta(1/N) \\ &(\text{variance } \sum_a \Theta(1/N^2) \cdot \Theta(1/N) = \Theta(1/N^2); \\ &h_0^3 \text{ entry scale } \Theta(1/\sqrt{N}) \text{ via (R2}\mu\text{-F1.6)}. \end{aligned} \quad (\text{R2}\mu\text{-F1.7})$$

### First backward pass and step-1 updates

**New intermediate scaling.** The intermediate quantity  $(W_0^{3,i})^\top (W_0^4)^\top$  has entries of variance  $\sigma_3^2 \|W_0^4\|^2 = \Theta(1) \cdot \Theta(1/N) = \Theta(1/N)$ , hence coordinate scale  $\Theta(1/\sqrt{N})$ .

**Per-layer gradients (with  $W_0^4/\Delta W^4$  split).**

$$\partial f_0 / \partial h_0^3 = (W_0^4)^\top \in \Theta(1/N). \quad (W_0^4 \text{ entries } \Theta(1/N); \Delta_t W^4 = 0 \text{ at } t = 0) \quad (\text{R2}\mu\text{-B1.1})$$

$$\partial f_0 / \partial h_0^{3,i} = (\phi_{i,0}/M)(W_0^4)^\top \in \Theta(1/(MN)) = \Theta(1/N^2). \quad (\phi_{i,0} \in \Theta(1)) \quad (\text{R2}\mu\text{-B1.2})$$

$$\begin{aligned} \partial f_0 / \partial h_0^{2,i} &= (\phi_{i,0}/M)(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/(M\sqrt{N})) = \Theta(1/N^{3/2}) \\ &((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N}) \text{ from new intermediate scaling above}). \end{aligned} \quad (\text{R2}\mu\text{-B1.3})$$

$$\begin{aligned}
 \partial f_0 / \partial \phi_{i,0} &= (1/M) \langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/N^{3/2}) \\
 \langle h_0^{3,i}, W_0^4 \rangle \text{ variance } &\sum_a \Theta(1) \cdot \Theta(1/N^2) = \Theta(1/N), \text{ coord } \Theta(1/\sqrt{N}); \\
 h_0^{3,i} \text{ entry scale } &\Theta(1) \text{ via (R2}\mu\text{-F1.5)).} \tag{R2}\mu\text{-B1.4}
 \end{aligned}$$

$$\begin{aligned}
 (\partial f_0 / \partial h_0^1)_{\text{exp}} &= (1/M) \sum_i \phi_{i,0} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N^{3/2}) \\
 (\text{at } t=0 \text{ only the all-init } &A_{4,1} \text{ piece is non-zero;} \\
 \text{per summand } (W_0^{2,i})^\top &\text{ acts on } (W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N}) \\
 \text{with } \|v\|^2 = N_e &\Theta(1/N) = \Theta(1/N); \\
 \text{variance } (1/N) \Theta(1/N) &= \Theta(1/N^2), \text{ per-}i \text{ coord } \Theta(1/N); \\
 \text{cross-}i \text{ CLT).} &\tag{R2}\mu\text{-B1.5}
 \end{aligned}$$

$$\begin{aligned}
 (\partial f_0 / \partial h_0^1)_{\text{router}} &= (1/M) Q_0^\top v, \quad v_i := \dot{\phi}_{i,0} \langle h_0^{3,i}, W_0^4 \rangle \\
 &\in \Theta(1/N^{3/2}) \\
 (v_i \in \Theta(1/\sqrt{N}) \text{ via (R2}\mu\text{-B1.4); } &\|v\|^2 = M \Theta(1/N) = \Theta(1); \\
 \text{cross-layer CLT giving } &Q_0^\top v \text{ coord } \Theta(1/\sqrt{N}); \\
 \mathbf{W}^{3,i}\text{-split } v_i = v_i^I &+ v_i^U \text{ where} \\
 v_i^I := \dot{\phi}_{i,0} \langle (W_0^{3,i}) h_0^{2,i}, &W_0^4 \rangle \\
 v_i^U := \dot{\phi}_{i,0} \langle (\Delta_0 W^{3,i}) h_0^{2,i}, &W_0^4 \rangle = 0 \text{ since } \Delta_0 W^{3,i} = 0; \\
 \text{so } (1/M) Q_0^\top v^I = \Theta(1/N^{3/2}) &\text{ and } (1/M) Q_0^\top v^U = 0). \tag{R2}\mu\text{-B1.6}
 \end{aligned}$$

$$\partial f_0 / \partial h_0^1 = (\partial f_0 / \partial h_0^1)_{\text{exp}} + (\partial f_0 / \partial h_0^1)_{\text{router}} \in \Theta(1/N^{3/2}). \tag{R2}\mu\text{-B1.7}$$

**Step-1 parameter updates.**

$$\begin{aligned}
 \Delta_1 W^4 &= -(\eta/N) \chi_0 (h_0^3)^\top \in \Theta(1/N^{3/2}) \\
 (h_0^3 \text{ entry scale } &\Theta(1/\sqrt{N}) \text{ via (R2}\mu\text{-F1.6); sub-leading vs } W_0^4 \in \Theta(1/N)). \tag{R2}\mu\text{-U1.4}
 \end{aligned}$$

$$\begin{aligned}
 \Delta_1 W^{3,i} &= -\eta_3 \chi_0 (\phi_{i,0}/M) (W_0^4)^\top (h_0^{2,i})^\top = -\eta N \chi_0 \phi_{i,0} (W_0^4)^\top (h_0^{2,i})^\top \\
 &\in \Theta(1) \text{ rank-1 along } (W_0^4)^\top \otimes h_0^{2,i} \\
 (\eta_3 = \eta M N; h_0^{2,i} \in &\Theta(1) \text{ via (R2}\mu\text{-F1.4);} \\
 \text{same scale as } W_0^{3,i} \in &\Theta(1)). \tag{R2}\mu\text{-U1.3}
 \end{aligned}$$

$$\begin{aligned}
 \Delta_1 W^{2,i} &= -(\eta \chi_0 \phi_{i,0}/N) (W_0^{3,i})^\top (W_0^4)^\top (h_0^1)^\top \in \Theta(1/N^{3/2}) \text{ rank-1} \\
 (\eta_2 = \eta M/N; (W_0^{3,i})^\top &(W_0^4)^\top \in \Theta(1/\sqrt{N}) \text{ from new intermediate scaling;} \\
 \text{sub-leading vs } W_0^{2,i} \in &\Theta(1/\sqrt{N})). \tag{R2}\mu\text{-U1.2}
 \end{aligned}$$

$$\begin{aligned} \Delta_1 W^1 &= -\eta_1 \chi_0 (\partial f_0 / \partial h_0^1) x^\top \in \Theta(1/\sqrt{N}) \\ &(\eta_1 = \eta N; \partial f_0 / \partial h_0^1 \in \Theta(1/N^{3/2}) \text{ via (R2}\mu\text{-B1.7)}). \end{aligned} \quad (\text{R2}\mu\text{-U1.1a})$$

$$\Delta_1 h^1 = \Delta_1 W^1 x \in \Theta(1/\sqrt{N}) \text{ aligned along } \partial f_0 / \partial h_0^1. \quad (\text{R2}\mu\text{-U1.1b})$$

$$\begin{aligned} \Delta_1 Q &= -\eta_Q \chi_0 (\partial f_0 / \partial \phi) (h_0^1)^\top \in \Theta(1/N^{3/2}). \quad (\partial f_0 / \partial \phi \in \Theta(1/N^{3/2}) \text{ via (R2}\mu\text{-B1.4)}) \\ &(\text{R2}\mu\text{-U1.Q}) \end{aligned}$$

### Second forward pass

$$h_1^1 = h_0^1 + \Delta_1 h^1 \in \Theta(1). \quad (h_0^1 \in \Theta(1) \text{ dominates } \Delta_1 h^1 \in \Theta(1/\sqrt{N}) \text{ via (R2}\mu\text{-U1.1b)}) \quad (\text{R2}\mu\text{-F2.1})$$

$$\psi_1 = \psi_0 + \Delta_1 \psi \in \Theta(1), \quad \phi_1 \in \Theta(1) \quad (\text{R2}\mu\text{-F2.2--R2}\mu\text{-F2.3})$$

$$(\Delta_1 \psi \in \Theta(1/\sqrt{N}) \text{ from } Q_0 \Delta_1 h^1 \text{ and } \Delta_1 Q h_0^1;$$

$$\text{scales by (R2}\mu\text{-U1.1b), (R2}\mu\text{-U1.Q);}$$

$$\text{cross } \Delta_1 Q \Delta_1 h^1 \in \Theta(1/N^{3/2}) \text{ sub-leading via } \partial f_0 / \partial h_0^1 \text{ via (R2}\mu\text{-B1.7)})$$

(R2 $\mu$ -F2.4)  $h_1^{2,i} = W_1^{2,i} h_1^1$ , four-piece decomposition:

$$h_1^{2,i} = h_0^{2,i} + W_0^{2,i} \Delta_1 h^1 + \Delta_1 W^{2,i} h_0^1 + \Delta_1 W^{2,i} \Delta_1 h^1. \quad (\text{R2}\mu\text{-F2.4})$$

$$\text{init: } W_0^{2,i} h_0^1 = h_0^{2,i} \in \Theta(1). \quad (\text{R2}\mu\text{-F1.4}) \quad (\text{R2}\mu\text{-F2.4a})$$

$$\text{prop: } W_0^{2,i} \Delta_1 h^1 \in \Theta(1/\sqrt{N}) \quad (\text{R2}\mu\text{-F2.4b})$$

$$(\sigma_2^2 \|\Delta_1 h^1\|^2 = (1/N) \Theta(1) = \Theta(1/N);$$

$$\|\Delta_1 h^1\|^2 \in \Theta(1) \text{ via (R2}\mu\text{-U1.1b)})$$

$$\text{eff: } \Delta_1 W^{2,i} h_0^1 = -(\eta \chi_0 \phi_{i,0} / N) \|h_0^1\|^2 (W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N}) \quad (\text{R2}\mu\text{-F2.4c})$$

((R2 $\mu$ -U1.2)substitution;

$$(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N}) \text{ from new intermediate scaling;}$$

$$\|h_0^1\|^2 \in \Theta(N))$$

$$\begin{aligned} \text{cross: } \Delta_1 W^{2,i} \Delta_1 h^1 &= -(\eta \chi_0 \phi_{i,0} / N) \langle h_0^1, \Delta_1 h^1 \rangle (W_0^{3,i})^\top (W_0^4)^\top \\ &\in \Theta(1/N^{3/2}) \end{aligned} \quad (\text{R2}\mu\text{-F2.4d})$$

$$((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N}) \text{ from new intermediate scaling;}$$

$$\langle h_0^1, \Delta_1 h^1 \rangle \in \Theta(1) \text{ random, derived via}$$

$$\partial f_0 / \partial h_0^1 \in \Theta(1/N^{3/2}) \text{ via (R2}\mu\text{-B1.7);}$$

sub-leading vs propagating and effective)

(R2 $\mu$ -F2.5)  $h_1^{3,i} = W_1^{3,i} h_1^{2,i}$ , four-piece decomposition:

$$h_1^{3,i} = h_0^{3,i} + W_0^{3,i} \Delta_1 h^{2,i} + \Delta_1 W^{3,i} h_0^{2,i} + \Delta_1 W^{3,i} \Delta_1 h^{2,i}. \quad (\text{R2}\mu\text{-F2.5})$$

$$\text{init: } h_0^{3,i} \in \Theta(1). \quad (\text{R2}\mu\text{-F1.5}) \quad (\text{R2}\mu\text{-F2.5a})$$

$$\text{prop: } W_0^{3,i} \Delta_1 h^{2,i} \in \Theta(1/\sqrt{N}) \quad (\text{R2}\mu\text{-F2.5b})$$

$$(\sigma_3^2 \|\Delta_1 h^{2,i}\|^2 = \Theta(1) \cdot \Theta(1/N) = \Theta(1/N);$$

$$\|\Delta_1 h^{2,i}\|^2 \in \Theta(1) \text{ via (R2}\mu\text{-F2.4)})$$

$$\text{eff: } \Delta_1 W^{3,i} h_0^{2,i} = -\eta N \chi_0 \phi_{i,0} \|h_0^{2,i}\|^2 (W_0^4)^\top \in \Theta(1) \text{ along } (W_0^4)^\top \quad (\text{R2}\mu\text{-F2.5c})$$

$$((\text{R2}\mu\text{-U1.3})\text{substitution; } \|h_0^{2,i}\|^2 \in \Theta(1) \text{ via (R2}\mu\text{-F1.4);}$$

$$(W_0^4)^\top \in \Theta(1/N))$$

$$\text{cross: } \Delta_1 W^{3,i} \Delta_1 h^{2,i} = -\eta N \chi_0 \phi_{i,0} \langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle (W_0^4)^\top \in \Theta(1/\sqrt{N}) \text{ along } (W_0^4)^\top \quad (\text{R2}\mu\text{-F2.5d})$$

((R2 $\mu$ -U1.3)substitution;

dominant piece of  $\Delta_1 h^{2,i} \propto (W_0^{3,i})^\top (W_0^4)^\top$  via (R2 $\mu$ -F2.4c),

so  $\langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle \sim -(\eta \chi_0 \phi/N) \|h_0^1\|^2 \langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$

random across  $i$ , via (R2 $\mu$ -B1.4);

random sign across  $i$ , sub-leading vs effective)

**Effective dominates:**  $\Delta_1 h^{3,i} \in \Theta(1)$  aligned along  $(W_0^4)^\top$  — the engine of feature learning at  $t = 1$ .

(R2 $\mu$ -F2.6)  $h_1^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D$ , where

$$A_1 := \frac{1}{M} \sum_i \phi_{i,1} h_0^{3,i},$$

$$A_{2,1} := \frac{1}{M} \sum_i \phi_{i,1} W_0^{3,i} W_0^{2,i} \Delta_1 h^1,$$

$$A_{2,2} := \frac{1}{M} \sum_i \phi_{i,1} W_0^{3,i} \Delta_1 W^{2,i} h_0^1,$$

$$A_3 := \frac{1}{M} \sum_i \phi_{i,1} \Delta_1 W^{3,i} h_0^{2,i},$$

$$D := \frac{1}{M} \sum_i \phi_{i,1} \Delta_1 W^{3,i} \Delta_1 h^{2,i}.$$

$$A_1 \in \Theta(1/\sqrt{N}). \quad (\text{cross-}i \text{ CLT on } \Theta(1)\text{-entry independent vectors}) \quad (\text{R2}\mu\text{-F2.6a})$$

$$A_{2,1} \in \Theta(1/(\sqrt{M}\sqrt{N})) = \Theta(1/N) \quad (\text{R2}\mu\text{-F2.6b})$$

(per summand chain CLT entries  $\Theta(1/\sqrt{N})$  random;

$$\|\Delta_1 h^1\|^2 \in \Theta(1) \text{ via (R2}\mu\text{-U1.1b); cross-}i \text{ CLT})$$

$$A_{2,2} = -(\eta \chi_0/N) \|h_0^1\|^2 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} W_0^{3,i} (W_0^{3,i})^\top (W_0^4)^\top$$

$$\in \Theta(1/N) \text{ in mostly-random directions} \quad (\text{R2}\mu\text{-F2.6c})$$

((R2 $\mu$ -U1.2)substitution;

$W_0^{3,i} (W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N})$  in mostly-random directions

by rank-deficient initialization Gram (§J.3.1);

$\|h_0^1\|^2 \in \Theta(N)$ ; per summand entries  $\Theta(1/\sqrt{N})$ , cross- $i$  CLT)

$$A_3 = -\eta N \chi_0 (W_0^4)^\top \left( \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \|h_0^{2,i}\|^2 \right) \in \Theta(1) \text{ along } (W_0^4)^\top \quad (\text{R2}\mu\text{-F2.6d})$$

((R2 $\mu$ -U1.3)substitution;  $\|h_0^{2,i}\|^2 \in \Theta(1)$  via (R2 $\mu$ -F1.4); LLN;  
 $(W_0^4)^\top \in \Theta(1/N)$ ; entry  $\eta N \cdot (1/N) \cdot \Theta(1) = \Theta(1)$ )

$$D = -\eta N \chi_0 (W_0^4)^\top \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle \in \Theta(1/N) \text{ along } (W_0^4)^\top \quad (\text{R2}\mu\text{-F2.6e})$$

((R2 $\mu$ -U1.3)substitution;

$\langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle \in \Theta(1/\sqrt{N})$  random across  $i$

via  $\langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$  via (R2 $\mu$ -B1.4);

cross- $i$  CLT in the random factor;

sub-leading vs  $A_3 \in \Theta(1)$ )

$$h_1^3 \in \Theta(1) \text{ along } (W_0^4)^\top, \text{ from } A_3 \text{ alone.} \quad (\text{R2}\mu\text{-F2.6})$$

$$f_1 = W_1^4 h_1^3 \approx \langle W_0^4, A_3 \rangle = -\eta N \chi_0 \|W_0^4\|^2 \left( \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \|h_0^{2,i}\|^2 \right) \in \Theta(1) \quad (\text{R2}\mu\text{-F2.7})$$

( $A_3$  aligned along  $(W_0^4)^\top$  via (R2 $\mu$ -F2.6d);

$\|W_0^4\|^2 = N \cdot (1/N^2) = 1/N$ ;  $\|h_0^{2,i}\|^2 \in \Theta(1)$  via (R2 $\mu$ -F1.4);

net  $\eta N \cdot \Theta(1) \cdot (1/N) = \Theta(1)$ )

## Second backward pass and step-2 updates

**New intermediate at  $t = 1$ .**

$$\begin{aligned} (\Delta_1 W^{3,i})^\top (W_1^4)^\top &\approx -\eta N \chi_0 \phi_{i,0} \|W_0^4\|^2 h_0^{2,i} \in \Theta(1) \text{ aligned } \|h_0^{2,i}\| \\ &((\text{R2}\mu\text{-U1.3})\text{rank-1 substitution; } \|W_0^4\|^2 = 1/N; \\ &h_0^{2,i} \in \Theta(1) \text{ via (R2}\mu\text{-F1.4); net } \eta N \cdot (1/N) = \Theta(1)) \end{aligned}$$

This  $\Theta(1)$  aligned-along- $h_0^{2,i}$  contribution from the rank-1 update is what drives  $\partial f_1 / \partial h_1^{2,i}$  to climb from  $\Theta(1/N^{3/2})$  at  $t = 0$  to  $\Theta(1/N)$  at  $t = 1$ .

**Per-layer gradients (with  $W_0^4 / \Delta W^4$  binary split).**

$$\partial f_1 / \partial h_1^3 = (W_0^4)^\top + (\Delta_1 W^4)^\top \in \Theta(1/N) \quad (\text{R2}\mu\text{-B2.1})$$

( $W_0^4 \in \Theta(1/N)$  dominates  $\Delta_1 W^4 \in \Theta(1/N^{3/2})$  via (R2 $\mu$ -U1.4))

$$\partial f_1 / \partial h_1^{3,i} = (\phi_{i,1}/M)(W_1^4)^\top \in \Theta(1/N^2) \quad (\text{R2}\mu\text{-B2.2})$$

(init piece  $\Theta(1/N^2)$ ;  $\Delta_1 W^4$  piece  $\Theta(1/N^{5/2})$  via (R2 $\mu$ -U1.4))

(R2 $\mu$ -B2.3)  $\partial f_1 / \partial h_1^{2,i} = (\phi_{i,1}/M)(W_1^{3,i})^\top (W_1^4)^\top$ , four-piece decomposition:

$$(\phi_{i,1}/M)(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N^{3/2}) \quad (\text{R2}\mu\text{-B2.3a})$$

( $(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/\sqrt{N})$  from new intermediate scaling at  $t = 0$ )

$$\begin{aligned}
 (\phi_{i,1}/M)(W_0^{3,i})^\top (\Delta_1 W^4)^\top &= -(\eta\chi_0\phi_{i,0}/(MN))(W_0^{3,i})^\top h_0^3 \\
 &\in \Theta(1/N^2) \tag{R2\mu-B2.3b}
 \end{aligned}$$

((R2\mu-U1.4)substitution;

$$(W_0^{3,i})^\top h_0^3 \in \Theta(1)$$

from  $j = i$  coherent contribution to the average)

$$(\phi_{i,1}/M)(\Delta_1 W^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ along } h_0^{2,i} \tag{R2\mu-B2.3c}$$

$$((\Delta_1 W^{3,i})^\top (W_0^4)^\top \in \Theta(1) \text{ along } h_0^{2,i}$$

from new intermediate at  $t = 1$  above; **dominant**)

$$(\phi_{i,1}/M)(\Delta_1 W^{3,i})^\top (\Delta_1 W^4)^\top \in \Theta(1/N^2) \tag{R2\mu-B2.3d}$$

(same alignment as (R2\mu-B2.3c);

$\Delta_1 W^4$  smaller by  $1/\sqrt{N}$  via (R2\mu-U1.4))

$$\partial f_1 / \partial h_1^{2,i} \in \Theta(1/N) \text{ along } h_0^{2,i}, \text{ from (R2\mu-B2.3c)}. \tag{R2\mu-B2.3}$$

$$\partial f_1 / \partial \phi_{i,1} = (1/M) \langle h_1^{3,i}, W_1^4 \rangle \in \Theta(1/N) \tag{R2\mu-B2.4}$$

(Decompose  $h_1^{3,i} = (W_0^{3,i}) h_1^{2,i} + (\Delta_1 W^{3,i}) h_1^{2,i}$  and  $W_1^4 = W_0^4 + \Delta_1 W^4$ ;

four-piece grid:

$$(\text{init} \cdot \text{init}) \langle (W_0^{3,i}) h_1^{2,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$$

via random IP  $\langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$  (R2\mu-B1.4);

$$(\text{init} \cdot \text{update}) \langle (W_0^{3,i}) h_1^{2,i}, \Delta_1 W^4 \rangle = -(\chi_0/N) \langle h_0^{3,i}, h_0^3 \rangle \in \Theta(1/N)$$

via  $\langle h_0^{3,i}, h_0^3 \rangle \in \Theta(1)$  ( $j = i$  diagonal  $(\phi_{i,0}/M) \|h_0^{3,i}\|^2 = \Theta(1)$  in  $\mu\text{P}$ );

$$(\text{update} \cdot \text{init}) \langle (\Delta_1 W^{3,i}) h_1^{2,i}, W_0^4 \rangle = -\eta N \chi_0 \phi_{i,0} \|h_0^{2,i}\|^2 \|W_0^4\|^2 \in \Theta(1)$$

via rank-1 alignment of  $\Delta_1 W^{3,i} h_0^{2,i}$  with  $(W_0^4)^\top$ ;  $\eta N \cdot \Theta(1) \cdot \Theta(1/N) = \Theta(1)$ ;

$$(\text{update} \cdot \text{update}) \langle (\Delta_1 W^{3,i}) h_1^{2,i}, \Delta_1 W^4 \rangle = \kappa \langle W_0^4, \Delta_1 W^4 \rangle \in \Theta(1/N)$$

via  $\kappa \in \Theta(N)$  (rank-1 prefactor) and  $\langle W_0^4, \Delta_1 W^4 \rangle \in \Theta(1/N^2)$ ;

/M:  $(\Theta(1/N^{3/2}), \Theta(1/N^2), \Theta(1/N), \Theta(1/N^2))$  respectively;

(update · init) dominates; total  $\Theta(1/N)$ )

**Expert pathway:**  $(\partial f_1 / \partial h_1^1)_{\text{exp}} = A_4 + A_5 + A_6 + E$ , **full 8-piece expansion.** Each of  $A_4, A_5, A_6, E$  splits into a .1 piece (using  $W_0^4$ ) and a .2 piece (using  $\Delta_1 W^4$ ):

$$A_{4.1} \in \Theta(1/N^{3/2}). \tag{R2\mu-B1.5} \tag{R2\mu-B2.5a.1}$$

$$A_{4.2} \in \Theta(1/N^2). \quad (\Delta_1 W^4 \text{ smaller by } 1/\sqrt{N} \text{ via (R2\mu-U1.4)}) \tag{R2\mu-B2.5a.2}$$

$$A_{5.1} = \frac{1}{M} \sum_i \phi_{i,1} (W_0^{2,i})^\top (\Delta_1 W^{3,i})^\top (W_0^4)^\top$$

(R2 $\mu$ -B2.5b.1)

$\in \Theta(1/N)$  in mostly-random directions

$((\Delta_1 W^{3,i})^\top (W_0^4)^\top) \in \Theta(1)$  along  $h_0^{2,i}$  from new intermediate at  $t = 1$ ;

$(W_0^{2,i})^\top h_0^{2,i} = (W_0^{2,i})^\top W_0^{2,i} h_0^1 \in \Theta(1/\sqrt{N})$  in mostly-random directions

by rank-deficient initialization Gram (§J.3.1);

coherent  $(N_e/N) h_0^1$  piece is  $\Theta(1/N)$  sub-leading;

cross- $i$  CLT:  $\Theta(1/\sqrt{N})/\sqrt{M} = \Theta(1/N)$

$$A_{5.2} \in \Theta(1/N^2). \quad (\text{same structure as (R2}\mu\text{-B2.5b.1); } \Delta_1 W^4 \text{ smaller by } 1/\sqrt{N} \text{ via (R2}\mu\text{-U1.4)})$$

(R2 $\mu$ -B2.5b.2)

$$A_{6.1} = \frac{1}{M} \sum_i \phi_{i,1} (\Delta_1 W^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top$$

$$= -(\eta\chi_0/N) h_0^1 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} S_i \in \Theta(1/N^2) \text{ along } h_0^1$$

(R2 $\mu$ -B2.5c.1)

((R2 $\mu$ -U1.2)substitution;  $S_i := W_0^4 W_0^{3,i} (W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N)$

quadratic form in  $W_0^4$  across rank- $N_e$  Gram, leading piece  $\|W_0^4\|^2$ ;

per summand entries  $(1/N) \cdot 1 \cdot (1/N) = \Theta(1/N^2)$  aligned along  $h_0^1$ )

$$A_{6.2} \in \Theta(1/N^3). \quad (\text{same structure as (R2}\mu\text{-B2.5c.1); } \Delta_1 W^4 \text{ smaller by } 1/\sqrt{N} \text{ via (R2}\mu\text{-U1.4)})$$

(R2 $\mu$ -B2.5c.2)

$$E_1 = \frac{1}{M} \sum_i \phi_{i,1} (\Delta_1 W^{2,i})^\top (\Delta_1 W^{3,i})^\top (W_0^4)^\top \in \Theta(1/N^2)$$

(R2 $\mu$ -B2.5d.1)

(combining (R2 $\mu$ -U1.2), (R2 $\mu$ -U1.3) substitutions;

per summand involves  $\langle W_0^4, h_0^{3,i} \rangle \in \Theta(1/\sqrt{N})$  random via (R2 $\mu$ -B1.4);

per summand entries  $\Theta(1/N^{3/2})$ ; cross- $i$  CLT)

$$E_2 \in \Theta(1/N^{5/2}). \quad (\text{same as (R2}\mu\text{-B2.5d.1); } \Delta_1 W^4 \text{ smaller by } 1/\sqrt{N} \text{ via (R2}\mu\text{-U1.4)})$$

(R2 $\mu$ -B2.5d.2)

$$(\partial f_1 / \partial h_1^1)_{\text{exp}} \in \Theta(1/N) \text{ in mostly-random directions, from } A_{5.1} \text{ alone.} \quad (\text{R2}\mu\text{-B2.5})$$

**Router pathway.** Decomposing  $\mathbf{v} = \mathbf{v}^{(0)} + \mathbf{v}^{(\Delta)}$  where  $\mathbf{v}_i^{(0)} = \dot{\phi}_{i,1} \langle h_1^{3,i}, W_0^4 \rangle \in \Theta(1)$  (coherent via the alignment in  $h_1^{3,i}$  from (R2 $\mu$ -F2.5c) / (R2 $\mu$ -B2.4)) and  $\mathbf{v}_i^{(\Delta)} \in \Theta(1/N)$  (via (R2 $\mu$ -U1.4)):

$$(1/M) Q_0^\top \mathbf{v}^{(0)} \in \Theta(1/N) \text{ random direction} \quad (\text{R2}\mu\text{-B2.6a})$$

( $Q_0$  indep. of  $\mathbf{v}^{(0)}$ ;  $\sigma_Q^2 = 1/N$ ;  $\|\mathbf{v}^{(0)}\|^2 \in \Theta(M)$ )

$$(1/M) Q_0^\top \mathbf{v}^{(\Delta)} \in \Theta(1/N^2). \quad (\text{R2}\mu\text{-B2.6b})$$

$$(1/M)\Delta_1 Q^\top \mathbf{v}^{(0)} \in \Theta(1/N^2) \text{ along } h_0^1 \quad (\text{R2}\mu\text{-B2.6c})$$

$(\Delta_1 Q \text{ rank-1 along } (h_0^1)^\top \text{ via (R2}\mu\text{-U1.Q);}$

$$\langle w_0, \mathbf{v}^{(0)} \rangle = (\dot{\phi}^2/M) \sum_i [r_i c_i + r_i^2] \in \Theta(1/N)$$

(both the cross- $i$  CLT piece  $\sum r_i c_i \sim 1$  and the LLN self-pairing

$\sum r_i^2 \sim M/N$  contribute  $\Theta(1)$  before the  $(1/M)$  prefactor);

$$\Delta_1 Q^\top \mathbf{v}^{(0)} = h_0^1 \eta_Q \chi_0 \langle w_0, \mathbf{v}^{(0)} \rangle \sim 1/N \text{ entries;}$$

after  $(1/M)$ :  $\Theta(1/N^2)$  — sub-leading vs (R2 $\mu$ -B2.6a))

$$(1/M)\Delta_1 Q^\top \mathbf{v}^{(\Delta)} \in \Theta(1/N^3). \quad (\text{R2}\mu\text{-B2.6d})$$

$(\partial f_1/\partial h_1^1)_{\text{router}} \in \Theta(1/N)$ , from (R2 $\mu$ -B2.6a) alone (other three pieces sub-leading).

(R2 $\mu$ -B2.6)

$W^{3,i}$ -imbalance propagates into the router pathway. The  $W^{3,i}$ -split of (R2 $\mu$ -B2.4),  $\mathbf{v} = \mathbf{v}^I + \mathbf{v}^U$  with  $\mathbf{v}_i^I = \dot{\phi}_{i,1} \langle (W_0^{3,i}) h_1^{2,i}, W_1^4 \rangle$  and  $\mathbf{v}_i^U = \dot{\phi}_{i,1} \langle (\Delta_1 W^{3,i}) h_1^{2,i}, W_1^4 \rangle$ , transmits as:  $\mathbf{v}_i^I \in \Theta(1/\sqrt{N})$  random across  $i$  (from  $\langle h_0^{3,i}, W_0^4 \rangle$ ), giving  $\|\mathbf{v}^I\|^2 \in \Theta(1)$  and  $(1/M)Q_t^\top \mathbf{v}^I \in \Theta(1/N^{3/2})$  via cross-layer CLT; whereas  $\mathbf{v}_i^U \in \Theta(1)$  coherent across  $i$  (rank-1 alignment with  $(W_0^4)^\top$ ), giving  $\|\mathbf{v}^U\|^2 \in \Theta(M)$  and  $(1/M)Q_t^\top \mathbf{v}^U \in \Theta(1/N)$ . The same  $\sqrt{N}$  deficit between  $W_0^{3,i}$ - and  $\Delta W^{3,i}$ -pieces is transmitted into  $\partial f_1/\partial h_1^1$  via the router pathway.

$$\partial f_1/\partial h_1^1 = (\partial f_1/\partial h_1^1)_{\text{exp}} + (\partial f_1/\partial h_1^1)_{\text{router}} \in \Theta(1/N). \quad (\text{R2}\mu\text{-B2.7})$$

### Step-2 parameter updates.

$$\Delta_2 W^4 = -(\chi_1/N)(h_1^3)^\top \in \Theta(1/N). \quad (h_1^3 \in \Theta(1) \text{ via (R2}\mu\text{-F2.6); now comparable to } W_0^4) \quad (\text{R2}\mu\text{-U2.4})$$

$$\Delta_2 W^{3,i} = -\eta N \chi_1 \phi_{i,1} (W_1^4)^\top (h_1^{2,i})^\top \in \Theta(1) \text{ rank-1 along } (W_1^4)^\top \otimes h_1^{2,i} \quad (\text{R2}\mu\text{-U2.3})$$

(same structure as (R2 $\mu$ -U1.3);

$$(W_1^4)^\top \in \Theta(1/N); h_1^{2,i} \in \Theta(1) \text{ via (R2}\mu\text{-F2.4)}$$

$$\Delta_2 W^{2,i} = -(\eta \chi_1 \phi_{i,1}/N)(W_1^{3,i})^\top (W_1^4)^\top (h_1^1)^\top \in \Theta(1/N) \text{ rank-1 along } h_0^{2,i} \otimes h_1^1 \quad (\text{R2}\mu\text{-U2.2})$$

$((W_1^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ along } h_0^{2,i} \text{ via new intermediate at } t = 1;$

$\sqrt{N}$  larger than  $\Delta_1 W^{2,i} \in \Theta(1/N^{3/2})$ )

$$\Delta_2 W^1 = -\eta_1 \chi_1 (\partial f_1/\partial h_1^1) x^\top \in \Theta(1) \quad (\text{R2}\mu\text{-U2.1a})$$

$$(\eta_1 = \eta N; \partial f_1/\partial h_1^1 \in \Theta(1/N) \text{ via (R2}\mu\text{-B2.7)})$$

$$\Delta_2 h^1 = \Delta_2 W^1 x \in \Theta(1) \text{ aligned along } h_0^1. \quad (\text{embedding feature-learns at } t = 2) \quad (\text{R2}\mu\text{-U2.1b})$$

$$\Delta_2 Q \in \Theta(1/N). \quad (\partial f_1/\partial \phi \in \Theta(1/N) \text{ via (R2}\mu\text{-B2.4)}) \quad (\text{R2}\mu\text{-U2.Q})$$

**Third forward pass** We compute activations at  $\theta^{(2)}$ . The cumulative embedding change  $\Delta h^1 := h_2^1 - h_0^1 = \Delta_1 h^1 + \Delta_2 h^1$ , with  $\Delta_1 h^1 \in \Theta(1/\sqrt{N})$  via (R2 $\mu$ -U1.1b) and  $\Delta_2 h^1 \in \Theta(1)$  aligned along  $h_0^1$  via (R2 $\mu$ -U2.1b); the latter dominates.

$$h_2^1 = h_0^1 + \Delta h^1 \in \Theta(1). \quad (\text{R2}\mu\text{-F3.1})$$

(R2 $\mu$ -F3.2)  $h_2^{2,i} = W_2^{2,i} h_2^1$ , four-piece decomposition with cumulative  $\Delta W^{2,i}$ :

$$\text{init: } W_0^{2,i} h_0^1 = h_0^{2,i} \in \Theta(1). \quad (\text{R2}\mu\text{-F1.4}) \quad (\text{R2}\mu\text{-F3.2a})$$

$$\text{prop: } W_0^{2,i} \Delta h^1 \in \Theta(1) \quad (\text{R2}\mu\text{-F3.2b})$$

$$(\sigma_2^2 \|\Delta h^1\|^2 = (1/N) \Theta(N) = \Theta(1);$$

$$\|\Delta h^1\|^2 \in \Theta(N) \text{ via } \Delta_2 h^1 \in \Theta(1) \text{ via (R2}\mu\text{-U2.1b)})$$

$$\text{eff: } \Delta_2 W^{2,i} h_0^1 = -(\eta \chi_1 \phi_{i,1}/N) (W_1^{3,i})^\top (W_1^4)^\top \langle h_1^1, h_0^1 \rangle \in \Theta(1) \text{ along } h_0^{2,i} \quad (\text{R2}\mu\text{-F3.2c})$$

((R2 $\mu$ -U2.2)substitution;

$$(W_1^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ along } h_0^{2,i} \text{ via new intermediate at } t = 1;$$

$$\langle h_1^1, h_0^1 \rangle \in \Theta(N) \text{ coherent})$$

$$\text{cross: } \Delta_2 W^{2,i} \Delta h^1 = -(\eta \chi_1 \phi_{i,1}/N) (W_1^{3,i})^\top (W_1^4)^\top \langle h_1^1, \Delta h^1 \rangle \in \Theta(1) \text{ along } h_0^{2,i} \quad (\text{R2}\mu\text{-F3.2d})$$

((R2 $\mu$ -U2.2)substitution;

$$(W_1^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ along } h_0^{2,i} \text{ via new intermediate at } t = 1;$$

$$\langle h_1^1, \Delta h^1 \rangle \in \Theta(N) \text{ coherent via } \Delta h^1 \parallel h_0^1 \text{ from (R2}\mu\text{-U2.1b);}$$

leading-scale contribution)

$\Delta h^{2,i} \in \Theta(1)$  entry-wise, with the coherent  $\Theta(1)$ -along- $h_0^{2,i}$  piece emerging from (R2 $\mu$ -F3.2c) + (R2 $\mu$ -F3.2d).

(R2 $\mu$ -F3.3)  $h_2^{3,i} = W_2^{3,i} h_2^{2,i}$ , four-piece decomposition with cumulative  $\Delta W^{3,i}$ :

$$\text{init: } h_0^{3,i} \in \Theta(1). \quad (\text{R2}\mu\text{-F1.5}) \quad (\text{R2}\mu\text{-F3.3a})$$

$$\text{prop: } W_0^{3,i} \Delta h^{2,i} \in \Theta(1) \quad (\text{R2}\mu\text{-F3.3b})$$

$$(\sigma_3^2 \|\Delta h^{2,i}\|^2 = \Theta(1) \cdot \Theta(1) = \Theta(1);$$

$$\|\Delta h^{2,i}\|^2 \in \Theta(1) \text{ via (R2}\mu\text{-F3.2); } N_e = \Theta(1))$$

$$\text{eff: } \Delta_2 W^{3,i} h_0^{2,i} = -\eta N \chi_1 \phi_{i,1} \|h_0^{2,i}\|^2 (W_1^4)^\top \in \Theta(1) \text{ along } (W_0^4)^\top \quad (\text{R2}\mu\text{-F3.3c})$$

((R2 $\mu$ -U2.3)substitution;  $\|h_0^{2,i}\|^2 \in \Theta(1)$  via (R2 $\mu$ -F1.4);

$$(W_1^4)^\top \in \Theta(1/N) \text{ dominated by } (W_0^4)^\top)$$

cross:  $\Delta_2 W^{3,i} \Delta h^{2,i}$  tracked as part of  $D'$  in (R2 $\mu$ -F3.4). (R2 $\mu$ -F3.3d)

$\Delta h^{3,i} \in \Theta(1)$  aligned along  $(W_0^4)^\top$ .

(R2 $\mu$ -F3.4)  $h_2^3 = A'_1 + A'_{2,1} + A'_{2,2} + A'_3 + D'$ , where

$$\begin{aligned} A'_1 &:= \frac{1}{M} \sum_i \phi_{i,2} h_0^{3,i}, & A'_{2,1} &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,i} W_0^{2,i} \Delta h^1, \\ A'_{2,2} &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,i} \Delta_2 W^{2,i} h_0^1, & A'_3 &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} h_0^{2,i}, \\ D' &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} \Delta h^{2,i}. \end{aligned}$$

$A'_1 \in \Theta(1/\sqrt{N})$ . (cross- $i$  CLT on  $\Theta(1)$ -entry independent vectors; same as (R2 $\mu$ -F2.6a)) (R2 $\mu$ -F3.4a)

$A'_{2,1} \in \Theta(1/\sqrt{N})$  (R2 $\mu$ -F3.4b)

(climbs from  $\Theta(1/N)$  at  $t = 1$  because  $\|\Delta h^1\|^2$   
grows from  $\Theta(1)$  to  $\Theta(N)$  via (R2 $\mu$ -U2.1b))

$A'_{2,2} \in \Theta(1/\sqrt{N})$  in mostly-random directions (R2 $\mu$ -F3.4c)

(same mechanism as  $A'_{2,1}$ ;  $\|\Delta h^1\|^2 \in \Theta(N)$  via (R2 $\mu$ -U2.1b);

rank-deficient initialization Gram (§J.3.1)

still does not collapse, so sub-leading)

$A'_3 = -\eta N \chi_1 (W_0^4)^\top \left( \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1} \|h_0^{2,i}\|^2 \right) \in \Theta(1)$  along  $(W_0^4)^\top$  (R2 $\mu$ -F3.4d)

(same structure as (R2 $\mu$ -F2.6d) with cumulative  $\Delta_2 W^{3,i}$ )

$D' = -\eta N \chi_1 (W_0^4)^\top \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1} \langle h_0^{2,i}, \Delta h^{2,i} \rangle$   
 $\in \Theta(1)$  along  $(W_0^4)^\top$  (R2 $\mu$ -F3.4e)

(climbs from  $\Theta(1/N)$  at  $t = 1$  to  $\Theta(1)$  at  $t = 2$ ;

$\langle h_0^{2,i}, \Delta h^{2,i} \rangle \in \Theta(1)$  coherent across  $i$  via (R2 $\mu$ -F3.2c)+(R2 $\mu$ -F3.2d),

replacing the  $t = 1$  random factor  $\langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$  via (R2 $\mu$ -B1.4);

LLN gives empirical average  $\Theta(1)$  rather than  $\Theta(1/\sqrt{M})$ )

$h_2^3 \in \Theta(1)$  along  $(W_0^4)^\top$ , from  $A'_3 + D'$ . (R2 $\mu$ -F3.4)

$f_2 = W_2^4 h_2^3 \in \Theta(1)$ . (R2 $\mu$ -F3.5)

( $A'_3 + D'$  aligned along  $(W_0^4)^\top$ ;  $\|W_0^4\|^2 = 1/N$ ; net  $\eta N \cdot \Theta(1) \cdot (1/N) = \Theta(1)$ )

**Third backward pass** The structure parallels (R2 $\mu$ -B2.1)–(R2 $\mu$ -B2.7) with cumulative  $\Delta W$ 's. The key change at  $t = 2$  vs  $t = 1$ :

$(\Delta W^4)^\top \in \Theta(1/N)$ , climbs from  $\Theta(1/N^{3/2})$  at  $t = 1$  (R2 $\mu$ -B3.1)

( $h_1^3 \in \Theta(1)$  via (R2 $\mu$ -F2.6);  $(\Delta_2 W^4)^\top \in \Theta(1/N)$  via (R2 $\mu$ -U2.4);

now comparable to  $(W_0^4)^\top \in \Theta(1/N)$ )

Combined with the alignment mechanism (driven by  $\langle W_1^4, h_1^{3,i} \rangle \in \Theta(1)$  coherent across  $i$  via (R2 $\mu$ -B2.4)), several pieces climb at  $t = 2$ :

**Router pathway  $\partial f_2/\partial\phi_{i,2}$ .**

$$\partial f_2/\partial\phi_{i,2} = (1/M)\langle h_2^{3,i}, W_2^4 \rangle \in \Theta(1/N), \quad (\text{R2}\mu\text{-B3.4})$$

four-piece grid (same decomposition as (R2 $\mu$ -B2.4) at  $t = 2$ ) :

(init · init)  $\langle (W_0^{3,i}) h_2^{2,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N})$  random, unchanged from  $t = 1$ ;

(init · update)  $\langle (W_0^{3,i}) h_2^{2,i}, \Delta_2 W^4 \rangle \in \Theta(1/\sqrt{N})$ , climbs from  $\Theta(1/N)$  at  $t = 1$   
via  $\langle h_0^{3,i}, A_3 \rangle \in \Theta(\sqrt{N})$  ( $A_3 \parallel (W_0^4)^\top$  from R2 $\mu$ -F2.6d)

followed by the  $1/N$  prefactor of  $\Delta_2 W^4$ ;

(update · init)  $\langle (\Delta_2 W^{3,i}) h_2^{2,i}, W_0^4 \rangle \in \Theta(1)$  via rank-1 alignment, unchanged;

(update · update)  $\langle (\Delta_2 W^{3,i}) h_2^{2,i}, \Delta_2 W^4 \rangle \in \Theta(1)$ , climbs from  $\Theta(1/N)$  at  $t = 1$   
via  $\langle W_0^4, h_1^3 \rangle \in \Theta(1)$  from  $A_3$

( $\kappa \in \Theta(N)$  rank-1 prefactor times  $\langle W_0^4, \Delta_2 W^4 \rangle \in \Theta(1/N)$ );

$/M$ :  $(\Theta(1/N^{3/2}), \Theta(1/N^{3/2}), \Theta(1/N), \Theta(1/N))$  respectively;

both update · \* pieces dominate; total  $\Theta(1/N)$ ;

the entire  $W_0^{3,i}$  row stays at  $\Theta(1/N^{3/2})$ :  $\sqrt{N}$  deficit preserved.

**Climbing pieces in expert pathway  $(\partial f_2/\partial h_2^1)_{\text{exp}}$ .**

$$A_{4.2} \in \Theta(1/N^{3/2}), \quad \text{climbs from } \Theta(1/N^2) \text{ at } t = 1 \quad (\text{R2}\mu\text{-B3.5a.2})$$

( $\sqrt{N}$  lift via (R2 $\mu$ -B3.1); still sub-leading vs  $A_{4.1}$ )

$$A_{5.2} \in \Theta(1/N), \quad \text{climbs from } \Theta(1/N^2) \text{ at } t = 1 \quad (\text{R2}\mu\text{-B3.5b.2})$$

(combined  $\sqrt{N}$  lift via (R2 $\mu$ -B3.1) and alignment-driven mechanism;

$\langle W_1^4, h_1^{3,i} \rangle \in \Theta(1)$  coherent via alignment of  $h_1^{3,i}$  along  $(W_0^4)^\top$   
from (R2 $\mu$ -F2.5c); joins  $A_{5.1}$  at leading scale)

$$A_{6.1} \in \Theta(1/N^2) \text{ at every step; } A_{6.2} \in \Theta(1/N^3) \text{ at } t = 1, \Theta(1/N^2) \text{ at } t = 2 \quad (\text{R2}\mu\text{-B3.5c})$$

( $A_{6.1} : S_i := W_0^4 \cdot M_W \cdot (W_0^4)^\top$  has trace-induced coherent piece

$\sigma_4^2 \cdot \text{tr}(M_W) = \Theta(1/N)$  at every step, LLN preserves  $\Theta(1/N^2)$ ;

$A_{6.2}$  at  $t = 1 : S_i$  random ( $\Delta_1 W^4$  indep of  $W_0^4$ ), cross- $i$  CLT gives  $1/\sqrt{N}$

combined with  $\Delta_1 W^4$  smaller by  $1/\sqrt{N}$  via (R2 $\mu$ -U1.4)  $\Rightarrow \Theta(1/N^3)$ ;

$A_{6.2}$  at  $t = 2 : \Delta W^4$  acquires alignment with  $W_0^4$  via  $h_1^3$

(FL piece in (R2 $\mu$ -F2.5c)), restoring trace-induced coherent piece in  $S_i$ ;

combined with  $\Delta W^4$  climb to  $\Theta(1/N)$  via (R2 $\mu$ -B3.1) gives  $\Theta(1/N^2)$ ;

both stay sub-leading vs  $A_{5.x}, E_x \in \Theta(1/N)$ )

$$E_1 \in \Theta(1/N), \quad \text{climbs from } \Theta(1/N^2) \text{ at } t = 1 \quad (\text{R2}\mu\text{-B3.5d.1})$$

(alignment-driven climb;

inner product structure inside  $E$  no longer CLT-cancels but LLN-accumulates

once  $\langle W_1^4, h_1^{3,i} \rangle \in \Theta(1)$  is coherent via (R2 $\mu$ -B2.4);

joins leading scale)

$$E_2 \in \Theta(1/N), \quad \text{climbs from } \Theta(1/N^{5/2}) \text{ at } t = 1 \quad (\text{R2}\mu\text{-B3.5d.2})$$

(combined alignment-driven climb and  $\sqrt{N}$  lift via (R2 $\mu$ -B3.1);

joins leading scale)

$$(\partial f_2 / \partial h_2^1)_{\text{exp}} \in \Theta(1/N), \quad \text{from } A_{5.1}, A_{5.2}, E_1, E_2. \quad (\text{R2}\mu\text{-B3.5})$$

### Router pathway.

$$(\partial f_2 / \partial h_2^1)_{\text{router}} \in \Theta(1/N), \quad \text{all four pieces at leading scale.} \quad (\text{R2}\mu\text{-B3.6})$$

The  $\mathbf{v}^{(\Delta)}$ -pieces lift via (R2 $\mu$ -B3.1); the  $\Delta_t Q^\top \mathbf{v}^{(0)}$  piece remains coherent via the rank-1 alignment of  $\Delta_t Q$  with  $\mathbf{v}^{(0)}$  (extending (R2 $\mu$ -B2.6c) to cumulative  $\Delta_t Q$ ).

$W^{3,i}$ -imbalance preserved at  $t = 2$ . Define  $v_i := \dot{\phi}_{i,2} \langle h_2^{3,i}, W_2^4 \rangle$  and split  $v_i = v_i^I + v_i^U$  where

$$\begin{aligned} v_i^I &:= \dot{\phi}_{i,2} \langle (W_0^{3,i}) h_2^{2,i}, W_2^4 \rangle, \\ v_i^U &:= \dot{\phi}_{i,2} \langle (\Delta_2 W^{3,i}) h_2^{2,i}, W_2^4 \rangle. \end{aligned}$$

$v_i^I \in \Theta(1/\sqrt{N})$  random across  $i$  (both  $(W_0^{3,i}, W_0^4)$  and  $(W_0^{3,i}, \Delta_2 W^4)$  entries are  $\Theta(1/\sqrt{N})$ , the latter via  $\langle h_0^{3,i}, A_3 \rangle \in \Theta(\sqrt{N})$  followed by the  $1/N$  prefactor of  $\Delta_2 W^4$ ).  $v_i^U \in \Theta(1)$  coherent. Cross-layer CLT then gives  $(1/M) Q_2^\top v^I \in \Theta(1/N^{3/2})$  (sub-leading) and  $(1/M) Q_2^\top v^U \in \Theta(1/N)$  (leading). The entire  $W_0^{3,i}$ -row of the router pathway sits a factor  $\sqrt{N}$  below the  $\Delta_2 W^{3,i}$ -row.

$$\partial f_2 / \partial h_2^1 \in \Theta(1/N), \quad (\text{R2}\mu\text{-B3.7})$$

with leading-scale contributions from  $A_{5.1}, A_{5.2}, E_1, E_2$  in the expert pathway and all four pieces in the router pathway.

### J.3.2. SUMMARY TABLES OF SIGNAL PROPAGATION FOR $\mu\text{P}$ IN REGIME II

Notation:  $\Delta_t W^\ell$  denotes the cumulative update  $W_t^\ell - W_0^\ell$ . Green = at the proper feature-learning scale at  $t = 2$ ; red = sub-leading throughout.

#### Forward.

Quantity	$t = 0$	$t = 1$	$t = 2$
$h_t^1 = h_0^1 + \Delta_t h^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^1 = W_0^1 x$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t h^1 = \Delta_t W^1 x$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
$\psi_t, \phi_t$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$h_t^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^{2,i} = W_0^{2,i} h_0^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{2,i} \Delta_t h^1$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
effective: $\Delta_t W^{2,i} h_0^1$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
cross: $\Delta_t W^{2,i} \Delta_t h^1$	0	$\Theta(1/N^{3/2})$	$\Theta(1)$
$h_t^{3,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^{3,i} = W_0^{3,i} h_0^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
effective: $\Delta_t W^{3,i} h_0^{2,i}$	0	$\Theta(1)$	$\Theta(1)$
cross: $\Delta_t W^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
$h_t^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D$	$\Theta(1/\sqrt{N})$	$\Theta(1)$	$\Theta(1)$
$A_1 = (1/M) \sum_i \phi_{i,t} h_0^{3,i}$	$\Theta(1/\sqrt{N})$	$\Theta(1/\sqrt{N})$	$\Theta(1/\sqrt{N})$
$A_{2,1} = (1/M) \sum_i \phi_{i,t} W_0^{3,i} W_0^{2,i} \Delta_t h^1$	0	$\Theta(1/N)$	$\Theta(1/\sqrt{N})$
$A_{2,2} = (1/M) \sum_i \phi_{i,t} W_0^{3,i} \Delta_t W^{2,i} h_0^1$	0	$\Theta(1/N)$	$\Theta(1/\sqrt{N})$
$A_3 = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} h_0^{2,i}$	0	$\Theta(1)$	$\Theta(1)$
$D = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1/N)$	$\Theta(1)$
$f_t = W_t^4 h_t^3$	$\Theta(1/N)$	$\Theta(1)$	$\Theta(1)$

**Backward.**

Quantity	$t = 0$	$t = 1$	$t = 2$
$\partial f_t / \partial h_t^3 = (W_0^4)^\top + (\Delta_t W^4)^\top$	$\Theta(1/N)$	$\Theta(1/N)$	$\Theta(1/N)$
init: $(W_0^4)^\top$	$\Theta(1/N)$	$\Theta(1/N)$	$\Theta(1/N)$
update: $(\Delta_t W^4)^\top$	0	$\Theta(1/N^{3/2})$	$\Theta(1/N)$
$\partial f_t / \partial h_t^{3,i} = (\phi_{i,t}/M)(W_t^4)^\top$	$\Theta(1/N^2)$	$\Theta(1/N^2)$	$\Theta(1/N^2)$
init: $(\phi_{i,t}/M)(W_0^4)^\top$	$\Theta(1/N^2)$	$\Theta(1/N^2)$	$\Theta(1/N^2)$
update: $(\phi_{i,t}/M)(\Delta_t W^4)^\top$	0	$\Theta(1/N^{5/2})$	$\Theta(1/N^2)$
$\partial f_t / \partial h_t^{2,i} = (\phi_{i,t}/M)(W_t^{3,i})^\top (W_t^4)^\top$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$
init-init: $(\phi_{i,t}/M)(W_0^{3,i})^\top (W_0^4)^\top$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$
init-update: $(\phi_{i,t}/M)(W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N^{3/2})$
update-init: $(\phi_{i,t}/M)(\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
update-update: $(\phi_{i,t}/M)(\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$\partial f_t / \partial \phi_{i,t} = (1/M) \langle h_t^{3,i}, W_t^4 \rangle$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$
init-init: $(1/M) \langle (W_0^{3,i}) h_t^{2,i}, W_0^4 \rangle$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$
init-update: $(1/M) \langle (W_0^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle$	0	$\Theta(1/N^2)$	$\Theta(1/N^{3/2})$
update-init: $(1/M) \langle (\Delta_t W^{3,i}) h_t^{2,i}, W_0^4 \rangle$	0	$\Theta(1/N)$	$\Theta(1/N)$
update-update: $(1/M) \langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$(\partial f_t / \partial h_t^1)_{\text{exp}}$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$
$A_{4.1} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$
$A_{4.2} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N^{3/2})$
$A_{5.1} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_{5.2} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$A_{6.1} = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N^2)$
$A_{6.2} = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^3)$	$\Theta(1/N^2)$
$E_1 = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$E_2 = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^{5/2})$	$\Theta(1/N)$
$(\partial f_t / \partial h_t^1)_{\text{router}} = (1/M) Q_t^\top [\dot{\phi}_{i,t} \langle h_t^{3,i}, W_t^4 \rangle]_i$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$
$(1/M) Q_0^\top [\dot{\phi}_{i,t} \langle (W_0^{3,i}) h_t^{2,i}, W_0^4 \rangle]_i$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$
$(1/M) Q_0^\top [\dot{\phi}_{i,t} \langle (W_0^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	$\Theta(1/N^2)$	$\Theta(1/N^{3/2})$
$(1/M) Q_0^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, W_0^4 \rangle]_i$	0	$\Theta(1/N)$	$\Theta(1/N)$
$(1/M) Q_0^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$(1/M) \Delta_t Q^\top [\dot{\phi}_{i,t} \langle (W_0^{3,i}) h_t^{2,i}, W_0^4 \rangle]_i$	0	$\Theta(1/N^2)$	$\Theta(1/N^2)$
$(1/M) \Delta_t Q^\top [\dot{\phi}_{i,t} \langle (W_0^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	$\Theta(1/N^3)$	$\Theta(1/N^2)$
$(1/M) \Delta_t Q^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, W_0^4 \rangle]_i$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$(1/M) \Delta_t Q^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	$\Theta(1/N^3)$	$\Theta(1/N)$
$\partial f_t / \partial h_t^1 = (\partial f_t / \partial h_t^1)_{\text{exp}} + (\partial f_t / \partial h_t^1)_{\text{router}}$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$

### J.3.3. DERIVING MSSP IN REGIME II

This subsection presents the MSSP-Regime-II derivation per Definition 5, with the user-specified modification  $W_0^4 = 0$ . MSSP-Regime-II differs from  $\mu\text{P}$ -Regime-II (Definition 4) in exactly two places: the variance of  $W_0^{3,i}$  is  $M/N_e$  (entries  $\Theta(\sqrt{N})$ ) rather than  $1/N_e$  (entries  $\Theta(1)$ ); and  $W_0^4 = 0$  rather than variance  $1/N^2$ . All other init variances, learning rates, gating, and standing assumptions are identical to  $\mu\text{P}$ -Regime-II.

Two structural consequences propagate through the derivation:

1. The boosted  $W_0^{3,i}$  variance gives  $h_0^{3,i}$  entries of scale  $\Theta(\sqrt{N})$  (vs  $\Theta(1)$  in  $\mu\text{P}$ ), which after  $1/M$ -aggregation with the cross- $i$  CLT gives  $h_0^3$  entries  $\Theta(1)$ . So the aggregated activation is at the leading scale at init, even though the prediction  $f_0 = 0$  vanishes.
2. Since  $W_0^4 = 0$ , every gradient at  $t = 0$  that carries a  $W_0^4$  factor is zero, so only  $W^4$  updates at  $t = 1$ . Hidden weights first move at  $t = 2$ . Three forward + three backward passes are needed to capture all post-update correlations.

Compared to  $\mu\text{P}$ -Regime-II: MSSP-Regime-II's larger  $W^{3,i}$  init means the bottleneck-layer expert weights start big and the per-expert hidden state  $h^{3,i}$  inherits that scale; updates to  $W^{3,i}$  are then sub-leading relative to its init (vs comparable in  $\mu\text{P}$ -Regime-II). The aggregated  $h^3$  still feature-learns at  $\Theta(1)$ , just through a different mechanism: the prediction reaches  $f_1 \in \Theta(1)$  at  $t = 1$  purely from the  $W^4 \cdot h_0^3$  contraction (since  $h_0^3 \in \Theta(1)$  entries,  $\|h_0^3\|^2 \in \Theta(N)$ ), with no hidden updates required.

#### Setup and standing assumptions

**Architecture.** Same as in §J.3.1.

**Initialization.** Per Definition 5 with  $W_0^4 = 0$ :

$$\begin{aligned} (W_0^1)_{ab} &\sim \mathcal{N}(0, 1/D), & (Q_0)_{ia} &\sim \mathcal{N}(0, 1/N), \\ (W_0^{2,i})_{ab} &\sim \mathcal{N}(0, 1/N), & (W_0^{3,i})_{ab} &\sim \mathcal{N}(0, M/N_e), \\ W_0^4 &= 0. \end{aligned}$$

All parameters drawn independently across  $i$ . Resulting entry scales:  $W_0^4 = 0$ ,  $W_0^{2,i} \in \Theta(1/\sqrt{N})$ ,  $W_0^{3,i} \in \Theta(\sqrt{M}) = \Theta(\sqrt{N})$ ,  $W_0^1 \in \Theta(1)$ ,  $Q_0 \in \Theta(1/\sqrt{N})$ .

**Learning rates (SGD).** Same as  $\mu\text{P}$ -Regime-II:  $\eta_1 = \eta N$ ,  $\eta_Q = \eta M/N$ ,  $\eta_2 = \eta M/N$ ,  $\eta_3 = \eta MN$ ,  $\eta_4 = \eta/N$ .

**Gating.** Sigmoid +  $1/M$  aggregation:  $\phi_{i,t}, \dot{\phi}_{i,t} \in \Theta(1)$ .

**Auxiliary scaling lemmas** We reuse CLT, LLN, and Lemma 10.

**Specialization of Lemma 10 to MSSP-Regime-II.** For  $W = W_0^{3,i} \in \mathbb{R}^{N \times N_e}$  with  $\sigma_W^2 = M/N_e$  and  $N_e = \Theta(1)$ :

- $W_0^{3,i} (W_0^{3,i})^\top \in \mathbb{R}^{N \times N}$  has rank  $\leq N_e = \Theta(1)$ . Diagonal mean  $n\sigma^2 = M = \Theta(N)$ , off-diagonal entries variance  $n\sigma^4 = M^2/N_e = \Theta(N^2)$ , coordinate scale  $\Theta(N)$ . **Every entry is  $\Theta(N)$**  — much larger than  $\mu\text{P}$ -Regime-II's  $\Theta(1)$  entries.

- $(W_0^{3,i})^\top W_0^{3,i} \in \mathbb{R}^{N_e \times N_e}$ , rank- $N_e$  full. Diagonal mean  $m\sigma^2 = NM/N_e = \Theta(N^2)$ , off-diagonal coord scale  $\Theta(\sigma^2 \sqrt{m}) = \Theta(N^{3/2})$ . So acting on a  $\Theta(1)$ -coordinate  $N_e$ -vector  $h_0^{2,i}$  gives entries of coordinate scale  $\Theta(N) \cdot \Theta(1) + \Theta(N^{3/2}) \cdot \Theta(1) = \Theta(N^2)$  at leading.

The other Gram products ( $W_0^{2,i}$  versions) are unchanged from  $\mu\text{P}$ -Regime-II.

### First forward pass

$$h_0^1 = W_0^1 x \in \Theta(1). \quad (\text{CLT}; \sigma_1^2 = 1/D, \|x\|^2 \in \Theta(D)) \quad (\text{R2s-F1.1})$$

$$\psi_0 = Q_0 h_0^1 \in \Theta(1). \quad (\text{CLT}; \sigma_Q^2 = 1/N, \|h_0^1\|^2 \in \Theta(N)) \quad (\text{R2s-F1.2})$$

$$\phi_{i,0} = \sigma(\psi_{i,0}) \in \Theta(1). \quad (\text{sigmoid bounded; gating assumption}) \quad (\text{R2s-F1.3})$$

$$\begin{aligned} h_0^{2,i} &= W_0^{2,i} h_0^1 \in \Theta(1), \quad \text{independent across } i \\ &(\text{CLT}; \sigma_2^2 = 1/N, \|h_0^1\|^2 \in \Theta(N); \\ &\|h_0^{2,i}\|^2 = N_e \cdot \Theta(1) = \Theta(1) \text{ since } N_e = \Theta(1)). \end{aligned} \quad (\text{R2s-F1.4})$$

$$\begin{aligned} h_0^{3,i} &= W_0^{3,i} h_0^{2,i} \in \Theta(\sqrt{N}), \quad \|h_0^{3,i}\|^2 \in \Theta(N^2) \quad (\text{boosted by MSSP variance}) \quad (\text{R2s-F1.5}) \\ &(\text{CLT}; \sigma_3^2 = M/N_e = \Theta(N); \|h_0^{2,i}\|^2 \in \Theta(1) \text{ via (R2s-F1.4);} \\ &\text{variance per entry } \Theta(N) \cdot \Theta(1) = \Theta(N)) \end{aligned}$$

$$\begin{aligned} h_0^3 &= (1/M) \sum_i \phi_{i,0} h_0^{3,i} \in \Theta(1) \quad (\text{R2s-F1.6}) \\ &(\text{cross-}i \text{ CLT on } h_0^{3,i} \text{ entries } \Theta(\sqrt{N}) \text{ via (R2s-F1.5);} \\ &\{h_0^{3,i}\} \text{ independent across } i; \\ &\text{variance per entry } (1/M^2) \sum_i \phi_{i,0}^2 \Theta(N) = \Theta(N/M) = \Theta(1); \\ &\sqrt{M} \text{ cross-}i \text{ CLT cancellation balances the } \sqrt{M} \text{ amplification of } h_0^{3,i}) \end{aligned}$$

$$f_0 = W_0^4 h_0^3 = 0. \quad (W_0^4 = 0 \text{ by MSSP-Regime-II standing assumption}) \quad (\text{R2s-F1.7})$$

**First backward pass and step-1 updates** By the chain rule, every hidden gradient at  $t = 0$  ( $\partial f_0 / \partial h_0^3$ ,  $\partial f_0 / \partial h_0^{3,i}$ ,  $\partial f_0 / \partial h_0^{2,i}$ ,  $\partial f_0 / \partial \phi_{i,0}$ , and both expert and router contributions to  $\partial f_0 / \partial h_0^1$ ) carries a  $W_0^4 = 0$  factor and is therefore zero. Hence  $W_1^\ell = W_0^\ell$  for  $\ell \in \{1, Q, (2, i), (3, i)\}$ ; only  $W^4$  is updated at  $t = 1$ .

### Step-1 parameter update.

$$\Delta_1 W^4 = -(\eta \chi_0 / N) (h_0^3)^\top \in \Theta(1/N). \quad (h_0^3 \in \Theta(1) \text{ via (R2s-F1.6)}) \quad (\text{R2s-U1.4})$$

**Second forward pass** Since the non-readout weights are unchanged at  $t = 1$ , every hidden activation is unchanged:  $h_1^\ell = h_0^\ell$  for  $\ell \in \{1, (2, i), (3, i), 3\}$ ,  $\psi_1 = \psi_0$ ,  $\phi_1 = \phi_0$ .

$$f_1 = W_1^4 h_0^3 = -\frac{\eta \chi_0}{N} \|h_0^3\|^2 \in \Theta(1). \quad (\|h_0^3\|^2 \in \Theta(N) \text{ via (R2s-F1.6)}) \quad (\text{R2s-F2.7})$$

### Second backward pass and step-2 updates

**New intermediate at  $t = 1$ .**

$$(W_0^{3,i})^\top h_0^3 \in \Theta(N) \text{ entry-wise}$$

$$\text{(using (R2s-F1.6) to expand } h_0^3 = (1/M) \sum_j \phi_{j,0} h_0^{3,j};$$

$$j = i \text{ piece } (\phi_{i,0}/M)(W_0^{3,i})^\top W_0^{3,i} h_0^{2,i} \text{ with}$$

$$(W_0^{3,i})^\top W_0^{3,i} h_0^{2,i} \in \Theta(N^2) \text{ from Lemma 10 specialization,}$$

$$\text{giving } (\phi_{i,0}/M) \cdot \Theta(N^2) = \Theta(N);$$

$$j \neq i \text{ pieces give cross-}j \text{ noise of comparable } \Theta(N) \text{ scale)}$$

$$\Rightarrow (W_0^{3,i})^\top (W_1^4)^\top = -(\eta\chi_0/N)(W_0^{3,i})^\top h_0^3 \in \Theta(1) \text{ entry-wise}$$

$$(W_1^4 \text{ via (R2s-U1.4)})$$

**Per-layer gradients.**

$$\partial f_1 / \partial h_1^3 = (W_1^4)^\top \in \Theta(1/N). \quad \text{(via (R2s-U1.4))} \quad \text{(R2s-B2.1)}$$

$$\partial f_1 / \partial h_1^{3,i} = (\phi_{i,0}/M)(W_1^4)^\top \in \Theta(1/N^2). \quad (\phi_{i,0} \in \Theta(1); M = \Theta(N)) \quad \text{(R2s-B2.2)}$$

$$\partial f_1 / \partial h_1^{2,i} = (\phi_{i,0}/M)(W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1/N) \quad \text{(R2s-B2.3)}$$

$$((W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ via the new intermediate above})$$

$$\partial f_1 / \partial \phi_{i,0} = (1/M) \langle h_0^{3,i}, W_1^4 \rangle = -(\eta\chi_0/(MN)) \langle h_0^{3,i}, h_0^3 \rangle \in \Theta(1/N) \quad \text{(R2s-B2.4)}$$

(Four-piece grid: only (init · update) is non-zero at  $t = 1$ ,

since  $W_0^4 = 0$  kills the \* · init column

and  $\Delta_1 W^{3,i} = 0$  (hidden weights frozen) kills the update · \* row;

(init · update)  $\langle (W_0^{3,1}) h_0^{2,i}, \Delta_1 W^4 \rangle \in \Theta(1)$  :

expanding  $h_0^3 = (1/M) \sum_j \phi_{j,0} h_0^{3,j}$  gives

$$\langle h_0^{3,i}, h_0^3 \rangle = (\phi_{i,0}/M) \|h_0^{3,i}\|^2 + (1/M) \sum_{j \neq i} \phi_{j,0} \langle h_0^{3,i}, h_0^{3,j} \rangle;$$

$$j = i \text{ diagonal term: } (\phi_{i,0}/M) \|h_0^{3,i}\|^2 = (1/M) \Theta(N^2) = \Theta(N)$$

using  $\|h_0^{3,i}\|^2 \in \Theta(N^2)$  via (R2s-F1.5);

$j \neq i$  off-diagonal terms give comparable  $\Theta(N)$  via cross- $j$  CLT;

so  $\langle h_0^{3,i}, h_0^3 \rangle \in \Theta(N)$ ;

substituting:  $\partial f_1 / \partial \phi_{i,0} = -(\eta\chi_0/(MN)) \cdot \Theta(N) = \Theta(1/M) = \Theta(1/N)$ ;

lifted from  $\mu\text{P-R2}$ 's (init · update) =  $\Theta(1/N^2)$  entry

by the boost  $\sigma_3^2 = M/N_e$ )

$$(\partial f_1 / \partial h_1^1)_{\text{exp}} = (1/M) \sum_i \phi_{i,0} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1/N) \quad (\text{R2s-B2.5})$$

$((W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1)$  via the new intermediate above;

$(W_0^{2,i})^\top$  acts on a  $\Theta(1)$  vector with  $\|v\|^2 = N_e \Theta(1) = \Theta(1)$ ;

variance per entry  $(1/N)\Theta(1) = \Theta(1/N)$ , coord  $\Theta(1/\sqrt{N})$  random per  $i$ ;

cross- $i$  CLT:  $\Theta(1/N)$ )

$$(\partial f_1 / \partial h_1^1)_{\text{router}} = (1/M) Q_0^\top v, \quad v_i := \dot{\phi}_{i,0} \langle h_0^{3,i}, W_1^4 \rangle \in \Theta(1/N) \quad (\text{R2s-B2.6})$$

$(v_i \in \Theta(1)$  via (R2s-B2.4);  $\|v\|^2 = M\Theta(1) = \Theta(M)$ ;

$Q_0^\top v$  entries variance  $(1/N)\Theta(M) = \Theta(1)$  via cross-layer CLT;

$(1/M)\Theta(1) = \Theta(1/N)$ ;

$W^{3,i}$ -split  $v_i = v_i^I + v_i^U$  where

$$v_i^I := \dot{\phi}_{i,0} \langle (W_0^{3,1}) h_0^{2,i}, W_1^4 \rangle$$

$$v_i^U := \dot{\phi}_{i,0} \langle (\Delta_1 W^{3,i}) h_0^{2,i}, W_1^4 \rangle = 0 \text{ since } \Delta_1 W^{3,i} = 0;$$

$v_i^I \in \Theta(1)$  coherent via boost (cf. R2s-B2.4);

$$(1/M) Q_0^\top v^I = \Theta(1/N) \text{ and } (1/M) Q_0^\top v^U = 0$$

$$\partial f_1 / \partial h_1^1 = (\partial f_1 / \partial h_1^1)_{\text{exp}} + (\partial f_1 / \partial h_1^1)_{\text{router}} \in \Theta(1/N). \quad (\text{R2s-B2.7})$$

## Step-2 parameter updates.

$$\Delta_2 W^4 = -(\eta \chi_1 / N) (h_0^3)^\top \in \Theta(1/N). \quad (h_0^3 \in \Theta(1) \text{ via (R2s-F1.6)}) \quad (\text{R2s-U2.4})$$

$$\Delta_2 W^{3,i} = -\eta N \chi_1 \phi_{i,0} (W_1^4)^\top (h_0^{2,i})^\top \in \Theta(1) \text{ rank-1} \quad (\text{R2s-U2.3})$$

$((W_1^4)^\top \in \Theta(1/N)$  via (R2s-U1.4);  $h_0^{2,i} \in \Theta(1)$  via (R2s-F1.4);

sub-leading vs  $W_0^{3,i} \in \Theta(\sqrt{N})$  via (R2s-F1.5))

$$\Delta_2 W^{2,i} = -(\eta \chi_1 \phi_{i,0} / N) (W_0^{3,i})^\top (W_1^4)^\top (h_0^1)^\top \in \Theta(1/N) \text{ rank-1} \quad (\text{R2s-U2.2})$$

$((W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1)$  via the new intermediate;

sub-leading vs  $W_0^{2,i} \in \Theta(1/\sqrt{N})$ )

$$\Delta_2 W^1 = -\eta_1 \chi_1 (\partial f_1 / \partial h_1^1) x^\top \in \Theta(1) \quad (\text{R2s-U2.1a})$$

$(\eta_1 = \eta N$ ;  $\partial f_1 / \partial h_1^1 \in \Theta(1/N)$  via (R2s-B2.7))

$$\Delta_2 h^1 = \Delta_2 W^1 x \in \Theta(1), \text{ aligned along } h_0^1. \quad (\text{R2s-U2.1b})$$

$$\Delta_2 Q \in \Theta(1/N). \quad (\partial f_1 / \partial \phi \in \Theta(1/N) \text{ via (R2s-B2.4)}) \quad (\text{R2s-U2.Q})$$

**Third forward pass**

$$h_2^1 = h_0^1 + \Delta_2 h^1 \in \Theta(1). \quad (\Delta_2 h^1 \in \Theta(1) \text{ via (R2s-U2.1b)}; h_0^1 \in \Theta(1)) \quad (\text{R2s-F3.1})$$

(R2s-F3.2)  $h_2^{2,i} = W_2^{2,i} h_2^1$ , four-piece decomposition:

$$h_2^{2,i} = h_0^{2,i} + W_0^{2,i} \Delta_2 h^1 + \Delta_2 W^{2,i} h_0^1 + \Delta_2 W^{2,i} \Delta_2 h^1. \quad (\text{R2s-F3.2})$$

$$\text{init: } h_0^{2,i} \in \Theta(1). \quad (\text{R2s-F1.4}) \quad (\text{R2s-F3.2a})$$

$$\text{prop: } W_0^{2,i} \Delta_2 h^1 \in \Theta(1)$$

$$(\sigma_2^2 \|\Delta_2 h^1\|^2 = (1/N)\Theta(N) = \Theta(1);$$

$$\|\Delta_2 h^1\|^2 \in \Theta(N) \text{ via (R2s-U2.1b)}). \quad (\text{R2s-F3.2b})$$

$$\text{eff: } \Delta_2 W^{2,i} h_0^1 = -(\eta\chi_1 \phi_{i,1}/N) \|h_0^1\|^2 (W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1)$$

((R2s-U2.2)substitution;

$$(W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ via the new intermediate at } t = 1;$$

$$\|h_0^1\|^2 \in \Theta(N)). \quad (\text{R2s-F3.2c})$$

$$\text{cross: } \Delta_2 W^{2,i} \Delta_2 h^1 = -(\eta\chi_1 \phi_{i,1}/N) \langle h_0^1, \Delta_2 h^1 \rangle (W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1)$$

$$((W_0^{3,i})^\top (W_1^4)^\top \in \Theta(1) \text{ via the new intermediate};$$

$$\langle h_0^1, \Delta_2 h^1 \rangle \in \Theta(N) \text{ coherent via (R2s-U2.1b)}). \quad (\text{R2s-F3.2d})$$

$\Delta_2 h^{2,i} \in \Theta(1)$  entry-wise.

(R2s-F3.3)  $h_2^{3,i} = W_2^{3,i} h_2^{2,i}$ , four-piece decomposition:

$$h_2^{3,i} = h_0^{3,i} + W_0^{3,i} \Delta_2 h^{2,i} + \Delta_2 W^{3,i} h_0^{2,i} + \Delta_2 W^{3,i} \Delta_2 h^{2,i}. \quad (\text{R2s-F3.3})$$

$$\text{init: } h_0^{3,i} \in \Theta(\sqrt{N}). \quad (\text{R2s-F1.5}) \quad (\text{R2s-F3.3a})$$

$$\text{prop: } W_0^{3,i} \Delta_2 h^{2,i} \in \Theta(\sqrt{N})$$

$$(\sigma_3^2 \|\Delta_2 h^{2,i}\|^2 = \Theta(N) \cdot \Theta(1) = \Theta(N);$$

$$\|\Delta_2 h^{2,i}\|^2 = N_e \Theta(1) = \Theta(1) \text{ via (R2s-F3.2)}; N_e = \Theta(1)). \quad (\text{R2s-F3.3b})$$

$$\text{eff: } \Delta_2 W^{3,i} h_0^{2,i} = -\eta N \chi_1 \phi_{i,1} \|h_0^{2,i}\|^2 (W_1^4)^\top \in \Theta(1)$$

((R2s-U2.3)substitution;

$$\|h_0^{2,i}\|^2 \in \Theta(1) \text{ via (R2s-F1.4)};$$

$$(W_1^4)^\top \in \Theta(1/N) \text{ via (R2s-U1.4)}). \quad (\text{R2s-F3.3c})$$

$$\text{cross: } \Delta_2 W^{3,i} \Delta_2 h^{2,i} \text{ tracked as part of } D \text{ in (R2s-F3.4)}. \quad (\text{R2s-F3.3d})$$

**Propagating dominates:**  $\Delta_2 h^{3,i} \in \Theta(\sqrt{N})$  — same scale as  $h_0^{3,i}$  at init. The effective piece  $\Theta(1)$  is sub-leading per-expert but contributes coherently to the aggregation in (R2s-F3.4d).

(R2s-F3.4)  $h_2^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D$ , where

$$\begin{aligned} A_1 &:= \frac{1}{M} \sum_i \phi_{i,2} h_0^{3,i}, & A_{2,1} &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,i} W_0^{2,i} \Delta_2 h^1, \\ A_{2,2} &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,i} \Delta_2 W^{2,i} h_0^1, & A_3 &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} h_0^{2,i}, \\ D &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} \Delta_2 h^{2,i}. \end{aligned}$$

$$A_1 \in \Theta(1). \quad (\text{cross-}i \text{ CLT, same calculation as (R2s-F1.6)}) \quad (\text{R2s-F3.4a})$$

$$A_{2,1} \in \Theta(1)$$

(per summand  $W_0^{3,i} W_0^{2,i} \Delta_2 h^1 \in \Theta(\sqrt{N})$  random;  
 $\|\Delta_2 h^1\|^2 \in \Theta(N)$  via (R2s-U2.1b);  
 cross- $i$  CLT:  $\Theta(\sqrt{N}/\sqrt{M}) = \Theta(1)$ ).

$$(\text{R2s-F3.4b})$$

$$A_{2,2} \in \Theta(1)$$

(per summand  $W_0^{3,i} [\Delta_2 W^{2,i} h_0^1] \in \Theta(\sqrt{N})$   
 since  $\Delta_2 W^{2,i} h_0^1 \in \Theta(1)$  via (R2s-F3.2c);  
 cross- $i$  CLT:  $\Theta(\sqrt{N}/\sqrt{M}) = \Theta(1)$ ).

$$(\text{R2s-F3.4c})$$

$$A_3 = -\eta N \chi_1 (W_1^4)^\top \left( \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1} \|h_0^{2,i}\|^2 \right) \in \Theta(1)$$

((R2s-U2.3)substitution;  
 $\|h_0^{2,i}\|^2 \in \Theta(1)$  via (R2s-F1.4); LLN;  
 $(W_1^4)^\top \in \Theta(1/N)$ ).

$$(\text{R2s-F3.4d})$$

$$D = -\eta N \chi_1 (W_1^4)^\top \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1} \langle h_0^{2,i}, \Delta_2 h^{2,i} \rangle \in \Theta(1)$$

((R2s-U2.3)substitution;  
 $\langle h_0^{2,i}, \Delta_2 h^{2,i} \rangle \in \Theta(1)$  coherent across  $i$ ,  
 from two pieces:  $\Delta_2 h^1 \parallel h_0^1$  via (R2s-U2.1b),  
 and the  $j = i$  Gram self-overlap in the effective piece (R2s-F3.2c);  
 LLN gives empirical average  $\Theta(1)$ ).

$$(\text{R2s-F3.4e})$$

$$h_2^3 \in \Theta(1) \text{ entry-wise.} \quad (\text{R2s-F3.4})$$

$$f_2 = W_2^4 h_2^3 = -(\eta(\chi_0 + \chi_1)/N) \langle h_0^3, h_2^3 \rangle \in \Theta(1)$$

( $W_2^4 = W_1^4 + \Delta_2 W^4$  via (R2s-U2.4);  
 $\langle h_0^3, h_2^3 \rangle \in \Theta(N)$  coherent, dominated by  $A_3 + D$  via (R2s-F3.4);  
 net  $\eta \cdot (1/N) \cdot \Theta(N) = \Theta(1)$ ).

$$(\text{R2s-F3.5})$$

**Third backward pass** The structure parallels (R2s-B2.1)–(R2s-B2.7) with cumulative weights at step 2. Since  $W_0^4 = 0$  in MSSP, the four “.1” pieces  $A_{4.1}, A_{5.1}, A_{6.1}, E_1$  — using  $W_0^4$  rather than  $\Delta W^4$  — are identically zero at every step; only the four “.2” pieces contribute. With  $W_2^4 = \Delta W^4$  now non-zero (via (R2s-U2.4)), all gradients are non-zero.

**Per-layer gradients.**

$$\partial f_2 / \partial h_2^3 = (W_2^4)^\top \in \Theta(1/N). \quad (\text{R2s-U2.4}) \quad (\text{R2s-B3.1})$$

$$\partial f_2 / \partial h_2^{3,i} = (\phi_{i,2}/M)(W_2^4)^\top \in \Theta(1/N^2). \quad (\phi_{i,2} \in \Theta(1); M = \Theta(N)) \quad (\text{R2s-B3.2})$$

$$\begin{aligned} \partial f_2 / \partial h_2^{2,i} &= (\phi_{i,2}/M)(W_2^{3,i})^\top (W_2^4)^\top \in \Theta(1/N) \\ &((W_0^{3,i})^\top (W_2^4)^\top \in \Theta(1) \text{ via the cumulative analog of the new intermediate at } t = 1; \\ &(\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1) \text{ via (R2s-U2.3)}). \end{aligned} \quad (\text{R2s-B3.3})$$

**Router pathway**  $\partial f_2 / \partial \phi_{i,2}$ .

$$\partial f_2 / \partial \phi_{i,2} = (1/M) \langle h_2^{3,i}, W_2^4 \rangle \in \Theta(1/N), \quad (\text{R2s-B3.4})$$

Four-piece grid: \* ·init column is identically zero ( $W_0^4 = 0$ );

(init · update)  $\langle (W_0^{3,1}) h_2^{2,i}, \Delta_2 W^4 \rangle \in \Theta(1)$ :

$$\langle h_0^{3,i}, W_2^4 \rangle = -\frac{\chi_0 + \chi_1}{N} \langle h_0^{3,i}, h_0^3 \rangle,$$

expanding  $h_0^3 = (1/M) \sum_j \phi_{j,0} h_0^{3,j}$  gives the diagonal  $j = i$  term

$$(\phi_{i,0}/M) \|h_0^{3,i}\|^2 = (1/M) \Theta(N^2) = \Theta(N) \text{ via (R2s-F1.5);}$$

$j \neq i$  off-diagonal terms give comparable  $\Theta(N)$ ;

$$\text{so } \langle h_0^{3,i}, W_2^4 \rangle \in \Theta(1);$$

(update · update)  $\langle (\Delta_2 W^{3,i}) h_2^{2,i}, \Delta_2 W^4 \rangle \in \Theta(1)$ :

$$\langle (\Delta_2 W^{3,i}) h_2^{2,i}, W_2^4 \rangle = \kappa \langle (W_1^4)^\top, W_2^4 \rangle$$

with  $\kappa = -\eta N \chi_1 \phi_{i,1} \|h_2^{2,i}\|^2 \in \Theta(N)$

$$\text{and } \langle W_1^4, W_2^4 \rangle = (\chi_0(\chi_0 + \chi_1)/N^2) \|h_0^3\|^2 \in \Theta(1/N);$$

$$\text{net } \Theta(N) \cdot \Theta(1/N) = \Theta(1);$$

substituting:  $\partial f_2 / \partial \phi_{i,2} = (1/M)(\Theta(1) + \Theta(1)) = \Theta(1/M) = \Theta(1/N)$ ;

lifted from  $\mu\text{P-R2}$ 's grid  $(\Theta(1/N^{3/2}), \Theta(1/N^{3/2}), \Theta(1/N), \Theta(1/N))$

to MSSP-R2's  $(0, \Theta(1/N), 0, \Theta(1/N))$  via  $W_0^4 = 0$  and  $\sigma_3^2 = M/N_e$ .

**Expert pathway:**  $(\partial f_2 / \partial h_2^1)_{\text{exp}} = A_{4.2} + A_{5.2} + A_{6.2} + E_2$ .

$$A_{4.2} = \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_2^4)^\top \in \Theta(1/N)$$

$((W_0^{3,i})^\top (W_2^4)^\top \in \Theta(1) \text{ via the cumulative analog of the new intermediate;}$

$(W_0^{2,i})^\top$  acts on  $\Theta(1)$  vector with  $\|v\|^2 = \Theta(1)$ ;

coord  $\Theta(1/\sqrt{N})$  random per  $i$ , cross- $i$  CLT:  $\Theta(1/N)$ ). (R2s-B3.5a)

$$\begin{aligned}
 A_{5.2} &= \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,i})^\top (\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1/N) \\
 &\quad \text{((R2s-U2.3)substitution gives } (\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1) \text{ along } h_0^{2,i}; \\
 &\quad (W_0^{2,i})^\top h_0^{2,i} \in \Theta(1/\sqrt{N}) \text{ random per } i; \\
 &\quad \text{cross-}i \text{ CLT: } \Theta(1/N)). \tag{R2s-B3.5b}
 \end{aligned}$$

$$\begin{aligned}
 A_{6.2} &= \frac{1}{M} \sum_i \phi_{i,2} (\Delta_2 W^{2,i})^\top (W_0^{3,i})^\top (W_2^4)^\top \in \Theta(1/N) \\
 &\quad \text{((R2s-U2.2)substitution;} \\
 &\quad \text{per summand } - (\eta\chi_1 \phi_{i,1}/N) h_0^1 \cdot S_i \text{ where} \\
 &\quad S_i := W_1^4 W_0^{3,i} (W_0^{3,i})^\top (W_2^4)^\top \in \Theta(1) \\
 &\quad \text{(Lemma 10 specialization: } W_0^{3,i} (W_0^{3,i})^\top \text{ entries } \Theta(N), \\
 &\quad \text{acting on } (W_2^4)^\top \in \Theta(1/N) \text{ gives } \Theta(\sqrt{N}); \\
 &\quad W_1^4 \text{ contraction yields } \Theta(1)); \\
 &\quad \text{per summand entries } \Theta(1/N), \text{ cross-}i \text{ avg: } \Theta(1/N)). \tag{R2s-B3.5c}
 \end{aligned}$$

$$\begin{aligned}
 E_2 &= \frac{1}{M} \sum_i \phi_{i,2} (\Delta_2 W^{2,i})^\top (\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1/N) \\
 &\quad \text{(combining (R2s-U2.2), (R2s-U2.3) substitutions;} \\
 &\quad (\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1) \text{ along } h_0^{2,i}; \\
 &\quad \langle W_1^4, h_0^{3,i} \rangle = -(\eta\chi_0/N) \langle h_0^3, h_0^{3,i} \rangle \in \Theta(1) \\
 &\quad \text{(coherent } j = i \text{ piece via } \|h_0^{3,i}\|^2 \in \Theta(N^2) \text{ from (R2s-F1.5);} \\
 &\quad \text{plus } j \neq i \text{ random of comparable scale);} \\
 &\quad \text{per summand entries } \Theta(1/N); \\
 &\quad \text{cross-}i \text{ LLN on coherent piece: } \Theta(1/N)). \tag{R2s-B3.5d}
 \end{aligned}$$

$$(\partial f_2 / \partial h_2^1)_{\text{exp}} \in \Theta(1/N), \text{ all four non-zero pieces at leading scale.} \tag{R2s-B3.5}$$

### Router pathway.

$$\begin{aligned}
 (\partial f_2 / \partial h_2^1)_{\text{router}} &= (1/M) Q_2^\top v, \quad v_i := \dot{\phi}_{i,2} \langle h_2^{3,i}, W_2^4 \rangle \in \Theta(1/N) \\
 &\quad (v_i \in \Theta(1) \text{ via (R2s-B3.4);} \\
 &\quad \text{by the analysis of (R2s-B2.6) extended to cumulative } Q_2 = Q_0 + \Delta_2 Q; \\
 &\quad Q_0 \text{ piece via cross-layer CLT;} \\
 &\quad \Delta_2 Q \text{ piece coherent via (R2s-U2.Q);} \\
 &\quad \text{both contribute at } \Theta(1/N)). \tag{R2s-B3.6}
 \end{aligned}$$

$W^{3,i}$ -balance transmits through the router pathway. The  $W^{3,i}$ -split of (R2s-B3.4),  $\mathbf{v} = \mathbf{v}^I + \mathbf{v}^U$  with  $\mathbf{v}_i^I = \dot{\phi}_{i,2} \langle (W_0^{3,1}) h_2^{2,i}, \Delta_2 W^4 \rangle$  and  $\mathbf{v}_i^U = \dot{\phi}_{i,2} \langle (\Delta_2 W^{3,i}) h_2^{2,i}, \Delta_2 W^4 \rangle$ , transmits as: both  $\mathbf{v}_i^I \in \Theta(1)$  and  $\mathbf{v}_i^U \in \Theta(1)$  coherent across  $i$  — the boost  $\sigma_3^2 = M/N_e$  promotes the  $W_0^{3,1}$ -piece's diagonal  $j = i$  contribution to  $\Theta(1)$  (cf. R2s-B3.4), matching the rank-1-aligned  $\Delta W^{3,i}$ -piece. So  $(1/M) Q_t^\top \mathbf{v}^I \in \Theta(1/N)$  and  $(1/M) Q_t^\top \mathbf{v}^U \in \Theta(1/N)$ : balanced, in contrast to the  $\sqrt{N}$  deficit under  $\mu\text{P-R2}$ .

$$\partial f_2 / \partial h_2^1 = (\partial f_2 / \partial h_2^1)_{\text{exp}} + (\partial f_2 / \partial h_2^1)_{\text{router}} \in \Theta(1/N). \tag{R2s-B3.7}$$

## J.3.4. SUMMARY TABLES OF SIGNAL PROPAGATION FOR MSSP IN REGIME II

**Forward features.**

Quantity	$t = 0$	$t = 1$	$t = 2$
$h_t^1$ (init)	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$\Delta_t h^1$	—	0 (no update)	$\Theta(1)$
effective $\Delta W^1 x$	—	0	$\Theta(1)$
$\psi_t$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$\Delta_t \psi$	—	0	$\Theta(1)$
effective $\Delta Q h_0^1$	—	0	$\Theta(1)$
propagating $Q_0 \Delta h^1$	—	0	$\Theta(1)$
$\phi_t$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$\Delta_t \phi$	—	0	$\Theta(1)$
$h_t^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$\Delta_t h^{2,i}$	—	0	$\Theta(1)$
effective $\Delta W^{2,i} h_0^1$	—	0	$\Theta(1)$
propagating $W_0^{2,i} \Delta h^1$	—	0	$\Theta(1)$
$h_t^{3,i}$	$\Theta(\sqrt{N})$	$\Theta(\sqrt{N})$	$\Theta(\sqrt{N})$
$\Delta_t h^{3,i}$	—	0	$\Theta(\sqrt{N})$
effective $\Delta W^{3,i} h_0^{2,i}$	—	0	$\Theta(1)$
propagating $W_0^{3,i} \Delta h^{2,i}$	—	0	$\Theta(\sqrt{N})$
$h_t^3$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$f_t$	0	$\Theta(1)$	$\Theta(1)$

**Forward aggregation:**  $h_t^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D$ .

Piece	$t = 0$	$t = 1$	$t = 2$
$A_1 = (1/M) \sum \phi h_0^{3,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$A_{2,1} = (1/M) \sum \phi W_0^{3,i} W_0^{2,i} \Delta h^1$	0	0	$\Theta(1)$
$A_{2,2} = (1/M) \sum \phi W_0^{3,i} \Delta W^{2,i} h_0^1$	0	0	$\Theta(1)$
$A_3 = (1/M) \sum \phi \Delta W^{3,i} h_0^{2,i}$	0	0	$\Theta(1)$
$D = (1/M) \sum \phi \Delta W^{3,i} \Delta h^{2,i}$	0	0	$\Theta(1)$
<b>Total <math>h_t^3</math></b>	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$

**Backward gradients (per-layer with  $W_0^4/\Delta W^4$  split).**

Piece	$t = 0$	$t = 1$	$t = 2$
$(W_0^4)^\top$	0	0	0
$(\Delta W^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$\partial f_t / \partial h_t^3$ <b>total</b>	0	$\Theta(1/N)$	$\Theta(1/N)$
$\partial f_t / \partial h_t^{3,i}$ <b>total</b>	0	$\Theta(1/N^2)$	$\Theta(1/N^2)$
$(\phi/M)(W_0^{3,i})^\top (W_0^4)^\top$	0	0	0
$(\phi/M)(W_0^{3,i})^\top (\Delta W^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$(\phi/M)(\Delta W^{3,i})^\top (W_0^4)^\top$	0	0	0
$(\phi/M)(\Delta W^{3,i})^\top (\Delta W^4)^\top$	0	0	$\Theta(1/N)$
$\partial f_t / \partial h_t^{2,i}$ <b>total</b>	0	$\Theta(1/N)$	$\Theta(1/N)$
init-init: $(1/M)\langle (W_0^{3,i}) h_t^{2,i}, W_0^4 \rangle$	0	0	0
init-update: $(1/M)\langle (W_0^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle$	0	$\Theta(1/N)$	$\Theta(1/N)$
update-init: $(1/M)\langle (\Delta_t W^{3,i}) h_t^{2,i}, W_0^4 \rangle$	0	0	0
update-update: $(1/M)\langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle$	0	0	$\Theta(1/N)$
$\partial f_t / \partial \phi_{i,t}$ <b>total</b>	0	$\Theta(1/N)$	$\Theta(1/N)$

**Expert pathway: 8-piece decomposition.**

Piece	$\Delta$ factors	$t = 0$	$t = 1$	$t = 2$
$A_{4.1}$ (no $\Delta W^4$ )	none	0	0	0
$A_{4.2}$ (with $\Delta W^4$ )	$W^4$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_{5.1}$ (no $\Delta W^4$ )	$W^{3,i}$	0	0	0
$A_{5.2}$ (with $\Delta W^4$ )	$W^{3,i}, W^4$	0	0	$\Theta(1/N)$
$A_{6.1}$ (no $\Delta W^4$ )	$W^{2,i}$	0	0	0
$A_{6.2}$ (with $\Delta W^4$ )	$W^{2,i}, W^4$	0	0	$\Theta(1/N)$
$E_1$ (no $\Delta W^4$ )	$W^{2,i}, W^{3,i}$	0	0	0
$E_2$ (with $\Delta W^4$ )	all three	0	0	$\Theta(1/N)$
<b>Expert total</b>	—	0	$\Theta(1/N)$	$\Theta(1/N)$

**Router pathway.**

Piece	$t = 0$	$t = 1$	$t = 2$
$(1/M)Q_t^\top [\dot{\phi}_{i,t} \langle h_t^{3,i}, \Delta_t W^4 \rangle]_i$ <b>(total)</b>	0	$\Theta(1/N)$	$\Theta(1/N)$
$(1/M)Q_0^\top [\dot{\phi}_{i,t} \langle (W_0^{3,1}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	$\Theta(1/N)$	$\Theta(1/N)$
$(1/M)Q_0^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	0	$\Theta(1/N)$
$(1/M)\Delta Q^\top [\dot{\phi}_{i,t} \langle (W_0^{3,1}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	0	$\Theta(1/N)$
$(1/M)\Delta Q^\top [\dot{\phi}_{i,t} \langle (\Delta_t W^{3,i}) h_t^{2,i}, \Delta_t W^4 \rangle]_i$	0	0	$\Theta(1/N)$

**Total backward.**

Quantity	$t = 0$	$t = 1$	$t = 2$
$\partial f_t / \partial h_t^1$	0	$\Theta(1/N)$	$\Theta(1/N)$

#### J.4. Scaling derivation for Regime III

In this section, we provide the heuristic scaling derivation for Regime III. Throughout we work under the proportional limit  $N \asymp N_e \asymp M \asymp K \rightarrow \infty$ .

##### Definition 6 ( $\mu\text{P}$ baseline, Regime III)

- **Initialization.** All parameters are drawn independently according to:

$$W_0^1 \sim \mathcal{N}(0, D^{-1}), \quad Q_0 \sim \mathcal{N}(0, N^{-1}), \quad W_0^{2,i} \sim \mathcal{N}(0, N^{-1}),$$

$$W_0^{3,i} \sim \mathcal{N}(0, N_e^{-1}), \quad W_0^4 \sim \mathcal{N}(0, N^{-2}), \quad \text{for } i \in [M].$$

- **SGD learning rates.**  $\eta_1 = \eta N$ ,  $\eta_Q = \eta$ ,  $\eta_2 = \eta_3 = \eta M$ ,  $\eta_4 = \eta N^{-1}$ .
- **Adam learning rates.**  $\eta_1 = \eta$ ,  $\eta_Q = \eta_2 = \eta_3 = \eta_4 = \eta N^{-1}$ .
- **Adam epsilon.**  $\epsilon_1 = \epsilon N^{-1}$ ,  $\epsilon_Q = \epsilon M^{-1}$ ,  $\epsilon_2 = \epsilon_3 = \epsilon N^{-1} M^{-1}$ ,  $\epsilon_4 = \epsilon$ .

**Definition 7 (Maximally Scale-Stable Parameterization (MSSP), Regime III)** *MSSP adopts the same per-layer initialization variances and learning rate scalings as the  $\mu\text{P}$  baseline (Definition 6), with the sole exception that expert weights are **shared across experts at initialization**: a single pair  $(W_0^2, W_0^3)$  is sampled and assigned to every expert, such that*

$$W_0^{2,i} = W_0^2 \quad \text{and} \quad W_0^{3,i} = W_0^3 \quad \text{for all } i = 1, \dots, M.$$

##### J.4.1. COMPOUND AGGREGATE LEMMAS IN REGIME III

The analyses in this section invoke three compound aggregate-level results stating how sums across  $M$  independent experts of products and Gram products of the per-expert weight matrices behave under proportional scaling. All three are exact moment computations under Gaussian initialization, with independence across experts and from the input vector  $v$  as the only inputs.

**Lemma 8 (Cross-layer product sum across experts)** *Let  $\{W^{2,i}\}_{i=1}^M$  have i.i.d. centered Gaussian entries of variance  $\sigma_2^2 = 1/N$  in  $\mathbb{R}^{N_e \times N}$ , and let  $\{W^{3,i}\}_{i=1}^M$  have i.i.d. centered Gaussian entries of variance  $\sigma_3^2 = 1/N_e$  in  $\mathbb{R}^{N \times N_e}$ , with all matrices mutually independent across  $i$  and across the two layers. Let  $c_i \in \Theta(1)$  be approximately independent of these weights, and let  $v \in \mathbb{R}^N$  have entries of coordinate scale  $\Theta(\alpha)$ , independent of all  $W^{2,i}, W^{3,i}$ . Then*

$$G := \frac{1}{M} \sum_{i=1}^M c_i (W^{2,i})^\top (W^{3,i})^\top \in \mathbb{R}^{N \times N}$$

*has zero mean, entries of coordinate scale  $\Theta(1/\sqrt{NM})$ , and  $Gv$  has entries of coordinate scale  $\Theta(\alpha/\sqrt{M}) = \Theta(\alpha/\sqrt{N})$ .*

**Proof**  $\mathbb{E}[G_{a,b}] = (1/M) \sum_i c_i \mathbb{E}[\sum_k W_{k,a}^{2,i} W_{b,k}^{3,i}] = 0$  by independence across the two layers. Variance:

$$\text{Var}(G_{a,b}) = \frac{1}{M^2} \sum_i c_i^2 \text{Var}\left(\sum_k W_{k,a}^{2,i} W_{b,k}^{3,i}\right) = \frac{1}{M^2} \sum_i c_i^2 \cdot N_e \sigma_2^2 \sigma_3^2 = \Theta\left(\frac{N_e \sigma_2^2 \sigma_3^2}{M}\right) = \Theta\left(\frac{1}{MN}\right).$$

Cross-entry covariances  $\mathbb{E}[G_{a,b_1}G_{a,b_2}]$  for  $b_1 \neq b_2$  vanish by independence, so  $\text{Var}((Gv)_a) = \sum_b \text{Var}(G_{a,b}) v_b^2 = N \cdot \Theta(1/(MN)) \cdot \alpha^2 = \Theta(\alpha^2/M)$ .  $\blacksquare$

**Lemma 9 (Gram concentration across experts)** *Let  $\{W^{a,i}\}_{i=1}^M$  have i.i.d. centered Gaussian entries of variance  $\sigma^2$  in  $\mathbb{R}^{m \times n}$ , independent across  $i$ . Let  $c_i \in \Theta(1)$  be approximately independent of the weights, with empirical mean  $\bar{c}$ . Let  $G^{a,i}$  denote either Gram  $(W^{a,i})^\top W^{a,i} \in \mathbb{R}^{n \times n}$  or  $W^{a,i}(W^{a,i})^\top \in \mathbb{R}^{m \times m}$ , and let  $v$  be a vector of compatible dimension at coordinate scale  $\Theta(\alpha)$ , independent of all  $W^{a,i}$ . Then*

$$\frac{1}{M} \sum_{i=1}^M c_i G^{a,i} v = \bar{c} (\text{tr-dim} \cdot \sigma^2) v + r, \quad r \text{ at coordinate scale } \Theta(\sigma^2 \sqrt{mn/M} \alpha),$$

where  $\text{tr-dim}$  is  $m$  for the inner Gram  $((W^{a,i})^\top W^{a,i})$  and  $n$  for the outer Gram  $(W^{a,i}(W^{a,i})^\top)$ . For the specific instances used in this section:

- *Inner Gram:*  $W^{2,i} \in \mathbb{R}^{N_e \times N}$  with  $\sigma_2^2 = 1/N$  gives leading constant  $m\sigma^2 = N_e/N$ ;  $W^{3,i} \in \mathbb{R}^{N \times N_e}$  with  $\sigma_3^2 = 1/N_e$  gives leading constant  $m\sigma^2 = N/N_e$ .
- *Outer Gram:* both layers give  $n\sigma^2 = 1$ .

All leading constants are  $\Theta(1)$  under proportional scaling, with random correction  $\Theta(\alpha/\sqrt{N})$ .

**Proof** Inner Gram case:  $\mathbb{E}[G_{j,k}] = (1/M) \sum_i c_i \mathbb{E}[\sum_a W_{a,j}^{a,i} W_{a,k}^{a,i}] = \bar{c} m\sigma^2 \delta_{j,k}$ , with off-diagonal variance  $\Theta(m\sigma^4/M)$  giving coordinate scale  $\Theta(\sigma^2 \sqrt{m/M})$  and off-diagonal contribution to  $(Gv)_j - \bar{c} m\sigma^2 v_j$  at coordinate scale  $\Theta(\sigma^2 \sqrt{mn/M} \alpha)$ . Outer Gram case is the same statement with the roles of  $m$  and  $n$  swapped, equivalently obtained by applying the inner-Gram result to  $W^\top$  (whose entries are i.i.d. Gaussian of the same variance with the dimension pair swapped).  $\blacksquare$

#### J.4.2. DERIVING $\mu\text{P}$ IN REGIME III

This subsection provides a self-contained derivation of the heuristic width scaling for the  $\mu\text{P}$  baseline (Definition 6) in Regime III. We invoke the compound aggregate lemmas of §J.4.1 for sums across independent experts of weight-product structures, and use CLT/LLN inline for single matrix-vector products and gating averages.

The principal structural consequence is the following. Although hidden weights move at  $t = 1$  (since  $W_0^4 \neq 0$  in  $\mu\text{P}$ ), the embedding update  $\Delta_1 h^1$  comes out at  $\Theta(1/\sqrt{N})$ , not  $\Theta(1)$ , because both the expert and router contributions to  $\partial f_0/\partial h_0^1$  cancel across independent experts via CLT. Full  $\Theta(1)$  feature learning at the embedding emerges at  $t = 2$ , when the rank-one updates  $\Delta_1 W^{2,i}$ ,  $\Delta_1 W^{3,i}$  inject coherent contributions ( $A_5$ ,  $A_6$ ) into the expert pathway gradient that align across experts along  $(W_0^4)^\top$  and so escape the cross- $i$  cancellation. We therefore analyse three forward and three backward passes.

##### First forward pass

$$h_0^1 = W_0^1 x \in \Theta(1). \quad (\text{CLT}; \sigma_1^2 = 1/D, \|x\|^2 \in \Theta(D)) \quad (\text{R3}\mu\text{-F1.1})$$

$$\psi_0 = Q_0 h_0^1 \in \Theta(1). \quad (\text{CLT}; \sigma_Q^2 = 1/N, \|h_0^1\|^2 \in \Theta(N)) \quad (\text{R3}\mu\text{-F1.2})$$

$$\phi_0 = \sigma(\psi_0) \in \Theta(1). \quad (\text{sigmoid bounded; gating assumption}) \quad (\text{R3}\mu\text{-F1.3})$$

$$h_0^{2,i} = W_0^{2,i} h_0^1 \in \Theta(1), \text{ indep. across } i. \quad (\text{CLT}; \sigma_2^2 = 1/N, \|h_0^1\|^2 \in \Theta(N)) \quad (\text{R3}\mu\text{-F1.4})$$

$$h_0^{3,i} = W_0^{3,i} h_0^{2,i} \in \Theta(1), \text{ indep. across } i. \quad (\text{CLT}; \sigma_3^2 = 1/N_e, \|h_0^{2,i}\|^2 \in \Theta(N_e)) \quad (\text{R3}\mu\text{-F1.5})$$

$$h_0^3 = (1/M) \sum_i \phi_{i,0} h_0^{3,i} \in \Theta(1/\sqrt{M}) = \Theta(1/\sqrt{N}). \quad (\text{cross-}i \text{ CLT}) \quad (\text{R3}\mu\text{-F1.6})$$

$$f_0 = W_0^4 h_0^3 = \langle W_0^4, h_0^3 \rangle \in \Theta(1/N). \quad (\text{CLT}; W_0^4 \text{ entries } \Theta(1/N) \text{ indep. of } h_0^3) \quad (\text{R3}\mu\text{-F1.7})$$

### First backward pass and step-1 updates

$$\partial f_0 / \partial h_0^3 = (W_0^4)^\top \in \Theta(1/N). \quad (\text{R3}\mu\text{-B1.1})$$

$$\partial f_0 / \partial h_0^{3,i} = (\phi_{i,0}/M)(W_0^4)^\top \in \Theta(1/(MN)). \quad (\text{R3}\mu\text{-B1.2})$$

$$\begin{aligned} \partial f_0 / \partial h_0^{2,i} &= (\phi_{i,0}/M)(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/(MN)). \quad (\text{R3}\mu\text{-B1.3}) \\ &((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ by CLT; } \sigma_3^2 = 1/N_e) \end{aligned}$$

$$\begin{aligned} \partial f_0 / \partial \phi_{i,0} &= (1/M) \langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/(M\sqrt{N})) = \Theta(1/N^{3/2}). \quad (\text{R3}\mu\text{-B1.4}) \\ &(\langle h_0^{3,i}, W_0^4 \rangle \in \Theta(1/\sqrt{N}) \text{ by CLT}) \end{aligned}$$

$$\begin{aligned} \left( \frac{\partial f_0}{\partial h_0^1} \right)_{\text{exp}} &= \frac{1}{M} \sum_i \phi_{i,0} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top \\ &\in \Theta(\sigma_2 \sigma_3 \sqrt{N_e} \| (W_0^4)^\top \| / \sqrt{M}) \text{ random direction} \quad (\text{Lem. 8; cross-}i \text{ CLT}) \\ &= \Theta(1/\sqrt{N} \cdot 1/\sqrt{N_e} \cdot \sqrt{N_e} \cdot 1/\sqrt{N} \cdot 1/\sqrt{M}) = \Theta(1/N^{3/2}). \quad (\text{R3}\mu\text{-B1.5}) \end{aligned}$$

$$\begin{aligned} \left( \frac{\partial f_0}{\partial h_0^1} \right)_{\text{router}} &= \frac{1}{M} Q_0^\top v, \quad v_i := \dot{\phi}_{i,0} \langle W_0^4, h_0^{3,i} \rangle \\ &\in \Theta(\sigma_Q \|v\|/M) \text{ random direction} \quad (Q_0 \text{ indep. of } v; \text{ cross-layer CLT}) \\ &= \Theta(1/\sqrt{N} \cdot 1 \cdot 1/N) = \Theta(1/N^{3/2}) \quad (\text{R3}\mu\text{-B1.6}) \\ &(\sigma_Q^2 = 1/N; v_i \in \Theta(1/\sqrt{N}) \text{ random across } i \Rightarrow \|v\|^2 \in \Theta(1)). \end{aligned}$$

$$\partial f_0 / \partial h_0^1 = (\partial f_0 / \partial h_0^1)_{\text{exp}} + (\partial f_0 / \partial h_0^1)_{\text{router}} \in \Theta(1/N^{3/2}). \quad (\text{R3}\mu\text{-B1.7})$$

**Step-1 parameter updates.**

$$\Delta_1 W^4 = -\frac{\chi_0}{N} (h_0^3)^\top \in \Theta(1/N^{3/2}). \quad (h_0^3 \in \Theta(1/\sqrt{N}) \text{ by (R3}\mu\text{-F1.6)}) \quad (\text{R3}\mu\text{-U1.4})$$

$$\Delta_1 W^{3,i} = -\eta \chi_0 \phi_{i,0} (W_0^4)^\top (h_0^{2,i})^\top \in \Theta(1/N). \quad (\eta_3 = \eta M) \quad (\text{R3}\mu\text{-U1.3})$$

$$\begin{aligned} \Delta_1 W^{2,i} &= -\eta \chi_0 \phi_{i,0} (W_0^{3,i})^\top (W_0^4)^\top (h_0^1)^\top \in \Theta(1/N). & (\text{R3}\mu\text{-U1.2}) \\ &((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ by (R3}\mu\text{-B1.3)}) \end{aligned}$$

$$\begin{aligned} \Delta_1 W^1 &= -\eta N \chi_0 (\partial f_0 / \partial h_0^1) x^\top \in \Theta(1/\sqrt{N}). & (\text{R3}\mu\text{-U1.1a}) \\ &(\eta_1 = \eta N, \partial f_0 / \partial h_0^1 \in \Theta(1/N^{3/2}) \text{ by (R3}\mu\text{-B1.7)}) \end{aligned}$$

$$\Delta_1 h^1 = \Delta_1 W^1 x \in \Theta(1/\sqrt{N}), \text{ sub-leading vs } h_0^1 \in \Theta(1). \quad (\text{R3}\mu\text{-U1.1b})$$

$$\begin{aligned} \Delta_1 Q &= -\eta_Q \chi_0 \text{diag}(\dot{\phi}_0) (\partial f_0 / \partial \phi) (h_0^1)^\top \in \Theta(1/N^{3/2}). & (\text{R3}\mu\text{-U1.Q}) \\ &(\partial f_0 / \partial \phi \in \Theta(1/N^{3/2}) \text{ by (R3}\mu\text{-B1.4); } \eta_Q = \eta) \end{aligned}$$

**Second forward pass**

$$h_1^1 = h_0^1 + \Delta_1 h^1 \in \Theta(1). \quad (h_0^1 \in \Theta(1) \text{ dominates } \Delta_1 h^1 \in \Theta(1/\sqrt{N}) \text{ by (R3}\mu\text{-U1.1b)}) \quad (\text{R3}\mu\text{-F2.1})$$

$$\begin{aligned} \psi_1 &= \psi_0 + \Delta_1 \psi \in \Theta(1), \quad \phi_1 \in \Theta(1). & (\text{R3}\mu\text{-F2.2--R3}\mu\text{-F2.3}) \\ &(\Delta_1 \psi \in \Theta(1/\sqrt{N}) \text{ from } Q_0 \Delta_1 h^1 \text{ and } \Delta_1 Q h_0^1; \text{ scales by (R3}\mu\text{-U1.1b), (R3}\mu\text{-U1.Q)}) \end{aligned}$$

(R3 $\mu$ -F2.4)  $h_1^{2,i} = W_1^{2,i} h_1^1$ , four-piece decomposition:

$$h_1^{2,i} = h_0^{2,i} + W_0^{2,i} \Delta_1 h^1 + \Delta_1 W^{2,i} h_0^1 + \Delta_1 W^{2,i} \Delta_1 h^1. \quad (\text{R3}\mu\text{-F2.4})$$

$$\text{init: } W_0^{2,i} h_0^1 = h_0^{2,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F1.4}) \quad (\text{R3}\mu\text{-F2.4a})$$

Op-norm of iid Gaussian  $W_0^{2,i} \Delta_1 h^1 \in \Theta(1/\sqrt{N})$  by (R3 $\mu$ -U1.1b):

$$\text{prop: } W_0^{2,i} \Delta_1 h^1 \in \Theta(1/\sqrt{N}). \quad (\text{R3}\mu\text{-F2.4b})$$

$$\begin{aligned} \text{eff: } \Delta_1 W^{2,i} h_0^1 &= -\eta \chi_0 \phi_{i,0} (W_0^{3,i})^\top (W_0^4)^\top \|h_0^1\|^2 \in \Theta(1). & (\text{R3}\mu\text{-F2.4c}) \\ &((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ by (R3}\mu\text{-B1.3); } \|h_0^1\|^2 \in \Theta(N)) \end{aligned}$$

$$\begin{aligned} \text{cross: } \Delta_1 W^{2,i} \Delta_1 h^1 &= -\eta \chi_0 \phi_{i,0} (W_0^{3,i})^\top (W_0^4)^\top \langle h_0^1, \Delta_1 h^1 \rangle \in \Theta(1/N). & (\text{R3}\mu\text{-F2.4d}) \\ &((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ by (R3}\mu\text{-B1.3); } \langle h_0^1, \Delta_1 h^1 \rangle \in \Theta(1) \text{ by CLT}) \end{aligned}$$

(R3 $\mu$ -F2.5)  $h_1^{3,i} = W_1^{3,i} h_1^{2,i}$ , four-piece decomposition:

$$h_1^{3,i} = h_0^{3,i} + W_0^{3,i} \Delta_1 h^{2,i} + \Delta_1 W^{3,i} h_0^{2,i} + \Delta_1 W^{3,i} \Delta_1 h^{2,i}. \quad (\text{R3}\mu\text{-F2.5})$$

$$\text{init: } h_0^{3,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F1.5}) \quad (\text{R3}\mu\text{-F2.5a})$$

$$\text{prop: } W_0^{3,i} \Delta_1 h^{2,i} \in \Theta(1). \quad (\text{op-norm of iid Gaussian } W_0^{3,i}; \Delta_1 h^{2,i} \in \Theta(1)) \quad (\text{R3}\mu\text{-F2.5b})$$

$$\begin{aligned} \text{eff: } \Delta_1 W^{3,i} h_0^{2,i} &= -\eta\chi_0 \phi_{i,0} (W_0^4)^\top \|h_0^{2,i}\|^2 \in \Theta(1). & (\text{R3}\mu\text{-F2.5c}) \\ &(\|h_0^{2,i}\|^2 \in \Theta(N_e); (W_0^4)^\top \in \Theta(1/N)) \end{aligned}$$

$$\text{cross: } \Delta_1 W^{3,i} \Delta_1 h^{2,i} \text{ tracked as part of } D \text{ in } (\text{R3}\mu\text{-F2.6}). \quad (\text{R3}\mu\text{-F2.5d})$$

$$(\text{R3}\mu\text{-F2.6}) h_1^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D, \text{ where}$$

$$\begin{aligned} A_1 &:= \frac{1}{M} \sum_i \phi_{i,1} h_0^{3,i}, & A_{2,1} &:= \frac{1}{M} \sum_i \phi_{i,1} W_0^{3,i} W_0^{2,i} \Delta_1 h^1, \\ A_{2,2} &:= \frac{1}{M} \sum_i \phi_{i,1} W_0^{3,i} \Delta_1 W^{2,i} h_0^1, & A_3 &:= \frac{1}{M} \sum_i \phi_{i,1} \Delta_1 W^{3,i} h_0^{2,i}, \\ D &:= \frac{1}{M} \sum_i \phi_{i,1} \Delta_1 W^{3,i} \Delta_1 h^{2,i}. \end{aligned}$$

$$A_1 \in \Theta(1/\sqrt{N}). \quad (\text{cross-}i \text{ CLT on } \Theta(1)\text{-entry independent vectors}) \quad (\text{R3}\mu\text{-F2.6a})$$

$$A_{2,1} \in \Theta(1/(\sqrt{N} \sqrt{M})) = \Theta(1/N). \quad (\text{R3}\mu\text{-F2.6b})$$

$$(\text{Lemma 8; } \Delta_1 h^1 \text{ entries } \Theta(1/\sqrt{N}) \text{ by } (\text{R3}\mu\text{-U1.1b}))$$

$$\begin{aligned} A_{2,2} &= -\eta\chi_0 \|h_0^1\|^2 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} W_0^{3,i} (W_0^{3,i})^\top (W_0^4)^\top \\ &\in \Theta(\|h_0^1\|^2 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0}) \cdot (W_0^4)^\top \text{ along } (W_0^4)^\top \quad (\text{Lem. 9}) \\ &= \Theta(N \cdot 1 \cdot 1/N) = \Theta(1) \text{ along } (W_0^4)^\top & (\text{R3}\mu\text{-F2.6c}) \\ &(\|h_0^1\|^2 \in \Theta(N); \text{LLN: } \frac{1}{M} \sum \phi\phi \in \Theta(1); (W_0^4)^\top \in \Theta(1/N)) \end{aligned}$$

$$\begin{aligned} A_3 &= -\eta\chi_0 (W_0^4)^\top \left( \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \|h_0^{2,i}\|^2 \right) \in \Theta(1) \text{ along } (W_0^4)^\top. & (\text{R3}\mu\text{-F2.6d}) \\ &(\|h_0^{2,i}\|^2 \in \Theta(N_e); \text{LLN}) \end{aligned}$$

$$\begin{aligned} D &= \frac{1}{M} \sum_i \phi_{i,1} \Delta_1 W^{3,i} \Delta_1 h^{2,i} = -\eta\chi_0 (W_0^4)^\top \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle \\ &\in \Theta\left(\frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle\right) \cdot (W_0^4)^\top \text{ along } (W_0^4)^\top \\ &(\text{cross-}i \text{ CLT: } \langle \cdot, \cdot \rangle \text{ random across } i, \text{ mean } 0) \\ &= \Theta(\sqrt{N}/\sqrt{M} \cdot 1/N) = \Theta(1/N) \text{ along } (W_0^4)^\top & (\text{R3}\mu\text{-F2.6e}) \\ &(\langle h_0^{2,i}, \Delta_1 h^{2,i} \rangle \in \Theta(\sqrt{N}); (W_0^4)^\top \in \Theta(1/N); M \in \Theta(N)) \end{aligned}$$

$$h_1^3 \in \Theta(1) \text{ along } (W_0^4)^\top, \text{ from } A_{2,2} + A_3. \quad (\text{R3}\mu\text{-F2.6})$$

$$W_0^4 h_1^3 = \langle W_0^4, A_{2,2} + A_3 \rangle + \dots \in \Theta(1). \quad (\text{R3}\mu\text{-F2.7a})$$

$$(A_{2,2}, A_3 \text{ aligned along } (W_0^4)^\top; \langle W_0^4, A_3 \rangle = -\eta\chi_0 \|W_0^4\|^2 \bar{c} N_e \in \Theta(1))$$

$$\Delta_1 W^4 h_1^3 = -(\chi_0/N) \langle h_0^3, h_1^3 \rangle \in \Theta(1/N). \quad (\text{R3}\mu\text{-F2.7b})$$

$$f_1 = W_1^4 h_1^3 \in \Theta(1). \quad (\text{R3}\mu\text{-F2.7})$$

**Second backward pass and step-2 updates**

$$\partial f_1 / \partial h_1^3 = (W_1^4)^\top = (W_0^4)^\top + (\Delta_1 W^4)^\top \in \Theta(1/N). \quad (\text{R3}\mu\text{-B2.1})$$

$$(W_0^4 \in \Theta(1/N) \text{ dominates } \Delta_1 W^4 \in \Theta(1/N^{3/2}) \text{ by } (\text{R3}\mu\text{-U1.4}))$$

$$\partial f_1 / \partial h_1^{3,i} = (\phi_{i,1}/M)(W_1^4)^\top \in \Theta(1/(MN)). \quad (\text{R3}\mu\text{-B2.2})$$

$$\partial f_1 / \partial h_1^{2,i} = (\phi_{i,1}/M)(W_1^{3,i})^\top (W_1^4)^\top \in \Theta(1/(MN)). \quad (\text{R3}\mu\text{-B2.3})$$

$$((W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N) \text{ by } (\text{R3}\mu\text{-B1.3}); (\Delta_1 W^{3,i})^\top W_0^4 \in \Theta(1/N);$$

$$\phi_{i,1} \in \Theta(1); M \in \Theta(N))$$

$$\partial f_1 / \partial \phi_{i,1} = (1/M) \langle h_1^{3,i}, W_1^4 \rangle \in \Theta(1/M). \quad (\text{R3}\mu\text{-B2.4})$$

$$(\langle h_1^{3,i}, W_0^4 \rangle = -\eta \chi_0 \phi_{i,0} \|h_0^{2,i}\|^2 \|W_0^4\|^2 \in \Theta(1) \text{ via eff. term of } (\text{R3}\mu\text{-F2.5});$$

$$\|h_0^{2,i}\|^2 \in \Theta(N_e); \|W_0^4\|^2 \in \Theta(1/N))$$

(R3μ-B2.5)  $(\partial f_1 / \partial h_1^1)_{\text{exp}} = A_4 + A_5 + A_6 + E$ , where

$$\begin{aligned} A_4 &:= \frac{1}{M} \sum_i \phi_{i,1} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_1^4)^\top, \\ A_5 &:= \frac{1}{M} \sum_i \phi_{i,1} (W_0^{2,i})^\top (\Delta_1 W^{3,i})^\top (W_1^4)^\top, \\ A_6 &:= \frac{1}{M} \sum_i \phi_{i,1} (\Delta_1 W^{2,i})^\top (W_0^{3,i})^\top (W_1^4)^\top, \\ E &:= \frac{1}{M} \sum_i \phi_{i,1} (\Delta_1 W^{2,i})^\top (\Delta_1 W^{3,i})^\top (W_1^4)^\top. \end{aligned}$$

$$A_4 \in \Theta(1/N^{3/2}). \quad (\text{Lemma 8}; v = (W_1^4)^\top \in \Theta(1/N)) \quad (\text{R3}\mu\text{-B2.5a})$$

$$\begin{aligned} A_5 &= -\eta \chi_0 \|W_0^4\|^2 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} (W_0^{2,i})^\top W_0^{2,i} h_0^1 \\ &\in \Theta(\|W_0^4\|^2 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0}) \cdot h_0^1 \text{ along } h_0^1 \quad (\text{Lem. 9}) \end{aligned}$$

$$= \Theta(1/N \cdot 1 \cdot 1) = \Theta(1/N) \text{ along } h_0^1 \quad (\text{R3}\mu\text{-B2.5b})$$

$$(\|W_0^4\|^2 \in \Theta(1/N); \text{LLN: } \frac{1}{M} \sum \phi \phi \in \Theta(1); h_0^1 \in \Theta(1))$$

$$\begin{aligned} A_6 &= -\eta \chi_0 h_0^1 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} W_0^4 W_0^{3,i} (W_0^{3,i})^\top (W_1^4)^\top \\ &\in \Theta(\frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0} \|W_0^4\|^2) \cdot h_0^1 \text{ along } h_0^1 \end{aligned}$$

$$(\text{Lem. 9: } W_0^{3,i} (W_0^{3,i})^\top (W_1^4)^\top \rightarrow (W_1^4)^\top; \langle W_0^4, (W_1^4)^\top \rangle \rightarrow \|W_0^4\|^2)$$

$$= \Theta(1 \cdot 1/N) = \Theta(1/N) \text{ along } h_0^1 \quad (\text{R3}\mu\text{-B2.5c})$$

$$(\|W_0^4\|^2 \in \Theta(1/N); \text{LLN}; h_0^1 \in \Theta(1))$$

$$\begin{aligned} E &= \frac{1}{M} \sum_i \phi_{i,1} (\Delta_1 W^{2,i})^\top (\Delta_1 W^{3,i})^\top (W_1^4)^\top \\ &= \eta^2 \chi_0^2 h_0^1 \frac{1}{M} \sum_i \phi_{i,1} \phi_{i,0}^2 \langle W_0^4, h_0^{3,i} \rangle (h_0^3)^\top (W_1^4)^\top \quad (\text{substituting } (\text{R3}\mu\text{-U1.2}), (\text{R3}\mu\text{-U1.3})) \\ &\in \Theta(\frac{1}{M} \sum_i \langle W_0^4, h_0^{3,i} \rangle) \cdot h_0^1 \text{ along } h_0^1 \quad (\text{cross-}i \text{ CLT: } \langle W_0^4, h_0^{3,i} \rangle \text{ random across } i, \text{ mean } 0) \\ &= \Theta(1/\sqrt{N} \cdot 1/\sqrt{M} \cdot \|h_0^3\|^2/N) = \Theta(1/N^2) \text{ along } h_0^1 \quad (\text{R3}\mu\text{-B2.5d}) \end{aligned}$$

$$(\langle W_0^4, h_0^{3,i} \rangle \in \Theta(1/\sqrt{N}) \text{ by } (\text{R3}\mu\text{-B1.4}); \|h_0^3\|^2 \in \Theta(1); M \in \Theta(N))$$

$$(\partial f_1/\partial h_1^1)_{\text{exp}} \in \Theta(1/N) \text{ along } h_0^1. \quad (\text{R3}\mu\text{-B2.5})$$

$$\begin{aligned} (\partial f_1/\partial h_1^1)_{\text{router}} &= \frac{1}{M} Q_0^\top v, \quad v_i := \dot{\phi}_{i,1} \langle h_1^{3,i}, W_1^4 \rangle \\ &\in \Theta(\sigma_Q \|v\|/M) \text{ random direction} \\ &\quad (Q_0 \text{ indep. of } v; \text{ cross-layer CLT; } \Delta_1 Q \text{ sub-leading by (R3}\mu\text{-U1.Q)}) \\ &= \Theta(1/\sqrt{N} \cdot \sqrt{M} \cdot 1/M) = \Theta(1/\sqrt{MN}) = \Theta(1/N) \quad (\text{R3}\mu\text{-B2.6}) \\ &\quad (\sigma_Q^2 = 1/N; v_i \in \Theta(1) \text{ via (R3}\mu\text{-B2.4)} \Rightarrow \|v\|^2 \in \Theta(M); M \in \Theta(N)) \end{aligned}$$

$$\partial f_1/\partial h_1^1 = (\partial f_1/\partial h_1^1)_{\text{exp}} + (\partial f_1/\partial h_1^1)_{\text{router}} \in \Theta(1/N), \quad (\text{R3}\mu\text{-B2.7})$$

the proper feature-learning rate, an order of magnitude larger than  $\partial f_0/\partial h_0^1$ .

### Step-2 parameter updates.

$$\Delta_2 W^4 = -\frac{\chi_1}{N} (h_1^3)^\top \in \Theta(1/N). \quad (h_1^3 \in \Theta(1) \text{ by (R3}\mu\text{-F2.6)}) \quad (\text{R3}\mu\text{-U2.4})$$

$$\begin{aligned} \Delta_2 W^{3,i}, \Delta_2 W^{2,i} &\in \Theta(1/N). \quad (\text{R3}\mu\text{-U2.3, R3}\mu\text{-U2.2}) \\ &\quad (\text{same structure as (R3}\mu\text{-U1.3), (R3}\mu\text{-U1.2) with cumulative weights}) \end{aligned}$$

$$\begin{aligned} \Delta_2 W^1 &= -\eta_1 \chi_1 (\partial f_1/\partial h_1^1) x^\top \in \Theta(1). \quad (\text{R3}\mu\text{-U2.1a}) \\ &\quad (\eta_1 = N, \partial f_1/\partial h_1^1 \in \Theta(1/N) \text{ by (R3}\mu\text{-B2.7)}) \end{aligned}$$

$$\Delta_2 h^1 = \Delta_2 W^1 x \in \Theta(1) \text{ aligned along } h_0^1. \quad (\text{embedding feature-learns at } t=2) \quad (\text{R3}\mu\text{-U2.1b})$$

$$\Delta_2 Q \in \Theta(1/N). \quad (\partial f_1/\partial \phi \in \Theta(1/N) \text{ by (R3}\mu\text{-B2.4)}) \quad (\text{R3}\mu\text{-U2.Q})$$

**Third forward pass** We compute activations at  $\theta^{(2)}$ . The cumulative embedding change  $\Delta h^1 := h_2^1 - h_0^1 = \Delta_1 h^1 + \Delta_2 h^1$ , with  $\Delta_1 h^1 \in \Theta(1/\sqrt{N})$  by (R3 $\mu$ -U1.1b) and  $\Delta_2 h^1 \in \Theta(1)$  aligned along  $h_0^1$  by (R3 $\mu$ -U2.1b).

$$h_2^1 = h_0^1 + \Delta h^1 \in \Theta(1). \quad (\text{R3}\mu\text{-F3.1})$$

(R3 $\mu$ -F3.2)  $h_2^{2,i} = W_2^{2,i} h_2^1$ , four-piece decomposition with cumulative  $\Delta W^{2,i}$ :

$$\text{init: } W_0^{2,i} h_0^1 = h_0^{2,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F1.4}) \quad (\text{R3}\mu\text{-F3.2a})$$

$$\begin{aligned} \text{prop: } W_0^{2,i} \Delta h^1 &\in \Theta(1). \quad (\text{op-norm of iid Gaussian } W_0^{2,i}; \Delta h^1 \in \Theta(1) \text{ by (R3}\mu\text{-U2.1b)}) \\ &\quad (\text{R3}\mu\text{-F3.2b}) \end{aligned}$$

$$\begin{aligned} \text{eff: } \Delta_t W^{2,i} h_0^1 &= -\eta \chi \phi(W_{t-1}^{3,i})^\top (W_{t-1}^4)^\top \langle h_{t-1}^1, h_0^1 \rangle \in \Theta(1). \quad (\text{R3}\mu\text{-F3.2c}) \\ &\quad ((W_{t-1}^{3,i})^\top (W_{t-1}^4)^\top \in \Theta(1/N) \text{ by (R3}\mu\text{-B1.3)/(R3}\mu\text{-B2.3);} \\ &\quad \langle h_{t-1}^1, h_0^1 \rangle \in \Theta(N) \text{ coherent via (R3}\mu\text{-U2.1b)}) \end{aligned}$$

$$\begin{aligned} \text{cross: } \Delta_t W^{2,i} \Delta h^1 &= -\eta \chi \phi(W_{t-1}^{3,i})^\top (W_{t-1}^4)^\top \langle h_{t-1}^1, \Delta h^1 \rangle \in \Theta(1). \quad (\text{R3}\mu\text{-F3.2d}) \\ &\quad ((W_{t-1}^{3,i})^\top (W_{t-1}^4)^\top \in \Theta(1/N) \text{ by (R3}\mu\text{-B1.3)/(R3}\mu\text{-B2.3);} \\ &\quad \langle h_{t-1}^1, \Delta h^1 \rangle \in \Theta(N) \text{ coherent via (R3}\mu\text{-U2.1b)}) \end{aligned}$$

$$h_2^{2,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F3.2})$$

(R3 $\mu$ -F3.3)  $h_2^{3,i} = W_2^{3,i} h_2^{2,i}$ , four-piece decomposition:

$$\text{init: } W_0^{3,i} h_0^{2,i} = h_0^{3,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F1.5}) \quad (\text{R3}\mu\text{-F3.3a})$$

$$\text{prop: } W_0^{3,i} \Delta h^{2,i} \in \Theta(1). \quad (\text{op-norm of iid Gaussian } W_0^{3,i}; \Delta h^{2,i} \in \Theta(1)) \quad (\text{R3}\mu\text{-F3.3b})$$

$$\text{eff: } \Delta_t W^{3,i} h_0^{2,i} = -\eta \chi \phi(W_{t-1}^4)^\top \langle h_{t-1}^{2,i}, h_0^{2,i} \rangle \in \Theta(1). \quad (\text{R3}\mu\text{-F3.3c})$$

$$((W_{t-1}^4)^\top \in \Theta(1/N); \langle h_{t-1}^{2,i}, h_0^{2,i} \rangle \in \Theta(N_e) \text{ coherent})$$

$$\text{cross: } \Delta W^{3,i} \Delta h^{2,i} \text{ tracked as part of } D' \text{ in (R3}\mu\text{-F3.4)}. \quad (\text{R3}\mu\text{-F3.3d})$$

$$h_2^{3,i} \in \Theta(1). \quad (\text{R3}\mu\text{-F3.3})$$

(R3 $\mu$ -F3.4)  $h_2^3 = A'_1 + A'_{2,1} + A'_{2,2} + A'_3 + D'$ :

$$A'_1 = (1/M) \sum_i \phi_{i,2} h_0^{3,i} \in \Theta(1/\sqrt{N}). \quad (\text{cross-}i \text{ CLT}) \quad (\text{R3}\mu\text{-F3.4a})$$

$$A'_{2,1} \in \Theta(1/\sqrt{M}) = \Theta(1/\sqrt{N}). \quad (\text{Lemma 8; } \Delta h^1 \in \Theta(1) \text{ by (R3}\mu\text{-U2.1b)}) \quad (\text{R3}\mu\text{-F3.4b})$$

$$\begin{aligned} A'_{2,2} &= -\eta \chi_1 \langle h_1^1, h_0^1 \rangle \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1} W_0^{3,i} (W_1^{3,i})^\top (W_1^4)^\top \\ &\in \Theta(\langle h_1^1, h_0^1 \rangle \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,1}) \cdot (W_0^4)^\top \text{ along } (W_0^4)^\top \quad (\text{Lem. 9}) \\ &= \Theta(N \cdot 1 \cdot 1/N) = \Theta(1) \text{ along } (W_0^4)^\top \\ &(\langle h_1^1, h_0^1 \rangle \in \Theta(N) \text{ coherent; LLN; } (W_0^4)^\top \in \Theta(1/N)). \end{aligned} \quad (\text{R3}\mu\text{-F3.4c})$$

$$A'_3 \in \Theta(1) \text{ along } (W_0^4)^\top. \quad (\Delta_2 W^{3,i} h_0^{2,i} \in \Theta(1) \text{ along } (W_1^4)^\top; \text{LLN}) \quad (\text{R3}\mu\text{-F3.4d})$$

$$D' \in \Theta(1) \text{ along } (W_0^4)^\top. \quad (\text{R3}\mu\text{-F3.4e})$$

$$(\langle h_1^{2,i}, \Delta h^{2,i} \rangle \in \Theta(N_e) \text{ coherent via alignment chain through (R3}\mu\text{-F3.2); LLN})$$

$$h_2^3 \in \Theta(1) \text{ along } (W_0^4)^\top. \quad (\text{R3}\mu\text{-F3.4})$$

(R3 $\mu$ -F3.5)  $f_2 = W_2^4 h_2^3$ , with  $W_2^4 = W_0^4 + \Delta_1 W^4 + \Delta_2 W^4$ :

$$W_0^4 h_2^3 \in \Theta(1). \quad (A'_{2,2} + A'_3 + D' = \kappa (W_0^4)^\top \text{ with } \kappa \in \Theta(N); \kappa \|W_0^4\|^2 \in \Theta(1)) \quad (\text{R3}\mu\text{-F3.5a})$$

$$\Delta_2 W^4 h_2^3 = -(\chi_1/N) \langle h_1^3, h_2^3 \rangle \in \Theta(1). \quad (\text{R3}\mu\text{-F3.5b})$$

$$(\langle h_1^3, h_2^3 \rangle \in \Theta(N) \text{ via aligned components along } (W_0^4)^\top;$$

$$h_1^3 \in \Theta(1) \text{ along } (W_0^4)^\top \text{ by (R3}\mu\text{-F2.6)})$$

$$\Delta_1 W^4 h_2^3 \in \Theta(1/\sqrt{N}). \quad (\Delta_1 W^4 \in \Theta(1/N^{3/2}) \text{ sub-leading by (R3}\mu\text{-U1.4)}) \quad (\text{R3}\mu\text{-F3.5c})$$

$$f_2 \in \Theta(1). \quad (\text{R3}\mu\text{-F3.5})$$

**Third backward pass** We compute gradients at  $\theta^{(2)}$ , with  $W_2^4 = W_0^4 + \Delta_1 W^4 + \Delta_2 W^4$  where  $W_0^4 \in \Theta(1/N)$ ,  $\Delta_1 W^4 \in \Theta(1/N^{3/2})$  by (R3 $\mu$ -U1.4), and  $\Delta_2 W^4 = -(\chi_1/N)(h_1^3)^\top \in \Theta(1/N)$  aligned along  $(h_1^3)^\top$  by (R3 $\mu$ -U2.4).

$$\partial f_2 / \partial h_2^3 = (W_2^4)^\top \in \Theta(1/N). \quad (\text{R3}\mu\text{-B3.1})$$

( $W_0^4 \in \Theta(1/N)$  and  $\Delta_2 W^4 \in \Theta(1/N)$  by (R3 $\mu$ -U2.4) both contribute;

$\Delta_1 W^4 \in \Theta(1/N^{3/2})$  by (R3 $\mu$ -U1.4) sub-leading)

$$\partial f_2 / \partial h_2^{3,i} = (\phi_{i,2}/M)(W_2^4)^\top \in \Theta(1/(MN)). \quad (\phi_{i,2} \in \Theta(1)) \quad (\text{R3}\mu\text{-B3.2})$$

$$\partial f_2 / \partial h_2^{2,i} = (\phi_{i,2}/M)(W_2^{3,i})^\top (W_2^4)^\top \in \Theta(1/(MN)). \quad (\text{R3}\mu\text{-B3.3})$$

( $(W_0^{3,i})^\top (W_0^4)^\top \in \Theta(1/N)$  via (R3 $\mu$ -B1.3);

$(\Delta W^{3,i})^\top (W_2^4)^\top \in \Theta(1/N)$  via Lem. 9 on aligned components)

$$\partial f_2 / \partial \phi_{i,2} = (1/M) \langle h_2^{3,i}, W_2^4 \rangle \in \Theta(1/M). \quad (\text{R3}\mu\text{-B3.4})$$

( $\langle h_2^{3,i}, W_2^4 \rangle \in \Theta(1)$  via the  $(W_0^4)^\top$ -aligned effective  
and propagating components of (R3 $\mu$ -F3.3))

(R3 $\mu$ -B3.5) Expert contribution. The chain through the experts gives

$$\left( \frac{\partial f_2}{\partial h_2^1} \right)_{\text{exp}} = \frac{1}{M} \sum_i \phi_{i,2} (W_2^{2,i})^\top (W_2^{3,i})^\top (W_2^4)^\top.$$

Expand  $W_2^{a,i} = W_0^{a,i} + \Delta W^{a,i}$  with cumulative update  $\Delta W^{a,i} = \Delta_1 W^{a,i} + \Delta_2 W^{a,i}$ . The four-term decomposition is:

$$\begin{aligned} A_4'' &:= \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_2^4)^\top, \\ A_5'' &:= \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,i})^\top (\Delta W^{3,i})^\top (W_2^4)^\top, \\ A_6'' &:= \frac{1}{M} \sum_i \phi_{i,2} (\Delta W^{2,i})^\top (W_0^{3,i})^\top (W_2^4)^\top, \\ E'' &:= \frac{1}{M} \sum_i \phi_{i,2} (\Delta W^{2,i})^\top (\Delta W^{3,i})^\top (W_2^4)^\top. \end{aligned}$$

$$A_4'' \in \Theta(1/N^{3/2}). \quad (\text{Lemma 8; } v = (W_2^4)^\top \in \Theta(1/N)) \quad (\text{R3}\mu\text{-B3.5a})$$

$$\begin{aligned} A_5'' &= -\eta \chi \|W_0^4\|^2 \frac{1}{M} \sum_i \phi_{i,2} \phi (W_0^{2,i})^\top W_0^{2,i} h_0^1 + \dots \\ &\in \Theta(\|W_0^4\|^2 \frac{1}{M} \sum_i \phi_{i,2} \phi) \cdot h_0^1 \text{ along } h_0^1 \\ &\quad (\text{Lem. 9; } \Delta_1 \text{ and } \Delta_2 \text{ contributions share leading direction}) \\ &= \Theta(1/N \cdot 1 \cdot 1) = \Theta(1/N) \text{ along } h_0^1 \\ &\quad (\|W_0^4\|^2 \in \Theta(1/N); \text{ LLN; } h_0^1 \in \Theta(1)) \end{aligned} \quad (\text{R3}\mu\text{-B3.5b})$$

$$A_6'' \in \Theta(1/N) \text{ along } h_0^1. \quad (\text{R3}\mu\text{-B3.5c})$$

(Lemma 9; LLN;  $(\Delta W^{2,i})^\top$  rank-1 with leading direction  $h_0^1 \otimes (W_0^4 W_0^{3,i})$ )

$$\begin{aligned} E'' &= \eta^2 \chi^2 h_0^1 \frac{1}{M} \sum_i \phi_{i,2} \phi^2 \langle W_t^4, h_t^{3,i} \rangle (h_t^3)^\top (W_2^4)^\top + \dots \\ &\in \Theta\left(\frac{1}{M} \sum_i \langle W_t^4, h_t^{3,i} \rangle \cdot \langle h_t^3, W_2^4 \rangle\right) \cdot h_0^1 \text{ along } h_0^1 \\ &\quad (\text{at } t=2: \langle W_t^4, h_t^{3,i} \rangle \in \Theta(1) \text{ coherent via eff. term of (R3}\mu\text{-F2.5); LLN}) \\ &= \Theta(1 \cdot 1) = \Theta(1/N) \text{ along } h_0^1 \quad (\text{R3}\mu\text{-B3.5d}) \\ &\quad (\text{after multiplying by } \langle h_t^3, W_2^4 \rangle \in \Theta(1) \text{ by (R3}\mu\text{-F3.5); prefactor } \eta^2 \|W_0^4\|^2 \in \Theta(1/N)) \\ (\partial f_2 / \partial h_2^1)_{\text{exp}} &= A_4'' + A_5'' + A_6'' + E'' \in \Theta(1/N) \text{ along } h_0^1. \quad (\text{R3}\mu\text{-B3.5}) \end{aligned}$$

$$\begin{aligned} (\partial f_2 / \partial h_2^1)_{\text{router}} &= \frac{1}{M} Q_0^\top v + \frac{1}{M} (\Delta_2 Q)^\top v, \quad v_i := \dot{\phi}_{i,2} \langle h_2^{3,i}, W_2^4 \rangle \\ &\in \Theta(\sigma_Q \|v\| / M) \text{ random direction} + \Theta(\langle \partial f_1 / \partial \phi, v \rangle / M) h_1^1 \text{ along } h_1^1 \\ &\quad (Q_0 \text{ piece via cross-layer CLT; } \Delta_2 Q \text{ piece coherent}) \\ &= \Theta(1/\sqrt{N} \cdot \sqrt{M} \cdot 1/M) = \Theta(1/\sqrt{MN}) = \Theta(1/N) \quad (\text{R3}\mu\text{-B3.6}) \\ &\quad (\sigma_Q^2 = 1/N; v_i \in \Theta(1) \text{ via (R3}\mu\text{-B2.4)} \Rightarrow \|v\|^2 \in \Theta(M); M \in \Theta(N)) \end{aligned}$$

$$\partial f_2 / \partial h_2^1 = (\partial f_2 / \partial h_2^1)_{\text{exp}} + (\partial f_2 / \partial h_2^1)_{\text{router}} \in \Theta(1/N), \quad (\text{R3}\mu\text{-B3.7})$$

dominated by  $A_5'' + A_6'' + E''$  along  $h_0^1$ ; same scale as  $\partial f_1 / \partial h_1^1$  at  $t=1$ , with  $E''$  now matching  $A_5'', A_6''$  via the alignment-coherence shift.

#### J.4.3. SUMMARY TABLE OF SIGNAL PROPAGATION FOR $\mu$ P IN REGIME III

Notation:  $\Delta_t W^\ell$  denotes the cumulative update  $W_t^\ell - W_0^\ell$ .

**Forward.**

Quantity	$t = 0$	$t = 1$	$t = 2$
$h_t^1 = h_0^1 + \Delta_t h^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^1 = W_0^1 x$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t h^1 = \Delta_t W^1 x$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
$\psi_t, \phi_t$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$h_t^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^{2,i} = W_0^{2,i} h_0^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{2,i} \Delta_t h^1$	0	$\Theta(1/\sqrt{N})$	$\Theta(1)$
effective: $\Delta_t W^{2,i} h_0^1$	0	$\Theta(1)$	$\Theta(1)$
cross: $\Delta_t W^{2,i} \Delta_t h^1$	0	$\Theta(1/N)$	$\Theta(1)$
$h_t^{3,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^{3,i} = W_0^{3,i} h_0^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t W^{3,i} h_0^{2,i}$	0	$\Theta(1)$	$\Theta(1)$
cross: $\Delta_t W^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1/N^{3/2})$	$\Theta(1)$
$h_t^3 = A_1 + A_{2,1} + A_{2,2} + A_3 + D$	$\Theta(1/\sqrt{N})$	$\Theta(1)$	$\Theta(1)$
$A_1 = (1/M) \sum_i \phi_{i,t} h_0^{3,i}$	$\Theta(1/\sqrt{N})$	$\Theta(1/\sqrt{N})$	$\Theta(1/\sqrt{N})$
$A_{2,1} = (1/M) \sum_i \phi_{i,t} W_0^{3,i} W_0^{2,i} \Delta_t h^1$	0	$\Theta(1/N)$	$\Theta(1/\sqrt{N})$
$A_{2,2} = (1/M) \sum_i \phi_{i,t} W_0^{3,i} \Delta_t W^{2,i} h_0^1$	0	$\Theta(1)$	$\Theta(1)$
$A_3 = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} h_0^{2,i}$	0	$\Theta(1)$	$\Theta(1)$
$D = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} \Delta_t h^{2,i}$	0	$\Theta(1/N)$	$\Theta(1)$
$f_t = W_t^4 h_t^3$	$\Theta(1/N)$	$\Theta(1)$	$\Theta(1)$
init: $W_0^4 h_0^3$	$\Theta(1/N)$	$\Theta(1/N)$	$\Theta(1/N)$
propagating: $W_0^4 \Delta_t h^3$	0	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t W^4 h_0^3$	0	$\Theta(1/N)$	$\Theta(1/N)$
cross: $\Delta_t W^4 \Delta_t h^3$	0	$\Theta(1/N)$	$\Theta(1)$

**Backward.**

Quantity	$t = 0$	$t = 1$	$t = 2$
$\partial f_t / \partial h_t^3 = (W_0^4)^\top + (\Delta_t W^4)^\top$	$\Theta(1/N)$	$\Theta(1/N)$	$\Theta(1/N)$
init: $(W_0^4)^\top$	$\Theta(1/N)$	$\Theta(1/N)$	$\Theta(1/N)$
update: $(\Delta_t W^4)^\top$	0	$\Theta(1/N^{3/2})$	$\Theta(1/N)$
$\partial f_t / \partial h_t^{3,i} = (\phi_{i,t}/M)(W_t^4)^\top$	$\Theta(1/(MN))$	$\Theta(1/(MN))$	$\Theta(1/(MN))$
init: $(\phi_{i,t}/M)(W_0^4)^\top$	$\Theta(1/(MN))$	$\Theta(1/(MN))$	$\Theta(1/(MN))$
update: $(\phi_{i,t}/M)(\Delta_t W^4)^\top$	0	$\Theta(1/(MN^{3/2}))$	$\Theta(1/(MN))$
$\partial f_t / \partial h_t^{2,i} = (\phi_{i,t}/M)(W_t^{3,i})^\top (W_t^4)^\top$	$\Theta(1/(MN))$	$\Theta(1/(MN))$	$\Theta(1/(MN))$
init-init: $(\phi_{i,t}/M)(W_0^{3,i})^\top (W_0^4)^\top$	$\Theta(1/(MN))$	$\Theta(1/(MN))$	$\Theta(1/(MN))$
init-update: $(\phi_{i,t}/M)(W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/(MN^{3/2}))$	$\Theta(1/(MN))$
update-init: $(\phi_{i,t}/M)(\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/(MN))$	$\Theta(1/(MN))$
update-update: $(\phi_{i,t}/M)(\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/(MN^2))$	$\Theta(1/(MN))$
$\partial f_t / \partial \phi_{i,t} = (1/M)\langle h_t^{3,i}, W_t^4 \rangle$	$\Theta(1/N^{3/2})$	$\Theta(1/M)$	$\Theta(1/M)$
init: $(1/M)\langle h_t^{3,i}, W_0^4 \rangle$	$\Theta(1/N^{3/2})$	$\Theta(1/M)$	$\Theta(1/M)$
update: $(1/M)\langle h_t^{3,i}, \Delta_t W^4 \rangle$	0	$\Theta(1/N^{3/2})$	$\Theta(1/M)$
$(\partial f_t / \partial h_t^1)_{\text{exp}}$	$\Theta(1/N^{3/2})$	$\Theta(1/N)$	$\Theta(1/N)$
$A_{4,1} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$	$\Theta(1/N^{3/2})$
$A_{4,2} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N^{3/2})$
$A_{5,1} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_{5,2} = (1/M) \sum_i \phi_{i,t} (W_0^{2,i})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$A_{6,1} = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (W_0^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_{6,2} = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (W_0^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$E_1 = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (\Delta_t W^{3,i})^\top (W_0^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$E_2 = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N^2)$	$\Theta(1/N)$
$(\partial f_t / \partial h_t^1)_{\text{router}} = (1/M) Q_t^\top v$	$\Theta(1/N^{3/2})$	$\Theta(1/\sqrt{MN})$	$\Theta(1/\sqrt{MN})$
init: $(1/M) Q_0^\top v$	$\Theta(1/N^{3/2})$	$\Theta(1/\sqrt{MN})$	$\Theta(1/\sqrt{MN})$
update: $(1/M) (\Delta_t Q)^\top v$	0	$\Theta(1/(MN))$	$\Theta(1/\sqrt{MN})$

**J.4.4. DERIVING MSSP IN REGIME III**

This subsection provides a self-contained, heuristic scaling analysis for the Maximally Scale-Stable Parameterization (MSSP, Definition 7) in Regime III.

**Auxiliary scaling lemma**

**Lemma 10 (Gram concentration)** *Let  $W \in \mathbb{R}^{m \times n}$  have i.i.d. centered entries of variance  $\sigma_W^2$ . Then*

$$\mathbb{E}[(W^\top W)_{jk}] = m\sigma_W^2 \delta_{jk}, \quad \mathbb{E}[(WW^\top)_{ij}] = n\sigma_W^2 \delta_{ij},$$

with off-diagonal fluctuations of order  $\sigma_W^2 \sqrt{m}$  (resp.  $\sigma_W^2 \sqrt{n}$ ). For any fixed  $v \in \mathbb{R}^n$  with  $\|v\|^2 \in \Theta(n)$  that is independent of  $W$ ,  $W^\top W v = m\sigma_W^2 v + r$  with  $r$  of coordinate scale  $\Theta(\sigma_W^2 \sqrt{mn})$ . The analogous statement holds for  $WW^\top$  acting on  $v' \in \mathbb{R}^m$  with  $\|v'\|^2 \in \Theta(m)$ .

**Proof [Sketch]**  $(W^\top W)_{jk} = \sum_a W_{aj}W_{ak}$  is a sum of  $m$  centered i.i.d. products with mean  $\sigma_W^2 \delta_{jk}$  and variance  $\sigma_W^4$  off-diagonal. The residual  $r = (W^\top W - m\sigma_W^2 I)v$  has entry-wise variance  $\sum_k m\sigma_W^4 v_k^2 = m\sigma_W^4 \|v\|^2$ , hence coordinate scale  $\sigma_W^2 \sqrt{mn}$ .  $\blacksquare$

**Specializations of Lemma 10 under proportional scaling.** Both Gram matrices of  $W_0^{2,1}$  and of  $W_0^{3,1}$  produce  $\Theta(1)$ -coordinate output when applied to  $\Theta(1)$ -coordinate input vectors (independent of or weakly correlated with the corresponding  $W$ ):

$$\begin{aligned} (W_0^{2,1})^\top W_0^{2,1} v &= (N_e/N) v + r, & W_0^{2,1} (W_0^{2,1})^\top v' &= v' + r', \\ (W_0^{3,1})^\top W_0^{3,1} u &= (N/N_e) u + r'', & W_0^{3,1} (W_0^{3,1})^\top u' &= u' + r''', \end{aligned}$$

where each residual is of the same  $\Theta(1)$  entry-wise order as the leading deterministic part. The Gram approximation is order-of-magnitude correct, with comparable random fluctuations.

### First forward pass

$$h_0^1 = W_0^1 x \in \Theta(1). \quad (\text{CLT}; \sigma_1^2 = 1/D, \|x\|^2 \in \Theta(D)) \quad (\text{R3s-F1.1})$$

$$\psi_0 = Q_0 h_0^1 \in \Theta(1). \quad (\text{CLT}; \sigma_Q^2 = 1/N, \|h_0^1\|^2 \in \Theta(N)) \quad (\text{R3s-F1.2})$$

$$\phi_0 = \sigma(\psi_0) \in \Theta(1). \quad (\text{sigmoid bounded; gating assumption}) \quad (\text{R3s-F1.3})$$

$$\begin{aligned} h_0^{2,i} = h_0^{2,1} &= W_0^{2,1} h_0^1 \in \Theta(1) \\ (\text{CLT}; \sigma_2^2 = 1/N, \|h_0^1\|^2 \in \Theta(N); h_0^{2,i} = h_0^{2,1} \text{ by MSSP}). & \quad (\text{R3s-F1.4}) \end{aligned}$$

$$\begin{aligned} h_0^{3,i} = h_0^{3,1} &= W_0^{3,1} h_0^{2,1} \in \Theta(1) \\ (\text{CLT}; \sigma_3^2 = 1/N_e, \|h_0^{2,1}\|^2 \in \Theta(N_e); h_0^{3,i} = h_0^{3,1} \text{ by MSSP}). & \quad (\text{R3s-F1.5}) \end{aligned}$$

$$\begin{aligned} h_0^3 &= (1/M) \sum_i \phi_{i,0} h_0^{3,1} = (\sum_k \phi_{k,0}/M) h_0^{3,1} \in \Theta(1) \\ ((\sum_k \phi_{k,0}/M) \in \Theta(1) \text{ by LLN}). & \quad (\text{R3s-F1.6}) \end{aligned}$$

$$f_0 = W_0^4 h_0^3 = 0. \quad (\text{standing assumption } W_0^4 = 0) \quad (\text{R3s-F1.7})$$

**First backward pass and step-1 updates** By the chain rule, every hidden gradient ( $\partial f_0/\partial h_0^3$ ,  $\partial f_0/\partial h_0^{3,i}$ ,  $\partial f_0/\partial h_0^{2,i}$ ,  $\partial f_0/\partial \phi_{i,0}$ ,  $\partial f_0/\partial h_0^1$  via either the expert or router path) carries a  $W_0^4 = 0$  factor and is therefore zero. Hence every non-readout weight is unchanged at  $t = 1$ :  $W_1^\ell = W_0^\ell$  for  $\ell \in \{1, Q, 2, 3\}$ .

### Step-1 parameter updates.

$$\Delta_1 W^4 = -\frac{\chi_0}{N} (h_0^3)^\top \in \Theta(1/N). \quad (\text{aligned along } (h_0^3)^\top) \quad (\text{R3s-U1.4})$$

**Second forward pass** Since the non-readout weights are unchanged at  $t = 1$ , every hidden activation is unchanged:  $h_1^\ell = h_0^\ell$  for  $\ell \in \{1, (2, i), (3, i), 3\}$ ,  $\psi_1 = \psi_0$ ,  $\phi_1 = \phi_0$ .

$$f_1 = W_1^4 h_0^3 = -\frac{\chi_0}{N} \|h_0^3\|^2 \in \Theta(1). \quad (\|h_0^3\|^2 \in \Theta(N) \text{ via (R3s-F1.6)}) \quad (\text{R3s-F2.7})$$

**Second backward pass and step-2 updates**

$$\partial f_1 / \partial h_1^3 = (W_1^4)^\top = (\Delta_1 W^4)^\top \in \Theta(1/N). \quad (\text{R3s-U1.4}); \quad W_0^4 = 0 \quad (\text{R3s-B2.1})$$

$$\partial f_1 / \partial h_1^{3,i} = (\phi_{i,0}/M)(W_1^4)^\top \in \Theta(1/(MN)). \quad (\phi_{i,0} \in \Theta(1)) \quad (\text{R3s-B2.2})$$

$$\begin{aligned} \partial f_1 / \partial h_1^{2,i} &= (\phi_{i,0}/M)(W_0^{3,1})^\top (W_1^4)^\top \in \Theta(1/(MN)) \\ ((W_0^{3,1})^\top (W_1^4)^\top &= -\frac{\chi_0}{MN} (\sum \phi)(W_0^{3,1})^\top W_0^{3,1} h_0^{2,1} \in \Theta(1/N) \text{ by Lem. 10;} \\ \phi_{i,0} &\in \Theta(1)). \end{aligned} \quad (\text{R3s-B2.3})$$

$$\begin{aligned} \partial f_1 / \partial \phi_{i,0} &= (1/M)(h_0^{3,1})^\top (W_1^4)^\top \in \Theta(1/M) \\ ((h_0^{3,1})^\top (W_1^4)^\top &= -\frac{\chi_0}{MN} (\sum \phi) \|h_0^{3,1}\|^2 \in \Theta(1) \text{ by LLN;} \\ \|h_0^{3,1}\|^2 &\in \Theta(N_e) \text{ via (R3s-F1.5); } N_e/N \in \Theta(1) \text{ in Regime III).} \end{aligned} \quad (\text{R3s-B2.4})$$

$$\begin{aligned} \left( \frac{\partial f_1}{\partial h_1^1} \right)_{\text{exp}} &= \left( \frac{1}{M} \sum \phi \right) (W_0^{2,1})^\top (W_0^{3,1})^\top (W_1^4)^\top \\ &\in \Theta\left(\frac{1}{M} \sum \phi\right) \cdot (W_0^{2,1})^\top h_0^{2,1} / N \text{ along } (W_0^{2,1})^\top h_0^{2,1} \\ &\quad (\text{Lem. 10 on } (W_0^{3,1})^\top h_0^3) \\ &= \Theta(1 \cdot 1/N) = \Theta(1/N) \\ &\quad (\text{LLN; } (W_0^{2,1})^\top h_0^{2,1} \in \Theta(1) \text{ by Lem. 10}). \end{aligned} \quad (\text{R3s-B2.5})$$

$$\begin{aligned} \left( \frac{\partial f_1}{\partial h_1^1} \right)_{\text{router}} &= \frac{1}{M} Q_0^\top v, \quad v_i := \dot{\phi}_{i,0} (h_0^{3,1})^\top (W_1^4)^\top \\ &\in \Theta(\sigma_Q \|v\|/M) \text{ random direction } (Q_0 \text{ indep. of } v; \text{ cross-layer CLT}) \\ &= \Theta(1/\sqrt{N} \cdot \sqrt{M} \cdot 1/M) = \Theta(1/\sqrt{MN}) \\ &\quad (\sigma_Q^2 = 1/N; v_i \in \Theta(1) \text{ via (R3s-B2.4)} \Rightarrow \|v\|^2 \in \Theta(M)). \end{aligned} \quad (\text{R3s-B2.6})$$

$$\partial f_1 / \partial h_1^1 = (\partial f_1 / \partial h_1^1)_{\text{exp}} + (\partial f_1 / \partial h_1^1)_{\text{router}} \in \Theta(1/N). \quad (\text{R3s-B2.7})$$

**Step-2 parameter updates.**

$$\begin{aligned} \Delta_2 W^4 &= -\frac{\chi_1}{N} (h_1^3)^\top \in \Theta(1/N), \quad W_2^4 = -\frac{\chi_0 + \chi_1}{N} (h_0^3)^\top \\ &\quad (h_1^3 = h_0^3). \end{aligned} \quad (\text{R3s-U2.4})$$

$$\Delta_2 W^{3,i} = -\frac{\chi_1 \phi_{i,0}}{N} h_0^3 (h_0^{2,1})^\top \in \Theta(1/N). \quad (\eta_3 = M; \text{R3s-B2.2}) \quad (\text{R3s-U2.3})$$

$$\Delta_2 W^{2,i} = -\frac{\chi_1 \phi_{i,0}}{N} (W_0^{3,1})^\top h_0^3 (h_0^1)^\top \in \Theta(1/N). \quad ((W_0^{3,1})^\top h_0^3 \in \Theta(1) \text{ via (R3s-B2.3)}) \quad (\text{R3s-U2.2})$$

$$\begin{aligned} \Delta_2 W^1 &= -\eta_1 \chi_1 (\partial f_1 / \partial h_1^1) x^\top \in \Theta(1) \\ &\quad (\eta_1 = N; \partial f_1 / \partial h_1^1 \in \Theta(1/N) \text{ via (R3s-B2.7)}). \end{aligned} \quad (\text{R3s-U2.1a})$$

$$\Delta_2 h^1 = \Delta_2 W^1 x \in \Theta(1) \text{ aligned along } \partial f_1 / \partial h_1^1. \quad (\text{R3s-U2.1b})$$

$$\begin{aligned} \Delta_2 Q &= -\eta_Q \chi_1 (\partial f_1 / \partial \phi) (h_1^1)^\top \in \Theta(1/M) = \Theta(1/N) \\ &\quad (\eta_Q = 1; \partial f_1 / \partial \phi \in \Theta(1/M) \text{ via (R3s-B2.4)}). \end{aligned} \quad (\text{R3s-U2.Q})$$

**Third forward pass** We compute activations at  $\theta^{(2)}$ , expanding each updated weight as  $W_2^\ell = W_0^\ell + \Delta_2 W^\ell$ .

$$h_2^1 = h_0^1 + \Delta_2 h^1 \in \Theta(1). \quad (\Delta_2 h^1 \in \Theta(1) \text{ via (R3s-U2.1b)}) \quad (\text{R3s-F3.1})$$

(R3s-F3.2)  $h_2^{2,i} = W_2^{2,i} h_2^1$ , four-piece decomposition:

$$h_2^{2,i} = h_0^{2,1} + W_0^{2,1} \Delta_2 h^1 + \Delta_2 W^{2,i} h_0^1 + \Delta_2 W^{2,i} \Delta_2 h^1. \quad (\text{R3s-F3.2})$$

$$\text{init: } W_0^{2,1} h_0^1 = h_0^{2,1} \in \Theta(1). \quad (\text{R3s-F1.4}) \quad (\text{R3s-F3.2a})$$

$$\text{prop: } W_0^{2,1} \Delta_2 h^1 \in \Theta(1)$$

$$(\text{op-norm of iid Gaussian } W_0^{2,1}; \Delta_2 h^1 \in \Theta(1) \text{ via (R3s-U2.1b)}). \quad (\text{R3s-F3.2b})$$

$$\text{eff: } \Delta_2 W^{2,i} h_0^1 = -\frac{\chi_1 \phi_{i,0}}{N} (W_0^{3,1})^\top h_0^3 \|h_0^1\|^2 \in \Theta(1)$$

$$((W_0^{3,1})^\top h_0^3 \in \Theta(1) \text{ via (R3s-B2.3)}; \|h_0^1\|^2 \in \Theta(N)). \quad (\text{R3s-F3.2c})$$

$$\text{cross: } \Delta_2 W^{2,i} \Delta_2 h^1 = -\frac{\chi_1 \phi_{i,0}}{N} (W_0^{3,1})^\top h_0^3 \langle h_0^1, \Delta_2 h^1 \rangle \in \Theta(1)$$

$$((W_0^{3,1})^\top h_0^3 \in \Theta(1) \text{ via (R3s-B2.3)}; \langle h_0^1, \Delta_2 h^1 \rangle \in \Theta(N) \text{ coherent via (R3s-U2.1b)}). \quad (\text{R3s-F3.2d})$$

(R3s-F3.3)  $h_2^{3,i} = W_2^{3,i} h_2^1$ , four-piece decomposition:

$$h_2^{3,i} = h_0^{3,1} + W_0^{3,1} \Delta_2 h^{2,i} + \Delta_2 W^{3,i} h_0^{2,1} + \Delta_2 W^{3,i} \Delta_2 h^{2,i}. \quad (\text{R3s-F3.3})$$

$$\text{init: } h_0^{3,1} \in \Theta(1). \quad (\text{R3s-F1.5}) \quad (\text{R3s-F3.3a})$$

$$\text{prop: } W_0^{3,1} \Delta_2 h^{2,i} \in \Theta(1)$$

$$(\text{op-norm of iid Gaussian } W_0^{3,1}; \Delta_2 h^{2,i} \in \Theta(1) \text{ via (R3s-F3.2)}). \quad (\text{R3s-F3.3b})$$

$$\text{eff: } \Delta_2 W^{3,i} h_0^{2,1} = -\frac{\chi_1 \phi_{i,0}}{N} h_0^3 \|h_0^{2,1}\|^2 \in \Theta(1)$$

$$(\|h_0^{2,1}\|^2 \in \Theta(N_e); h_0^3 \in \Theta(1) \text{ via (R3s-F1.6)}). \quad (\text{R3s-F3.3c})$$

$$\text{cross: } \Delta_2 W^{3,i} \Delta_2 h^{2,i} \text{ tracked as part of } D \text{ in (R3s-F3.4)}. \quad (\text{R3s-F3.3d})$$

(R3s-F3.4)  $h_2^3 = A_1 + A_2 + A_3 + D$ , where

$$\begin{aligned} A_1 &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,1} h_0^{2,1}, & A_2 &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} h_0^{2,1}, \\ A_3 &:= \frac{1}{M} \sum_i \phi_{i,2} W_0^{3,1} \Delta_2 h^{2,i}, & D &:= \frac{1}{M} \sum_i \phi_{i,2} \Delta_2 W^{3,i} \Delta_2 h^{2,i}. \end{aligned}$$

$$A_1 = (\sum \phi_{\cdot,2}/M) h_0^{3,1} \in \Theta(1). \quad (i\text{-independent matrix; LLN}) \quad (\text{R3s-F3.4a})$$

$$A_2 = -\frac{\chi_1 \|h_0^{2,1}\|^2}{N} \left( \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0} \right) h_0^3 \in \Theta(1) \text{ along } h_0^3 \\ (\|h_0^{2,1}\|^2/N \in \Theta(1); \text{ LLN}). \quad (\text{R3s-F3.4b})$$

$$A_3 \in \Theta(1) \\ (\text{three-piece expansion of } \Delta_2 h^{2,i} \text{ from (R3s-F3.2);} \\ \text{Lemma 10 on } h_0^3, h_0^1; \langle h_0^1, \Delta_2 h^1 \rangle \in \Theta(N); \text{ LLN}). \quad (\text{R3s-F3.4c})$$

$$D = -\frac{\chi_1}{N} h_0^3 \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0} \langle h_0^{2,1}, \Delta_2 h^{2,i} \rangle \\ \in \Theta\left(\frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0} \langle h_0^{2,1}, \Delta_2 h^{2,i} \rangle / N\right) \cdot h_0^3 \text{ along } h_0^3 \\ (\text{Lem. 10; LLN}) \\ = \Theta(N/N) = \Theta(1) \text{ along } h_0^3 \\ (\langle h_0^{2,1}, \Delta_2 h^{2,i} \rangle \in \Theta(N) \text{ coherent via (R3s-U2.1b); } h_0^3 \in \Theta(1)). \quad (\text{R3s-F3.4d})$$

$$h_2^3 = A_1 + A_2 + A_3 + D \in \Theta(1). \quad (\text{R3s-F3.4})$$

$$f_2 = W_2^4 h_2^3 = -\frac{\chi_0 + \chi_1}{N} \langle h_0^3, h_2^3 \rangle \in \Theta(1) \\ (\langle h_0^3, A_1 \rangle = (\sum \phi_{\cdot,2}/M) \langle h_0^3, h_0^{3,1} \rangle \in \Theta(N) \text{ coherent via (R3s-F3.4a);} \\ \langle h_0^3, A_2 + A_3 + D \rangle \in \Theta(N) \text{ via aligned components (R3s-F3.4b)–(R3s-F3.4d)}). \quad (\text{R3s-F3.5})$$

### Third backward pass

$$\partial f_2 / \partial h_2^3 = (W_2^4)^\top = -\frac{\chi_0 + \chi_1}{N} h_0^3 \in \Theta(1/N). \quad ((\text{R3s-U2.4}); h_0^3 \in \Theta(1)) \quad (\text{R3s-B3.1})$$

$$\partial f_2 / \partial h_2^{3,i} = (\phi_{i,2}/M) (W_2^4)^\top \in \Theta(1/(MN)). \quad (\phi_{i,2} \in \Theta(1)) \quad (\text{R3s-B3.2})$$

$$\partial f_2 / \partial h_2^{2,i} = (\phi_{i,2}/M) (W_2^{3,i})^\top (W_2^4)^\top \in \Theta(1/(MN)) \\ ((W_0^{3,1})^\top (W_2^4)^\top \in \Theta(1/N) \text{ via (R3s-B2.3);} \\ (\Delta_2 W^{3,i})^\top (W_2^4)^\top \in \Theta(1/N) \text{ along } h_0^{2,1} \text{ using } \|h_0^3\|^2 \in \Theta(N) \text{ via (R3s-F1.6);} \\ \phi_{i,2} \in \Theta(1)). \quad (\text{R3s-B3.3})$$

$$\partial f_2 / \partial \phi_{i,2} = (1/M) (h_2^{3,i})^\top (W_2^4)^\top \in \Theta(1/M) \\ ((h_2^{3,i})^\top (W_2^4)^\top = -\frac{\chi_0 + \chi_1}{N} \langle h_2^{3,i}, h_0^3 \rangle \in \Theta(1); \\ \langle h_2^{3,i}, h_0^3 \rangle \in \Theta(N) \text{ coherent via (R3s-F3.3)}). \quad (\text{R3s-B3.4})$$

(R3s-B3.5)  $(\partial f_2/\partial h_2^1)_{\text{exp}} = \tilde{A}_4 + \tilde{A}_5 + \tilde{A}_6 + \tilde{E}$ , where

$$\begin{aligned}\tilde{A}_4 &:= \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,1})^\top (W_0^{3,1})^\top (W_2^4)^\top, \\ \tilde{A}_5 &:= \frac{1}{M} \sum_i \phi_{i,2} (W_0^{2,1})^\top (\Delta_2 W^{3,i})^\top (W_2^4)^\top, \\ \tilde{A}_6 &:= \frac{1}{M} \sum_i \phi_{i,2} (\Delta_2 W^{2,i})^\top (W_0^{3,1})^\top (W_2^4)^\top, \\ \tilde{E} &:= \frac{1}{M} \sum_i \phi_{i,2} (\Delta_2 W^{2,i})^\top (\Delta_2 W^{3,i})^\top (W_2^4)^\top.\end{aligned}$$

$$\begin{aligned}\tilde{A}_4 &= -\frac{\chi_0 + \chi_1}{N} (\sum \phi_{\cdot,2}/M) (W_0^{2,1})^\top (W_0^{3,1})^\top h_0^3 \\ &\in \Theta((W_0^{2,1})^\top (W_0^{3,1})^\top h_0^3/N) \\ &\quad (i\text{-indep. matrix under MSSP; Lem. 10; LLN}) \\ &= \Theta(1 \cdot 1/N) = \Theta(1/N) \\ &\quad ((W_0^{2,1})^\top (W_0^{3,1})^\top h_0^3 \in \Theta(1) \text{ via (R3s-B2.3)}). \tag{R3s-B3.5a}\end{aligned}$$

$$\begin{aligned}\tilde{A}_5 &= \frac{\chi_1(\chi_0 + \chi_1) \|h_0^3\|^2}{N^2} \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0} (W_0^{2,1})^\top h_0^{2,1} \\ &\in \Theta(\|h_0^3\|^2 \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0}/N^2) \cdot h_0^1 \text{ along } h_0^1 \\ &\quad (\text{Lem. 10: } (W_0^{2,1})^\top h_0^{2,1} \rightarrow h_0^1) \\ &= \Theta(N \cdot 1 \cdot 1/N^2) = \Theta(1/N) \text{ along } h_0^1 \\ &\quad (\|h_0^3\|^2 \in \Theta(N) \text{ via (R3s-F1.6); LLN; } h_0^1 \in \Theta(1)). \tag{R3s-B3.5b}\end{aligned}$$

$$\begin{aligned}\tilde{A}_6 &= -\frac{\chi_1}{N} h_0^1 \frac{1}{M} \sum_i \phi_{i,0} S_i, \quad S_i := (h_0^3)^\top W_0^{3,1} (W_0^{3,1})^\top (W_2^4)^\top \\ &\in \Theta(\frac{1}{M} \sum_i \phi_{i,0} S_i/N) \cdot h_0^1 \text{ along } h_0^1 \\ &\quad (\text{Lem. 10: } S_i \in \Theta(1); \text{ LLN}) \\ &= \Theta(1 \cdot 1/N) = \Theta(1/N) \text{ along } h_0^1 \\ &\quad (S_i \in \Theta(1); h_0^1 \in \Theta(1)). \tag{R3s-B3.5c}\end{aligned}$$

$$\begin{aligned}\tilde{E} &= -\frac{\chi_1^2(\chi_0 + \chi_1) \|h_0^3\|^2}{N^3} h_0^1 \frac{1}{M} \sum_i \phi_{i,2} \phi_{i,0}^2 \langle h_0^3, h_0^{3,1} \rangle \\ &\in \Theta(\langle h_0^3, h_0^{3,1} \rangle \|h_0^3\|^2/N^3) \cdot h_0^1 \text{ along } h_0^1 \\ &\quad (\text{LLN; } \langle h_0^3, h_0^{3,1} \rangle \text{ coherent under MSSP}) \\ &= \Theta(N \cdot N/N^3) = \Theta(1/N) \text{ along } h_0^1 \\ &\quad (\langle h_0^3, h_0^{3,1} \rangle \in \Theta(N) \text{ coherent via (R3s-F1.6); } \|h_0^3\|^2 \in \Theta(N) \text{ via (R3s-F1.6)}). \tag{R3s-B3.5d}\end{aligned}$$

$$(\partial f_2/\partial h_2^1)_{\text{exp}} \in \Theta(1/N) \text{ along } h_0^1. \tag{R3s-B3.5}$$

$$\begin{aligned}
 (\partial f_2 / \partial h_2^1)_{\text{router}} &= \frac{1}{M} Q_0^\top v + \frac{1}{M} (\Delta_2 Q)^\top v, \quad v_i := \dot{\phi}_{i,2} (h_2^{3,i})^\top (W_2^4)^\top \\
 &\in \Theta(\sigma_Q \|v\| / M) \text{ random direction} + \Theta(\langle \partial f_1 / \partial \phi, v \rangle / M) h_1^1 \text{ along } h_1^1 \\
 &\quad (Q_0 \text{ piece via cross-layer CLT; } \Delta_2 Q \text{ piece coherent under MSSP}) \\
 &= \Theta(1/\sqrt{N} \cdot \sqrt{M} \cdot 1/M) = \Theta(1/\sqrt{MN}) \\
 &\quad (\sigma_Q^2 = 1/N; \|v\|^2 \in \Theta(M) \text{ via } v_i \in \Theta(1) \text{ from (R3s-B3.4);}) \\
 &\quad \langle \partial f_1 / \partial \phi, v \rangle \in \Theta(1) \text{ via MSSP coherence).} \tag{R3s-B3.6}
 \end{aligned}$$

$$\partial f_2 / \partial h_2^1 = (\partial f_2 / \partial h_2^1)_{\text{exp}} + (\partial f_2 / \partial h_2^1)_{\text{router}} \in \Theta(1/N). \tag{R3s-B3.7}$$

#### J.4.5. SUMMARY TABLE OF SIGNAL PROPAGATION FOR MSSP IN REGIME III

Notation:  $\Delta_t W^\ell$  denotes the cumulative update  $W_t^\ell - W_0^\ell$ .

##### Forward.

Quantity	$t = 0$	$t = 1$	$t = 2$
$h_t^1 = h_0^1 + \Delta_t h^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $h_0^1 = W_0^1 x$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t h^1 = \Delta_t W^1 x$	0	0	$\Theta(1)$
$\psi_t, \phi_t$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$h_t^{2,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $W_0^{2,1} h_0^1$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{2,1} \Delta_t h^1$	0	0	$\Theta(1)$
effective: $\Delta_t W^{2,i} h_0^1$	0	0	$\Theta(1)$
cross: $\Delta_t W^{2,i} \Delta_t h^1$	0	0	$\Theta(1)$
$h_t^{3,i}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
init: $W_0^{3,1} h_0^{2,1}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
propagating: $W_0^{3,1} \Delta_t h^{2,i}$	0	0	$\Theta(1)$
effective: $\Delta_t W^{3,i} h_0^{2,1}$	0	0	$\Theta(1)$
cross: $\Delta_t W^{3,i} \Delta_t h^{2,i}$	0	0	$\Theta(1)$
$h_t^3 = A_1 + A_2 + A_3 + D$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$A_1 = (1/M) \sum_i \phi_{i,t} W_0^{3,1} h_0^{2,1}$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$
$A_2 = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} h_0^{2,1}$	0	0	$\Theta(1)$
$A_3 = (1/M) \sum_i \phi_{i,t} W_0^{3,1} \Delta_t h^{2,i}$	0	0	$\Theta(1)$
$D = (1/M) \sum_i \phi_{i,t} \Delta_t W^{3,i} \Delta_t h^{2,i}$	0	0	$\Theta(1)$
$f_t = W_t^4 h_t^3$	0	$\Theta(1)$	$\Theta(1)$
effective: $\Delta_t W^4 h_0^3$	0	$\Theta(1)$	$\Theta(1)$
cross: $\Delta_t W^4 \Delta_t h^3$	0	0	$\Theta(1)$

##### Backward.

Quantity	$t = 0$	$t = 1$	$t = 2$
$\partial f_t / \partial h_t^3 = (\Delta_t W^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$\partial f_t / \partial h_t^{3,i} = (\phi_{i,t}/M)(\Delta_t W^4)^\top$	0	$\Theta(1/(MN))$	$\Theta(1/(MN))$
$\partial f_t / \partial h_t^{2,i} = (\phi_{i,t}/M)(W_t^{3,i})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/(MN))$	$\Theta(1/(MN))$
init: $(\phi_{i,t}/M)(W_0^{3,1})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/(MN))$	$\Theta(1/(MN))$
update: $(\phi_{i,t}/M)(\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	0	$\Theta(1/(MN))$
$\partial f_t / \partial \phi_{i,t} = (1/M)\langle h_t^{3,i}, \Delta_t W^4 \rangle$	0	$\Theta(1/M)$	$\Theta(1/M)$
$(\partial f_t / \partial h_t^1)_{\text{exp}}$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_4 = (1/M) \sum_i \phi_{i,t} (W_0^{2,1})^\top (W_0^{3,1})^\top (\Delta_t W^4)^\top$	0	$\Theta(1/N)$	$\Theta(1/N)$
$A_5 = (1/M) \sum_i \phi_{i,t} (W_0^{2,1})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	0	$\Theta(1/N)$
$A_6 = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (W_0^{3,1})^\top (\Delta_t W^4)^\top$	0	0	$\Theta(1/N)$
$E = (1/M) \sum_i \phi_{i,t} (\Delta_t W^{2,i})^\top (\Delta_t W^{3,i})^\top (\Delta_t W^4)^\top$	0	0	$\Theta(1/N)$
$(\partial f_t / \partial h_t^1)_{\text{router}} = (1/M) Q_t^\top v$	0	$\Theta(1/\sqrt{MN})$	$\Theta(1/\sqrt{MN})$
init: $(1/M) Q_0^\top v$	0	$\Theta(1/\sqrt{MN})$	$\Theta(1/\sqrt{MN})$
update: $(1/M) (\Delta_t Q)^\top v$	0	0	$\Theta(1/\sqrt{MN})$

## K. DMFT Analysis for Regime I

### K.1. Setup

We use the following architecture:

$$\begin{aligned}
 h_\mu^1 &= W^1 x_\mu \\
 \psi_\mu &= \frac{1}{N} Q \sigma(h_\mu^1) \\
 \phi_\mu &= \text{softmax}(\psi_\mu) \\
 h_{\mu,i}^2 &= \frac{1}{\sqrt{N}} W_i^2 \sigma(h_\mu^1), \\
 h_{\mu,i}^{\ell+1} &= \frac{1}{\sqrt{N}} W_i^{\ell+1} \sigma(h_{\mu,i}^\ell), \quad \ell \in \{2, 3, \dots, L-1\} \\
 h_\mu^L &= \sum_{i=1}^M \phi_\mu^i \cdot h_{\mu,i}^L \\
 h_\mu^{L+1} &= \frac{1}{\sqrt{N}} w^{L+1 \top} h_\mu^L \\
 f_\mu &= \frac{1}{\gamma} h_\mu^{L+1} \\
 w_\alpha^{L+1}(0) &\sim \mathcal{N}(0, 1) \\
 W_{\alpha\beta,i}^\ell(0) &\sim \mathcal{N}(0, 1) \text{ for } \ell \in \{2, 3, \dots, L\}, i \in \{1, 2, \dots, M\} \\
 W_{\alpha\beta}^1(0) &\sim \mathcal{N}(0, 1) \\
 Q_{\alpha\beta}(0) &\sim \mathcal{N}(0, 1) \\
 \eta &= \eta_0 \gamma^2 \\
 \gamma &= \gamma_0 \sqrt{N} \\
 \eta_0, \gamma_0 &\sim O(1)
 \end{aligned} \tag{K.1}$$

$x_\mu \in \mathbb{R}^D$  is the embedding vector of input  $\mu \in \{1, 2, \dots, P\}$ ,  $\psi_\mu \in \mathbb{R}^M$  gives the router's preferences for each of the  $M$  experts,  $\phi_\mu^i$  denotes the  $i^{\text{th}}$  component of  $\phi_\mu$ , for  $i \in \{1, 2, \dots, M\}$ ,  $\sigma$  is an element-wise nonlinearity,  $Q \in \mathbb{R}^{M \times N}$ ,  $w^4 \in \mathbb{R}^N$ ,  $W_i^2, W_i^3 \in \mathbb{R}^{N \times N}$ ,  $W^1 \in \mathbb{R}^{N \times D}$ , the preactivations  $h$  are in  $\mathbb{R}^N$ , except for  $h_\mu^4$ , which is a scalar. We define  $\text{softmax} : \mathbb{R}^p \rightarrow (0, 1)^p$ ,  $p > 1$  as the function which takes a tuple  $\mathbf{z} = (z_1, z_2, \dots, z_p) \in \mathbb{R}^p$ , and computes each component of the vector  $\text{softmax}(\mathbf{z}) \in (0, 1)^p$  with  $[\text{softmax}(\mathbf{z})]_i = \frac{e^{\beta z_i}}{\sum_{j=1}^p e^{\beta z_j}}$ , where  $T = \frac{1}{\beta}$  is the temperature of the softmax. We will be interested in the low temperature limit  $\beta \rightarrow \infty$ , where softmax implements Top-1.

We denote for brevity  $\boldsymbol{\theta} = \text{Vec}\{W^1, W_i^\ell, w^{L+1}, Q\}_{i,\ell}$ , for  $i \in \{1, 2, \dots, M\}$  and  $\ell \in \{2, 3, \dots, L\}$ , and define the gradients

$$\begin{aligned}
 g_\mu^l &= \sqrt{N} \frac{\partial h_\mu^{L+1}}{\partial h_\mu^l}, \quad l \in \text{Vec}\{1, (2, i), (3, i), \dots, (L, i), L, L+1\}_i \\
 g_\mu^\phi &= \frac{1}{\sqrt{N}} \frac{\partial h_\mu^{L+1}}{\partial \phi_\mu}
 \end{aligned} \tag{K.2}$$

$g^1$  has two terms, as we must track the gradient flowing both through the router, and through the experts. Using K.4, the components of  $g_\mu^\phi$  are

$$(g_\mu^\phi)_i = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^{L+1}}{\partial \phi_\mu^i} = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^{L+1}}{\partial h_\mu^L} \cdot \frac{\partial h_\mu^L}{\partial \phi_\mu^i} = \frac{1}{N} \mathbf{g}^L \cdot \mathbf{h}_{\mu,i}^L = O(1). \tag{K.3}$$

In addition, we define the following *pre-gradient fields* which are defined as the gradients with respect to the activations instead of the pre-activations. For linear layers, they clearly reduce to the gradient fields.

$$\begin{aligned}
 g_\mu^L &= z_\mu^L = w^{L+1} \\
 g_\mu^{L,i} &= z_{\mu,i}^L = g_\mu^L \phi_\mu^i \\
 g_{\mu,i}^\ell &= \dot{\sigma}(h_{\mu,i}^\ell) \odot z_{\mu,i}^\ell, \quad z_{\mu,i}^\ell = \frac{1}{\sqrt{N}} W_i^{\ell+1\top} g_{\mu,i}^{\ell+1}, \quad \ell \in \{2, 3, 4, \dots, L-1\}, i \in \{1, 2, \dots, M\} \\
 (g_\mu^\phi)_i &= (z_\mu^\phi)_i = \frac{1}{N} g_\mu^L \cdot h_{\mu,i}^L, \quad i \in \{1, 2, \dots, M\} \\
 g_\mu^1 &= \dot{\sigma}(h_\mu^1) \odot z_\mu^1, \quad z_\mu^1 = \sum_{i=1}^M \frac{1}{\sqrt{N}} W_i^{2\top} g_{\mu,i}^2 + Q^\top g_\mu^\phi
 \end{aligned} \tag{K.4}$$

For convenience, we define  $\tilde{z}_\mu^\phi := Q^\top g_\mu^\phi$  and  $\tilde{z}_{\mu,i}^1 := \frac{1}{\sqrt{N}} W_i^{2\top} g_{\mu,i}^2$ . Then we have  $z_\mu^1 = \sum_{i=1}^M \tilde{z}_{\mu,i}^1$  and  $z_\mu^\phi = \tilde{z}_\mu^\phi$ .

We define also the following feature and gradient kernels (The gradient kernels may be considered to be defined in terms of the z-objects via K.4) :

$$\begin{aligned}
 \Phi_{\mu\nu}^0 &:= x_\mu \cdot x_\nu, \\
 \Phi_{\mu\nu}^1(s, t) &:= \frac{1}{N} \sigma(h_\mu^1(s)) \cdot \sigma(h_\nu^1(t)), \\
 \Phi_{\mu\nu,i}^\ell(s, t) &:= \frac{1}{N} \sigma(h_{\mu,i}^\ell(s)) \cdot \sigma(h_{\nu,i}^\ell(t)) \quad \ell \in \{2, 3, \dots, L-1\}, \\
 \Phi_{\mu\nu}^L(s, t) &:= \frac{1}{N} h_\mu^L(s) \cdot h_\nu^L(t) \\
 G_{\mu\nu}^1(s, t) &:= \frac{1}{N} g_\mu^1(s) \cdot g_\nu^1(t), \\
 G_{\mu\nu,i}^\ell(s, t) &:= \frac{1}{N} g_{\mu,i}^\ell(s) \cdot g_{\nu,i}^\ell(t) \quad \ell \in \{2, 3, \dots, L\} \\
 G_{\mu\nu}^\phi(s, t) &:= g_\mu^\phi(s) g_\nu^\phi(t)
 \end{aligned} \tag{K.5}$$

All of these kernels are normalised to be order unity in  $N$ .

## K.2. Dynamics

We train using gradient flow with learning rate  $\eta$  on the loss function

$$\mathcal{L} = \frac{1}{P} \sum_{\mu=1}^P \ell(f_{\mu}, y_{\mu}) \quad (\text{K.6})$$

This induces the following dynamics:

$$\begin{aligned} \frac{d\theta}{dt} &= \frac{\eta_0 \gamma}{P} \sum_{\mu} \Delta_{\mu} \frac{\partial h_{\mu}^{L+1}}{\partial \theta} \\ \Delta_{\mu} &= -\frac{\partial \mathcal{L}}{\partial f_{\mu}} \end{aligned} \quad (\text{K.7})$$

In particular, for a MSE loss  $\mathcal{L} = \frac{1}{P} \sum_{\nu} (y_{\nu} - f_{\nu})^2$ , we have  $\Delta_{\nu} = 2(y_{\nu} - f_{\nu})$ . The logits update as

$$\frac{df_{\mu}(t)}{dt} = \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{d\theta}{dt} = \frac{\eta}{P} \sum_{\alpha} \Delta_{\alpha} K_{\mu\alpha}^{\text{NTK}}(t, t) \quad (\text{K.8})$$

$$\text{where } K_{\mu\alpha}^{\text{NTK}}(t, s) \equiv \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{\partial f_{\alpha}(s)}{\partial \theta} \text{ is the neural tangent kernel.} \quad (\text{K.9})$$

We can derive from this definition and the repeated use of chain rule that

$$\gamma^2 K_{\mu\nu}^{\text{NTK}}(t, t) = G_{\mu\nu}^1 \Phi_{\mu\nu}^0 + \sum_{i=1}^M [G_{\mu\nu,i}^2 \Phi_{\mu\nu}^1 + \sum_{\ell=3}^L G_{\mu\nu,i}^{\ell} \Phi_{\mu\nu,i}^{\ell-1}] + \Phi_{\mu\nu}^L + G_{\mu\nu}^{\phi} \Phi_{\mu\nu}^1 \quad (\text{K.10})$$

Hence

$$\begin{aligned} \frac{df_{\mu}(t)}{dt} &= \frac{\eta_0}{P} \sum_{\nu=1}^P \left[ G_{\mu\nu}^1(t, t) \Phi_{\mu\nu}^0 + \sum_{i=1}^M [G_{\mu\nu,i}^2(t, t) \Phi_{\mu\nu}^1(t, t) + \sum_{\ell=3}^L G_{\mu\nu,i}^{\ell}(t, t) \Phi_{\mu\nu,i}^{\ell-1}(t, t)] + \Phi_{\mu\nu}^L(t, t) \right. \\ &\quad \left. + G_{\mu\nu}^{\phi}(t, t) \Phi_{\mu\nu}^1(t, t) \right] \Delta_{\nu}(t) \end{aligned} \quad (\text{K.11})$$

Using the learning dynamics [K.7](#), we find for the weight matrices

$$\frac{dW^1(t)}{dt} = \frac{\eta_0 \gamma_0}{P} \sum_{\mu} \Delta_{\mu} g_{\mu}^1(t) x_{\mu}^{\top} \quad (\text{K.12})$$

$$\frac{dW_i^2(t)}{dt} = \frac{\eta_0 \gamma_0}{\sqrt{N} P} \sum_{\mu} \Delta_{\mu} g_{\mu,i}^2 \sigma(h_{\mu}^1)^{\top} \quad (\text{K.13})$$

$$\frac{dW_i^{\ell}(t)}{dt} = \frac{\eta_0 \gamma_0}{\sqrt{N} P} \sum_{\mu} \Delta_{\mu} g_{\mu,i}^{\ell} \sigma(h_{\mu,i}^{\ell-1})^{\top}, \quad \ell \in \{3, 4, \dots, L\} \quad (\text{K.14})$$

$$\frac{dw^{L+1}(t)}{dt} = \frac{\eta_0\gamma_0\sqrt{N}}{P} \sum_{\mu} \Delta_{\mu} \frac{\partial h_{\mu}^{L+1}}{\partial w^{L+1}} = \frac{\eta_0\gamma_0}{P} \sqrt{N} \sum_{\mu} \Delta_{\mu} \frac{1}{\sqrt{N}} h_{\mu}^L = \frac{\gamma_0\eta_0}{P} \sum_{\mu} \Delta_{\mu} h_{\mu}^L \quad (\text{K.15})$$

$$\underbrace{\frac{dQ(t)}{dt}}_{\in \mathbb{R}^{M \times N}} = \frac{\eta_0\gamma_0}{P} \sum_{\mu=1}^P \Delta_{\mu}(t) \underbrace{g_{\mu}^{\phi}(t)}_{\in \mathbb{R}^{M \times 1}} \underbrace{[\sigma(h_{\mu}^1(t))]}_{\in \mathbb{R}^{1 \times N}}^{\top} \quad (\text{K.16})$$

We can now integrate these expressions and use the definitions of the pre-activations from [K.1](#) to obtain:

$$\begin{aligned} W^1(t) &= W^1(0) + \int_0^t ds \frac{\eta_0\gamma_0}{P} \sum_{\mu} \Delta_{\mu} g_{\mu}^1(s) \cdot x_{\mu} \\ \implies h_{\mu}^1(t) &= W^1(0)x_{\mu} + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_{\nu} g_{\nu}^1(s) \Phi_{\mu\nu}^0 \\ W_i^2(t) &= W_i^2(0) + \int_0^t ds \frac{\eta_0\gamma_0}{\sqrt{N}P} \sum_{\mu} \Delta_{\mu} g_{\mu,i}^2(s) \sigma(h_{\mu}^1(s))^{\top} \\ \implies h_{\mu,i}^2(t) &= \frac{1}{\sqrt{N}} W_i^2(0) \sigma(h_{\mu}^1(t)) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_{\nu} g_{\nu,i}^2(s) \Phi_{\mu\nu}^1(s, t) \\ W_i^{\ell}(t) &= W_i^{\ell}(0) + \frac{\eta_0\gamma_0}{\sqrt{N}P} \int_0^t ds \sum_{\mu} \Delta_{\mu} g_{\mu,i}^{\ell}(s) \sigma(h_{\mu}^{\ell-1})^{\top} \\ \implies h_{\mu,i}^{\ell}(t) &= \frac{1}{\sqrt{N}} W_i^{\ell}(0) \sigma(h_{\mu,i}^{\ell-1}(t)) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_{\nu} g_{\nu,i}^{\ell}(s) \Phi_{\mu\nu,i}^{\ell-1}(s, t), \quad \ell \in \{3, 4, \dots, L\} \\ Q(t) &= Q(0) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_{\mu=1}^P \Delta_{\mu}(s) g_{\mu}^{\phi}(s) (h_{\mu}^1(s))^{\top} \\ \implies \psi_{\mu}(t) &= \frac{1}{N} Q(0) \sigma(h_{\mu}^1(t)) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_{\nu=1}^P \Delta_{\nu}(s) g_{\nu}^{\phi}(s) \Phi_{\mu\nu}^1(s, t) \end{aligned} \quad (\text{K.17})$$

This is useful, because we hope to average over the weights at initialization to obtain expressions in terms of only features. These expressions [K.17](#) motivate the definition of the following fields:

$$\begin{aligned} \chi_{\mu}^1 &= W^1(0)x_{\mu} \in \mathbb{R}^N \\ \chi_{\mu,i}^2 &= \frac{1}{\sqrt{N}} W_i^2(0) \sigma(h_{\mu}^1(t)) \in \mathbb{R}^N \\ \chi_{\mu,i}^{\ell} &= \frac{1}{\sqrt{N}} W_i^{\ell}(0) \sigma(h_{\mu,i}^{\ell-1}(t)) \in \mathbb{R}^N \quad \ell \in \{3, 4, \dots, L\} \\ \chi_{\mu}^{\phi} &= \frac{1}{N} Q(0) \sigma(h_{\mu}^1(t)) \in \mathbb{R}^M \end{aligned} \quad (\text{K.18})$$

We repeat the analysis of K.17 for the gradient fields which are defined as  $z_\mu^\ell = \frac{1}{\sqrt{N}} W^{\ell+1}(t)^\top g_\mu^{\ell+1}(t)$ , where  $\ell$  indexes the layer, and may involve an expert index too, and  $z_\mu^\phi(t) = Q(t)^\top g_\mu^\phi(t)$ : In particular, we have  $\forall i \in \{1, 2, \dots, M\}$

$$\begin{aligned}
 z_\mu^\phi(t) &= \xi_\mu^\phi(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu=1}^P \Delta_\nu(s) \sigma(h_\nu^1(s)) G_{\mu\nu}^\phi(t, s) \\
 z_{\mu,i}^1(t) &= \xi_{\mu,i}^1(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu \sigma(h_\nu^1(s)) G_{\mu\nu,i}^2(s, t) \\
 z_{\mu,i}^\ell(t) &= \xi_{\mu,i}^\ell(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu \sigma(h_{\nu,i}^\ell(s)) G_{\mu\nu,i}^{\ell+1}(s, t) \quad \ell \in \{2, 3, \dots, L-1\} \\
 z_{\mu,i}^L &= \phi_{\mu,i}^L z_\mu^L \\
 z_\mu^L(t) &= w_\mu^{L+1}(t) = \xi_\mu^L(t) + \frac{\gamma_0 \eta_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu h_\nu^L,
 \end{aligned} \tag{K.19}$$

where we define the following fields:

$$\begin{aligned}
 \xi_\mu^\phi(t) &= Q(0)^\top g_\mu^\phi(t) \\
 \xi_{\mu,i}^\ell(t) &= \frac{1}{\sqrt{N}} W_i^{\ell+1}(0)^\top g_{\mu,i}^{\ell+1}(t) \quad \ell \in \{1, 2, \dots, L-1\} \\
 \xi_\mu^L &= w^{L+1}(0) \sim \mathcal{N}(0, 1)
 \end{aligned} \tag{K.20}$$

Note that  $\xi_\mu^L$  is in fact time-independent, and equal to the initialization of  $w^{L+1}$ . We know that this object has entries distributed as  $\mathcal{N}(0, 1)$ , and we will also recover this from the DMFT for completeness.

Note also that the final gradient is  $g_\mu^4 = \sqrt{N}$ , so we can think of the final equation in K.19 as involving the identity kernel.

Finally, we write

$$\begin{aligned}
 [g^\phi]_i &= \frac{1}{N} \sum_{n=1}^N g_n^L [h_i^L]_n = \frac{1}{N^{3/2}} \sum_{n,m=1}^N g_n^L [W_i^L]_{nm} \sigma([h_i^{L-1}]_m) \\
 [g_\mu^\phi(t)]_i &= \xi_{\mu,i}^{g^\phi}(t) \\
 &+ \frac{1}{N^2} \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu(s) \sum_{n,m=1}^N g_{\mu n}^L(t) g_{\nu n}^L(s) \phi_\nu^i(s) \sigma([h_{\mu i}^{L-1}(t)]_m) \sigma([h_{\nu i}^{L-1}(s)]_m) \\
 &= \xi_{\mu,i}^{g^\phi}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu(s) G_{\mu\nu}^L(t, s) \Phi_{\mu\nu}^{L-1}(t, s) \phi_\nu^i(s)
 \end{aligned} \tag{K.21}$$

Where we define the field

$$\xi_{\mu,i}^{g^\phi}(t) = \frac{1}{N^{3/2}} g_\mu^{L\top}(t) W_i^L(0) \sigma(h_{\mu i}^{L-1}(t)). \tag{K.22}$$

### K.3. Mean field theory

If we can show that the kernels defined above concentrate, then our DMFT will be determined by the stochastic fields  $\mathcal{F} = \{\chi_\mu^1, \chi_{\mu,i}^2, \chi_{\mu,i}^3, \dots, \chi_{\mu,i}^L, \chi_\mu^\phi, \xi_{\mu,i}^1, \xi_{\mu,i}^2, \dots, \xi_{\mu,i}^{L-1}, \xi_\mu^L, \xi_\mu^\phi\}_{i \in \{1,2,\dots,M\}}$ . Therefore to construct the DMFT, we must compute the moment generating functional of these stochastic processes.

For brevity, we define, alongside the set  $\mathcal{F}$  of fields, the set  $\mathcal{S}$  of sources, which contains a corresponding source  $j$  (or  $v$ ) for every  $\chi$  (or  $\xi$ ) in  $\mathcal{F}$ . So  $\mathcal{S} = \{j_\mu^1, j_{\mu,i}^2, j_{\mu,i}^3, \dots, j_{\mu,i}^L, j_\mu^\phi, v_{\mu,i}^1, v_{\mu,i}^2, \dots, v_{\mu,i}^{L-1}, v_\mu^L, v_\mu^\phi\}_{i \in \{1,2,\dots,M\}}$ . This notation is constructed so that we can write  $\mathcal{F} \cdot \mathcal{S}$  to denote  $j_\mu^1(t) \cdot \chi_\mu^1(t) + j_{\mu,i}^2(t) \cdot \chi_{\mu,i}^2(t) + \dots + v_\mu^\phi(t) \cdot \xi_\mu^\phi(t)$ . Note that  $v_\mu^\phi \in \mathbb{R}^M$ , but all the other dummy variables are in  $\mathbb{R}^N$ . We also define  $\theta_0 = \text{Vec}\{W^1(0), W_i^2(0), \dots, W_i^L(0), w^{L+1}(0), Q(0)\}_{i \in \{1,\dots,M\}}$ .

The MGF with which we are concerned is therefore

$$Z[\{j, v\}] \left\langle \exp \left( \sum_\mu \int_0^\infty dt \mathcal{F} \cdot \mathcal{S} \right) \right\rangle_{\theta_0} \quad (\text{K.23})$$

$v_\mu^\phi \in \mathbb{R}^M$ , but all the other dummy variables are in  $\mathbb{R}^N$ .

In order to integrate, we enforce definitions of our random fields using resolutions of the identity to obtain:

$$\begin{aligned} Z = & \left\langle \frac{1}{(2\pi)^{8N}} \prod_{\mu=1}^P \prod_{i=1}^M \int_0^\infty dt \left\{ \exp \left( \sum_\mu \int_0^\infty dt \mathcal{F} \cdot \mathcal{S} \right) + \right. \\ & \int d\mathcal{F} \exp \left[ i \left( \hat{\chi}_\mu^1 \cdot (\chi_\mu^1 - W^1(0)x_\mu) + \hat{\chi}_{\mu,i}^2 \cdot \left( \chi_{\mu,i}^2 - \frac{1}{\sqrt{N}} W_i^2(0) \sigma(h_\mu^1(t)) \right) \right. \right. \\ & + \sum_{\ell=3}^L \hat{\chi}_{\mu,i}^\ell \cdot \left( \chi_{\mu,i}^\ell - \frac{1}{\sqrt{N}} W_i^\ell(0) \sigma(h_{\mu,i}^{\ell-1}(t)) \right) + \sum_{\ell=1}^{L-1} \hat{\xi}_{\mu,i}^\ell \cdot \left( \xi_{\mu,i}^\ell - \frac{1}{\sqrt{N}} W_i^{\ell+1}(0)^\top g_{\mu,i}^{\ell+1}(t) \right) \\ & + \hat{\xi}_\mu^L \cdot \left( \xi_\mu^L - w^{L+1}(0)^\top \right) + \hat{\chi}_\mu^\phi \cdot \left( \chi_\mu^\phi - \frac{1}{N} Q(0) \sigma(h_\mu^1(t)) \right) + \hat{\xi}_\mu^\phi \cdot \left( \xi_\mu^\phi - Q(0)^\top g_\mu^\phi(t) \right) \\ & \left. \left. \left. + \hat{\xi}^{g^\phi} \cdot \left( \xi_\mu^{g^\phi} - \frac{1}{N^{3/2}} g_\mu^{L\top} W_i^L(0) \sigma(h_{\mu,i}^{L-1}) \right) \right) \right] \right\} \right\rangle_{\theta_0} \quad (\text{K.24}) \end{aligned}$$

We then integrate using the fact that for a Gaussian variable  $z \sim \mathcal{N}(0, \sigma^2)$ ,  $\langle e^{-ia \cdot z} \rangle_z \propto e^{-\frac{1}{2} \sigma^2 |a|^2}$  to obtain

$$\begin{aligned}
 Z \propto & \prod_{\mu=1}^P \prod_{\nu=1}^P \prod_{i=1}^M \int_0^\infty dt \int dt \int d\mathcal{F} \\
 & \times \exp \left\{ -\frac{1}{2} \sum_{\mu,\nu=1}^P \int_0^\infty dt \int_0^\infty ds \left[ \hat{\chi}_\mu^1(t) \cdot \hat{\chi}_\nu^1(s) \Phi_{\mu\nu}^0(t, s) \right. \right. \\
 & + (\hat{\chi}_{\mu,i}^2(t) \cdot \hat{\chi}_{\nu,i}^2(s) + \frac{1}{N} \hat{\chi}_\mu^\phi(t) \cdot \hat{\chi}_\nu^\phi(s)) \Phi_{\mu\nu}^1(t, s) + \sum_{\ell=3}^L \hat{\chi}_{\mu,i}^\ell(s) \cdot \hat{\chi}_{\nu,i}^\ell(t) \Phi_{\mu\nu,i}^{\ell-1}(s, t) \\
 & \left. + \sum_{\ell=2}^L \hat{\xi}_{\mu,i}^\ell(s) \cdot \hat{\xi}_{\nu,i}^\ell(t) G_{\mu\nu,i}^\ell(s, t) + \hat{\xi}_\mu^{L+1}(s) \cdot \hat{\xi}_\nu^{L+1}(t) + \hat{\xi}_\mu^\phi(s) \cdot \hat{\xi}_\nu^\phi(t) G_{\mu\nu}^\phi(t, s) \right] \\
 & - i \sum_{\mu,\nu=1}^P \int_0^\infty dt \int_0^\infty ds \left[ \sum_{\ell=2}^L \hat{\chi}_{\mu,i}^\ell(t) \cdot g_{\nu,i}^\ell(t) A_{\mu\nu,i}^\ell(t, t) + \hat{\chi}_\mu^\phi(t) \cdot g_\nu^\phi(t) A_{\mu\nu}^\phi(t, t) \right] \\
 & + \sum_{\mu=1}^P \int_0^\infty dt \int_0^\infty ds \left[ \sum_{\ell=2}^{L-1} \left( (j_{\mu,i}^\ell(t) + i \hat{\chi}_{\mu,i}^\ell(t)) \cdot \chi_{\mu,i}^\ell(t) + (v_{\mu,i}^\ell(t) + i \hat{\xi}_{\mu,i}^\ell) \cdot \xi_{\mu,i}^\ell(t) \right) \right. \\
 & + (j_\mu^1(t) + i \hat{\chi}_\mu^1(t)) \cdot \chi_\mu^1(t) + (v_\mu^L + i \hat{\xi}_\mu^L(t)) \cdot \xi_\mu^L(t) \\
 & + (j_{\mu,i}^L(t) + i \hat{\chi}_{\mu,i}^L(t)) \cdot \chi_{\mu,i}^L(t) + (v_{\mu,i}^1(t) + i \hat{\xi}_{\mu,i}^1(t)) \cdot \xi_{\mu,i}^1(t) + (v_\mu^\phi + i \hat{\xi}_\mu^\phi(t)) \cdot \xi_\mu^\phi(t) \\
 & \left. + (j_\mu^\phi + i \hat{\chi}_\mu^\phi(t)) \cdot \chi_\mu^\phi(t) + (v_\mu^{g^\phi} + i \hat{\xi}_\mu^{g^\phi}(t)) \cdot \xi_\mu^{g^\phi}(t) \right] \left. \right\}. \tag{K.25}
 \end{aligned}$$

Where we have defined the kernels

$$A_{\mu\nu,i}^\ell(s, t) = \begin{cases} -\frac{i}{N} \hat{\xi}_{\mu,i}^1(s) \cdot \sigma(h_\nu^1(t)) & \ell = 2 \\ -\frac{i}{N} \hat{\xi}_{\mu,i}^{\ell-1}(s) \cdot \sigma(h_{\nu,i}^{\ell-1}(t)) & \ell \in \{3, 4, \dots, L\} \end{cases} \tag{K.26}$$

$$A_{\mu\nu}^\phi(s, t) = -\frac{i}{N} \hat{\xi}_\mu^\phi(s) \cdot \sigma(h_\nu^1(t)) \tag{K.27}$$

There are also three terms arising from the coupling of  $\xi^{g^\phi}$  to itself and to the forwards and backwards fields at layer L. These vanish as  $\frac{1}{N}$ , so have been excluded from K.25 for brevity.

The  $\Phi$  and  $G$  kernels are required for computing the evolution of the function output, via the NTK. The kernels  $A$  are not involved in the NTK, but rather arise from the coupling of the fields across a single layer's initial weight matrix. We now enforce the definitions of all three types of kernels using integral representations of Dirac delta-functions. In particular, for each pair  $\mu, \nu$  of samples, each expert  $i$  and each pair  $t, s$  of times, we multiply by

$$1 = \int \frac{d\Phi_{\mu\nu}^0(t, s) d\hat{\Phi}_{\mu\nu}^0(t, s)}{2\pi i} \exp \left[ \hat{\Phi}_{\mu\nu}^0(t, s) (\Phi_{\mu\nu}^0(t, s) - x_\mu \cdot x_\nu) \right] \tag{K.28}$$

$$1 = \int \frac{d\Phi_{\mu\nu}^1(t, s) d\hat{\Phi}_{\mu\nu}^1(t, s)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^1(t, s) (N\Phi_{\mu\nu}^1(t, s) - \sigma(h_\mu^1(t)) \cdot \sigma(h_\nu^1(s))) \right] \tag{K.29}$$

$$1 = \int \frac{d\Phi_{\mu\nu,i}^\ell(t,s) d\hat{\Phi}_{\mu\nu,i}^\ell(t,s)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu,i}^\ell(t,s) (N\Phi_{\mu\nu,i}^\ell(t,s) - \sigma(h_{\mu,i}^\ell(t)) \cdot \sigma(h_{\nu,i}^\ell(s))) \right] \quad (\text{K.30})$$

$\ell \in \{2, 3, \dots, L-1\}$

$$1 = \int \frac{d\Phi_{\mu\nu}^L(t,s) d\hat{\Phi}_{\mu\nu}^L(t,s)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^L(t,s) (N\Phi_{\mu\nu}^L(t,s) - h_\mu^L(t) \cdot h_\nu^L(s)) \right] \quad (\text{K.31})$$

$$1 = \int \frac{dG_{\mu\nu}^1(t,s) d\hat{G}_{\mu\nu}^1(t,s)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^1(t,s) (NG_{\mu\nu}^1(t,s) - g_\mu^1(t) \cdot g_\nu^1(s)) \right] \quad (\text{K.32})$$

$$1 = \int \frac{dG_{\mu\nu,i}^\ell(t,s) d\hat{G}_{\mu\nu,i}^\ell(t,s)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu,i}^\ell(t,s) (NG_{\mu\nu,i}^\ell(t,s) - g_{\mu,i}^\ell(t) \cdot g_{\nu,i}^\ell(s)) \right] \quad (\text{K.33})$$

$\ell \in \{2, 3, \dots, L\}$

$$1 = \int \frac{dA_{\mu\nu,i}^2(t,s) dB_{\mu\nu,i}^2(t,s)}{2\pi i N^{-1}} \exp \left[ -B_{\mu\nu,i}^2(t,s) (NA_{\mu\nu,i}^2(t,s) + i\hat{\xi}_i^1(t) \cdot \sigma(h_\nu^1(s))) \right] \quad (\text{K.34})$$

$$1 = \int \frac{dA_{\mu\nu}^\phi(t,s) dB_{\mu\nu}^\phi(t,s)}{2\pi i N^{-1}} \exp \left[ -B_{\mu\nu}^\phi(t,s) (NA_{\mu\nu}^\phi(t,s) + i\hat{\xi}_\mu^\phi(t) \cdot \sigma(h_\nu^1(s))) \right] \quad (\text{K.35})$$

$$1 = \int \frac{dA_{\mu\nu,i}^\ell(t,s) dB_{\mu\nu,i}^\ell(t,s)}{2\pi i N^{-1}} \exp \left[ -B_{\mu\nu,i}^\ell(t,s) (NA_{\mu\nu,i}^\ell(t,s) + i\hat{\xi}_{\mu,i}^{\ell-1}(t) \cdot \sigma(h_{\nu,i}^{\ell-1}(s))) \right] \quad (\text{K.36})$$

$\ell \in \{3, 4, \dots, L\}$

$$1 = \int \frac{dG_{\mu\nu}^\phi(t,s) d\hat{G}_{\mu\nu}^\phi(t,s)}{2\pi i} \exp \left[ \hat{G}_{\mu\nu}^\phi(t,s) (G_{\mu\nu}^\phi(t,s) + g_\mu^\phi(t) \cdot g_\nu^\phi(t)) \right] \quad (\text{K.37})$$

We call the conjugate kernel of  $A$   $B$ , rather than  $\hat{A}$ . This is because it will turn out to be equal to  $\hat{\chi} \cdot g$ , which appears in the gradient stream update equation in the final DMFT.

Multiplying in all these factors of unity, we arrive at a partition function which factorises over each of the  $N$  sites in each hidden layer. There are two things to note here:

1. We want to treat  $g, h$  as scalars, so each of the resolutions of the identity [K.28](#) to [K.37](#) above will be duplicated  $N$  times leading to an extensive factor  $N$ , which pulls through to the front of the exponential, since all the neurons in a given layer are indistinguishable in the limit, so have the same  $\Phi$  or  $G$ . That is, we can write for instance

$$1 = \int \frac{dG_{\mu\nu}^1(t,s) d\hat{G}_{\mu\nu}^1(t,s)}{2\pi i N^{-1}} \exp \left[ N\hat{G}_{\mu\nu}^1(t,s) (G_{\mu\nu}^1(t,s) - g_\mu^1(t)g_\nu^1(s)) \right] \quad (\text{K.38})$$

Where  $g_\mu^1 \in \mathbb{R}$  are the components of the former vector of the same name, which are considered statistically indistinguishable in the limit. The exception to this however is [K.37](#), since  $g_\mu^\phi \in \mathbb{R}^M$ .  $M$  is not being taken to the infinite limit, so we cannot consider the  $M$  components of  $g_\mu^\phi$  to be indistinguishable. Rather we must sum over them:

$$1 = \int \frac{dG_{\mu\nu}^\phi(t,s) d\hat{G}_{\mu\nu}^\phi(t,s)}{2\pi i} \exp \left[ \sum_{\beta=1}^M \hat{G}_{\mu\nu}^\phi(t,s) (G_{\mu\nu}^\phi(t,s) + g_{\mu,\beta}^\phi(t)g_{\nu,\beta}^\phi(t)) \right]. \quad (\text{K.39})$$

2. Each term in the argument of the exponential in K.25 is a dot product of two vectors in  $\mathbb{R}^N$ , so contains  $N$  terms which are identical in the limit. We can therefore bring out a factor  $N$  here too. Once again, the exception is the term  $(j_\mu^\phi(t) + i\hat{\chi}_\mu^\phi(t))\chi_\mu^\phi(t)$ , which is repeated  $M$  times. Therefore when we bring a factor  $N$  out front of the action, we must include a prefactor  $\frac{M}{N}$  with this term.

Define the set  $\mathcal{K}$  of all order parameters. This includes all the kernels defined above .

$$Z \propto \int \prod_{\mu,\nu,t,s,i} d\mathcal{K} \exp\left(N S[\{\hat{\Phi}, \Phi, \hat{G}, G, A, B\}]\right) \quad (\text{K.40})$$

Where the DMFT action  $S$  is  $O(1)$ , and has the form

$$\begin{aligned} S[\{\hat{\Phi}, \Phi, \hat{G}, G, A, B\}] = & \sum_{\mu,\nu} \int_0^\infty dt \int_0^\infty ds \left[ \hat{G}_{\mu\nu}^1(t,s) G_{\mu\nu}^1(t,s) + \sum_{\ell=2}^L \hat{G}_{\mu\nu,i}^\ell(t,s) G_{\mu\nu,i}^\ell(t,s) \right. \\ & + \hat{\Phi}_{\mu\nu}^1(t,s) \Phi_{\mu\nu}^1(t,s) + \sum_{\ell=2}^{L-1} \hat{\Phi}_{\mu\nu,i}^\ell(t,s) \Phi_{\mu\nu,i}^\ell(t,s) + \hat{\Phi}_{\mu\nu}^L(t,s) \Phi_{\mu\nu}^L(t,s) \\ & \left. - \sum_{L=2}^L B_{\mu\nu,i}^\ell(t,s) A_{\mu\nu,i}^\ell(t,s) + \frac{1}{N} \hat{G}_{\mu\nu}^\phi(t,s) G_{\mu\nu}^\phi(t,s) - B_{\mu\nu}^\phi(t,s) A_{\mu\nu}^\phi(t,s) \right] \\ & + \frac{1}{N} \sum_{\alpha=1}^N \ln \mathcal{Z}_N[\Phi, \hat{\Phi}, G, \hat{G}, A, B, j_\alpha, v_\alpha] + \frac{1}{N} \sum_{\beta=1}^M \ln \mathcal{Z}_M[\Phi^1, \hat{G}^\phi, j_\beta^\phi] \end{aligned} \quad (\text{K.41})$$

such that it consists of inner products of the order parameters  $\{\Phi, G, A\}$  and their duals  $\{\hat{\Phi}, \hat{G}, B\}$ , in addition to the single-site MGFs

$$\begin{aligned}
 & \mathcal{Z}_N[\{\Phi, \hat{\Phi}, G, \hat{G}, A, B, j, v\}] \\
 &= \prod_{\mu, \nu, i, t} \int_0^\infty dt \int d\mathcal{F} \\
 & \times \exp \left( + \sum_{\mu=1}^P \int_0^\infty dt \int_0^\infty ds \left[ \sum_{\ell=2}^L \left( (j_{\mu,i}^\ell(t) + i\hat{\chi}_{\mu,i}^\ell(t)) \cdot \chi_{\mu,i}^\ell(t) + (v_{\mu,i}^\ell(t) + i\hat{\xi}_{\mu,i}^\ell(t)) \cdot \xi_{\mu,i}^\ell(t) \right) \right. \right. \\
 & + (j_\mu^1(t) + i\hat{\chi}_\mu^1(t)) \cdot \chi_\mu^1(t) + (v_\mu^{L+1} + i\hat{\xi}_\mu^{L+1}(t)) \cdot \xi_\mu^{L+1}(t) + (j_\mu^\phi(t) + i\hat{\chi}_\mu^\phi(t)) \cdot \chi_\mu^\phi(t) \\
 & \left. \left. + (v_\mu^\phi(t) + i\hat{\xi}_\mu^\phi(t)) \cdot \xi_\mu^\phi(t) \right] \right) \\
 & \times \exp \left( - \frac{1}{2} \sum_{\mu, \nu, i} \int_0^\infty dt \int_0^\infty ds \left[ \hat{\chi}_\mu^1(t) \cdot \hat{\chi}_\nu^1(s) \Phi_{\mu\nu}^0(t, s) + (\hat{\chi}_{\mu,i}^2(t) \cdot \hat{\chi}_{\nu,i}^2(s) \right. \right. \\
 & + \frac{1}{N} \hat{\chi}_\mu^\phi(t) \cdot \hat{\chi}_\mu^\phi(s) \Phi_{\mu\nu}^1(t, s) + \sum_{\ell=3}^L \hat{\chi}_{\mu,i}^\ell(s) \cdot \hat{\chi}_{\nu,i}^\ell(t) \Phi_{\mu\nu,i}^{\ell-1}(s, t) + \sum_{\ell=2}^L \hat{\xi}_{\mu,i}^\ell(s) \cdot \hat{\xi}_{\nu,i}^\ell(t) G_{\mu\nu,i}^\ell(s, t) \\
 & \left. \left. + \hat{\xi}_\mu^{L+1}(s) \cdot \hat{\xi}_\nu^{L+1}(t) + \hat{\xi}_\mu^\phi(t) \hat{\xi}_\nu^\phi(s) G_{\mu\nu}^\phi(t, s) \right] \right) \\
 & - i \sum_{\mu, \nu=1}^P \int_0^\infty dt \int_0^\infty ds \sum_{\ell=2}^L \hat{\chi}_{\mu,i}^\ell(t) \cdot g_{\nu,i}^\ell(t) A_{\mu\nu,i}^\ell(t, t) \\
 & + i \sum_{\mu} \int_0^\infty dt \left[ \hat{h}_\mu^3(t) h_\mu^3(t) - \sum_{i=1}^M \phi_\mu^i(t) \hat{h}_\mu^3(t) h_\mu^3(t) \right] \\
 & \times \exp \left( - \sum_{\mu, \nu, i} \int_0^\infty dt \int_0^\infty dt \left[ x_\mu x_\nu \hat{\Phi}_{\mu\nu}^0 + \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \hat{\Phi}_{\mu\nu}^1(t, s) \right. \right. \\
 & + \sum_{\ell=2}^{L-1} \sigma(h_{\mu,i}^\ell(t)) \sigma(h_{\nu,i}^\ell(s)) \hat{\Phi}_{\mu\nu,i}^\ell(t, s) + h_\mu^L(t) h_\nu^L(s) \hat{\Phi}_{\mu\nu}^L(t, s) + g_\mu^1(t) g_\nu^1(s) \hat{G}_{\mu\nu}^1(t, s) \\
 & + \sum_{\ell=2}^L g_{\mu,i}^\ell(t) g_{\nu,i}^\ell(s) \hat{G}_{\mu\nu,i}^\ell(t, s) + \hat{\xi}_{\mu,i}^2(t) \sigma(h_\nu^1(s)) B_{\mu\nu,i}^2(t, s) \\
 & \left. \left. + i \sum_{\ell=3}^L \hat{\xi}_{\mu,i}^\ell(t) \sigma(h_{\nu,i}^{\ell-1}(s)) B_{\mu\nu,i}^\ell(t, s) + i \hat{\xi}_\mu^\phi(t) \sigma(h_\nu^1(s)) B_{\mu\nu}^\phi(t, s) \right] \right)
 \end{aligned} \tag{K.42}$$

And

$$\begin{aligned}
 & \mathcal{Z}_M[\{\Phi^1, \hat{G}^\phi, j_\beta^\phi\}] \\
 &= \prod_{\mu, \nu, t} \int_0^\infty dt \int d\hat{\chi}_{\mu, \beta}^\phi(t) d\chi_{\mu, \beta}^\phi(t) \exp\left( (j_{\mu, \beta}^\phi(t) + i\hat{\chi}_{\mu, \beta}^\phi(t))\chi_\mu^\phi(t) + (v_\mu^{g^\phi} + i\hat{\xi}_\mu^{g^\phi}(t)) \cdot \xi_\mu^{g^\phi}(t) \right. \\
 & \quad \left. - \frac{1}{2} \frac{1}{N} \hat{\chi}_{\mu, \beta}^\phi(t) \hat{\chi}_{\nu, \beta}^\phi(s) \Phi_{\mu\nu}^1(t, s) - g_{\mu, \beta}^\phi g_{\nu, \beta}^\phi \hat{G}_{\mu\nu}^\phi(t, s) - i\hat{\chi}_{\mu, \beta}^\phi(t) g_{\nu, \beta}^\phi(t) A_{\mu\nu}^\phi(t, t) \right)
 \end{aligned} \tag{K.43}$$

The sum over  $\beta$  is an average over sites (rows of vectors in  $\mathbb{R}^M$ ). We regard  $h_\mu(t)$  and  $g_\mu(t)$  as functions of  $\chi$  and  $\xi$ .

It is manifest in the form of [K.40](#) that the actual state of the system in the limit  $N \rightarrow \infty$  will be a saddle point of the action. This is the point where  $\delta S = 0$  for *any* variation of our many order parameters. That is, we impose the following  $\forall t, s, \mu, \nu, i$ :

$$\begin{aligned}
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu}^1(t, s)} &= \hat{G}_{\mu\nu}^1(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu}^1(t, s)} = \hat{G}_{\mu\nu}^1(t, s) - 0 = 0 \\
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu, i}^\ell(t, s)} &= \hat{G}_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu, i}^\ell(t, s)} = \hat{G}_{\mu\nu, i}^\ell(t, s) - \frac{1}{2} \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\xi}_{\mu, i}^\ell(t) \hat{\xi}_{\nu, i}^\ell(s) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu}^1(t, s)} &= \hat{\Phi}_{\mu\nu}^1(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu}^1(t, s)} = \hat{\Phi}_{\mu\nu}^1(t, s) - \frac{1}{2} \left( \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\chi}_{\mu, i}^2(t) \hat{\chi}_{\nu, i}^2(s) \rangle_\alpha \right. \\
 &\quad \left. + \frac{1}{N} \frac{1}{N} \sum_{\beta=1}^M \langle \hat{\chi}_\mu^\phi(t) \hat{\chi}_\nu^\phi(s) \rangle_\beta \right) = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu, i}^\ell(t, s)} &= \hat{\Phi}_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu, i}^\ell(t, s)} = \hat{\Phi}_{\mu\nu, i}^\ell(t, s) - \frac{1}{2} \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\chi}_{\mu, i}^\ell(t) \hat{\chi}_{\nu, i}^\ell(s) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu}^L(t, s)} &= \hat{\Phi}_{\mu\nu}^L(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu}^L(t, s)} = \hat{\Phi}_{\mu\nu}^L(t, s) - 0 = 0 \\
 \frac{\delta \mathcal{S}}{\delta A_{\mu\nu, i}^\ell(t, s)} &= -B_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta A_{\mu\nu, i}^\ell(t, s)} = -B_{\mu\nu, i}^\ell(t, s) - i \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\chi}_{\mu, i}^\ell(t) g_{\nu, i}^\ell(s) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta A_{\mu\nu}^\phi(t, s)} &= -B_{\mu\nu}^\phi(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta A_{\mu\nu}^\phi(t, s)} = -B_{\mu\nu}^\phi(t, s) - i \frac{1}{N} \sum_{\beta=1}^M \langle \hat{\chi}_\mu^\phi(t) g_\nu^\phi(s) \rangle_\beta = 0 \\
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu}^\phi(t, s)} &= \frac{1}{N} \hat{G}_{\mu\nu}^\phi(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu}^\phi(t, s)} = \frac{1}{N} \hat{G}_{\mu\nu}^\phi(t, s) - \frac{1}{2} \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\xi}_\mu^\phi(t) \hat{\xi}_\nu^\phi(s) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu}^1(t, s)} &= \Phi_{\mu\nu}^1(t, s) - \frac{1}{N} \sum_{\alpha=1}^N \langle \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu, i}^\ell(t, s)} &= \Phi_{\mu\nu, i}^\ell(t, s) - \frac{1}{N} \sum_{\alpha=1}^N \langle \sigma(h_{\mu, i}^\ell(t)) \sigma(h_{\nu, i}^\ell(s)) \rangle_\alpha = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu}^L(t, s)} &= \Phi_{\mu\nu}^L(t, s) - \frac{1}{N} \sum_{\alpha=1}^N \langle h_\mu^L(t) h_\nu^L(s) \rangle_\alpha = 0
 \end{aligned}$$

(K.44)

$$\begin{aligned}
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu}^1(t, s)} &= G_{\mu\nu}^1(t, s) - \frac{1}{N} \sum_{\alpha=1}^N \langle g_{\mu}^1(t) g_{\nu}^1(s) \rangle_{\alpha} = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu, i}^{\ell}(t, s)} &= G_{\mu\nu, i}^{\ell}(t, s) - \frac{1}{N} \sum_{\alpha=1}^N \langle g_{\mu, i}^{\ell}(t) g_{\nu, i}^{\ell}(s) \rangle_{\alpha} = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu, i}^2(t, s)} &= -A_{\mu\nu, i}^2(t, s) - i \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\xi}_{\mu, i}^2(t) \sigma(h_{\nu}^1(s)) \rangle_{\alpha} = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu, i}^{\ell}(t, s)} &= -A_{\mu\nu, i}^{\ell}(t, s) - i \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\xi}_{\mu, i}^{\ell}(t) \sigma(h_{\nu, i}^{\ell-1}(s)) \rangle_{\alpha} = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu}^{\phi}(t, s)} &= -A_{\mu\nu}^{\phi}(t, s) - i \frac{1}{N} \sum_{\alpha=1}^N \langle \hat{\xi}_{\mu}^{\phi}(t) \sigma(h_{\nu}^1(s)) \rangle_{\alpha} = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu}^{\phi}(t, s)} &= \frac{1}{N} G_{\mu\nu}^{\phi}(t, s) - \frac{1}{N} \sum_{\beta=1}^M \langle g_{\mu}^{\phi}(t) g_{\nu}^{\phi}(s) \rangle_{\beta} = 0
 \end{aligned} \tag{K.45}$$

Here the average  $\langle \cdot \rangle_{\alpha}$  denotes the  $\alpha^{\text{th}}$  single-site average of an observable  $\mathcal{O}(\{\chi, \xi, u\})$ , defined as:

$$\langle \mathcal{O}(\{\chi, \xi\}) \rangle_{\alpha} = \frac{1}{\mathcal{Z}[j_{\alpha}, v_{\alpha}]} \int \prod_{\mu} \prod_{\text{layers}} d\chi_{\mu}^{\ell} d\xi_{\mu}^{\ell} \exp(-\mathcal{H}[\{\chi, \xi\}, \{j, v\}]) \mathcal{O}(\{\chi, \xi\}) \tag{K.46}$$

Where  $\mathcal{H}$  is the logarithm of the integrand in K.42 or K.43.

At zero source  $j, v \rightarrow 0$ , all single-site averages over  $N$  terms are equivalent, and we can write (for instance)  $\Phi_{\mu\nu}^1 = \langle \sigma(h_{\mu}^1) \sigma(h_{\nu}^1) \rangle$ , where  $\langle \cdot \rangle$  is the average over single-site distributions for  $j, v \rightarrow 0$ . The sum  $\frac{1}{N} \frac{1}{N} \sum_{\beta=1}^M \langle \hat{\chi}_{\mu}^{\phi}(t) \hat{\chi}_{\nu}^{\phi}(s) \rangle_{\beta} \rightarrow 0$  in the limit, but notably,  $\frac{1}{N^2} \sum_{\beta=1}^M \langle g_{\mu}^{\phi}(t) g_{\nu}^{\phi}(s) \rangle_{\beta}$  remains

$O(1)$  but doesn't concentrate. We have

$$\begin{aligned}
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu}^1(t, s)} &= \hat{G}_{\mu\nu}^1(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu}^1(t, s)} = \hat{G}_{\mu\nu}^1(t, s) - 0 = 0 \\
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu, i}^\ell(t, s)} &= \hat{G}_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu, i}^\ell(t, s)} = \hat{G}_{\mu\nu, i}^\ell(t, s) - \frac{1}{2} \langle \hat{\xi}_{\mu, i}^\ell(t) \hat{\xi}_{\nu, i}^\ell(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu}^0} &= \hat{\Phi}_{\mu\nu}^0 - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu}^0} = \hat{\Phi}_{\mu\nu}^0 - \frac{1}{2N} \langle \hat{\chi}_\mu^1(t) \hat{\chi}_\nu^1(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu}^1(t, s)} &= \hat{\Phi}_{\mu\nu}^1(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu}^1(t, s)} = \hat{\Phi}_{\mu\nu}^1(t, s) - \frac{1}{2} \langle \hat{\chi}_{\mu, i}^2(t) \hat{\chi}_{\nu, i}^2(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu, i}^\ell(t, s)} &= \hat{\Phi}_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu, i}^\ell(t, s)} = \hat{\Phi}_{\mu\nu, i}^\ell(t, s) - \frac{1}{2} \langle \hat{\chi}_{\mu, i}^\ell(t) \hat{\chi}_{\nu, i}^\ell(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \Phi_{\mu\nu}^L(t, s)} &= \hat{\Phi}_{\mu\nu}^L(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta \Phi_{\mu\nu}^L(t, s)} = \hat{\Phi}_{\mu\nu}^L(t, s) - 0 = 0 \\
 \frac{\delta \mathcal{S}}{\delta A_{\mu\nu, i}^\ell(t, s)} &= -B_{\mu\nu, i}^\ell(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta A_{\mu\nu, i}^\ell(t, s)} = -B_{\mu\nu, i}^\ell(t, s) - i \langle \hat{\chi}_{\mu, i}^\ell(t) g_{\nu, i}^\ell(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta A_{\mu\nu}^\phi(t, s)} &= -B_{\mu\nu}^\phi(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta A_{\mu\nu}^\phi(t, s)} = -B_{\mu\nu}^\phi(t, s) = 0 \\
 \frac{\delta \mathcal{S}}{\delta G_{\mu\nu}^\phi(t, s)} &= \frac{1}{N} \hat{G}_{\mu\nu}^\phi(t, s) - \frac{1}{\mathcal{Z}} \frac{\delta \mathcal{Z}}{\delta G_{\mu\nu}^\phi(t, s)} = \frac{1}{N} \hat{G}_{\mu\nu}^\phi(t, s) - \frac{1}{2} \langle \hat{\xi}_\mu^\phi(t) \hat{\xi}_\nu^\phi(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu}^1(t, s)} &= \Phi_{\mu\nu}^1(t, s) - \langle \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu, i}^\ell(t, s)} &= \Phi_{\mu\nu, i}^\ell(t, s) - \langle \sigma(h_{\mu, i}^\ell(t)) \sigma(h_{\nu, i}^\ell(s)) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{\Phi}_{\mu\nu}^L(t, s)} &= \Phi_{\mu\nu}^L(t, s) - \langle h_\mu^L(t) h_\nu^L(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu}^1(t, s)} &= G_{\mu\nu}^1(t, s) - \langle g_\mu^1(t) g_\nu^1(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu, i}^\ell(t, s)} &= G_{\mu\nu, i}^\ell(t, s) - \langle g_{\mu, i}^\ell(t) g_{\nu, i}^\ell(s) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu, i}^2(t, s)} &= -A_{\mu\nu, i}^2(t, s) - i \langle \hat{\xi}_{\mu, i}^2(t) \sigma(h_\nu^1(s)) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu, i}^\ell(t, s)} &= -A_{\mu\nu, i}^\ell(t, s) - i \langle \hat{\xi}_{\mu, i}^\ell(t) \sigma(h_{\nu, i}^{\ell-1}(s)) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta B_{\mu\nu}^\phi(t, s)} &= -A_{\mu\nu}^\phi(t, s) - i \langle \hat{\xi}_\mu^\phi(t) \sigma(h_\nu^1(s)) \rangle = 0 \\
 \frac{\delta \mathcal{S}}{\delta \hat{G}_{\mu\nu}^\phi(t, s)} &= G_{\mu\nu}^\phi(t, s) - \sum_{\beta=1}^M \langle g_\mu^\phi(t) g_\nu^\phi(s) \rangle_\beta = 0
 \end{aligned} \tag{K.47}$$

#### K.4. Averages of dual variables

We see that the dual variables  $\{\hat{\Phi}, \hat{G}, B\}$  are given in terms of such correlators as  $\langle \hat{\chi} \hat{\chi} \rangle$ . We will now show that these vanish, finding along the way expressions for the fields  $\{\chi, \xi\}$ , as well as the kernels  $\{A, B\}$ .

It will serve brevity to work with a vectorised notation. To wit: let  $\chi_i^\ell = \text{Vec}\{\chi_{\mu,i}^\ell(t)\}_{\mu \in \{1, \dots, P\}, t \in \mathbb{R}^+}$  denote the vectorization of the stochastic field over different samples and times, with analogous objects defined for the other fields. We denote the dot product between such objects  $\mathbf{a} \cdot \mathbf{b} = \sum_{\mu=1}^P \int_0^\infty a_\mu(t) b_\mu(t)$ . A similar procedure is applied to matrices by defining  $\Phi = \text{Mat}\{\Phi_{\mu\nu}\}_{\mu\nu \in \{1, \dots, P\}, ts \in \mathbb{R}^+}$ , with the appropriate matrix product defined as  $[\mathbf{A}\mathbf{b}]_{\mu,t} = \int_0^t ds \frac{1}{P} \sum_{\nu=1}^P A_{\mu\nu}(t, s) b_\nu(s)$ .

With this notation in place, we can write  $\langle \hat{\chi}_i^\ell \hat{\chi}_i^\ell \rangle$  for  $\ell \in \{3, \dots, L-1\}$  in terms of the moment generating function for  $\hat{\chi}_i^\ell$ , and the dummy field  $\mathbf{w}$

$$\begin{aligned}
 \langle \hat{\chi}_i^\ell \hat{\chi}_i^\ell \rangle &= - \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \left\langle \exp(i \mathbf{w} \cdot \hat{\chi}_i^\ell) \right\rangle \Big|_{\mathbf{w}=0} \\
 &= - \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \left[ \frac{1}{\mathcal{Z}} \int d\hat{\chi}_i^\ell d\chi_i^\ell \exp(i \mathbf{w} \cdot \hat{\chi}_i^\ell) \exp(-\mathcal{H}) \right] \Big|_{\mathbf{w}=0} \\
 &= - \frac{1}{\mathcal{Z}} \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \int d\hat{\chi}_i^\ell d\chi_i^\ell \exp\left(-\frac{1}{2} \hat{\chi}_i^{\ell\top} \Phi_i^{\ell-1} \hat{\chi}_i^\ell + i(-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell) \cdot \hat{\chi}_i^\ell\right) \Big|_{\mathbf{w}=0} \\
 &= - \frac{1}{\mathcal{Z}} \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \int d\chi_i^\ell \exp\left(-\frac{1}{2} (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top [\Phi_i^{\ell-1}]^{-1} (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell)\right) \Big|_{\mathbf{w}=0} \\
 &= \frac{1}{\mathcal{Z}} \int d\chi_i^\ell \exp\left(-\frac{1}{2} (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top [\Phi_i^{\ell-1}]^{-1} (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell)\right) \\
 &\quad \times \left( [\Phi_i^{\ell-1}]^{-1} - [\Phi_i^{\ell-1}]^{-1} (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell) (-\chi_i^\ell + \mathbf{w} + \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top [\Phi_i^{\ell-1}]^{-1} \right) \Big|_{\mathbf{w}=0} \\
 &= [\Phi_i^{\ell-1}]^{-1} - [\Phi_i^{\ell-1}]^{-1} \left\langle (\chi_i^\ell - \mathbf{A}_i^\ell \mathbf{g}_i^\ell) (\chi_i^\ell - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top \right\rangle [\Phi_i^{\ell-1}]^{-1}
 \end{aligned} \tag{K.48}$$

Where we use the definition K.46 of the average, and retain only the relevant parts of the integral. Through similar derivations, we obtain also the following:

$$\langle \hat{\chi}_i^2 \hat{\chi}_i^2 \rangle = [\Phi^1]^{-1} - [\Phi^1]^{-1} \left\langle (\chi_i^2 - \mathbf{A}_i^2 \mathbf{g}_i^2) (\chi_i^2 - \mathbf{A}_i^2 \mathbf{g}_i^2)^\top \right\rangle [\Phi^1]^{-1} \tag{K.49}$$

$$\langle \hat{\chi}^1 \hat{\chi}^1 \rangle = [\Phi^0]^{-1} - [\Phi^0]^{-1} \left\langle \chi^1 \chi^{1\top} \right\rangle [\Phi^0]^{-1} \tag{K.50}$$

$$\langle \hat{\xi}_i^\ell \hat{\xi}_i^\ell \rangle = [\mathbf{G}_i^\ell]^{-1} - [\mathbf{G}_i^\ell]^{-1} \left\langle (\xi_i^\ell - \mathbf{B}_i^\ell \mathbf{g}_i^\ell) (\xi_i^\ell - \mathbf{B}_i^\ell \mathbf{g}_i^\ell)^\top \right\rangle [\mathbf{G}_i^\ell]^{-1} \quad \ell \in \{2, 3, \dots, L\} \tag{K.51}$$

$$\langle \hat{\xi}^\phi \hat{\xi}^\phi \rangle = [\mathbf{G}^\phi]^{-1} - [\mathbf{G}^\phi]^{-1} \left\langle (\xi^\phi - \mathbf{B}^\phi \mathbf{g}^\phi) (\xi^\phi - \mathbf{B}^\phi \mathbf{g}^\phi)^\top \right\rangle [\mathbf{G}^\phi]^{-1} \tag{K.52}$$

$$\begin{aligned}
 -i \langle \hat{\chi}_i^\ell \mathbf{g}_i^\ell \rangle &= \frac{\partial}{\partial \mathbf{w}} \left\langle \exp(-i \mathbf{w} \cdot \hat{\chi}_i^\ell) \mathbf{g}_i^{\ell\top} \right\rangle \Big|_{\mathbf{w}=0} \\
 &= \frac{1}{Z} \frac{\partial}{\partial \mathbf{w}} \int d\hat{\chi}_i^\ell d\chi_i^\ell \exp\left(-\frac{1}{2} \hat{\chi}_i^{\ell\top} \Phi_i^{\ell-1} \hat{\chi}_i^\ell + i(\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell) \cdot \hat{\chi}_i^\ell\right) \mathbf{g}_i^{\ell\top} \Big|_{\mathbf{w}=0} \\
 &= \frac{1}{Z} \frac{\partial}{\partial \mathbf{w}} \int d\chi_i^\ell \exp\left(-\frac{1}{2} (\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top [\Phi_i^{\ell-1}]^{-1} (\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)\right) \mathbf{g}_i^{\ell\top} \Big|_{\mathbf{w}=0} \\
 &= -\frac{1}{Z} \int d\chi_i^\ell \exp\left(-\frac{1}{2} (\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)^\top [\Phi_i^{\ell-1}]^{-1} (\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)\right) \\
 &\quad \times [\Phi_i^{\ell-1}]^{-1} (\chi_i^\ell + \mathbf{w} - \mathbf{A}_i^\ell \mathbf{g}_i^\ell) \mathbf{g}_i^{\ell\top} \Big|_{\mathbf{w}=0} \\
 &= [\Phi_i^{\ell-1}]^{-1} \langle (\chi_i^\ell - \mathbf{A}_i^\ell \mathbf{g}_i^\ell) \mathbf{g}_i^{\ell\top} \rangle \quad \ell \in \{2, 3, \dots, L\}
 \end{aligned} \tag{K.53}$$

$$-i \langle \hat{\xi}_i^\ell \sigma(\mathbf{h}_i^{\ell-1}) \rangle = [\mathbf{G}_i^\ell]^{-1} \langle (\xi_i^\ell - \mathbf{B}_i^\ell \sigma(\mathbf{h}_i^{\ell-1})) \sigma(\mathbf{h}_i^{\ell-1}) \rangle \quad \ell \in \{3, 4, \dots, L\} \tag{K.54}$$

$$-i \langle \hat{\xi}^\phi \sigma(\mathbf{h}^1) \rangle = [\mathbf{G}^\phi]^{-1} \langle (\xi^\phi - \mathbf{B}^\phi \sigma(\mathbf{h}^1)) \sigma(\mathbf{h}^1) \rangle \tag{K.55}$$

$$-i \langle \hat{\xi}^{2,i} \sigma(\mathbf{h}^1) \rangle = [\mathbf{G}^{2,i}]^{-1} \langle (\xi^{2,i} - \mathbf{B}_i^2 \sigma(\mathbf{h}^1)) \sigma(\mathbf{h}^1) \rangle \tag{K.56}$$

Now the Hubbard Trick [24] states that

$$\exp\left(-\frac{1}{2} \mathbf{x}^\top A \mathbf{x}\right) = \int_{\mathbb{R}^d} \frac{d\mathbf{u}}{(2\pi)^{d/2} \sqrt{\det A}} \exp\left(-\frac{1}{2} \mathbf{u}^\top A^{-1} \mathbf{u} - i \mathbf{u} \cdot \mathbf{x}\right) = \langle \exp(-i \mathbf{u} \cdot \mathbf{x}) \rangle_{\mathbf{u} \sim \mathcal{N}(0, A)}. \tag{K.57}$$

This allows us to rewrite the quadratic terms in our single-site MGF as follows:

$$\begin{aligned}
 &\exp\left(-\frac{1}{2} \sum_{\mu\nu} \int_0^\infty dt \int_0^\infty ds \hat{\chi}_\mu^1(t) \cdot \hat{\chi}_\nu^1(s) \Phi_{\mu\nu}^0(t, s)\right) = \\
 &= \left\langle \exp\left(-i \sum_\mu \int_0^\infty dt \alpha_\mu^1(t) \hat{\chi}_\mu^1(t)\right) \right\rangle_{\alpha^1 \sim \mathcal{GP}(0, \Phi^0)}
 \end{aligned} \tag{K.58}$$

$$\begin{aligned}
 &\exp\left(-\frac{1}{2} \sum_{\mu\nu i} \int_0^\infty dt \int_0^\infty ds \hat{\chi}_{\mu,i}^2(t) \cdot \hat{\chi}_{\nu,i}^2(s) \Phi_{\mu\nu}^1(t, s)\right) = \\
 &= \left\langle \exp\left(-i \sum_{\mu i} \int_0^\infty dt \alpha_{\mu,i}^2(t) \hat{\chi}_{\mu,i}^2(t)\right) \right\rangle_{\{\alpha_i^2\} \sim \mathcal{GP}(0, \Phi^1)}
 \end{aligned} \tag{K.59}$$

$$\begin{aligned}
 &\exp\left(-\frac{1}{2} \sum_{\mu\nu i} \int_0^\infty dt \int_0^\infty ds \hat{\chi}_{\mu,i}^\ell(t) \cdot \hat{\chi}_{\nu,i}^\ell(s) \Phi_{\mu\nu,i}^{\ell-1}(t, s)\right) = \\
 &= \left\langle \exp\left(-i \sum_{\mu i} \int_0^\infty dt \alpha_{\mu,i}^\ell(t) \hat{\chi}_{\mu,i}^\ell(t)\right) \right\rangle_{\{\alpha_i^\ell\} \sim \mathcal{GP}(0, \Phi_i^{\ell-1})} \quad \ell \in \{3, 4, \dots, L\}
 \end{aligned} \tag{K.60}$$

$$\begin{aligned} & \exp\left(-\frac{1}{2} \sum_{\mu\nu i} \int_0^\infty dt \int_0^\infty ds \hat{\xi}_{\mu,i}^\ell(t) \cdot \hat{\xi}_{\nu,i}^\ell(s) G_{\mu\nu,i}^\ell(t,s)\right) = \\ & = \left\langle \exp\left(-i \sum_{\mu i} \int_0^\infty dt \beta_{\mu,i}^\ell(t) \hat{\xi}_{\mu,i}^\ell(t)\right) \right\rangle_{\{\beta_i^\ell\} \sim \mathcal{GP}(0, G_i^\ell)} \quad \ell \in \{2, 3, 4, \dots, L\} \end{aligned} \quad (\text{K.61})$$

$$\begin{aligned} & \exp\left(-\frac{1}{2} \sum_{\mu\nu i} \int_0^\infty dt \int_0^\infty ds \hat{\xi}_\mu^{L+1}(t) \cdot \hat{\xi}_\nu^{L+1}(s)\right) = \\ & = \left\langle \exp\left(-i \sum_{\mu i} \int_0^\infty dt \beta_\mu^{L+1}(t) \hat{\xi}_\mu^{L+1}(t)\right) \right\rangle_{\{\beta^{L+1}\} \sim \mathcal{GP}(0,1)} \end{aligned} \quad (\text{K.62})$$

$$\exp\left(-\frac{1}{2} \sum_{\mu\nu i} \int_0^\infty dt \int_0^\infty ds \hat{\xi}_\mu^\phi(t) \cdot \hat{\xi}_\nu^\phi(s)\right) = \left\langle \exp\left(-i \sum_{\mu i} \int_0^\infty dt \beta_\mu^\phi(t) \hat{\xi}_\mu^\phi(t)\right) \right\rangle_{\{\beta^\phi\} \sim \mathcal{GP}(0, G^\phi)} \quad (\text{K.63})$$

We can now easily integrate over all the  $\{\hat{\chi}, \hat{\xi}\}$  variables, since the argument of the exponential in  $\mathcal{Z}$  has been linearised with respect to them all. Doing so yields the following delta functions:

$$\int \prod_{\mu,t} \frac{d\hat{\chi}_\mu^1(t)}{2\pi} \exp(i \hat{\chi}^1 \cdot (\boldsymbol{\chi}^1 - \boldsymbol{\alpha}^1)) = \delta(\boldsymbol{\chi}^1 - \boldsymbol{\alpha}^1) \quad (\text{K.64})$$

$$\int \prod_{\mu,t} \frac{d\hat{\chi}_{\mu,i}^\ell(t)}{2\pi} \exp[i \hat{\chi}_i^\ell \cdot (\boldsymbol{\chi}_i^\ell - \boldsymbol{\alpha}_i^\ell - \mathbf{A}_i^\ell \mathbf{g}_i^\ell)] = \delta(\boldsymbol{\chi}_i^\ell - \boldsymbol{\alpha}_i^\ell - \mathbf{A}_i^\ell \mathbf{g}_i^\ell) \quad \ell \in \{2, 3, \dots, L\} \quad (\text{K.65})$$

$$\int \prod_{\mu,t} \frac{d\hat{\xi}_{\mu,i}^{2,i}(t)}{2\pi} \exp(i \hat{\xi}_i^{2,i} \cdot (\boldsymbol{\xi}_i^{2,i} - \boldsymbol{\beta}_i^{2,i} - \mathbf{B}_i^{2,i\top} \boldsymbol{\sigma}(\mathbf{h}^1))) = \delta(\boldsymbol{\xi}_i^{2,i} - \boldsymbol{\beta}_i^{2,i} - \mathbf{B}_i^{2,i\top} \boldsymbol{\sigma}(\mathbf{h}^1)) \quad (\text{K.66})$$

$$\int \prod_{\mu,t} \frac{d\hat{\xi}_{\mu,i}^\ell(t)}{2\pi} \exp(i \hat{\xi}_i^\ell \cdot (\boldsymbol{\xi}_i^\ell - \boldsymbol{\beta}_i^\ell - \mathbf{B}_i^{\ell\top} \boldsymbol{\sigma}(\mathbf{h}_i^{\ell-1}))) = \delta(\boldsymbol{\xi}_i^\ell - \boldsymbol{\beta}_i^\ell - \mathbf{B}_i^{\ell\top} \boldsymbol{\sigma}(\mathbf{h}_i^{\ell-1})) \quad \ell \in \{3, 4, \dots, L\} \quad (\text{K.67})$$

$$\int \prod_{\mu,t} \frac{d\hat{\xi}_\mu^\phi(t)}{2\pi} \exp(i \hat{\xi}^\phi \cdot (\boldsymbol{\xi}^\phi - \boldsymbol{\beta}^\phi - \mathbf{B}^{\phi\top} \boldsymbol{\sigma}(\mathbf{h}^1))) = \delta(\boldsymbol{\xi}^\phi - \boldsymbol{\beta}^\phi - \mathbf{B}^{\phi\top} \boldsymbol{\sigma}(\mathbf{h}^1)) \quad (\text{K.68})$$

$$\int \prod_{\mu,t} \frac{d\hat{\xi}_\mu^{L+1}(t)}{2\pi} \exp(i \hat{\xi}^{L+1} \cdot (\boldsymbol{\xi}^{L+1} - \boldsymbol{\beta}^{L+1})) = \delta(\boldsymbol{\xi}^{L+1} - \boldsymbol{\beta}^{L+1}) \quad (\text{K.69})$$

Using these in the expressions K.49 to K.56, we note that all four classes of correlators vanish, since the average term collapses to the covariance matrix of the respective Gaussian process  $\{\alpha, \beta\}$ , which cancels with its inverse, giving zero overall. We summarise this as follows (where the expression is implied to hold for all allowed arguments, superscripts and subscripts):

$$\hat{G} = \hat{\Phi} = \langle \hat{\chi} \hat{\chi} \rangle = \langle \hat{\xi} \hat{\xi} \rangle = 0 \quad (\text{K.70})$$

### K.5. The coupling kernels A and B

What of  $A$  and its dual  $B$ ? Well inserting K.66 into K.54, and inserting this in turn into the relevant saddle point equation, we obtain for  $\ell \in \{3, 4, \dots, L\}$

$$\begin{aligned} A_{\mu\nu,i}^\ell(t, s) &= -i \langle \hat{\xi}_{\nu,i}^\ell(t) \sigma(h_{\mu,i}^{\ell-1}(s)) \rangle \\ &= [G_i^\ell]^{-1} \langle (\xi_i^\ell - B_i^\ell \sigma(h_i^{\ell-1})) \sigma(h_i^{\ell-1}) \rangle \\ &= \left\langle \frac{\partial \sigma(h_{\mu,i}^{\ell-1}(t))}{\partial \beta_{\nu,i}^{\ell\top}(s)} \right\rangle \end{aligned} \quad (\text{K.71})$$

In the final line, we use Stein's lemma [48] which states that for a normally distributed random variable  $X$  with expectation  $\mu$  and variance  $\sigma^2$ , and differentiable function  $g$ ,  $\mathbb{E}[g(X)(X - \mu)] = \sigma^2 \mathbb{E}[g'(X)]$ .

Similarly:

$$A_{\mu\nu,i}^2(t, s) = \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \beta_{\nu,i}^{2\top}(s)} \right\rangle \quad (\text{K.72})$$

$$A_{\mu\nu}^\phi(t, s) = \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \beta_\nu^{\phi\top}(s)} \right\rangle \quad (\text{K.73})$$

$$B_{\mu\nu,i}^\ell(t, s) = \left\langle \frac{\partial g_{\mu,i}^\ell(t)}{\partial \alpha_{\nu,i}^{\ell\top}(s)} \right\rangle \quad (\text{K.74})$$

Note that the saddle equations tell us that  $B_{\mu\nu}^\phi(t, s) = 0$ .

### K.6. Final DMFT

To remedy the asymmetry manifest in equations K.66 and K.67, we redefine  $B$  as its transpose. We also rescale  $A$  to  $\gamma_0 \eta_0 A$ , and  $B$  to  $\gamma_0 \eta_0 B$ . This makes it clear that the non-Gaussian corrections to  $h$  and  $z$  are  $O(\gamma_0 \eta_0)$ . We also replace instances of  $g$  with their corresponding  $z$  using K.4.

We then have the following complete self-consistent DMFT equations:

$$\begin{aligned}
 \alpha^1(t) &\sim \mathcal{GP}(0, \Phi^0), & \alpha_i^\ell(t) &\sim \begin{cases} \mathcal{GP}(0, \Phi^1(t, s)) & \ell = 2 \\ \mathcal{GP}(0, \Phi_i^{\ell-1}(t, s)) & \ell \in \{3, 4, \dots, L\} \end{cases} \\
 \beta_i^\ell(t) &\sim \mathcal{GP}(0, G_i^{\ell+1}(t, s)) & \ell \in \{1, 2, \dots, L-1\}, & \beta^L(t) &\sim \mathcal{N}(0, 1), \\
 \beta^\phi(t) &\sim \mathcal{GP}(0, G^\phi(t, s)) & \Phi_{\mu\nu}^0 = \langle x_\mu x_\nu \rangle, & \Phi_{\mu\nu}^1(t, s) &= \langle \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \rangle, \\
 \Phi_{\mu\nu, i}^\ell(t, s) &= \langle \sigma(h_{\mu, i}^\ell(t)) \sigma(h_{\nu, i}^\ell(s)) \rangle & \ell \in \{2, 3, \dots, L-1\}, \\
 \Phi_{\mu\nu}^L(t, s) &= \left\langle \sum_{i, j=1}^M \phi_\mu^i(t) \phi_\nu^j(s) h_{\mu, i}^L(t) h_{\nu, j}^L(s) \right\rangle, \\
 G_{\mu\nu}^1(t, s) &= \langle [\dot{\sigma}(h_\mu^1(t)) \odot z_\mu^1(t)] [\dot{\sigma}(h_\nu^1(s)) \odot z_\nu^1(s)] \rangle \\
 G_{\mu\nu, i}^\ell(t, s) &= \langle [\dot{\sigma}(h_{\mu, i}^\ell(t)) \odot z_{\mu, i}^\ell(t)] [\dot{\sigma}(h_{\nu, i}^\ell(s)) \odot z_{\nu, i}^\ell(s)] \rangle & \ell \in \{2, 3, \dots, L-1\}, \\
 G_{\mu\nu, i}^L(t, s) &= \langle z_{\mu, i}^L(t) z_{\nu, i}^L(s) \rangle \\
 G_{\mu\nu}^\phi(t, s) &= \sum_{i=1}^M \langle h_{\mu, i}^L(t) z_\mu^L(t) \rangle_i \langle h_{\nu, i}^L(s) z_\nu^L(s) \rangle_i \\
 A_{\mu\nu, i}^2(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \beta_{\nu, i}^{2\top}(s)} \right\rangle, & A_{\mu\nu, i}^\ell(t, s) &= \left\langle \frac{\partial \sigma(h_{\mu, i}^{\ell-1}(t))}{\partial \beta_{\nu, i}^{\ell\top}(s)} \right\rangle & \ell \in \{3, 4, \dots, L\}, \\
 A_{\mu\nu}^\phi(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \beta_{\nu}^{\phi\top}(s)} \right\rangle \\
 B_{\mu\nu, i}^\ell(t, s) &= \left\langle \frac{\partial (\dot{\sigma}(h_{\mu, i}^\ell(t)) z_{\mu, i}^\ell(t))}{\partial \alpha_{\nu, i}^{\ell\top}(s)} \right\rangle & \ell \in \{2, 3, \dots, L-1\}, & B_{\mu\nu, i}^L(t, s) &= \left\langle \frac{\partial z_{\mu, i}^L(t)}{\partial \alpha_{\nu, i}^{L\top}(s)} \right\rangle \\
 \tilde{z}_{\mu, i}^1(t) &= \beta_{\mu, i}^1(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} [B_{\mu\nu, i}^2(s, t) + \Delta_\nu(s) G_{\mu\nu, i}^2(s, t)] \sigma(h_\nu^1(s)) \\
 \tilde{z}_\mu^\phi(t) &= \beta_\mu^\phi(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} [\Delta_\nu(s) G_{\mu\nu}^\phi(s, t)] \sigma(h_\nu^1(s)), \\
 z_\mu^1(t) &= \sum_{i=1}^M \tilde{z}_{\mu, i}^1(t) + \tilde{z}_\mu^\phi \\
 z_{\mu, i}^\ell(t) &= \beta_{\mu, i}^\ell(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} [B_{\mu\nu, i}^{\ell+1}(s, t) + \Delta_\nu(s) G_{\mu\nu, i}^{\ell+1}(s, t)] \sigma(h_{\nu, i}^\ell(s)) \\
 & & \ell \in \{2, 3, \dots, L-1\} \\
 z_\mu^L(t) &= \beta_\mu^L(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu} \Delta_\nu(s) h_\nu^L(s), & z_{\mu, i}^L(t) &= \phi^i(t) z_\mu^L(t)
 \end{aligned} \tag{K.75}$$

$$\begin{aligned}
 h_\mu^1(t) &= \alpha_\mu^1(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu \Phi_{\mu\nu}^0 (\dot{\sigma}(h_\mu^1) \odot z_\mu^1(s)), \\
 h_{\mu,i}^2(t) &= \alpha_{\mu,i}^2(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [A_{\mu\nu,i}^2(s,t) + \Delta_\nu(s)\Phi_{\mu\nu}^1(s,t)] (\dot{\sigma}(h_{\nu,i}^2(s)) \odot z_{\nu,i}^2(s)), \\
 h_{\mu,i}^\ell(t) &= \alpha_{\mu,i}^\ell(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [A_{\mu\nu,i}^\ell(s,t) + \Delta_\nu(s)\Phi_{\mu\nu,i}^{\ell-1}(s,t)] (\dot{\sigma}(h_{\nu,i}^\ell(s)) \odot z_{\nu,i}^\ell(s)) \\
 &\hspace{25em} \ell \in \{3, 4, \dots, L-1\}, \\
 h_{\mu,i}^L(t) &= \alpha_{\mu,i}^L(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [A_{\mu\nu,i}^L(s,t) + \Delta_\nu(s)\Phi_{\mu\nu,i}^{L-1}(s,t)] z_{\nu,i}^L(s), \\
 h_\mu^L &= \sum_{i=1}^M \phi_\mu^i \cdot h_{\mu,i}^L, \quad \phi_\mu(t) = \text{softmax}(\psi_\mu(t)) \\
 \frac{df_\mu(t)}{dt} &= \frac{\eta_0}{P} \sum_{\nu=1}^P \left[ G_{\mu\nu}^1(t,t)\Phi_{\mu\nu}^0 + \sum_{i=1}^M [G_{\mu\nu,i}^2(t,t)\Phi_{\mu\nu}^1(t,t) + \sum_{\ell=3}^L G_{\mu\nu,i}^\ell(t,t)\Phi_{\mu\nu,i}^{\ell-1}(t,t)] \right. \\
 &\quad \left. + \Phi_{\mu\nu}^L(t,t) + G_{\mu\nu}^\phi(t,t)\Phi_{\mu\nu}^1(t,t) \right] \Delta_\nu(t) \\
 \psi_\mu^i(t) &= \frac{\gamma_0\eta_0}{P} \int_0^t ds \sum_{\nu=1}^P [A_{\nu\mu}^\phi(s,t) + \Delta_\nu(s)\Phi_{\nu\mu}^1(s,t)] [g_\nu^\phi(s)]_i, \\
 [g_\mu^\phi(t)]_i &= \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\mu\nu}^L(t,s) \Phi_{\mu\nu}^{L-1}(t,s) \phi_\nu^i(s)
 \end{aligned} \tag{K.76}$$

## L. DMFT Analysis for Regime II (MSSP)

### L.1. Architectural definitions

Here we distinguish the external width  $N$  and the internal width  $N_e$ , maintaining a standard continuous MoE architecture, but now with  $Q \in \mathbb{R}^{M \times N}$  and  $W_i^{2,\text{in}} \in \mathbb{R}^{N_e \times N}$ . All other architectural objects live in the standard spaces and the total number of experts is  $M$ .

We are concerned with the thermodynamic limit where  $M, N \rightarrow \infty$  at a fixed ratio. Specifically, we define  $\kappa = \frac{M}{N}$ , where  $\kappa$  is a constant of order unity in  $N$ . The internal expert width  $N_e$  will also remain strictly order unity in this setting. The forward pass, following the *MSSP*, is defined as:

$$\begin{aligned}
 \mathbb{R}^N \ni h_\mu^1 &= \frac{1}{\sqrt{D}} W^1 x_\mu & \mathbb{R}^M \ni \psi_\mu &= \frac{1}{\sqrt{N}} Q \sigma(h_\mu^1) \\
 \mathbb{R}^{N_e} \ni h_{\mu,i}^{2,\text{in}} &= \frac{1}{\sqrt{N}} W_i^{2,\text{i}} \sigma(h_\mu^1) & \mathbb{R}^M \ni \phi_\mu &= \text{softmax}(\psi_\mu) \\
 \mathbb{R}^N \ni h_{\mu,i}^{2,\text{o}} &= \sqrt{\frac{\kappa N}{N_e}} W_i^{2,\text{o}} \sigma(h_{\mu,i}^{2,\text{in}}) & \mathbb{R}^N \ni h_\mu^3 &= \sum_{i=1}^M \phi_\mu^i h_{\mu,i}^{2,\text{o}} \\
 \mathbb{R} \ni h_\mu^4 &= \frac{1}{\sqrt{N}} w^{4\text{T}} h_\mu^3 & f_\mu &= \frac{1}{\gamma} h_\mu^4
 \end{aligned} \tag{L.1}$$

### L.2. Learning rates and initialization

To ensure that the network evolves dynamically in the infinite-width limit without gradients vanishing or exploding, we apply specific learning rate scalings:

$$\eta = \eta_0 \gamma^2 \quad \eta_Q = \eta_0 \gamma^2 \kappa \quad \gamma = \gamma_0 \sqrt{N} \quad \eta_0, \gamma_0 \sim O(1) \tag{L.2}$$

The network weights are initialized from standard Gaussian distributions:

$$w_\alpha^4(0), [W_i^{2,\text{out}}(0)]_{\alpha\beta}, [W_i^{2,\text{in}}(0)]_{\alpha\beta}, W_{\alpha\beta}^1(0), Q_{\alpha\beta}(0) \sim \mathcal{N}(0, 1) \tag{L.3}$$

For brevity in our derivations, we denote the collection of all network parameters as  $\theta = \text{Vec}\{W^1, W_i^{2,\text{in}}, W_i^{2,\text{out}}, w^4, Q\}$ .

Due to the normalization constraint of the softmax operator over  $M$  experts, the routing probabilities scale as  $\phi^i \sim O(\frac{1}{\kappa N})$ . It is mathematically tidier to formulate the DMFT using variables which remain  $O(1)$  in the limit. We therefore define and track the rescaled routing variables:

$$\tilde{\phi}_\mu^i := \kappa N \phi_\mu^i \tag{L.4}$$

### L.3. Gradient definitions

We mathematically *define* the pre-activation gradients, ensuring they contain the correct scaling factors to remain finite in the  $N \rightarrow \infty$  limit:

$$\begin{aligned}
 g_\mu^1 &:= \sqrt{N} \frac{\partial h_\mu^4}{\partial h_\mu^1} & g_{\mu,i}^{2,\text{in}} &:= \kappa \sqrt{N} \frac{\partial h_\mu^4}{\partial h_{\mu,i}^{2,\text{in}}} \\
 g_{\mu,i}^{2,\text{out}} &:= \kappa N^{\frac{3}{2}} \frac{\partial h_\mu^4}{\partial h_{\mu,i}^{2,\text{o}}} & g_\mu^3 &:= \sqrt{N} \frac{\partial h_\mu^4}{\partial h_\mu^3} = w^4 = z^3 \\
 g_\mu^\phi &:= \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial \phi_\mu} & z_{\mu,i}^{2,\text{out}} &:= \kappa N \phi_\mu^i g_\mu^3 = \tilde{\phi}_\mu^i g_\mu^3 \\
 z_{\mu,i}^{2,\text{in}} &:= \sqrt{\frac{\kappa}{N_e N}} W_i^{2,\text{out} \top} g_{\mu,i}^{2,\text{out}} & \tilde{z}_\mu^\phi &:= \frac{1}{\kappa \sqrt{N}} Q^\top g^{1,\phi} \\
 \tilde{z}_{\mu,i}^1 &:= \frac{1}{\kappa \sqrt{N}} W_i^{2,\text{in} \top} g_{\mu,i}^{2,\text{in}} & z^1 &:= \frac{1}{\kappa \sqrt{N}} Q^\top g^{1,\phi} + \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^M W_i^{2,\text{in} \top} g_{\mu,i}^{2,\text{in}} \\
 & & &= \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^M \tilde{z}_{\mu,i}^1 + \tilde{z}_\mu^{1,\phi}
 \end{aligned} \tag{L.5}$$

Note that the components of the router gradient  $g_\mu^\phi$  naturally scale as  $O(1)$ :

$$(g_\mu^\phi)_i = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial \phi_\mu^i} = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial h_\mu^3} \cdot \frac{\partial h_\mu^3}{\partial \phi_\mu^i} = \frac{1}{N} g^3 \cdot h_{\mu,i}^{2,\text{out}} = O(1). \tag{L.6}$$

The backward pass through the router dictates that the gradient at layer 1 involves the quantity  $g^{1,\phi}$ , defined component-wise to handle the Jacobian of the softmax:

$$g_k^{1,\phi} := \kappa N \sum_{i=1}^{\kappa N} g_i^\phi \phi^i (\delta_{ik} - \phi^k) = \sum_{i=1}^{\kappa N} g_i^\phi \tilde{\phi}^i (\delta_{ik} - \frac{\tilde{\phi}^k}{\kappa N}) \tag{L.7}$$

#### L.4. DMFT kernels and order parameters

We define the following macroscopic *kernels*, all of which are  $O_N(1)$  and will serve as the fundamental order parameters in the dynamics:

$$\begin{aligned}
 \Phi_{\mu\nu}^0 &= \frac{1}{D} x_\mu \cdot x_\nu & \Phi_{\mu\nu}^1(s, t) &= \frac{1}{N} \sigma(h_\mu^1(s)) \cdot \sigma(h_\nu^1(t)) \\
 \Phi_{\mu\nu,i}^{2,\text{in}}(s, t) &= \frac{1}{N_e} \sigma(h_{\mu,i}^{2,\text{in}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) & \Phi_{\mu\nu,i}^{2,\text{o}}(s, t) &= \frac{1}{N} \sigma(h_{\mu,i}^{2,\text{out}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{out}}(t)) \\
 \Phi_{\mu\nu}^3(s, t) &= \frac{1}{N} h_\mu^3(s) \cdot h_\nu^3(t) & G_{\mu\nu}^1(s, t) &= \frac{1}{N} g_\mu^1(s) \cdot g_\nu^1(t) \\
 G_{\mu\nu,i}^{2,\text{in}}(s, t) &= \frac{1}{N_e} g_{\mu,i}^{2,\text{in}}(s) \cdot g_{\nu,i}^{2,\text{in}}(t) & G_{\mu\nu,i}^{2,\text{out}}(s, t) &= \frac{1}{N} g_{\mu,i}^{2,\text{out}}(s) \cdot g_{\nu,i}^{2,\text{out}}(t) \\
 G_{\mu\nu}^{1,\phi}(s, t) &= \frac{1}{\kappa N} g_\mu^{1,\phi}(s) \cdot g_\nu^{1,\phi}(t) & G_{\mu\nu}^3(s, t) &= \frac{1}{N} g_\mu^3(s) \cdot g_\nu^3(t)
 \end{aligned} \tag{L.8}$$

With the exception of the input data Gram matrix  $\Phi^0$ , all of these kernels are dynamically evolving macroscopic variables.

### L.5. Learning Dynamics and the Neural Tangent Kernel

We train the network using continuous-time gradient flow. The variable learning rate scheme defined previously is necessary so that the activation updates within the experts and for the router scores do not vanish or explode as  $N, M \rightarrow \infty$ .

Taking a standard empirical risk minimization loss of the form:

$$\mathcal{L} = \frac{1}{P} \sum_{\mu=1}^P \ell(f_{\mu}, y_{\mu}) \quad (\text{L.9})$$

The network parameters evolve according to:

$$\begin{aligned} \frac{d\theta}{dt} &= \frac{\eta}{P\gamma} \sum_{\mu} \Delta_{\mu} \frac{\partial h_{\mu}^4}{\partial \theta} \\ \Delta_{\mu} &= -\frac{\partial \mathcal{L}}{\partial f_{\mu}} \end{aligned} \quad (\text{L.10})$$

For a Mean Squared Error (MSE) loss where  $\mathcal{L} = \frac{1}{P} \sum_{\nu} (y_{\nu} - f_{\nu})^2$ , the error signal (or residual) simplifies to  $\Delta_{\nu} = 2(y_{\nu} - f_{\nu})$ .

The logits update dynamically via the chain rule, driven by the Neural Tangent Kernel (NTK):

$$\frac{df_{\mu}(t)}{dt} = \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{d\theta}{dt} = \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{\eta_{\theta}}{P} \sum_{\alpha} \Delta_{\alpha} \frac{\partial f_{\alpha}(t)}{\partial \theta} = \frac{\eta}{P} \sum_{\alpha} \Delta_{\alpha} K_{\mu\alpha}^{\text{NTK}}(t, t) \quad (\text{L.11})$$

Where the infinite-width NTK is defined as:

$$K_{\mu\alpha}^{\text{NTK}}(t, s) \equiv \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{\partial f_{\alpha}(s)}{\partial \theta} \quad (\text{L.12})$$

### L.6. Decomposition of the MoE NTK

To explicitly compute the NTK, we require that the magnitude  $\eta_{\theta} \frac{\partial f_{\mu}(t)}{\partial \theta_i} \cdot \frac{\partial f_{\alpha}(s)}{\partial \theta_i}$  remains order 1 for each parameter block  $\theta_i$ . We evaluate this constraint block by block.

For the expert input layer  $W^{2,\text{in}}$ :

$$\begin{aligned} & \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,j=1}^{N_e} \sum_{k=1}^N \frac{\partial h^4}{\partial [h_i^{2,\text{in}}]_l} \frac{\partial [h_i^{2,\text{in}}]_l}{\partial W_{jk,i}^2} \frac{\partial h^4}{\partial [h_i^{2,\text{in}}]_m} \frac{\partial [h_i^{2,\text{in}}]_m}{\partial W_{jk,i}^2} = \\ &= \eta_0 N \sum_{i=1}^{\kappa N} \sum_{l,m,j=1}^{N_e} \sum_{k=1}^N \frac{[g_i^{2,\text{in}}]_l}{\kappa\sqrt{N}} \frac{[g_i^{2,\text{in}}]_m}{\kappa\sqrt{N}} \frac{\sigma(h_k^1)}{\sqrt{N}} \frac{\sigma(h_k^1)}{\sqrt{N}} \delta_{mj} \delta_{lj} \\ &= \eta_0 \frac{N_e}{\kappa} \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \left[ \frac{1}{N_e} \sum_{j=1}^{N_e} [g_i^{2,\text{in}}]_j [g_i^{2,\text{in}}]_j \right] \left[ \frac{1}{N} \sum_{k=1}^N \sigma(h_k^1)^2 \right] \\ &= \eta_0 \frac{N_e}{\kappa} \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_i^{2,\text{in}} \Phi^1 = \eta_0 \frac{N_e}{\kappa} \bar{G}^2 \Phi^1 \end{aligned} \quad (\text{L.13})$$

For the expert output layer  $W^{2,\text{out}}$ :

$$\begin{aligned}
 & \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,k=1}^{N_e} \sum_{j=1}^N \frac{\partial h^4}{\partial [h_i^{2,\text{out}}]_l} \frac{\partial [h_i^{2,\text{out}}]_l}{\partial W_{jk}^{2,\text{out}}} \frac{\partial h^4}{\partial [h_i^{2,\text{out}}]_m} \frac{\partial [h_i^{2,\text{out}}]_m}{\partial W_{jk}^{2,\text{out}}} \\
 &= \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,k=1}^{N_e} \sum_{j=1}^N \frac{[g_i^{2,\text{out}}]_l}{\kappa N^{3/2}} \frac{[g_i^{2,\text{out}}]_m}{\kappa N^{3/2}} \frac{\sqrt{\kappa N}}{\sqrt{N_e}} \sigma([h_i^{2,\text{in}}]_k) \frac{\sqrt{\kappa N}}{\sqrt{N_e}} \sigma([h_i^{2,\text{in}}]_k) \delta_{mj} \delta_{lj} \\
 &= \eta_0 \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_i^{2,\text{out}} \Phi_i^{2,\text{in}}
 \end{aligned} \tag{L.14}$$

For the shared base layer  $W^1$ :

$$\eta_0 \sum_{k=1}^D \sum_{j,m=1}^N \frac{\partial h^4}{\partial h_n^1} \frac{\partial h_n^1}{\partial W_{jk}^1} \frac{\partial h^4}{\partial h_m^1} \frac{\partial h_m^1}{\partial W_{jk}^1} = \eta_0 \sum_{k=1}^D \sum_{j,m,n=1}^N \frac{g_n^1}{\sqrt{N}} \frac{g_m^1}{\sqrt{N}} \delta_{nj} \delta_{mj} \frac{x_k}{\sqrt{D}} \frac{x_k}{\sqrt{D}} = \eta_0 G^1 \Phi^0$$
(L.15)

The router weight matrix  $Q$  involves the complex Jacobian of the softmax operator. Integrating this out, we find:

$$\eta_0 \kappa \sum_{j=1}^{\kappa N} \sum_{k=1}^N \frac{\partial h^4}{\partial Q_{jk}} \frac{\partial h^4}{\partial Q_{jk}} = \eta_0 \kappa N \sum_{j=1}^{\kappa N} \sum_{i,i'=1}^{\kappa N} \Phi^1 g_i^\phi g_{i'}^\phi \phi^i \phi^{i'} (\delta_{ij} - \phi^j) (\delta_{i'j} - \phi^j) = \eta_0 \Phi^1 G^{1,\phi}$$
(L.16)

Aggregating these contributions yields the full macroscopic evolution equation:

$$\begin{aligned}
 \frac{df_\mu(t)}{dt} &= \frac{\eta_0}{P} \sum_\nu \Delta_\nu \left[ G_{\mu\nu}^1(t,t) \Phi_{\mu\nu}^0(t,t) + \frac{N_e}{\kappa} \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_{\mu\nu,i}^{2,\text{in}}(t,t) \Phi_{\mu\nu}^1(t,t) \right. \\
 &\quad \left. + \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_{\mu\nu,i}^{2,\text{out}}(t,t) \Phi_{\mu\nu,i}^{2,\text{in}}(t,t) + \Phi_{\mu\nu}^3(t,t) + \Phi_{\mu\nu}^1(t,t) G_{\mu\nu}^{1,\phi}(t,t) \right]
 \end{aligned} \tag{L.17}$$

## L.7. Evolution of weights, preactivations, and pregradients

The updates to the router weights are:

$$\begin{aligned}
 \frac{dQ_{jk}}{dt} &= \frac{\eta}{P} \sum_\mu \Delta_\mu \frac{\partial f_\mu}{\partial Q_{jk}} = \frac{\eta_0 \gamma \kappa \sqrt{N}}{P} \sum_\mu \Delta_\mu \frac{\partial h_\mu^4}{\partial Q_{jk}} = \frac{\eta_0 \gamma_0 \kappa N}{P} \sum_\mu \Delta_\mu \frac{\partial h_\mu^4}{\partial Q_{jk}} \\
 &= \frac{\eta_0 \gamma_0}{P \sqrt{N}} \sum_\mu \Delta_\mu \left[ \kappa N \sum_{a=1}^M g_a^\phi \phi_a (\delta_{aj} - \phi_j) \right] \sigma(h_k^1) \\
 &= \frac{\eta_0 \gamma_0}{P \sqrt{N}} \sum_\mu \Delta_\mu g_j^{1,\phi} \sigma(h_k^1)
 \end{aligned} \tag{L.18}$$

We obtain thus:

$$\begin{aligned}
 \psi_\mu(t) &= \frac{1}{\sqrt{N}} Q(t) \sigma(h_\mu^1(t)) \\
 &= \chi_\mu^\phi(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) \Phi_{\mu\nu}^1(s, t) g_\nu^{1,\phi}(t) \\
 \tilde{z}_\mu^{1,\phi} &= \tilde{\xi}_\mu^{1,\phi} + \frac{\eta_0 \gamma_0}{P} \int_0^\infty dt \sum_\nu \Delta_\nu G_{\nu\mu}^{1,\phi} \sigma(h_{\nu j}^1)
 \end{aligned} \tag{L.19}$$

The other preactivation and pregradient updates follow similarly:

$$\begin{aligned}
 \frac{dW^1(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{g_\mu^1(t)}{\sqrt{N}} \frac{1}{\sqrt{D}} x_\mu \\
 \Rightarrow h_\mu^1(t) &= \frac{1}{\sqrt{D}} W^1(0) x_\mu + \int_0^t ds \frac{\eta_0 \gamma_0}{P} \sum_\nu \Delta_\nu(s) g_\nu^1(s) \Phi_{\mu\nu}^0(t, s) \\
 \frac{dW_i^{2,\text{in}}(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{g_{\mu,i}^{2,\text{in}}(t)}{\kappa \sqrt{N}} \frac{1}{\sqrt{N}} \sigma(h_\mu^1(t)) \\
 \Rightarrow h_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{in}}(0) \sigma(h_\mu^1(t)) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) g_{\nu,i}^{2,\text{in}}(s) \Phi_{\mu\nu}^1(t, s) \\
 \frac{dW_i^{2,\text{out}}(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{g_{\mu,i}^{2,\text{out}}(t)}{N^{3/2} \kappa} \frac{\sqrt{\kappa N}}{\sqrt{N_e}} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \\
 &= \frac{\eta_0 \gamma_0}{P \sqrt{\kappa}} \frac{1}{\sqrt{N N_e}} \sum_\mu \Delta_\mu(t) g_i^{2,\text{out}}(t) \sigma(h_i^{2,\text{in}}(t)) \\
 \Rightarrow h_{\mu,i}^{2,\text{out}}(t) &= \frac{\sqrt{\kappa N}}{\sqrt{N_e}} W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) g_{\nu,i}^{2,\text{out}}(s) \Phi_{\mu\nu,i}^{2,\text{in}}(t, s) \\
 \frac{dw^4(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{1}{\sqrt{N}} h_\mu^3(t) = \frac{\eta_0 \gamma_0}{P} \sum_\mu \Delta_\mu(t) h_\mu^3(t)
 \end{aligned} \tag{L.20}$$

Using the same expressions for the evolution of the weights to derive expressions for the pregradients:

$$\begin{aligned}
 z_\mu^3(t) &= g_\mu^3(t) = w_\mu^4(t) = \xi_\mu^3(0) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 z_{\mu,i}^{2,\text{out}}(t) &= \kappa g_{\mu,i}^{2,\text{out}}(t) = \tilde{\phi}_\mu^i(t) g_\nu^3(t) = \frac{\eta_0 \gamma_0}{P} \tilde{\phi}_\mu^i(t) \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 z_{\mu,i}^{2,\text{in}}(t) &= \sqrt{\frac{\kappa}{NN_e}} W_i^{2,\text{out}\top}(t) g_{\mu,i}^{2,\text{out}}(t) \\
 &= \xi_{\mu,i}^{2,\text{in}}(t) + \frac{\eta_0 \gamma_0}{PN_e} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\nu\mu,i}^{2,\text{out}}(s, t) \sigma(h_{\nu,i}^{2,\text{in}}(s)) \\
 z_{\mu,i}^1(t) &= W_i^{2,\text{in}\top}(t) g_i^{2,\text{in}}(t) \\
 &= \tilde{\xi}_{\mu,i}^1(t) + \frac{\eta_0 \gamma_0 N_e}{\sqrt{NP}} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\nu\mu,i}^{2,\text{in}}(s, t) \sigma(h_\nu^1(s))
 \end{aligned} \tag{L.21}$$

Finally for the router gradient:

$$\begin{aligned}
 g_i^{1,\phi} &= \frac{1}{\kappa N} \sum_{m=1}^{\kappa N} g_m^\phi \tilde{\phi}^m (\kappa N \delta_{mi} - \tilde{\phi}^i) \\
 &= \xi_i^{g^{1,\phi}} + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu G^3 \Phi_i^{2,\text{in}} \tilde{\phi}^i \tilde{\phi}^i - \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu G^3 \bar{\Phi}^{2,\text{in}} \tilde{\phi}^i.
 \end{aligned} \tag{L.22}$$

### L.8. Stochastic initial fields

To isolate the purely deterministic part of the trajectory from the random initialization, we define a set of initial stochastic fields  $\mathcal{F}$ . These fields encapsulate the randomness of the initial weights:

$$\mathcal{F} = \{\chi_\mu^1(t), \chi_{\mu,i}^{2,\text{in}}(t), \chi_{\mu,i}^{2,\text{out}}(t), \chi_\mu^\phi(t), \tilde{\xi}_{\mu,i}^1(t), \xi_{\mu,i}^{2,\text{in}}(t), \xi_\mu^3(t), \xi_\mu^\phi(t), \xi_{\mu i}^{g^\phi}(t)\}_{i \in \{1, \dots, M\}, \mu \in \{1, \dots, P\}} \tag{L.23}$$

$$\begin{aligned}
 \chi_\mu^1 &= \frac{1}{\sqrt{D}} W^1(0) x_\mu \\
 \chi_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{in}}(0) \sigma(h_\mu^1(t)) \\
 \chi_{\mu,i}^{2,\text{out}}(t) &= \sqrt{\frac{\kappa N}{N_e}} W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) \\
 \chi_\mu^\phi(t) &= \frac{1}{\sqrt{N}} Q(0) \sigma(h_\mu^1(t)) \\
 \tilde{\xi}_{\mu,i}^1(t) &= (W_i^{2,\text{in}}(0))^\top g_{\mu,i}^{2,\text{in}}(t) \\
 \xi_{\mu,i}^{2,\text{in}}(t) &= \sqrt{\frac{\kappa}{N N_e}} (W_i^{2,\text{out}}(0))^\top g_{\mu,i}^{2,\text{out}}(t) \\
 \xi_\mu^{1,\phi}(t) &= \frac{1}{\kappa \sqrt{N}} Q(0)^\top g_\mu^{1,\phi}(t) \\
 \xi_{\mu i}^{g^\phi}(t) &= \frac{1}{\sqrt{N_e N}} g_\mu^{3\top}(t) W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t))
 \end{aligned} \tag{L.24}$$

A critical step in rendering the MoE partition function computationally tractable is distinguishing between expert-local fields (which are specific to an expert  $i$ ) and global fields (which impact the shared base layer or the final aggregated output). This distinction allows the massive partition function to factorize over the experts.

We define the router-averaged stochastic fields to capture the macroscopic effect of the local processes:

$$\bar{\chi}_\mu^3(t) = \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \tilde{\phi}_\mu^i(t) \chi_{\mu,i}^{2,\text{out}}(t) \tag{L.25}$$

$$\bar{\xi}_\mu^1(t) = \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^{\kappa N} \tilde{\xi}_{\mu,i}^1(t) \tag{L.26}$$

Substituting these, the global aggregated output  $h_\mu^3$  and the gradient arriving at the base layer  $z_\mu^1$  can be expressed entirely in terms of order parameters that are smooth over the ensemble of experts:

$$h_\mu^3(t) = \bar{\chi}_\mu^3(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) z_\nu^3(s) \tag{L.27}$$

$$z_\mu^1(t) = \bar{\xi}_\mu^1(t) + \xi_\mu^{1,\phi}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) \left[ \frac{1}{\sqrt{N}} \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) + G_{\mu\nu}^{1,\phi}(s, t) \right] \sigma(h_\nu^1(s)) \tag{L.28}$$

### L.9. Deriving the DMFT Action

We formulate the DMFT by writing the moment-generating function for the system's trajectories and performing a disorder average over the initial weights  $\theta_0$ .

$$\begin{aligned}
 Z \propto & \left\langle \int d\mathcal{F} \exp \left( \sum_{\mu} \int_0^{\infty} dt i \left[ \hat{\chi}_{\mu}^1 \cdot \left( \chi_{\mu}^1 - \frac{1}{\sqrt{D}} W^1(0) x_{\mu} \right) \right. \right. \right. \\
 & + \sum_{i=1}^{\kappa N} \hat{\chi}_{\mu,i}^{2,\text{in}} \cdot \left( \chi_{\mu,i}^{2,\text{in}} - \frac{1}{\sqrt{N}} W_i^{2,\text{in}}(0) \sigma(h_{\mu}^1(t)) \right) \\
 & + \hat{\chi}_{\mu}^3 \cdot \left( \bar{\chi}_{\mu}^3 - \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \tilde{\phi}_{\mu}^i(t) \sqrt{\frac{\kappa N}{N_e}} W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) \right) \\
 & + \hat{\xi}_{\mu}^1 \cdot \left( \bar{\xi}_{\mu}^1 - \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^{\kappa N} W_i^{2,\text{in}}(0)^{\top} g_{\mu,i}^{2,\text{in}}(t) \right) \\
 & + \sum_{i=1}^{\kappa N} \hat{\xi}_{\mu,i}^{2,\text{in}} \cdot \left( \xi_{\mu,i}^{2,\text{in}} - \sqrt{\frac{\kappa}{N N_e}} W_i^{2,\text{out}}(0)^{\top} g_{\mu,i}^{2,\text{out}}(t) \right) + \hat{\xi}_{\mu}^3 \cdot \left( \xi_{\mu}^3 - w^4(0)^{\top} \right) \\
 & + \hat{\chi}_{\mu}^{\phi} \cdot \left( \chi_{\mu}^{\phi} - \frac{1}{\sqrt{N}} Q(0) \sigma(h_{\mu}^1(t)) \right) + \hat{\xi}_{\mu}^{1,\phi} \cdot \left( \xi_{\mu}^{1,\phi} - \frac{1}{\kappa \sqrt{N}} Q(0)^{\top} g_{\mu}^{1,\phi}(t) \right) \\
 & \left. + \sum_{m=1}^{\kappa N} \hat{\xi}_{\mu,m}^{g^{1,\phi}}(t) \left( \xi_{\mu,m}^{g^{1,\phi}}(t) - \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \frac{1}{\sqrt{N_e} N} g_{\mu}^{3T}(t) W_i^{2,\text{out}} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \tilde{\phi}_{\mu}^i(t) \left( \kappa N \delta_{im} - \tilde{\phi}_{\mu}^m(t) \right) \right) \right] \right) \\
 & \times \exp \left( \sum_{\mu} \int_0^{\infty} dt \left[ j_{\mu}^1(t) \cdot \chi_{\mu}^1(t) + \sum_{i=1}^M j_{\mu,i}^{2,\text{in}}(t) \cdot \chi_{\mu,i}^{2,\text{in}}(t) + \bar{j}_{\mu}^3(t) \cdot \bar{\chi}_{\mu}^3(t) + \sum_{i=1}^M v_{\mu,i}^{2,\text{in}}(t) \cdot \xi_{\mu,i}^{2,\text{in}}(t) \right. \right. \\
 & \left. \left. + \bar{v}_{\mu}^1(t) \cdot \bar{\xi}_{\mu}^1(t) + v_{\mu}^3 \cdot \xi_{\mu}^3(t) + j_{\mu}^{\phi}(t) \cdot \chi_{\mu}^{\phi}(t) + v_{\mu}^{1,\phi}(t) \cdot \xi_{\mu}^{1,\phi}(t) + v_{\mu}^{g^{1,\phi}}(t) \cdot \xi_{\mu}^{g^{1,\phi}}(t) \right] \right) \right\rangle_{\theta_0}. \tag{L.29}
 \end{aligned}$$

The resulting partition function obtained after the integration over the initial weights  $\theta_0$  is:

$$\begin{aligned}
 Z \propto \int d\mathcal{F} \exp \left\{ -\frac{1}{2} \sum_{\mu\nu} \int_0^\infty dt \int_0^\infty ds \left[ \hat{\chi}_\mu^1(t) \cdot \hat{\chi}_\nu^1(s) \Phi_{\mu\nu}^0(t, s) + \left( \sum_{i=1}^M \hat{\chi}_{\mu,i}^{2,\text{in}}(t) \cdot \hat{\chi}_{\nu,i}^{2,\text{in}}(s) \right. \right. \right. \\
 + \hat{\chi}_\mu^\phi(t) \cdot \hat{\chi}_\nu^\phi(s) \Phi_{\mu\nu}^1(t, s) + \frac{N_e}{\kappa} \hat{\xi}_\mu^1(s) \cdot \hat{\xi}_\nu^1(t) \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) \\
 + \frac{\kappa}{N_e} \sum_{i=1}^M \hat{\xi}_{\mu,i}^{2,\text{in}}(s) \cdot \hat{\xi}_{\nu,i}^{2,\text{in}}(t) G_{\mu\nu}^3(s, t) \tilde{\phi}_\mu^i(s) \tilde{\phi}_\nu^i(t) \\
 + \hat{\xi}_\mu^3(s) \cdot \hat{\xi}_\nu^3(t) + \frac{1}{\kappa} \hat{\xi}_\mu^{1,\phi}(s) \cdot \hat{\xi}_\nu^{1,\phi}(t) G_{\mu\nu}^{1,\phi}(t, s) \\
 + \kappa \sum_{i=1}^{\kappa N} \hat{\xi}_{\mu,i}^{g^{1,\phi}}(s) \hat{\xi}_{\nu,i}^{g^{1,\phi}}(t) \Phi_{\mu\nu,i}^{2,\text{in}}(s, t) G_{\mu\nu}^3(s, t) \tilde{\phi}_\mu^i(s) \tilde{\phi}_\nu^i(t) \\
 \left. \left. \left. + \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) \left( \hat{\chi}_\mu^3(s) \cdot \hat{\chi}_\nu^3(t) \right) \right] - i \sum_{\mu,\nu} \int_0^\infty dt \int_0^\infty ds \left[ \bar{A}_{\mu\nu}^1(s, t) \frac{1}{\kappa} \sum_{i=1}^{\kappa N} \hat{\chi}_{\mu,i}^{2,\text{in}}(s) \cdot g_{\nu,i}^{2,\text{in}}(t) \right. \right. \\
 + \frac{1}{\kappa} \hat{\chi}_\mu^\phi(s) \cdot g_{\nu}^{1,\phi}(t) A_{\mu\nu}^\phi(s, t) + \frac{1}{N_e} \sum_{i=1}^{\kappa N} \left( \hat{\xi}_{\mu,i}^{2,\text{in}}(s) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) \right) \tilde{\phi}_\mu^i(s) \tilde{\phi}_\nu^i(t) \bar{B}_{\mu\nu}^{2,\text{in}}(s, t) \\
 + \sum_{i=1}^{\kappa N} \hat{\xi}_{\mu,i}^{g^{1,\phi}}(t) \tilde{\phi}_\nu^i(s) \bar{A}_{\mu\nu,i}^{\tilde{\phi}}(t, s) \left. \right] + \sum_{\mu} \int_0^\infty dt \left[ j_\mu^1(t) \cdot \chi_\mu^1(t) + \sum_{i=1}^M j_{\mu,i}^{2,\text{in}}(t) \cdot \chi_{\mu,i}^{2,\text{in}}(t) + \bar{j}_\mu^3(t) \cdot \bar{\chi}_\mu^3(t) \right. \\
 \left. \left. + \sum_{i=1}^M v_{\mu,i}^{2,\text{in}}(t) \cdot \xi_{\mu,i}^{2,\text{in}}(t) + \bar{v}_\mu^1(t) \cdot \bar{\xi}_\mu^1(t) + v_\mu^3 \cdot \xi_\mu^3(t) + j_\mu^\phi(t) \cdot \chi_\mu^\phi(t) + v_\mu^{1,\phi}(t) \cdot \xi_\mu^{1,\phi}(t) \right] \right\} \quad (\text{L.30})
 \end{aligned}$$

Where we have introduced the  $O(1)$  kernels (also called response functions)

$$\begin{aligned}
 \bar{B}_{\mu\nu}^{2,\text{in}}(s, t) &= -\frac{i}{N} \hat{\chi}_\mu^3(s) \cdot g_\nu^3(t) \\
 \bar{A}_{\mu\nu}^1(s, t) &= -\frac{i}{N} \hat{\xi}_\mu^1(s) \cdot \sigma(h_\nu^1(t)) \\
 A_{\mu\nu}^{1,\phi}(s, t) &= -\frac{i}{N} \hat{\xi}_\mu^{1,\phi}(s) \cdot \sigma(h_\nu^1(t)) \\
 \bar{A}_{\mu\nu,i}^{\tilde{\phi}}(t, s) &= \kappa \bar{\Phi}_{\mu\nu,i}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s) \tilde{\phi}_\mu^i(t) A_{\mu\nu}^{g^{1,\phi}}(t, s) + \kappa A_{\mu\nu,i}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s) \tilde{\phi}_\mu^i(t) \\
 &\quad + \Phi_{\mu\nu,i}^{2,\text{in}}(t, s) \tilde{\phi}_\mu^i(t) \bar{B}_{\mu\nu}^{2,\text{in}}(t, s) \quad + \frac{1}{2} \kappa \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s) A_{\mu\nu}^{g^{1,\phi}}(t, s) \\
 &\quad + \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, s) \bar{B}_{\mu\nu}^{2,\text{in}}(t, s) + \kappa \bar{A}_{\mu\nu}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s)
 \end{aligned} \quad (\text{L.31})$$

We must constrain the solutions to be physical, so for each  $i, \mu, \nu, s, t$  we multiply in the following resolutions of the identity to enforce the definitions of the order parameters.

$$1 = \int \frac{d\Phi_{\mu\nu}^1(s, t) d\hat{\Phi}_{\mu\nu}^1(s, t)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^1(s, t) \left( N \Phi_{\mu\nu}^1(s, t) - \sigma(h_\mu^1(s)) \cdot \sigma(h_\nu^1(t)) \right) \right] \quad (\text{L.32})$$

$$1 = \int \frac{d\Phi_{\mu\nu,i}^{2,\text{in}}(s,t) d\hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(s,t)}{2\pi i N_e^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(s,t) \left( N_e \Phi_{\mu\nu,i}^{2,\text{in}}(s,t) - \sigma(h_{\mu,i}^{2,\text{in}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) \right) \right] \quad (\text{L.33})$$

$$1 = \int \frac{dG_{\mu\nu}^1(s,t) d\hat{G}_{\mu\nu}^1(s,t)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^1(s,t) \left( N G_{\mu\nu}^1(s,t) - g_{\mu}^1(s) \cdot g_{\nu}^1(t) \right) \right] \quad (\text{L.34})$$

$$1 = \int \frac{dG_{\mu\nu,i}^{2,\text{in}}(s,t) d\hat{G}_{\mu\nu,i}^{2,\text{in}}(s,t)}{2\pi i N_e^{-1}} \exp \left[ \hat{G}_{\mu\nu,i}^{2,\text{in}}(s,t) \left( N_e G_{\mu\nu,i}^{2,\text{in}}(s,t) - g_{\mu,i}^2(s) \cdot g_{\nu,i}^2(t) \right) \right] \quad (\text{L.35})$$

$$1 = \int \frac{dG_{\mu\nu}^3(s,t) d\hat{G}_{\mu\nu}^3(s,t)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^3(s,t) \left( N G_{\mu\nu}^3(s,t) - g_{\mu}^3(s) \cdot g_{\nu}^3(t) \right) \right] \quad (\text{L.36})$$

$$1 = \int \frac{dG_{\mu\nu}^{1,\phi}(s,t) d\hat{G}_{\mu\nu}^{1,\phi}(s,t)}{2\pi i \kappa N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^{1,\phi}(s,t) \left( \kappa N G_{\mu\nu}^{1,\phi}(s,t) - g_{\mu}^{1,\phi}(s) \cdot g_{\nu}^{1,\phi}(t) \right) \right] \quad (\text{L.37})$$

$$1 = \int \frac{d\bar{B}_{\mu\nu}^{2,\text{in}}(s,t) d\bar{A}_{\mu\nu}^{2,\text{in}}(s,t)}{2\pi i N^{-1}} \exp \left[ -\bar{A}_{\mu\nu}^{2,\text{in}}(s,t) \left( N \bar{B}_{\mu\nu}^{2,\text{in}}(s,t) + i \hat{\chi}_{\mu}^3(s) \cdot g_{\nu}^3(t) \right) \right] \quad (\text{L.38})$$

$$1 = \int \frac{d\bar{A}_{\mu\nu}^1(s,t) d\bar{B}_{\mu\nu}^1(s,t)}{2\pi i N^{-1}} \exp \left[ -\bar{B}_{\mu\nu}^1(s,t) \left( N \bar{A}_{\mu\nu}^1(s,t) + i \hat{\xi}_{\mu}^1(s) \cdot \sigma(h_{\nu}^1(t)) \right) \right] \quad (\text{L.39})$$

$$1 = \int \frac{dA_{\mu\nu}^{1,\phi}(s,t) dB_{\mu\nu}^{1,\phi}(s,t)}{2\pi i N^{-1}} \exp \left[ -B_{\mu\nu}^{1,\phi}(s,t) \left( N A_{\mu\nu}^{1,\phi}(s,t) + i \hat{\xi}_{\mu}^{1,\phi}(s) \cdot \sigma(h_{\nu}^1(t)) \right) \right] \quad (\text{L.40})$$

$$1 = \int \frac{dA_{\mu\nu}^{g^{1,\phi}}(s,t) dB_{\mu\nu}^{g^{1,\phi}}(s,t)}{2\pi i (\kappa N)^{-1}} \exp \left[ -B_{\mu\nu}^{g^{1,\phi}}(s,t) \left( \kappa N A_{\mu\nu}^{g^{1,\phi}}(s,t) + i \hat{\xi}_{\mu}^{g^{1,\phi}}(s) \cdot \tilde{\phi}_{\nu}(t) \right) \right] \quad (\text{L.41})$$

$$1 = \int \frac{dA_{\mu\nu,i}^{2,\text{in}}(s,t) dB_{\mu\nu,i}^{2,\text{in}}(s,t)}{2\pi i N_e^{-1}} \exp \left[ B_{\mu\nu,i}^{2,\text{in}}(s,t) \left( N_e A_{\mu\nu,i}^{2,\text{in}}(s,t) + i \hat{\xi}_{\mu,i}^{2,\text{in}}(s) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) \right) \right] \quad (\text{L.42})$$

$$1 = \int \frac{d\Phi_{\mu\nu}^3(t,s) d\hat{\Phi}_{\mu\nu}^3(t,s)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^3(t,s) \left( N \Phi_{\mu\nu}^3(t,s) - h_{\mu}^3(t) \cdot h_{\nu}^3(s) \right) \right] \quad (\text{L.43})$$

$$1 = \int \frac{d\bar{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) d\hat{\bar{\Phi}}_{\mu\nu}^{2,\text{in}}(t,s)}{2\pi i} \exp \left[ \hat{\bar{\Phi}}_{\mu\nu}^{2,\text{in}}(t,s) \left( \kappa N \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) - \sum_{i=1}^{\kappa N} \tilde{\phi}_{\mu}^i(s) \tilde{\phi}_{\nu}^i(t) \Phi_{\mu\nu,i}^{2,\text{in}}(s,t) \right) \right] \quad (\text{L.44})$$

$$1 = \int \frac{d\bar{G}_{\mu\nu}^{2,\text{in}}(t,s) d\hat{\bar{G}}_{\mu\nu}^{2,\text{in}}(t,s)}{2\pi i (\kappa N)^{-1}} \exp \left[ \hat{\bar{G}}_{\mu\nu}^{2,\text{in}}(t,s) \left( \kappa N \bar{G}_{\mu\nu}^{2,\text{in}}(t,s) - \sum_{i=1}^{\kappa N} G_{\mu\nu,i}^{2,\text{in}}(s,t) \right) \right] \quad (\text{L.45})$$

### L.10. Softmax

We need to encode the behaviour of softmax in the DMFT. Recall that  $\phi^i = \text{softmax}(\psi^i)$ , so  $\tilde{\phi}^i = \kappa N \phi^i = \kappa N \text{softmax}(\psi^i)$ . That is,

$$\tilde{\phi}_\mu^i(t) = \kappa N \phi_\mu^i(t) = \frac{e^{\psi_\mu^i(t)}}{\frac{1}{\kappa N} \sum_{j=1}^{\kappa N} e^{\psi_\mu^j(t)}} \quad (\text{L.46})$$

This invites the definition for each input  $\mu$  of the *global* order parameter

$$\mathcal{S}_\mu(t) := \frac{1}{\kappa N} \sum_{j=1}^{\kappa N} e^{\psi_\mu^j(t)}, \quad (\text{L.47})$$

which we enforce via the Fourier-transformed delta function

$$1 = \int \frac{d\mathcal{S}_\mu(t) d\hat{\mathcal{S}}_\mu(t)}{2\pi i (\kappa N)^{-1}} \exp \left[ \hat{\mathcal{S}}_\mu(t) \left( \kappa N \mathcal{S}_\mu(t) - \sum_{i=1}^{\kappa N} e^{\psi_\mu^i(t)} \right) \right] \quad (\text{L.48})$$

### L.11. Partition function

Define the set of all kernels and conjugates (indexed for time and feature, although this is omitted for brevity in L.49):

$$\begin{aligned} \mathcal{K} = \{ & \Phi^1, \hat{\Phi}^1, \Phi_i^{2,\text{in}}, \hat{\Phi}_i^{2,\text{in}}, \bar{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \Phi^3, \hat{\Phi}^3, G^1, \hat{G}^1, \bar{G}^{2,\text{in}}, \hat{G}^{2,\text{in}}, G_i^{2,\text{in}}, \hat{G}_i^{2,\text{in}}, \\ & G^3, \hat{G}^3, G^{1,\phi}, \hat{G}^{1,\phi}, \bar{A}^1, \bar{B}^1, \bar{B}^{2,\text{in}}, \bar{A}^{2,\text{in}}, A^{1,\phi}, B^{1,\phi}, A^{g^{1,\phi}}, B^{g^{1,\phi}}, B_i^{2,\text{in}}, A_i^{2,\text{in}} \} \end{aligned} \quad (\text{L.49})$$

We can partition these order parameters into the set  $\mathcal{K}_{global}$  of global order parameters, and the set  $\mathcal{K}_{exp-local}$  of expert-local order parameters.

$$\begin{aligned} \mathcal{K}_{global} = \{ & \Phi^1, \hat{\Phi}^1, \bar{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \Phi^3, \hat{\Phi}^3, G^1, \hat{G}^1, \bar{G}^{2,\text{in}}, \hat{G}^{2,\text{in}}, G^{1,\phi}, \hat{G}^{1,\phi}, G^3, \hat{G}^3, \\ & A^1, \bar{B}^1, \bar{B}^{2,\text{in}}, \bar{A}^{2,\text{in}}, A^{1,\phi}, B^{1,\phi}, A^{g^{1,\phi}}, B^{g^{1,\phi}} \} \end{aligned} \quad (\text{L.50})$$

$$\mathcal{K}_{exp-local} = \{ \Phi_i^{2,\text{in}}, \hat{\Phi}_i^{2,\text{in}}, \Phi_i^{2,\text{out}}, \hat{\Phi}_i^{2,\text{out}}, G_i^{2,\text{in}}, \hat{G}_i^{2,\text{in}}, B_i^{2,\text{in}}, A_i^{2,\text{in}} \} \quad (\text{L.51})$$

$$Z \propto \int \left( \prod_{\mu,\nu} \prod_{t,s} d\mathcal{K}_{global} \right) \exp \left( N S[\mathcal{K}] \right), \quad \text{where} \quad (\text{L.52})$$

$$\begin{aligned}
 S[\mathcal{K}] = & \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\Phi}_{\mu\nu}^1(t,s) \Phi_{\mu\nu}^1(t,s) + \hat{G}_{\mu\nu}^1(t,s) G_{\mu\nu}^1(t,s) + \kappa \hat{G}_{\mu\nu}^{1,\phi}(t,s) G_{\mu\nu}^{1,\phi}(t,s) \right. \\
 & - \bar{A}_{\mu\nu}^{2,\text{in}}(t,s) \bar{B}_{\mu\nu}^{2,\text{in}}(t,s) - \bar{B}_{\mu\nu}^1(t,s) \bar{A}_{\mu\nu}^1(t,s) - B_{\mu\nu}^{1,\phi}(t,s) A_{\mu\nu}^{1,\phi}(t,s) - \kappa B_{\mu\nu}^{g^{1,\phi}}(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) \\
 & + \hat{\Phi}_{\mu\nu}^3(t,s) \Phi_{\mu\nu}^3(t,s) + \kappa \hat{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) + \kappa \hat{G}_{\mu\nu}^{2,\text{in}}(t,s) \bar{G}_{\mu\nu}^{2,\text{in}}(t,s) \\
 & \left. + \kappa \hat{\mathcal{S}}_{\mu}(t) \mathcal{S}_{\nu}(s) + \hat{G}_{\mu\nu}^3(t,s) G_{\mu\nu}^3(t,s) \right] \\
 & + \frac{1}{N} \sum_{n=1}^N \ln \mathcal{Z}_N^{\text{global}} [j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3] + \frac{\kappa}{\kappa N} \sum_{i=1}^{\kappa N} \ln Z^{\text{local}} [j_i^{2,\text{in}}, v_i^{2,\text{in}}, j_i^{3,\text{out}}, j_i^{\phi}]
 \end{aligned} \tag{L.53}$$

$$\begin{aligned}
 \ln \mathcal{Z}_N^{\text{global}} [j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3] = & -\frac{1}{2} \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\chi}_{\mu}^1(t) \hat{\chi}_{\nu}^1(s) \Phi_{\mu\nu}^0(t,s) + \hat{\xi}_{\mu}^3(t) \hat{\xi}_{\nu}^3(s) \right. \\
 & + \frac{1}{\kappa} \hat{\xi}_{\mu}^{1,\phi}(t) \hat{\xi}_{\nu}^{1,\phi}(s) G_{\mu\nu}^{1,\phi}(t,s) + \frac{N_e}{\kappa} \hat{\xi}_{\mu}^1(t) \hat{\xi}_{\nu}^1(s) \bar{G}_{\mu\nu}^{2,\text{in}}(t,s) + \hat{\chi}_{\mu}^3(t) \hat{\chi}_{\nu}^3(s) \bar{\Phi}^{2,\text{in}} \\
 & - \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\Phi}_{\mu\nu}^1(t,s) \sigma(h_{\mu}^1(t)) \sigma(h_{\nu}^1(s)) \right. \\
 & + \hat{G}_{\mu\nu}^1(t,s) g_{\mu}^1(t) g_{\nu}^1(s) + \hat{\Phi}_{\mu\nu}^3(t,s) h_{\mu}^3(t) h_{\nu}^3(s) + \hat{G}_{\mu\nu}^3(t,s) g_{\mu}^3(t) g_{\nu}^3(s) \\
 & - i \sum_{\mu,\nu} \int dt \int ds \left[ B_{\mu\nu}^{1,\phi}(t,s) \hat{\xi}_{\mu}^{1,\phi}(t) \sigma(h_{\nu}^1(s)) + \bar{B}_{\mu\nu}^1(t,s) \hat{\xi}_{\mu}^1(t) \sigma(h_{\nu}^1(s)) \right. \\
 & + \bar{A}_{\mu\nu}^{2,\text{in}}(t,s) \hat{\chi}_{\mu}^3(t) g_{\nu}^3(s) \\
 & + \sum_{\mu} \int dt \left[ (v_{\mu,n}^3 + i \hat{\xi}_{\mu,n}^3(t)) \xi_{\mu,n}^3(t) + (v_{\mu,n}^{1,\phi} + i \hat{\xi}_{\mu,n}^{1,\phi}(t)) \xi_{\mu,n}^{1,\phi}(t) + (j_{\mu,n}^1(t) \right. \\
 & \left. + i \hat{\chi}_{\mu,n}^1(t)) \chi_{\mu,n}^1(t) + (\bar{j}_{\mu,n}^3 + i \hat{\chi}_{\mu,n}^3(t)) \bar{\chi}_{\mu,n}^3(t) + (\bar{v}_{\mu,n}^1 + i \hat{\xi}_{\mu,n}^1(t)) \bar{\xi}_{\mu,n}^1(t) \left. \right]
 \end{aligned} \tag{L.54}$$

$$\begin{aligned}
 \mathcal{Z}^{\text{local}} [j_i^2, v_i^2, j_i^{3,i}, j_i^{\phi}] = & \int \left( \prod_{\mu\nu} \prod_{t,s} d\mathcal{K}_{\text{exp-loc}} \right) \times \exp \left( \sum_{\mu\nu} \int_0^{\infty} dt \int_0^{\infty} ds \left[ N_e \hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(t,s) \Phi_{\mu\nu,i}^{2,\text{in}}(t,s) \right. \right. \\
 & \left. \left. + N_e \hat{G}_{\mu\nu,i}^{2,\text{in}}(t,s) G_{\mu\nu,i}^{2,\text{in}}(t,s) + N_e B_{\mu\nu,i}^{2,\text{in}}(t,s) A_{\mu\nu,i}^{2,\text{in}}(t,s) \right] \right) \\
 & \times \frac{1}{N} \prod_{j=1}^{N_e} \mathcal{Z}_{N_e} \left[ [j_i^{2,\text{in}}]_j, [v_i^{2,\text{in}}]_j \right] \times \mathcal{Z}_M [j_i^{\phi}, v_i^{g^{1,\phi}}]
 \end{aligned} \tag{L.55}$$

$$\begin{aligned}
 \mathcal{Z}_{N_e j} \left[ [j_i^{2,\text{in}}]_j, [v_i^{2,\text{in}}]_j \right] &:= \int d\mathcal{F} \exp \left\{ -\frac{1}{2} \sum_{\mu,\nu} \int dt \int ds \left[ [\hat{\chi}_{\mu,i}^{2,\text{in}}(t)]_j [\hat{\chi}_{\nu,i}^{2,\text{in}}(s)]_j \Phi_{\mu\nu}^1(t,s) \right. \right. \\
 &+ \frac{\kappa}{N_e} [\hat{\xi}_{\mu,i}^{2,\text{in}}(t)]_j [\hat{\xi}_{\nu,i}^{2,\text{in}}(s)]_j G_{\mu\nu}^3(s,t) \tilde{\phi}_\mu^i(s) \tilde{\phi}_\nu^i(t) \\
 &+ \frac{1}{N_e} G_{\mu\nu}^3(s,t) \sigma([h_{\mu,i}^{2,\text{in}}(s)]_j) \sigma([h_{\nu,i}^{2,\text{in}}(t)]_j) \hat{\xi}_{\mu,i}^{g^\phi}(t) \hat{\xi}_{\nu,i}^{g^\phi}(s) \\
 &- i \sum_{\mu,\nu} \int dt \int ds \left[ \frac{1}{\kappa} [\hat{\chi}_{\mu,i}^{2,\text{in}}(t)]_j [g_{\nu,i}^{2,\text{in}}(s)]_j \bar{A}_{\mu\nu}^1(t,s) \right. \\
 &- \left. B_{\mu\nu,i}^{2,\text{in}}(t,s) [\hat{\xi}_{\mu,i}^{2,\text{in}}(t)] \sigma([h_{\nu,i}^{2,\text{in}}(s)]_j) \right] \\
 &- \sum_{\mu,\nu} \int dt \int ds [\hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(t,s) \sigma([h_{\mu,i}^{2,\text{in}}(t)]_j) \sigma([h_{\nu,i}^{2,\text{in}}(s)]_j) \\
 &+ \hat{G}_{\mu\nu,i}^{2,\text{in}}(t,s) [g_{\mu,i}^{2,\text{in}}(t)]_j [g_{\nu,i}^{2,\text{in}}(s)]_j \\
 &+ \sum_{\mu} \int dt \left( ([j_{\mu,i}^{2,\text{in}}(t)]_j + i [\hat{\chi}_{\mu,i}^{2,\text{in}}(t)]_j) [\chi_{\mu,i}^{2,\text{in}}(t)]_j \right. \\
 &\left. + ([v_{\mu,i}^{2,\text{in}}(t)]_j + i [\hat{\xi}_{\mu,i}^{2,\text{in}}(t)]_j) [\xi_{\mu,i}^{2,\text{in}}(t)]_j \right) \left. \right\}. \tag{L.56}
 \end{aligned}$$

$$\begin{aligned}
 \mathcal{Z}_M [j_m^\phi, v_m^{g^{1,\phi}}] &:= \int d\mathcal{F} \exp \left\{ -\frac{1}{2} \sum_{\mu,\nu} \int dt \int ds [\hat{\chi}_{\mu m}^\phi(t) \hat{\chi}_{\nu m}^\phi(s) \Phi_{\mu\nu}^1(t,s) \right. \\
 &+ \kappa \hat{\xi}_{\mu,m}^{g^{1,\phi}}(s) \hat{\xi}_{\nu,m}^{g^{1,\phi}}(t) \Phi_{\mu\nu,m}^{2,\text{in}}(s,t) G_{\mu\nu}^3(s,t) \tilde{\phi}_\mu^m(s) \tilde{\phi}_\nu^m(t) \\
 &- i \sum_{\mu,\nu} \int dt \int ds \left[ \frac{1}{\kappa} \hat{\chi}_{\mu m}^\phi(t) g_{\nu}^{1,\phi}(s) A_{\mu\nu}^{1,\phi}(t,s) + \hat{\xi}_{\mu m}^{g^{1,\phi}}(t) \tilde{\phi}_\nu^m(s) [B_{\mu\nu}^{g^{1,\phi}}(t,s) + \bar{A}_{\mu\nu,m}^{\tilde{\phi}}(t,s)] \right] \\
 &+ \sum_{\mu,\nu} \int_0^\infty dt \int_0^\infty ds [A_{\mu\nu,m}^{2,\text{in}}(s,t) \tilde{\phi}_\mu^m(s) \tilde{\phi}_\nu^m(t) \bar{B}_{\mu\nu}^{2,\text{in}}(s,t) \\
 &- \sum_{\mu,\nu} \int dt \int ds [\hat{G}_{\mu\nu}^{1,\phi}(t,s) g_{\mu m}^{1,\phi}(t) g_{\nu m}^{1,\phi}(s) + \hat{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \tilde{\phi}_\mu^m(t) \tilde{\phi}_\nu^m(s) \Phi_{\mu\nu,m}^{2,\text{in}}(t,s) \\
 &+ \hat{G}_{\mu\nu}^{2,\text{in}}(t,s) G_{\mu\nu,m}^{2,\text{in}}(t,s) \\
 &- \sum_{\mu} \int dt \hat{S}_\mu(t) e^{\psi_\mu^m(t)} + \sum_{\mu} \int dt (j_{\mu m}^\phi(t) + i \hat{\chi}_{\mu m}^\phi(t)) \chi_{\mu m}^\phi(t) + (v_{\mu m}^{g^\phi}(t) + i \hat{\xi}_{\mu m}^{g^{1,\phi}}(t)) \xi_{\mu m}^{g^{1,\phi}}(t) \left. \right\}. \tag{L.57}
 \end{aligned}$$

## L.12. Saddle point approximation

To write down the saddle point equations, we define first the single-site distributions.

We can write

$$\mathcal{Z}_N^{global}[j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3] = \int d\mathcal{F} \exp \left\{ -\mathcal{H}_N^{global}[\{\chi_{\mu n}^1, \bar{\chi}_{\mu n}^3, \bar{\xi}_{\mu n}^1, \xi_{\mu n}^{1,\phi}, \xi_{\mu n}^3, j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3\}_{\mu}] \right\} \quad (\text{L.58})$$

$$\mathcal{Z}_M[j_m^\phi] = \int d\mathcal{F} \exp \left\{ -\mathcal{H}_M[\{\chi_{\mu,m}^\phi, j_{\mu,m}^\phi\}_{\mu}] \right\} \quad (\text{L.59})$$

Where  $\mathcal{H}$  in each case is the logarithm of the integrand of the corresponding  $\mathcal{Z}$ . We can then define for each  $\mathcal{Z} \in \{\mathcal{Z}_N^{global}, \mathcal{Z}_{N_e j}, \mathcal{Z}_M\}$  and the corresponding  $\mathcal{H}$  the average

$$\langle \mathcal{O}(\{\chi, \xi\}) \rangle_{\mathcal{Z}} = \frac{1}{\mathcal{Z}} \int \prod_{\mu} d\mathcal{F} \exp(-\mathcal{H}[\{\chi, \xi\}, \{j, v\}]) \mathcal{O}(\{\chi, \xi\}) \quad (\text{L.60})$$

With this apparatus in place, we can take saddle equations. We treat expert-local kernels as microvariables which are implicitly defined in terms of  $\chi, \xi$ , and so do not take saddle equations of them at this point.

### L.13. Saddle-Point Equations

Since the integrand in eqn.L.52 has the form  $e^{NS[\mathcal{K}]}$ , in the  $\lim_N \rightarrow \infty$  we can use the saddle point approximation, which consist in finding the set of equations that lead to a stationary action  $S[\mathcal{K}]$ .

We note the following:

- At zero source, all single site averages  $\langle \rangle_{\mathcal{Z}_N^{global}[j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3]}$  are equivalent, so we may write  $\langle \rangle_{\mathcal{Z}_N^{global}}$  for the average over the single-site distributions when  $j^1, \bar{j}^3, \bar{v}^1, v^{1,\phi}, v^3 \rightarrow 0$ .
- *Conditional on the set  $\mathcal{K}_{global}$  of global kernels,  $\mathcal{Z}^{local}$  factorises over experts, and so we can write  $\langle \rangle_{\mathcal{Z}_{N_e j} | \mathcal{K}_{global}}$  for the conditional average over the distribution defined by  $\mathcal{Z}_{N_e j} \left[ \left[ j_i^{2,\text{in}} \right]_j, \left[ v_i^{2,\text{in}} \right]_j \right]$  as  $j^{2,\text{in}}, v^{2,\text{in}} \rightarrow 0$ , and  $\langle \rangle_{\mathcal{Z}_M | \mathcal{K}_{global}}$  for the conditional average over the distribution defined by  $\mathcal{Z}_M[j_m^\phi, v_m^{g^\phi}]$  as  $j^\phi, v^{g^{1,\phi}} \rightarrow 0$ .*
- Expert-local variables follow single-site processes, so we can drop expert indices
- By a derivation similar to that in sec.K, we can prove that conjugate kernels defined as covariances between conjugate fields vanish, since they have no physical meaning and can not influence the dynamics in addition to imposing constraint when introducing kernel definitions.

The global kernels are given by their expectation values over the single-site global distribution  $\mathcal{Z}_N^{global}$ :

$$\begin{aligned}
 \Phi_{\mu\nu}^1(s, t) &= \langle \sigma(h_\mu^1(s)) \sigma(h_\nu^1(t)) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 \Phi_{\mu\nu}^3(s, t) &= \langle h_\mu^3(s) h_\nu^3(t) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 G_{\mu\nu}^1(s, t) &= \langle g_\mu^1(s) g_\nu^1(t) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 G_{\mu\nu}^3(s, t) &= \langle g_\mu^3(s) g_\nu^3(t) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 A_{\mu\nu}^{1,\phi}(s, t) &= -i \langle \hat{\xi}_\mu^{1,\phi}(s) \sigma(h_\nu^1(t)) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 \bar{A}_{\mu\nu}^1(s, t) &= -i \langle \hat{\xi}_\mu^1(t) \sigma(h_\nu^1(s)) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 \bar{B}_{\mu\nu}^{2,\text{in}}(s, t) &= -i \langle \hat{\chi}_\mu^3(t) g_\nu^3(s) \rangle_{\mathcal{Z}_N^{\text{global}}}
 \end{aligned} \tag{L.61}$$

The router kernels require expectations over the expert ensemble distribution  $\mathcal{Z}_M$ , conditioned on the global fields:

$$\begin{aligned}
 \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \Phi_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle G_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 G_{\mu\nu}^{1,\phi}(s, t) &= \langle g_\mu^{1,\phi}(s) g_\nu^{1,\phi}(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 \mathcal{S}_\mu(t) &= \langle e^{\psi_\mu(t)} \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}
 \end{aligned} \tag{L.62}$$

#### L.14. Expert-local order parameters

The local kernels are determined by independent processes within each expert's internal dimensions:

$$\begin{aligned}
 \Phi_{\mu\nu}^{2,\text{in}}(s, t) &= \frac{1}{N_e} \sum_{j=1}^{N_e} \sigma([h_\mu^{2,\text{in}}(t)]_j) \sigma([h_\nu^{2,\text{in}}(s)]_j) \\
 G_{\mu\nu}^{2,\text{in}}(s, t) &= \frac{1}{N_e} \sum_{j=1}^{N_e} [g_\mu^{2,\text{in}}(t)]_j [g_\nu^{2,\text{in}}(s)]_j \\
 A_{\mu\nu}^{2,\text{in}}(s, t) &= i \kappa \sum_{j=1}^{N_e} [\hat{\xi}_\mu^{2,\text{in}}(s)]_j \sigma([h_\nu^{2,\text{in}}]_j) \\
 B_{\mu\nu}^{2,\text{in}}(s, t) &= -\frac{1}{N_e} \langle \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \bar{B}_{\mu\nu}^{2,\text{in}}(s, t)
 \end{aligned} \tag{L.63}$$

All non-physical conjugate kernels identically vanish at the saddle point:  $\hat{\Phi}^3 = \hat{\Phi}^1 = \hat{G}^1 = \hat{\Phi}^{2,\text{in}} = \dots = 0$ .

### L.15. Hubbard-Stratonovich transformation

The Hubbard-Stratonovich trick allows us to rewrite the quadratic terms over conjugate fields as Gaussian integrals over linear conjugate fields:

$$\exp\left(-\frac{1}{2} \hat{\mathbf{x}}^\top A \hat{\mathbf{x}}\right) = \int_{\mathbb{R}^d} \frac{d\mathbf{u}}{(2\pi)^{d/2} \sqrt{\det A}} \exp\left(-\frac{1}{2} \mathbf{u}^\top A^{-1} \mathbf{u} - i \mathbf{u} \cdot \hat{\mathbf{x}}\right) = \left\langle \exp(-i \mathbf{u} \cdot \hat{\mathbf{x}}) \right\rangle_{\mathbf{u} \sim \mathcal{N}(0, A)}. \quad (\text{L.64})$$

where  $\hat{\mathbf{x}}$  is a generic conjugate fields in the partition function above. Using Stein's lemma, (integration by parts on a Gaussian variable), we can reformulate the definitions of each of the kernels A and B as response functions:

$$\begin{aligned} A_{\mu\nu,i}^\ell(t, s) &= -i \left\langle \hat{\xi}_{\nu,i}^\ell(t) \sigma(h_{\mu,i}^{\ell-1}(s)) \right\rangle_{\beta_{\nu,i}^\ell(s)} \\ &= [G_i^\ell]^{-1} \left\langle (\xi_i^\ell - B_i^\ell \sigma(h_i^{\ell-1})) \sigma(h_i^{\ell-1}) \right\rangle_{\beta_{\nu,i}^\ell(s)} \\ &= \left\langle \frac{\partial \sigma(h_{\mu,i}^{\ell-1}(t))}{\partial \beta_{\nu,i}^{\ell\top}(s)} \right\rangle_{\beta_{\nu,i}^\ell(s)} \end{aligned} \quad (\text{L.65})$$

It easier now integrate over all the  $\hat{\mathbf{x}}$ 's, since the argument of the exponential in  $\mathcal{Z}$  has been linearised with respect to them all. Doing so yields delta functions that give us the final DMFT dynamics.

### L.16. DMFT Dynamics

Piecing together the self-consistent order parameters and the exact integrations over the Hubbard-Stratonovich fields, the fully rigorous DMFT description of the infinite-width MoE is summarized below.

$$\begin{aligned} \alpha^1(t) &\sim \mathcal{N}(0, \Phi^0) & \alpha^\phi(t) &\sim \mathcal{GP}(0, \Phi^1) & [\alpha^{2,\text{in}}]_j &\sim \mathcal{GP}(0, \Phi^1) \\ \bar{\alpha}^3(t) &\sim \mathcal{GP}(0, \bar{\Phi}^{2,\text{in}}) & \beta^{g^{1,\phi}} &\sim \mathcal{GP}(0, \kappa \Phi^{2,\text{in}} \tilde{\phi} \tilde{\phi} G^3) \\ \bar{\beta}^1(t) &\sim \mathcal{GP}(0, \frac{N_e}{\kappa} \bar{G}^{2,\text{in}}) & [\beta^{2,\text{in}}]_j &\sim \mathcal{GP}(0, \frac{\kappa}{N_e} G^3 \tilde{\phi} \tilde{\phi}) \\ \beta^3(t) &\sim \mathcal{N}(0, \mathbb{1}) & \tilde{\beta}^{1,\phi}(t) &\sim \mathcal{GP}(0, \frac{1}{\kappa} G^{1,\phi}) \end{aligned} \quad (\text{L.66})$$

$$\begin{aligned}
 g_\mu^{1,\phi}(t) &= \beta_\mu^{1,\phi}(t) + \frac{\eta_0\gamma_0}{\kappa P} \int_0^t ds \sum_\nu [\Delta_\nu(s) G_{\mu\nu}^3(s,t) [\Phi_{\mu\nu}^{2,\text{in}}(s,t) \tilde{\phi}_\nu(t) - \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s,t)] \\
 &\quad + \tilde{A}^{\tilde{\phi}}] \tilde{\phi}_\mu(s) \\
 z_\mu^1(t) &= \bar{\beta}_\mu^1(t) + \tilde{\beta}_\mu^{1,\phi}(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [B_{\nu\mu}^{1,\phi}(s,t) + \bar{B}_{\mu\nu}^1(s,t) \\
 &\quad + \Delta_\nu(s) G_{\mu\nu}^{1,\phi}(s,t)] \sigma(h_\nu^1(s)) \\
 [z_\mu^{2,\text{in}}(t)]_j &= [\beta_\mu^{2,\text{in}}(t)]_j + \frac{\eta_0\gamma_0}{PN_e} \int_0^t ds \sum_\nu [\Delta_\nu(s) G_{\mu\nu}^3(s,t) \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) + B_{\mu\nu}^{2,\text{in}}(t,s)] \sigma([h_\nu^{2,\text{in}}(s)]_j) \\
 z_\mu^3(t) &= \beta_\mu^3(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 h_\mu^1(t) &= \alpha_\mu^1(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu \Phi_{\mu\nu}^0(z_\mu^1(s) \dot{\sigma}(h_\mu^1(s))) \\
 [h_\mu^{2,\text{in}}(t)]_j &= \alpha_\mu^{2,\text{in}}(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [\Delta_\nu(s) \Phi_{\mu\nu}^1(t,s) \\
 &\quad + \frac{1}{\kappa} \bar{A}_{\mu\nu}^1(s,t)] \left( \dot{\sigma}([h_\mu^{2,\text{in}}(t)]_j) \odot [z_\nu^{2,\text{in}}(s)]_j \right) \\
 h_\mu^3(t) &= \bar{\alpha}_\mu^3(t) + \frac{\eta_0\gamma_0}{P} \int_0^t ds \sum_\nu [\bar{A}_{\mu\nu}^{2,\text{in}}(t,s) + \Delta_\nu(s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s,t)] z_\nu^3(s)
 \end{aligned}$$

(L.67)

**M. DMFT Analysis for Regime II ( $\mu\text{P}$ )**

For  $\mu\text{P}$  in Regime II, an analogous derivation as in the previous section, which we omit for conciseness, yields the following set of self-consistent DMFT equations:

$$\begin{aligned}
 \alpha^1(t) &\sim \mathcal{N}(0, \Phi^0) & \alpha^\phi(t) &\sim \mathcal{GP}(0, \Phi^1) & [\alpha^{2,\text{in}}]_j &\sim \mathcal{GP}(0, \Phi^1) \\
 \bar{\beta}^1(t) &\sim \mathcal{GP}(0, \frac{N_e}{\kappa} \bar{G}^{2,\text{in}}) & \beta^3(t) &\sim \mathcal{N}(0, \mathbb{I}) & \tilde{\beta}^{1,\phi}(t) &\sim \mathcal{GP}(0, \frac{1}{\kappa} G^{1,\phi}) \\
 \Phi_{\mu\nu}^0 &= \frac{1}{D} x_\mu x_\nu^\top, & \Phi_{\mu\nu}^1(t, s) &= \langle \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \rangle_{\mathcal{Z}_N^{\text{global}}}, \\
 \Phi_{\mu\nu}^{2,\text{in}}(s, t) &= \frac{1}{N_e} \sum_{j=1}^{N_e} \langle \sigma([h_\mu^{2,\text{in}}(t)]_j) \sigma([h_\nu^{2,\text{in}}(s)]_j) \rangle_{\mathcal{Z}_{N_e j} | \mathcal{K}_{\text{global}}} \\
 \Phi_{\mu\nu}^3(s, t) &= \langle h_\mu^3(t) h_\nu^3(s) \rangle_{\mathcal{Z}_N^{\text{global}}}, & G_{\mu\nu}^1(t, s) &= \langle [\dot{\sigma}(h_\mu^1(t)) \odot z_\mu^1(t)] [\dot{\sigma}(h_\nu^1(s)) \odot z_\nu^1(s)] \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \Phi_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} & \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle G_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 G_{\mu\nu}^{2,\text{in}}(s, t) &= \frac{1}{N_e} \sum_{j=1}^{N_e} \langle \left[ \dot{\sigma}([h_\mu^{2,\text{in}}(t)]_j) \odot [z_\mu^{2,\text{in}}(t)]_j \right] \left[ \dot{\sigma}([h_\nu^{2,\text{in}}(s)]_j) \odot [z_\nu^{2,\text{in}}(s)]_j \right] \rangle_{\mathcal{Z}_{N_e j} | \mathcal{K}_{\text{global}}} \\
 G_{\mu\nu}^{1,\phi}(s, t) &= \langle g_\mu^{1,\phi}(s) g_\nu^{1,\phi}(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}, & G_{\mu\nu}^3(s, t) &= \langle z_\mu^3(s) z_\nu^3(t) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 A_{\mu\nu}^1(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \tilde{\beta}_\nu^{1\top}(s)} \right\rangle_{\mathcal{Z}_N^{\text{global}}}, & \bar{B}_{\mu\nu}^1(t, s) &= \sum_{j=1}^{N_e} \left\langle \frac{\partial [g_\mu^{2,\text{in}}(t)]_j}{\partial [\alpha_\nu^{2,\text{in}\top}(s)]_j} \right\rangle_{\mathcal{Z}_{N_e j} | \mathcal{K}_{\text{global}}} \\
 A_{\mu\nu}^{1,\phi}(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \tilde{\beta}_\nu^{1,\phi\top}(s)} \right\rangle_{\mathcal{Z}_N^{\text{global}}}, & B_{\mu\nu}^{1,\phi}(t, s) &= \left\langle \frac{\partial g_\mu^{1,\phi}(t)}{\partial \alpha_\nu^{\phi\top}(s)} \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 g_\mu^{1,\phi}(t) &= \frac{\eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\mu\nu}^3(s, t) \left[ \Phi_{\mu\nu}^{2,\text{in}}(s, t) \tilde{\phi}_\nu(t) - \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) \right] \tilde{\phi}_\mu(s) \\
 z_\mu^1(t) &= \bar{\beta}_\mu^1(t) + \tilde{\beta}_\mu^{1,\phi}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left\{ B_{\nu\mu}^{1,\phi}(s, t) + \bar{B}_{\mu\nu}^1(s, t) + \Delta_\nu(s) \left[ \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) + G_{\mu\nu}^{1,\phi}(s, t) \right] \right\} \sigma(h_\nu^1(s)) \\
 [z_\mu^{2,\text{in}}(t)]_j &= \frac{\eta_0 \gamma_0}{P \sqrt{N_e}} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\mu\nu}^3(s, t) \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \sigma([h_\nu^{2,\text{in}}(s)]_j), \quad j \in \{1, 2, \dots, N_e\} \\
 z_\mu^3(t) &= \beta_\mu^3(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 h_\mu^1(t) &= \alpha_\mu^1(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu \Phi_{\mu\nu}^0(z_\mu^1(s) \dot{\sigma}(h_\mu^1(s))) \\
 [h_\mu^{2,\text{in}}(t)]_j &= \alpha_\mu^{2,\text{in}}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left[ \Delta_\nu(s) \Phi_{\mu\nu}^1(t, s) + \frac{1}{\kappa} \bar{A}_{\mu\nu}^1(s, t) \right] \left( \dot{\sigma}([h_\mu^{2,\text{in}}(t)]_j) \odot [z_\nu^{2,\text{in}}(s)]_j \right), \quad \forall j \\
 h_\mu^3(t) &= \frac{\eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) z_\nu^3(s) \\
 \tilde{\phi}_\mu(t) &= \frac{e^{\hat{\psi}_\mu(t)}}{\mathcal{S}_\mu(t)}, & \mathcal{S}_\mu(t) &= \langle e^{\hat{\psi}_\mu(t)} \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}, & \hat{\psi}_\mu(t) &= m_\mu(t) \psi_\mu(t), \\
 m_\mu(t) &= \mathbf{1}(\psi_\mu(t) - \tau_\mu(t)) & \text{where } \tau_\mu(t) &\text{ is chosen such that } \rho = \langle m_\mu(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 \frac{df_\mu(t)}{dt} &= \frac{\eta_0}{P} \sum_\nu \Delta_\nu(t) \left[ G_{\mu\nu}^1(t, t) \Phi_{\mu\nu}^0(t, t) + \left[ \frac{N_e}{\kappa} \bar{G}_{\mu\nu}^{2,\text{in}}(t, t) + G_{\mu\nu}^{1,\phi}(t, t) \right] \Phi_{\mu\nu}^1(t, t) \right. \\
 &\quad \left. + \kappa G_{\mu\nu}^3(t, t) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, t) + \Phi_{\mu\nu}^3(t, t) \right]
 \end{aligned}$$

## N. DMFT Analysis for Regime III (MSSP)

### N.1. Architectural Definitions

In this section we are considering the case in which  $M, N, N_e \rightarrow \infty$  at a fixed ratio. Specifically, we define  $\kappa = \frac{M}{N}, \iota = \frac{N_e}{N}$ , where  $\kappa, \iota$  are order one in  $N$ . The forward pass, following *MSSP*, is defined as:

$$\begin{aligned}
 \mathbb{R}^N \ni h_\mu^1 &= \frac{1}{\sqrt{D}} W^1 x_\mu & \mathbb{R}^M \ni \psi_\mu &= \frac{1}{\sqrt{N}} Q \sigma(h_\mu^1) \\
 \mathbb{R}^{N_e} \ni h_{\mu,i}^{2,\text{in}} &= \frac{1}{\sqrt{N}} W_i^{2,\text{i}} \sigma(h_\mu^1) & \mathbb{R}^M \ni \phi_\mu &= \text{softmax}(\psi_\mu) \\
 \mathbb{R}^N \ni h_{\mu,i}^{2,\text{out}} &= \frac{1}{\sqrt{N}} W_i^{2,\text{o}} \sigma(h_{\mu,i}^{2,\text{in}}) & \mathbb{R}^N \ni h_\mu^3 &= \sum_{i=1}^M \phi_\mu^i h_{\mu,i}^{2,\text{o}} \\
 \mathbb{R} \ni h_\mu^4 &= \frac{1}{\sqrt{N}} w^{4\text{T}} h_\mu^3 & f_\mu &= \frac{1}{\gamma} h_\mu^4
 \end{aligned} \tag{N.1}$$

### N.2. Learning Rates and Initialization

To ensure that the network evolves dynamically in the infinite-width limit without gradients vanishing or exploding, we apply specific learning rate scalings:

$$\eta = \eta_0 \gamma^2 \quad \eta_Q = \eta_0 \gamma^2 \kappa \quad \eta_E = \eta_0 \gamma^2 N \quad \gamma = \gamma_0 \sqrt{N} \quad \eta_0, \gamma_0 \sim O(1) \tag{N.2}$$

The network weights are initialized from standard Gaussian distributions:

$$w_\alpha^4(0), [W_i^{2,\text{out}}(0)]_{\alpha\beta}, [W_i^{2,\text{in}}(0)]_{\alpha\beta}, W_{\alpha\beta}^1(0), Q_{\alpha\beta}(0) \sim \mathcal{N}(0, 1) \tag{N.3}$$

but noticing that for each expert's matrices  $W_i^{2,\text{in}}, W_i^{2,\text{out}}$  we share the same initial weights  $W_i^{2,\text{in}}(0) = W^{2,\text{in}}(0), W_i^{2,\text{out}}(0) = W^{2,\text{out}}(0)$  For brevity in our derivations, we denote the collection of all network parameters as  $\theta = \text{Vec}\{W^1, W_i^{2,\text{in}}, W_i^{2,\text{out}}, w^4, Q\}$ .

Due to the normalization constraint of the softmax operator over  $M$  experts, the routing probabilities scale as  $\phi^i \sim O(\frac{1}{\kappa N})$ . It is mathematically tidier to formulate the DMFT using variables which remain  $O(1)$  in the limit. We therefore define and track the rescaled routing variables:

$$\tilde{\phi}_\mu^i := \kappa N \phi_\mu^i \tag{N.4}$$

### N.3. Gradient Definitions

We mathematically *define* the pre-activation gradients, ensuring they contain the correct scaling factors to remain finite in the  $N \rightarrow \infty$  limit:

$$\begin{aligned}
 g_\mu^1 &:= \sqrt{N} \frac{\partial h_\mu^4}{\partial h_\mu^1} & g_{\mu,i}^{2,\text{in}} &:= \kappa N^{\frac{3}{2}} \frac{\partial h_\mu^4}{\partial h_{\mu,i}^{2,\text{in}}} \\
 g_{\mu,i}^{2,\text{out}} &:= \kappa N^{\frac{3}{2}} \frac{\partial h_\mu^4}{\partial h_{\mu,i}^{2,\text{o}}} & g_\mu^3 &:= \sqrt{N} \frac{\partial h_\mu^4}{\partial h_\mu^3} = w^4 = z^3 \\
 g_\mu^\phi &:= \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial \phi_\mu} & z_{\mu,i}^{2,\text{out}} &:= \kappa N \phi_\mu^i g_\mu^3 = \tilde{\phi}_\mu^i g_\mu^3 \\
 z_{\mu,i}^{2,\text{in}} &:= \frac{1}{\sqrt{N}} W_i^{2,\text{out}\top} g_{\mu,i}^{2,\text{out}} & \tilde{z}_\mu^\phi &:= \frac{1}{\kappa \sqrt{N}} Q^\top g^{1,\phi} \\
 \tilde{z}_{\mu,i}^1 &:= \frac{1}{\kappa \sqrt{N}} W_i^{2,\text{in}\top} g_{\mu,i}^{2,\text{in}} & z^1 &:= \frac{1}{\kappa \sqrt{N}} Q^\top g^{1,\phi} + \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^M W_i^{2,\text{in}\top} g_{\mu,i}^{2,\text{in}} \\
 & & &= \frac{1}{\kappa \sqrt{N}} \sum_{i=1}^M \tilde{z}_{\mu,i}^1 + \tilde{z}_\mu^{1,\phi}
 \end{aligned} \tag{N.5}$$

Note that the components of the router gradient  $g_\mu^\phi$  naturally scale as  $O(1)$ :

$$(g_\mu^\phi)_i = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial \phi_\mu^i} = \frac{1}{\sqrt{N}} \frac{\partial h_\mu^4}{\partial h_\mu^3} \cdot \frac{\partial h_\mu^3}{\partial \phi_\mu^i} = \frac{1}{N} \mathbf{g}^3 \cdot \mathbf{h}_{\mu,i}^{2,\text{out}} = O(1). \tag{N.6}$$

The backward pass through the router dictates that the gradient at layer 1 involves the quantity  $g^{1,\phi}$ , defined component-wise to handle the Jacobian of the softmax:

$$g_k^{1,\phi} := \kappa N \sum_{i=1}^{\kappa N} g_i^\phi \phi^i (\delta_{ik} - \phi^k) = \sum_{i=1}^{\kappa N} g_i^\phi \tilde{\phi}^i (\delta_{ik} - \frac{\tilde{\phi}^k}{\kappa N}) \tag{N.7}$$

#### N.4. DMFT Kernels and Order Parameters

We define the following macroscopic *kernels*, all of which are  $O_N(1)$  and will serve as the fundamental order parameters in the dynamics:

$$\begin{aligned}
 \Phi_{\mu\nu}^0 &= \frac{1}{D} x_\mu \cdot x_\nu & \Phi_{\mu\nu}^1(s, t) &= \frac{1}{N} \sigma(h_\mu^1(s)) \cdot \sigma(h_\nu^1(t)) \\
 \Phi_{\mu\nu,i}^{2,\text{in}}(s, t) &= \frac{1}{\iota N} \sigma(h_{\mu,i}^{2,\text{in}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) & \Phi_{\mu\nu,i}^{2,\text{o}}(s, t) &= \frac{1}{N} \sigma(h_{\mu,i}^{2,\text{out}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{out}}(t)) \\
 \Phi_{\mu\nu}^3(s, t) &= \frac{1}{N} h_\mu^3(s) \cdot h_\nu^3(t) & G_{\mu\nu}^1(s, t) &= \frac{1}{N} g_\mu^1(s) \cdot g_\nu^1(t) \\
 G_{\mu\nu,i}^{2,\text{in}}(s, t) &= \frac{1}{\iota N} g_{\mu,i}^{2,\text{in}}(s) \cdot g_{\nu,i}^{2,\text{in}}(t) & G_{\mu\nu,i}^{2,\text{out}}(s, t) &= \frac{1}{N} g_{\mu,i}^{2,\text{out}}(s) \cdot g_{\nu,i}^{2,\text{out}}(t) \\
 G_{\mu\nu}^{1,\phi}(s, t) &= \frac{1}{\kappa N} g_\mu^{1,\phi}(s) \cdot g_\nu^{1,\phi}(t) & G_{\mu\nu}^3(s, t) &= \frac{1}{N} g_\mu^3(s) \cdot g_\nu^3(t)
 \end{aligned} \tag{N.8}$$

With the exception of the input data Gram matrix  $\Phi^0$ , all of these kernels are dynamically evolving macroscopic variables.

### N.5. Learning Dynamics and the Neural Tangent Kernel

We train the network using continuous-time gradient flow. The variable learning rate scheme defined previously is necessary so that the activation updates within the experts and for the router scores do not vanish or explode as  $N, M \rightarrow \infty$ .

Taking a standard empirical risk minimization loss of the form:

$$\mathcal{L} = \frac{1}{P} \sum_{\mu=1}^P \ell(f_{\mu}, y_{\mu}) \quad (\text{N.9})$$

The network parameters evolve according to:

$$\begin{aligned} \frac{d\theta}{dt} &= \frac{\eta}{P\gamma} \sum_{\mu} \Delta_{\mu} \frac{\partial h_{\mu}^4}{\partial \theta} \\ \Delta_{\mu} &= -\frac{\partial \mathcal{L}}{\partial f_{\mu}} \end{aligned} \quad (\text{N.10})$$

For a Mean Squared Error (MSE) loss where  $\mathcal{L} = \frac{1}{P} \sum_{\nu} (y_{\nu} - f_{\nu})^2$ , the error signal (or residual) simplifies to  $\Delta_{\nu} = 2(y_{\nu} - f_{\nu})$ .

The logits update dynamically via the chain rule, driven by the Neural Tangent Kernel (NTK):

$$\frac{df_{\mu}(t)}{dt} = \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{d\theta}{dt} = \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{\eta_{\theta}}{P} \sum_{\alpha} \Delta_{\alpha} \frac{\partial f_{\alpha}(t)}{\partial \theta} = \frac{\eta}{P} \sum_{\alpha} \Delta_{\alpha} K_{\mu\alpha}^{\text{NTK}}(t, t) \quad (\text{N.11})$$

Where the infinite-width NTK is defined as:

$$K_{\mu\alpha}^{\text{NTK}}(t, s) \equiv \frac{\partial f_{\mu}(t)}{\partial \theta} \cdot \frac{\partial f_{\alpha}(s)}{\partial \theta} \quad (\text{N.12})$$

### N.6. Decomposition of the MoE NTK

To explicitly compute the NTK, we require that the magnitude  $\eta_{\theta} \frac{\partial f_{\mu}(t)}{\partial \theta_i} \cdot \frac{\partial f_{\alpha}(s)}{\partial \theta_i}$  remains order 1 for each parameter block  $\theta_i$ . We evaluate this constraint block by block.

For the expert input layer  $W^{2,\text{in}}$ :

$$\begin{aligned} & \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,j=1}^{N_e} \sum_{k=1}^N \frac{\partial h^4}{\partial [h_i^{2,\text{in}}]_l} \frac{\partial [h_i^{2,\text{in}}]_l}{\partial W_{jk,i}^2} \frac{\partial h^4}{\partial [h_i^{2,\text{in}}]_m} \frac{\partial [h_i^{2,\text{in}}]_m}{\partial W_{jk,i}^2} = \\ &= \eta_0 N \sum_{i=1}^{\kappa N} \sum_{l,m,j=1}^{N_e} \sum_{k=1}^N \frac{[g_i^{2,\text{in}}]_l}{\kappa N^{3/2}} \frac{[g_i^{2,\text{in}}]_m}{\kappa N^{3/2}} \frac{\sigma(h_k^1)}{\sqrt{N}} \frac{\sigma(h_k^1)}{\sqrt{N}} \delta_{mj} \delta_{lj} \\ &= \eta_0 \frac{\iota}{\kappa} \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \left[ \frac{1}{N_e} \sum_{j=1}^{N_e} [g_i^{2,\text{in}}]_j [g_i^{2,\text{in}}]_j \right] \left[ \frac{1}{N} \sum_{k=1}^N \sigma(h_k^1)^2 \right] \\ &= \eta_0 \frac{\iota}{\kappa} \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_i^{2,\text{in}} \Phi^1 = \eta_0 \frac{\iota}{\kappa} \bar{G}^2 \Phi^1 \end{aligned} \quad (\text{N.13})$$

For the expert output layer  $W^{2,\text{out}}$ :

$$\begin{aligned}
 & \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,k=1}^{N_e} \sum_{j=1}^N \frac{\partial h^4}{\partial [h_i^{2,\text{out}}]_l} \frac{\partial [h_i^{2,\text{out}}]_l}{\partial W_{jk}^{2,\text{out}}} \frac{\partial h^4}{\partial [h_i^{2,\text{out}}]_m} \frac{\partial [h_i^{2,\text{out}}]_m}{\partial W_{jk}^{2,\text{out}}} \\
 &= \eta_0 \sum_{i=1}^{\kappa N} \sum_{l,m,k=1}^{N_e} \sum_{j=1}^N \frac{[g_i^{2,\text{out}}]_l}{\kappa N^{3/2}} \frac{[g_i^{2,\text{out}}]_m}{\kappa N^{3/2}} \frac{\sqrt{\kappa N}}{\sqrt{N_e}} \sigma([h_i^{2,\text{in}}]_k) \frac{\sqrt{\kappa N}}{\sqrt{N_e}} \sigma([h_i^{2,\text{in}}]_k) \delta_{mj} \delta_{lj} \\
 &= \eta_0 \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} G_i^{2,\text{out}} \Phi_i^{2,\text{in}}
 \end{aligned} \tag{N.14}$$

For the shared base layer  $W^1$ :

$$\eta_0 \sum_{k=1}^D \sum_{j,m=1}^N \frac{\partial h^4}{\partial h_n^1} \frac{\partial h_n^1}{\partial W_{jk}^1} \frac{\partial h^4}{\partial h_m^1} \frac{\partial h_m^1}{\partial W_{jk}^1} = \eta_0 \sum_{k=1}^D \sum_{j,m,n=1}^N \frac{g_n^1}{\sqrt{N}} \frac{g_m^1}{\sqrt{N}} \delta_{nj} \delta_{mj} \frac{x_k}{\sqrt{D}} \frac{x_k}{\sqrt{D}} = \eta_0 G^1 \Phi^0 \tag{N.15}$$

The router weight matrix  $Q$  involves the complex Jacobian of the softmax operator. Integrating this out, we find:

$$\eta_0 \kappa \sum_{j=1}^{\kappa N} \sum_{k=1}^N \frac{\partial h^4}{\partial Q_{jk}} \frac{\partial h^4}{\partial Q_{jk}} = \eta_0 \kappa N \sum_{j=1}^{\kappa N} \sum_{i,i'=1}^{\kappa N} \Phi^1 g_i^\phi g_{i'}^\phi \phi^i \phi^{i'} (\delta_{ij} - \phi^j) (\delta_{i'j} - \phi^j) = \eta_0 \Phi^1 G^{1,\phi} \tag{N.16}$$

Aggregating these contributions yields the full macroscopic evolution equation:

$$\begin{aligned}
 \frac{df_\mu(t)}{dt} &= \frac{\eta_0}{P} \sum_\nu \Delta_\nu \left[ G_{\mu\nu}^1(t,t) \Phi_{\mu\nu}^0(t,t) + \frac{l}{\kappa} \left[ G_{\mu\nu,i}^{2,\text{in}} + G_{\mu\nu}^{1,\phi}(t,t) \right] (t,t) \Phi_{\mu\nu}^1(t,t) \right. \\
 &\quad \left. + \frac{l}{\kappa} G_{\mu\nu}^3(t,t) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t,t) + \Phi_{\mu\nu}^3(t,t) \right]
 \end{aligned} \tag{N.17}$$

## N.7. Evolution of weights, preactivations, and pregradients

The updates to the router weights are:

$$\begin{aligned}
 \frac{dQ_{jk}}{dt} &= \frac{\eta}{P} \sum_\mu \Delta_\mu \frac{\partial f_\mu}{\partial Q_{jk}} = \frac{\eta_0 \gamma \kappa \sqrt{N}}{P} \sum_\mu \Delta_\mu \frac{\partial h_\mu^4}{\partial Q_{jk}} = \frac{\eta_0 \gamma_0 \kappa N}{P} \sum_\mu \Delta_\mu \frac{\partial h_\mu^4}{\partial Q_{jk}} \\
 &= \frac{\eta_0 \gamma_0}{P \sqrt{N}} \sum_\mu \Delta_\mu \left[ \kappa N \sum_{a=1}^M g_a^\phi \phi_a (\delta_{aj} - \phi_j) \right] \sigma(h_k^1) \\
 &= \frac{\eta_0 \gamma_0}{P \sqrt{N}} \sum_\mu \Delta_\mu g_j^{1,\phi} \sigma(h_k^1)
 \end{aligned} \tag{N.18}$$

We obtain thus:

$$\begin{aligned}
 \psi_\mu(t) &= \frac{1}{\sqrt{N}} Q(t) \sigma(h_\mu^1(t)) \\
 &= \chi_\mu^\phi(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) \Phi_{\mu\nu}^1(s, t) g_\nu^{1,\phi}(t) \\
 \tilde{z}_\mu^{1,\phi} &= \tilde{\xi}_\mu^{1,\phi} + \frac{\eta_0 \gamma_0}{P} \int_0^\infty dt \sum_\nu \Delta_\nu G_{\nu\mu}^{1,\phi} \sigma(h_{\nu j}^1)
 \end{aligned} \tag{N.19}$$

The other preactivation and pregradient updates follow similarly:

$$\begin{aligned}
 \frac{dW^1(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{g_\mu^1(t)}{\sqrt{N}} \frac{1}{\sqrt{D}} x_\mu \\
 \Rightarrow h_\mu^1(t) &= \frac{1}{\sqrt{D}} W^1(0) x_\mu + \int_0^t ds \frac{\eta_0 \gamma_0}{P} \sum_\nu \Delta_\nu(s) g_\nu^1(s) \Phi_{\mu\nu}^0(t, s) \\
 \frac{dW_i^{2,\text{in}}(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N} N}{P} \sum_\mu \Delta_\mu(t) \frac{g_{\mu,i}^{2,\text{in}}(t)}{\kappa \sqrt{N}} \frac{1}{N^{3/2}} \sigma(h_\mu^1(t)) \\
 \Rightarrow h_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{in}}(0) \sigma(h_\mu^1(t)) + \frac{\eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) g_{\nu,i}^{2,\text{in}}(s) \Phi_{\mu\nu}^1(t, s) \\
 \frac{dW_i^{2,\text{out}}(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N} N}{P} \sum_\mu \Delta_\mu(t) \frac{g_{\mu,i}^{2,\text{out}}(t)}{N^{3/2} \kappa} \frac{1}{\sqrt{N}} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \\
 \Rightarrow h_{\mu,i}^{2,\text{out}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) + \frac{\eta_0 \gamma_0 t}{P \kappa} \int_0^t ds \sum_\nu \Delta_\nu(s) g_{\nu,i}^{2,\text{out}}(s) \Phi_{\mu\nu,i}^{2,\text{in}}(t, s) \\
 \frac{dw^4(t)}{dt} &= \frac{\eta_0 \gamma_0 \sqrt{N}}{P} \sum_\mu \Delta_\mu(t) \frac{1}{\sqrt{N}} h_\mu^3(t) = \frac{\eta_0 \gamma_0}{P} \sum_\mu \Delta_\mu(t) h_\mu^3(t)
 \end{aligned} \tag{N.20}$$

Using the same expressions for the evolution of the weights to derive expressions for the pregradients:

$$\begin{aligned}
 z_\mu^3(t) &= g_\mu^3(t) = w_\mu^4(t) = \xi_\mu^3(0) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 z_{\mu,i}^{2,\text{out}}(t) &= \kappa g_{\mu,i}^{2,\text{out}}(t) = \tilde{\phi}_\mu^i(t) g_\nu^3(t) = \frac{\eta_0 \gamma_0}{P} \tilde{\phi}_\mu^i(t) \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s) \\
 z_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{out}\top}(t) g_{\mu,i}^{2,\text{out}}(t) \\
 &= \xi_{\mu,i}^{2,\text{in}}(t) + \frac{\eta_0 \gamma_0}{P\kappa} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\nu\mu,i}^{2,\text{out}}(s,t) \sigma(h_{\nu,i}^{2,\text{in}}(s)) \\
 z_{\mu,i}^1(t) &= W_i^{2,\text{in}\top}(t) g_i^{2,\text{in}}(t) \\
 &= \tilde{\xi}_{\mu,i}^1(t) + \frac{\eta_0 \gamma_0 t}{P\kappa} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\nu\mu,i}^{2,\text{in}}(s,t) \sigma(h_\nu^1(s))
 \end{aligned} \tag{N.21}$$

Finally for the router gradient:

$$\begin{aligned}
 g_i^{1,\phi} &= \frac{1}{\kappa N} \sum_{m=1}^{\kappa N} g_m^\phi \tilde{\phi}^m (\kappa N \delta_{mi} - \tilde{\phi}^i) \\
 &= \xi_i^{g^{1,\phi}} + \frac{\eta_0 \gamma_0 t}{P\kappa} \int_0^t ds \sum_\nu \Delta_\nu G_\nu^3 \Phi_i^{2,\text{in}} \tilde{\phi}^i \tilde{\phi}^i - \frac{\eta_0 \gamma_0 t}{P\kappa} \int_0^t ds \sum_\nu \Delta_\nu G_\nu^3 \bar{\Phi}^{2,\text{in}} \tilde{\phi}^i.
 \end{aligned} \tag{N.22}$$

### N.8. Stochastic Initial Fields

To isolate the purely deterministic part of the trajectory from the random initialization, we define a set of initial stochastic fields  $\mathcal{F}$ . These fields encapsulate the randomness of the initial weights:

$$\mathcal{F} = \{\chi_\mu^1(t), \chi_{\mu,i}^{2,\text{in}}(t), \chi_{\mu,i}^{2,\text{out}}(t), \chi_\mu^\phi(t), \tilde{\xi}_{\mu,i}^1(t), \xi_{\mu,i}^{2,\text{in}}(t), \xi_\mu^3(t), \xi_\mu^\phi(t), \xi_{\mu i}^{g^\phi}(t)\}_{i \in \{1, \dots, M\}, \mu \in \{1, \dots, P\}} \tag{N.23}$$

$$\begin{aligned}
 \chi_\mu^1 &= \frac{1}{\sqrt{D}} W^1(0) x_\mu \\
 \chi_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{in}}(0) \sigma(h_\mu^1(t)) \\
 \chi_{\mu,i}^{2,\text{out}}(t) &= \frac{1}{\sqrt{N}} W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) \\
 \chi_\mu^\phi(t) &= \frac{1}{\sqrt{N}} Q(0) \sigma(h_\mu^1(t)) \\
 \tilde{\xi}_{\mu,i}^1(t) &= \frac{1}{\sqrt{N}} (W_i^{2,\text{in}}(0))^\top g_{\mu,i}^{2,\text{in}}(t) \\
 \xi_{\mu,i}^{2,\text{in}}(t) &= \frac{1}{\sqrt{N}} (W_i^{2,\text{out}}(0))^\top g_{\mu,i}^{2,\text{out}}(t) \\
 \xi_\mu^{1,\phi}(t) &= \frac{1}{\kappa\sqrt{N}} Q(0)^\top g_\mu^{1,\phi}(t) \\
 \xi_{\mu i}^{g^\phi}(t) &= \frac{1}{N^{3/2}} g_\mu^{3\top}(t) W_i^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t))
 \end{aligned} \tag{N.24}$$

A critical step in rendering the MoE partition function computationally tractable is distinguishing between expert-local fields (which are specific to an expert  $i$ ) and global fields (which impact the shared base layer or the final aggregated output). This distinction allows the massive partition function to factorize over the experts.

We define the router-averaged stochastic fields to capture the macroscopic effect of the local processes:

$$\bar{\chi}_\mu^3(t) = \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \tilde{\phi}_\mu^i(t) \chi_{\mu,i}^{2,\text{out}}(t) \tag{N.25}$$

$$\bar{\xi}_\mu^1(t) = \frac{1}{\kappa\sqrt{N}} \sum_{i=1}^{\kappa N} \tilde{\xi}_{\mu,i}^1(t) \tag{N.26}$$

Substituting these, the global aggregated output  $h_\mu^3$  and the gradient arriving at the base layer  $z_\mu^1$  can be expressed entirely in terms of order parameters that are smooth over the ensemble of experts:

$$h_\mu^3(t) = \bar{\chi}_\mu^3(t) + \frac{\eta_0 \gamma_0 \ell}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) z_\nu^3(s) \tag{N.27}$$

$$z_\mu^1(t) = \bar{\xi}_\mu^1(t) + \xi_\mu^{1,\phi}(t) + \frac{\eta_0 \gamma_0 \ell}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) \left[ \frac{1}{\sqrt{N}} \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) + G_{\mu\nu}^{1,\phi}(s, t) \right] \sigma(h_\nu^1(s)) \tag{N.28}$$

### N.9. Deriving the DMFT Action

We formulate the DMFT by writing the moment-generating function for the system's trajectories and performing a disorder average over the initial weights  $\theta_0$ .

$$\begin{aligned}
 Z \propto & \left\langle \int d\mathcal{F} \exp \left( \sum_{\mu} \int_0^{\infty} dt i \left[ \hat{\chi}_{\mu}^1 \cdot \left( \chi_{\mu}^1 - \frac{1}{\sqrt{D}} W^1(0) x_{\mu} \right) + \hat{\chi}_{\mu}^{2,\text{in}} \cdot \left( \chi_{\mu}^{2,\text{in}} - \frac{1}{\sqrt{N}} W^{2,\text{in}}(0) \sigma(h_{\mu}^1(t)) \right) \right. \right. \\
 & + \hat{\chi}_{\mu}^3 \cdot \left( \bar{\chi}_{\mu}^3 - \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \tilde{\phi}_{\mu}^i(t) \frac{1}{\sqrt{N}} W^{2,\text{out}}(0) \sigma(h_{\mu,i}^{2,\text{in}}(t)) \right) + \hat{\xi}_{\mu}^1 \cdot \left( \bar{\xi}_{\mu}^1 - \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \frac{1}{\sqrt{N}} W^{2,\text{in}}(0)^{\top} g_{\mu,i}^{2,\text{in}}(t) \right) \\
 & + \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \hat{\xi}_{\mu,i}^{2,\text{in}} \cdot \left( \xi_{\mu,i}^{2,\text{in}} - \frac{1}{\sqrt{N}} W^{2,\text{out}}(0)^{\top} g_{\mu,i}^{2,\text{out}}(t) \right) + \hat{\xi}_{\mu}^3 \cdot \left( \xi_{\mu}^3 - w^4(0)^{\top} \right) \\
 & + \hat{\chi}_{\mu}^{\phi} \cdot \left( \chi_{\mu}^{\phi} - \frac{1}{\sqrt{N}} Q(0) \sigma(h_{\mu}^1(t)) \right) + \hat{\xi}_{\mu}^{1,\phi} \cdot \left( \xi_{\mu}^{1,\phi} - \frac{1}{\kappa \sqrt{N}} Q(0)^{\top} g_{\mu}^{1,\phi}(t) \right) \\
 & \left. + \sum_{m=1}^{\kappa N} \hat{\xi}_{\mu,m}^{g^{1,\phi}}(t) \left( \xi_{\mu,m}^{g^{1,\phi}}(t) - \frac{1}{\kappa N} \sum_{i=1}^{\kappa N} \frac{1}{N^{3/2}} g_{\mu}^{3T}(t) W^{2,\text{out}} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \tilde{\phi}_{\mu}^i(t) \left( \kappa N \delta_{im} - \tilde{\phi}_{\mu}^m(t) \right) \right) \right] \Bigg) \\
 & \times \exp \left( \sum_{\mu} \int_0^{\infty} dt \left[ j_{\mu}^1(t) \cdot \chi_{\mu}^1(t) + j_{\mu}^{2,\text{in}}(t) \cdot \chi_{\mu}^{2,\text{in}}(t) + \bar{j}_{\mu}^3(t) \cdot \bar{\chi}_{\mu}^3(t) + \frac{1}{\kappa N} \sum_{i=1}^M v_{\mu,i}^{2,\text{in}}(t) \cdot \xi_{\mu,i}^{2,\text{in}}(t) \right. \right. \\
 & \left. \left. + \bar{v}_{\mu}^1(t) \cdot \bar{\xi}_{\mu}^1(t) + v_{\mu}^3 \cdot \xi_{\mu}^3(t) + j_{\mu}^{\phi}(t) \cdot \chi_{\mu}^{\phi}(t) + v_{\mu}^{1,\phi}(t) \cdot \xi_{\mu}^{1,\phi}(t) + v_{\mu}^{g^{1,\phi}}(t) \cdot \xi_{\mu}^{g^{1,\phi}}(t) \right] \right) \Bigg\rangle_{\theta_0}.
 \end{aligned} \tag{N.29}$$

Given the shared initialisation across experts we now have terms coupling different experts (e.g.  $\sigma(h_{\mu}^{2,\text{in}})_i \sigma(h_{\nu}^{2,\text{in}})_j \tilde{\phi}_{\mu}^i \tilde{\phi}_{\nu}^j$ ). In order to obtain a stable set of order parameters that scale coherently with  $N$  we define the following set of "averaged fields":

$$\begin{aligned}
 1 &= \int \frac{d\hat{g}_{\mu}^{2,\text{in}}(t)}{2\pi} \exp \left[ \hat{g}_{\mu}^{2,\text{in}}(t) \left( \kappa N \bar{g}_{\mu}^{2,\text{in}}(t) - \sum_{m=1}^{\kappa N} g_{\mu,m}^{2,\text{in}}(t) \right) \right] \\
 1 &= \int \frac{d\hat{\sigma}(h_{\mu}^{2,\text{in}}(t))}{2\pi} \exp \left[ \hat{\sigma}(h_{\mu}^{2,\text{in}}(t)) \left( \kappa N \bar{\sigma}(h_{\mu}^{2,\text{in}}(t)) - \sum_{i=1}^{\kappa N} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \tilde{\phi}^i(t) \right) \right] \\
 1 &= \int \frac{d\hat{\sigma}^{\phi}(h_{\mu}^{2,\text{in}}(t))}{2\pi} \exp \left[ \hat{\sigma}^{\phi}(h_{\mu}^{2,\text{in}}(t)) \left( \kappa N \bar{\sigma}^{\phi}(h_{\mu}^{2,\text{in}}(t)) - \sum_{i=1}^{\kappa N} \sigma(h_{\mu,i}^{2,\text{in}}(t)) \tilde{\phi}^i(t) \hat{\xi}_i^{g^{1,\phi}}(t) \right) \right] \\
 1 &= \int \frac{d\hat{\xi}_{\mu}^{2,\text{in}}(t)}{2\pi} \exp \left[ \hat{\xi}_{\mu}^{2,\text{in}}(t) \left( \kappa N \bar{\xi}_{\mu}^{2,\text{in}}(t) - \sum_{i=1}^{\kappa N} \xi_{\mu,i}^{2,\text{in}}(t) \tilde{\phi}^i(t) \right) \right]
 \end{aligned} \tag{N.30}$$

We must constrain the solutions to be physical, so for each  $i, \mu, \nu, s, t$  we multiply in the following resolutions of the identity:

$$1 = \int \frac{d\hat{\Phi}_{\mu\nu}^1(s, t) d\hat{\Phi}_{\mu\nu}^1(s, t)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^1(s, t) \left( N \Phi_{\mu\nu}^1(s, t) - \sigma(h_{\mu}^1(s)) \cdot \sigma(h_{\nu}^1(t)) \right) \right] \tag{N.31}$$

$$1 = \int \frac{dG_{\mu\nu}^1(s, t) d\hat{G}_{\mu\nu}^1(s, t)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^1(s, t) \left( N G_{\mu\nu}^1(s, t) - g_{\mu}^1(s) \cdot g_{\nu}^1(t) \right) \right] \tag{N.32}$$

$$1 = \int \frac{dG_{\mu\nu}^3(s, t) d\hat{G}_{\mu\nu}^3(s, t)}{2\pi i N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^3(s, t) \left( NG_{\mu\nu}^3(s, t) - g_\mu^3(s) \cdot g_\nu^3(t) \right) \right] \quad (\text{N.33})$$

$$1 = \int \frac{dG_{\mu\nu}^{1,\phi}(s, t) d\hat{G}_{\mu\nu}^{1,\phi}(s, t)}{2\pi i \kappa N^{-1}} \exp \left[ \hat{G}_{\mu\nu}^{1,\phi}(s, t) \left( \kappa N G_{\mu\nu}^{1,\phi}(s, t) - g_\mu^{1,\phi}(s) \cdot g_\nu^{1,\phi}(t) \right) \right] \quad (\text{N.34})$$

$$1 = \int \frac{d\Phi_{\mu\nu,i}^{2,\text{in}}(s, t) d\hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(s, t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(s, t) \left( \iota N \Phi_{\mu\nu,i}^{2,\text{in}}(s, t) - \sigma(h_{\mu,i}^{2,\text{in}}(s)) \cdot \sigma(h_{\nu,i}^{2,\text{in}}(t)) \right) \right] \quad (\text{N.35})$$

$$1 = \int \frac{dG_{\mu\nu,i}^{2,\text{in}}(s, t) d\hat{G}_{\mu\nu,i}^{2,\text{in}}(s, t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{G}_{\mu\nu,i}^{2,\text{in}}(s, t) \left( \iota N G_{\mu\nu,i}^{2,\text{in}}(s, t) - g_{\mu,i}^2(s) \cdot g_{\nu,i}^2(t) \right) \right] \quad (\text{N.36})$$

$$1 = \int \frac{d\bar{A}_{\mu\nu}^1(s, t) d\bar{B}_{\mu\nu}^1(s, t)}{2\pi i N^{-1}} \exp \left[ -\bar{B}_{\mu\nu}^1(s, t) \left( N\bar{A}_{\mu\nu}^1(s, t) + i\hat{\xi}_\mu^1(s) \cdot \sigma(h_\nu^1(t)) \right) \right] \quad (\text{N.37})$$

$$1 = \int \frac{dA_{\mu\nu}^{1,\phi}(s, t) dB_{\mu\nu}^{1,\phi}(s, t)}{2\pi i N^{-1}} \exp \left[ -B_{\mu\nu}^{1,\phi}(s, t) \left( NA_{\mu\nu}^{1,\phi}(s, t) + i\hat{\xi}_\mu^{1,\phi}(s) \cdot \sigma(h_\nu^1(t)) \right) \right] \quad (\text{N.38})$$

$$1 = \int \frac{dA_{\mu\nu}^{g^{1,\phi}}(s, t) dB_{\mu\nu}^{g^{1,\phi}}(s, t)}{2\pi i (\kappa N)^{-1}} \exp \left[ -B_{\mu\nu}^{g^{1,\phi}}(s, t) \left( \kappa N A_{\mu\nu}^{g^{1,\phi}}(s, t) + i\hat{\xi}_\mu^{g^{1,\phi}}(s) \cdot \tilde{\phi}_\nu(t) \right) \right] \quad (\text{N.39})$$

$$1 = \int \frac{d\Phi_{\mu\nu}^3(t, s) d\hat{\Phi}_{\mu\nu}^3(t, s)}{2\pi i N^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^3(t, s) \left( N\Phi_{\mu\nu}^3(t, s) - h_\mu^3(t) \cdot h_\nu^3(s) \right) \right] \quad (\text{N.40})$$

$$1 = \int \frac{d\bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, s) d\hat{\bar{\Phi}}_{\mu\nu}^{2,\text{in}}(t, s)}{2\pi i} \exp \left[ \hat{\bar{\Phi}}_{\mu\nu}^{2,\text{in}}(t, s) \left( \kappa N \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, s) - \sum_{i=1}^{\kappa N} \tilde{\phi}_\mu^i(s) \Phi_{\mu\nu,i}^{2,\text{in}}(s, t) \right) \right] \quad (\text{N.41})$$

$$1 = \int \frac{d\bar{G}_{\mu\nu}^{2,\text{in}}(t, s) d\hat{\bar{G}}_{\mu\nu}^{2,\text{in}}(t, s)}{2\pi i (\kappa N)^{-1}} \exp \left[ \hat{\bar{G}}_{\mu\nu}^{2,\text{in}}(t, s) \left( \kappa N \bar{G}_{\mu\nu}^{2,\text{in}}(t, s) - \sum_{i=1}^{\kappa N} G_{\mu\nu,i}^{2,\text{in}}(s, t) \right) \right] \quad (\text{N.42})$$

$$1 = \int \frac{d\tilde{G}_{\mu\nu}^{2,\text{in}}(s, t) d\hat{\tilde{G}}_{\mu\nu}^{2,\text{in}}(s, t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\tilde{G}}_{\mu\nu}^{2,\text{in}}(s, t) \left( \iota N \tilde{G}_{\mu\nu}^{2,\text{in}}(s, t) - \bar{g}_\mu^{2,\text{in}}(s) \cdot \bar{g}_\nu^{2,\text{in}}(t) \right) \right] \quad (\text{N.43})$$

$$1 = \int \frac{d\tilde{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) d\hat{\tilde{\Phi}}_{\mu\nu}^{2,\text{in}}(s, t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\tilde{\Phi}}_{\mu\nu}^{2,\text{in}}(s, t) \left( \iota N \tilde{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) - \bar{\sigma}_\mu^{2,\text{in}}(s) \cdot \bar{\sigma}_\nu^{2,\text{in}}(t) \right) \right] \quad (\text{N.44})$$

$$1 = \int \frac{d\tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi}(s, t) d\hat{\tilde{\Phi}}_{\mu\nu}^{2,\text{in},\phi}(s, t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\tilde{\Phi}}_{\mu\nu}^{2,\text{in},\phi}(s, t) \left( \iota N \tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi}(s, t) - \bar{\sigma}_\mu^{2,\text{in},\phi}(s) \cdot \bar{\sigma}_\nu^{2,\text{in},\phi}(t) \right) \right] \quad (\text{N.45})$$

$$1 = \int \frac{d\hat{\Xi}_{\mu\nu}^{2,\text{in}}(s,t) d\hat{\Xi}_{\mu\nu}^{2,\text{in}}(s,t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\Xi}_{\mu\nu}^{2,\text{in}}(s,t) \left( \iota N \hat{\Xi}_{\mu\nu}^{2,\text{in}}(s,t) - \hat{\xi}_{\mu}^{2,\text{in}}(s) \cdot \hat{\xi}_{\nu}^{2,\text{in}}(t) \right) \right] \quad (\text{N.46})$$

$$1 = \int \frac{d\hat{\Phi}_{\mu\nu}^{\phi\phi,2,\text{in}}(s,t) d\hat{\Phi}_{\mu\nu}^{\phi\phi,2,\text{in}}(s,t)}{2\pi i (\iota N)^{-1}} \exp \left[ \hat{\Phi}_{\mu\nu}^{\phi\phi,2,\text{in}}(s,t) \left( \iota N \hat{\Phi}_{\mu\nu}^{\phi\phi,2,\text{in}}(s,t) - \bar{\sigma}_{\mu}^{\phi,2,\text{in}}(s) \cdot \bar{\sigma}_{\nu}^{\phi,2,\text{in}}(t) \right) \right] \quad (\text{N.47})$$

$$1 = \int \frac{d\tilde{A}_{\mu\nu}^{2,\text{in}}(s,t) d\tilde{B}_{\mu\nu}^{2,\text{in}}(s,t)}{2\pi i (\iota N)^{-1}} \exp \left[ -\tilde{B}_{\mu\nu}^{2,\text{in}}(s,t) \left( \iota N \tilde{A}_{\mu\nu}^{2,\text{in}}(s,t) - \bar{\sigma}_{\mu}^{2,\text{in}}(s) \cdot \bar{\xi}_{\nu}^{2,\text{in}}(t) \right) \right] \quad (\text{N.48})$$

$$1 = \int \frac{d\tilde{A}_{\mu\nu}^{\phi,2,\text{in}}(s,t) d\tilde{B}_{\mu\nu}^{\phi,2,\text{in}}(s,t)}{2\pi i (\iota N)^{-1}} \exp \left[ -\tilde{B}_{\mu\nu}^{\phi,2,\text{in}}(s,t) \left( \iota N \tilde{A}_{\mu\nu}^{\phi,2,\text{in}}(s,t) - \bar{\sigma}_{\mu}^{\phi,2,\text{in}}(s) \cdot \bar{\xi}_{\nu}^{2,\text{in}}(t) \right) \right] \quad (\text{N.49})$$

## N.10. Softmax

As for the previous regimes we enforce the softmax layer order parameter via the Fourier-transformed delta function

$$1 = \int \frac{d\mathcal{S}_{\mu}(t) d\hat{\mathcal{S}}_{\mu}(t)}{2\pi i (\kappa N)^{-1}} \exp \left[ \hat{\mathcal{S}}_{\mu}(t) \left( \kappa N \mathcal{S}_{\mu}(t) - \sum_{i=1}^{\kappa N} e^{\psi_{\mu}^i(t)} \right) \right] \quad (\text{N.50})$$

## N.11. Partition function

Define the set of all kernels and conjugates (indexed for time and feature, although this is omitted for brevity in N.51):

$$\begin{aligned} \mathcal{K} = \{ & \Phi^1, \hat{\Phi}^1, \bar{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \Phi^3, \hat{\Phi}^3, G^1, \hat{G}^1, \bar{G}^{2,\text{in}}, \hat{G}^{2,\text{in}}, \Phi_i^{2,\text{in}}, \hat{\Phi}_i^{2,\text{in}}, G_i^{2,\text{in}}, \hat{G}_i^{2,\text{in}}, \\ & G^3, \hat{G}^3, G^{1,\phi}, \hat{G}^{1,\phi}, \bar{A}^1, \bar{B}^1, \bar{B}^{2,\text{in}}, A^{1,\phi}, B^{1,\phi}, A^{g^{1,\phi}}, B^{g^{1,\phi}}, B^3, A^3, \\ & \tilde{G}^{2,\text{in}}, \tilde{G}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \tilde{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in},\phi}, \tilde{\Phi}^{2,\text{in},\phi}, \tilde{B}^{2,\text{in}}, \tilde{A}^{2,\text{in}}, \tilde{B}^{2,\text{in},\phi}, \tilde{A}^{2,\text{in},\phi}, \hat{\Xi}^{2,\text{in}}, \hat{\Xi}^{2,\text{in}}, \hat{\Phi}^{\phi\phi,2,\text{in}}, \Phi^{\phi\phi,2,\text{in}} \} \end{aligned} \quad (\text{N.51})$$

We can partition these order parameters into the set  $\mathcal{K}_{\text{global}}$  of global order parameters, the set  $\mathcal{K}_{\text{exp-loc}}$  of expert-local ones, together with the shared fields  $\mathcal{F}_j^{2,\text{in,sh}}$ , the router fields  $\mathcal{F}_i^{\phi}$ , the experts fields  $\mathcal{F}_{i,j}^{2,\text{in}}, \mathcal{F}_n^{\text{global-site}}$ :

$$\begin{aligned} \mathcal{K}_{\text{global}} = \{ & \Phi^1, \hat{\Phi}^1, \bar{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \Phi^3, \hat{\Phi}^3, G^1, \hat{G}^1, \bar{G}^{2,\text{in}}, \hat{G}^{2,\text{in}}, G^{1,\phi}, \hat{G}^{1,\phi}, G^3, \hat{G}^3, \bar{A}^1, \bar{B}^1, \bar{A}^{2,\text{in}}, A^{1,\phi}, \\ & B^{1,\phi}, A^{g^{1,\phi}}, B^{g^{1,\phi}}, \tilde{G}^{2,\text{in}}, \hat{G}^{2,\text{in}}, \tilde{\Phi}^{2,\text{in}}, \hat{\Phi}^{2,\text{in}}, \tilde{\Phi}^{2,\text{in},\phi}, \hat{\Phi}^{2,\text{in},\phi}, \tilde{B}^{2,\text{in}}, \\ & \tilde{A}^{2,\text{in}}, \tilde{B}^{2,\text{in},\phi}, \tilde{A}^{2,\text{in},\phi}, \hat{\Xi}^{2,\text{in}}, \hat{\Xi}^{2,\text{in}}, \Phi^{\phi\phi,2,\text{in}}, \hat{\Phi}^{\phi\phi,2,\text{in}}, S, \hat{S} \}. \end{aligned} \quad (\text{N.52})$$

$$\mathcal{K}_{\text{exp-loc},i} = \{ \Phi_i^{2,\text{in}}, \hat{\Phi}_i^{2,\text{in}}, G_i^{2,\text{in}}, \hat{G}_i^{2,\text{in}} \}. \quad (\text{N.53})$$

$$\mathcal{F}_j^{2,\text{in,sh}} = \{\hat{\chi}_{\mu,j}^{2,\text{in}}(t), \bar{g}_{\mu,j}^{2,\text{in}}(t), \hat{g}_{\mu,j}^{2,\text{in}}(t), \bar{\sigma}_{\mu,j}^{2,\text{in}}(t), \hat{\sigma}_{\mu,j}^{2,\text{in}}(t), \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t), \hat{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t), \bar{\xi}_{\mu,j}^{2,\text{in}}(t), \hat{\xi}_{\mu,j}^{2,\text{in}}(t)\}_{\mu,t}.$$

$$d\mathcal{F}_j^{2,\text{in,sh}} := \prod_{\mu,t} d\hat{\chi}_{\mu,j}^{2,\text{in}}(t) d\bar{g}_{\mu,j}^{2,\text{in}}(t) d\hat{g}_{\mu,j}^{2,\text{in}}(t) \times d\bar{\sigma}_{\mu,j}^{2,\text{in}}(t) d\hat{\sigma}_{\mu,j}^{2,\text{in}}(t) d\bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) d\hat{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \times d\bar{\xi}_{\mu,j}^{2,\text{in}}(t) d\hat{\xi}_{\mu,j}^{2,\text{in}}(t). \quad (\text{N.54})$$

$$\mathcal{F}_i^\phi = \{\chi_{\mu i}^\phi(t), \hat{\chi}_{\mu i}^\phi(t), g_{\mu i}^{1,\phi}(t), \xi_{\mu i}^{g1,\phi}(t), \hat{\xi}_{\mu i}^{g1,\phi}(t), \tilde{\phi}_\mu^i(t), \psi_\mu^i(t)\}_{\mu,t}. \quad (\text{N.55})$$

$$\mathcal{F}_{i,j}^{2,\text{in}} = \{\chi_{\mu,i,j}^{2,\text{in}}(t), h_{\mu,i,j}^{2,\text{in}}(t), g_{\mu,i,j}^{2,\text{in}}(t), \xi_{\mu,i,j}^{2,\text{in}}(t), \hat{\xi}_{\mu,i,j}^{2,\text{in}}(t)\}_{\mu,t}. \quad (\text{N.56})$$

$$\begin{aligned} \mathcal{F}_n^{\text{global-site}} = & \{h_{\mu,n}^1(t), g_{\mu,n}^1(t), \chi_{\mu,n}^1(t), \hat{\chi}_{\mu,n}^1(t), h_{\mu,n}^3(t), g_{\mu,n}^3(t), \xi_{\mu,n}^3(t), \hat{\xi}_{\mu,n}^3(t), \\ & \xi_{\mu,n}^{1,\phi}(t), \hat{\xi}_{\mu,n}^{1,\phi}(t), \bar{\chi}_{\mu,n}^3(t), \hat{\chi}_{\mu,n}^3(t), \bar{\xi}_{\mu,n}^1(t), \hat{\xi}_{\mu,n}^1(t)\}_{\mu,t}. \end{aligned}$$

Then the partition function can be written as:

$$\begin{aligned} Z \propto & \int \left( \prod_{\mu,\nu} \prod_{t,s} d\mathcal{K}_{\text{global}} \right) \exp \{N \mathcal{S}_{\text{global}}[\mathcal{K}_{\text{global}}]\} \times \prod_{n=1}^N \mathcal{Z}_N^{\text{global}} [j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3] \\ & \times \int \left[ \prod_{j=1}^{N_e} d\mathcal{F}_j^{2,\text{in,sh}} \right] \times \left[ \exp \left\{ \sum_{j=1}^{N_e} \mathcal{S}_j^{2,\text{in,sh}} \right\} \right. \\ & \left. \times \prod_{i=1}^{\kappa N} \mathcal{Z}_i^{\text{local}} \left[ j_i^{2,\text{in}}, v_i^{2,\text{in}}, j_i^{3,\text{out}}, j_i^\phi, v_i^{g1,\phi}; \{\mathcal{F}_j^{2,\text{in,sh}}\}_{j=1}^{N_e} \right] \right]. \end{aligned}$$

with the global action defined as:

$$\begin{aligned}
 \mathcal{S}_{\text{global}} = & \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\Phi}_{\mu\nu}^1(t,s) \Phi_{\mu\nu}^1(t,s) + \hat{G}_{\mu\nu}^1(t,s) G_{\mu\nu}^1(t,s) + \kappa \hat{G}_{\mu\nu}^{1,\phi}(t,s) G_{\mu\nu}^{1,\phi}(t,s) \right. \\
 & - \bar{B}_{\mu\nu}^1(t,s) \bar{A}_{\mu\nu}^1(t,s) - B_{\mu\nu}^{1,\phi}(t,s) A_{\mu\nu}^{1,\phi}(t,s) + \hat{\Phi}_{\mu\nu}^3(t,s) \Phi_{\mu\nu}^3(t,s) \\
 & + \kappa \hat{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \\
 & + \kappa \hat{G}_{\mu\nu}^{2,\text{in}}(t,s) \bar{G}_{\mu\nu}^{2,\text{in}}(t,s) + \kappa \hat{S}_\mu(t) S_\nu(s) + \hat{G}_{\mu\nu}^3(t,s) G_{\mu\nu}^3(t,s) \\
 & - \kappa B_{\mu\nu}^{g^{1,\phi}}(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) + \iota \hat{G}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \tilde{G}_{\mu\nu}^{2,\text{in}}(t,s) \\
 & + \iota \hat{\Phi}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \tilde{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) + \iota \hat{\Phi}_{\mu\nu}^{\hat{2},\text{in},\phi}(t,s) \tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi}(t,s) \\
 & - \iota \tilde{B}_{\mu\nu}^{2,\text{in}}(t,s) \tilde{A}_{\mu\nu}^{2,\text{in}}(t,s) - \iota \tilde{B}_{\mu\nu}^{2,\text{in},\phi}(t,s) \tilde{A}_{\mu\nu}^{2,\text{in},\phi}(t,s) \\
 & + \iota \hat{\Xi}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \hat{\Xi}_{\mu\nu}^{2,\text{in}}(t,s) + \iota \hat{\Phi}_{\mu\nu}^{\phi\phi,2,\text{in}}(t,s) \Phi_{\mu\nu}^{\phi\phi,2,\text{in}}(t,s) \\
 & + i\kappa \bar{A}_{\mu\nu}^{2,\text{in}}(t,s) G_{\mu\nu}^3(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) - \iota \frac{1}{2} \hat{\Xi}_{\mu\nu}^{2,\text{in}} G^3 \\
 & + \kappa A_{\mu\nu}^{g^{1,\phi}}(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) G_{\mu\nu}^3(t,s) \tilde{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \\
 & + i\kappa \iota \tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi}(t,s) G_{\mu\nu}^3(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) - \frac{\iota}{2} \Phi_{\mu\nu}^{\phi\phi,2,\text{in}}(t,s) G_{\mu\nu}^3(t,s) \\
 & \left. - \iota \kappa \tilde{A}_{\mu\nu}^{\phi,2,\text{in}}(t,s) G_{\mu\nu}^3(t,s) + i\kappa \tilde{A}_{\mu\nu}^{2,\text{in}}(t,s) G_{\mu\nu}^3(t,s) A_{\mu\nu}^{g^{1,\phi}}(t,s) \right].
 \end{aligned}$$

and the single-site global action as:

$$\mathcal{Z}_N^{\text{global}}[j_n^1, \bar{j}_n^3, \bar{v}_n^1, v_n^{1,\phi}, v_n^3] := \int d\mathcal{K}_n^{\text{global-site}} \exp \left\{ \mathcal{S}_n^{\text{global-site}} \right\}.$$

$$\begin{aligned}
 \mathcal{S}_n^{\text{global-site}} = & -\frac{1}{2} \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\chi}_{\mu,n}^1(t) \hat{\chi}_{\nu,n}^1(s) \Phi_{\mu\nu}^0(t,s) + \hat{\xi}_{\mu,n}^3(t) \hat{\xi}_{\nu,n}^3(s) \right. \\
 & \left. + \hat{\xi}_{\mu,n}^{1,\phi}(t) \hat{\xi}_{\nu,n}^{1,\phi}(s) G_{\mu\nu}^{1,\phi}(t,s) \right] \\
 & - \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\Phi}_{\mu\nu}^1(t,s) \sigma(h_{\mu,n}^1(t)) \sigma(h_{\nu,n}^1(s)) + \hat{G}_{\mu\nu}^1(t,s) g_{\mu,n}^1(t) g_{\nu,n}^1(s) \right. \\
 & + \hat{\Phi}_{\mu\nu}^3(t,s) h_{\mu,n}^3(t) h_{\nu,n}^3(s) + \hat{G}_{\mu\nu}^3(t,s) g_{\mu,n}^3(t) g_{\nu,n}^3(s) \\
 & - \frac{1}{2} \tilde{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \hat{\chi}_{\mu,n}^3(t) \hat{\chi}_{\nu,n}^3(s) + \frac{1}{2} \tilde{G}_{\mu\nu}^{2,\text{in}}(t,s) \hat{\xi}_{\mu,n}^1(t) \hat{\xi}_{\nu,n}^1(s) \\
 & \left. - \iota\kappa \tilde{\Phi}_{\mu\nu}^{2,\text{in}} \hat{\chi}^3 g^3 A^{g^{1,\phi}} + \iota\kappa \tilde{A}_{\mu\nu}^{2,\text{in}} \hat{\chi}^3 g^3 \right] \\
 & - i \sum_{\mu,\nu} \int dt \int ds \left[ B_{\mu\nu}^{1,\phi}(t,s) \hat{\xi}_{\mu,n}^{1,\phi}(t) \sigma(h_{\nu,n}^1(s)) + \bar{B}_{\mu\nu}^1(t,s) \hat{\xi}_{\mu,n}^1(t) \sigma(h_{\nu,n}^1(s)) \right. \\
 & \left. + \iota\kappa \tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi} \hat{\chi}^3 g^3 \right] \\
 & + \sum_{\mu} \int dt \left[ (v_{\mu,n}^3(t) + i\hat{\xi}_{\mu,n}^3(t)) \xi_{\mu,n}^3(t) + (v_{\mu,n}^{1,\phi}(t) + i\hat{\xi}_{\mu,n}^{1,\phi}(t)) \xi_{\mu,n}^{1,\phi}(t) \right. \\
 & + (j_{\mu,n}^1(t) + i\hat{\chi}_{\mu,n}^1(t)) \chi_{\mu,n}^1(t) + (\bar{j}_{\mu,n}^3(t) + i\hat{\chi}_{\mu,n}^3(t)) \bar{\chi}_{\mu,n}^3(t) \\
 & \left. + (\bar{v}_{\mu,n}^1(t) + i\hat{\xi}_{\mu,n}^1(t)) \bar{\xi}_{\mu,n}^1(t) \right].
 \end{aligned}$$

$$\begin{aligned}
 \mathcal{S}_j^{2,\text{in,sh}} = & -\frac{1}{2} \sum_{\mu,\nu} \int dt \int ds N \left[ \hat{\chi}_{\mu,j}^{2,\text{in}}(t) \hat{\chi}_{\nu,j}^{2,\text{in}}(s) \Phi_{\mu\nu}^1(t,s) - \hat{G}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \bar{g}_{\mu,j}^{2,\text{in}}(t) \bar{g}_{\nu,j}^{2,\text{in}}(s) \right. \\
 & - \hat{\Phi}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \bar{\sigma}_{\mu,j}^{2,\text{in}}(t) \bar{\sigma}_{\nu,j}^{2,\text{in}}(s) - \hat{\Phi}_{\mu\nu}^{\hat{2},\text{in},\phi}(t,s) \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \bar{\sigma}_{\nu,j}^{2,\text{in}}(s) \\
 & - \hat{\Xi}_{\mu\nu}^{\hat{2},\text{in}}(t,s) \hat{\xi}_{\mu,j}^{\hat{2},\text{in}}(t) \hat{\xi}_{\nu,j}^{\hat{2},\text{in}}(s) \\
 & - \hat{\Phi}_{\mu\nu}^{\phi,2,\text{in}}(t,s) \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \bar{\sigma}_{\nu,j}^{\phi,2,\text{in}}(s) - \tilde{B}_{\mu\nu}^{2,\text{in}}(t,s) \bar{\sigma}_{\mu,j}^{2,\text{in}}(t) \hat{\xi}_{\nu,j}^{\hat{2},\text{in}}(s) \\
 & \left. - \tilde{B}_{\mu\nu}^{\phi,2,\text{in}}(t,s) \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \hat{\xi}_{\nu,j}^{\hat{2},\text{in}}(s) \right] \\
 & - i \sum_{\mu} \int dt \left[ \hat{g}_{\mu,j}^{2,\text{in}}(t) \bar{g}_{\mu,j}^{2,\text{in}}(t) + \hat{\sigma}_{\mu,j}^{2,\text{in}}(t) \bar{\sigma}_{\mu,j}^{2,\text{in}}(t) \right. \\
 & \left. + \hat{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) + \hat{\xi}_{\mu,j}^{\hat{2},\text{in}}(t) \hat{\xi}_{\mu,j}^{\hat{2},\text{in}}(t) \right] \\
 & - i \sum_{\mu,\nu} \int dt \int ds \left[ -\frac{1}{2} \Phi_{\mu\nu}^1(t,s) \hat{\chi}_{\mu,j}^{2,\text{in}}(t) \hat{\chi}_{\nu,j}^{2,\text{in}}(s) \right].
 \end{aligned}$$

$$\begin{aligned}
 \mathcal{Z}_i^{\text{local}} \left[ j_i^{2,\text{in}}, v_i^{2,\text{in}}, j_i^{3,\text{out}}, j_i^{\phi}, v_i^{g^{1,\phi}}; \{\mathcal{K}_j^{2,\text{in,sh}}\}_{j=1}^{N_e} \right] & := \int d\mathcal{K}_{\text{exp-loc},i} \\
 \int d\mathcal{K}_i^{\phi} \times \exp \left\{ \mathcal{S}_i^{\text{exp-loc}} + \mathcal{S}_i^{\phi} \right\} \times \frac{\iota}{\iota N} \prod_{j=1}^{N_e} \left[ \int d\mathcal{K}_{i,j}^{2,\text{in}} \exp \left\{ \mathcal{S}_{i,j}^{2,\text{in}} \right\} \right].
 \end{aligned}$$

$$\mathcal{S}_i^{\text{exp-loc}} = \sum_{\mu,\nu} \int_0^\infty dt \int_0^\infty ds \left[ \iota \hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(t,s) \Phi_{\mu\nu,i}^{2,\text{in}}(t,s) + \iota \hat{G}_{\mu\nu,i}^{2,\text{in}}(t,s) G_{\mu\nu,i}^{2,\text{in}}(t,s) \right]. \quad (\text{N.57})$$

$$\begin{aligned} \mathcal{S}_i^\phi &= \sum_{\mu} \int_0^\infty dt g_{\mu i}^{1,\phi}(t) g_{\mu i}^{1,\phi}(t) \hat{G}_{\mu\nu}^{1,\phi} - \frac{1}{2} \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\chi}_{\mu i}^\phi(t) \hat{\chi}_{\nu i}^\phi(s) \Phi_{\mu\nu}^1(t,s) \right] \\ &\quad - i \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\chi}_{\mu i}^\phi(t) g_{\nu i}^{1,\phi}(s) A_{\mu\nu}^{1,\phi}(t,s) - B_{\mu\nu}^{g^{1,\phi}}(t,s) \hat{\xi}_{\mu i}^{g^{1,\phi}}(t) \tilde{\phi}_\nu^i(s) \right] \\ &\quad - \sum_{\mu,\nu} \int dt \int ds \left[ \hat{G}_{\mu\nu}^{1,\phi}(t,s) g_{\mu i}^{1,\phi}(t) g_{\nu i}^{1,\phi}(s) + \hat{\Phi}_{\mu\nu}^{2,\text{in}}(t,s) \tilde{\phi}_\mu^i(t) \tilde{\phi}_\nu^i(s) \Phi_{\mu\nu,i}^{2,\text{in}}(t,s) \right] \\ &\quad + \hat{G}_{\mu\nu}^{2,\text{in}}(t,s) G_{\mu\nu,i}^{2,\text{in}}(t,s) \\ &\quad - \sum_{\mu} \int dt \hat{S}_\mu(t) e^{\psi_\mu^i(t)} + \sum_{\mu} \int dt \left[ (j_{\mu i}^\phi(t) + i \hat{\chi}_{\mu i}^\phi(t)) \chi_{\mu i}^\phi(t) + (v_{\mu i}^{g^{1,\phi}}(t) + i \hat{\xi}_{\mu i}^{g^{1,\phi}}(t)) \xi_{\mu i}^{g^{1,\phi}}(t) \right]. \end{aligned}$$

$$\begin{aligned} \mathcal{S}_{i,j}^{2,\text{in}} &= -i \sum_{\mu,\nu} \int dt \int ds \left[ \frac{1}{\kappa} \hat{\chi}_{\mu,j}^{2,\text{in}}(t) g_{\nu,i,j}^{2,\text{in}}(s) \bar{A}_{\mu\nu}^1(t,s) \right] \\ &\quad - i \sum_{\mu} \int dt \left[ \hat{g}_{\mu,j}^{2,\text{in}}(t) g_{\mu,i,j}^{2,\text{in}}(t) \hat{\sigma}_{\mu,j}^{2,\text{in}}(t) \sigma \left( h_{\mu,i,j}^{2,\text{in}}(t) \right) \tilde{\phi}_\mu^i(t) \right. \\ &\quad \quad \left. - \hat{\sigma}_{\mu,j}^{\phi,2,\text{in}}(t) \sigma \left( h_{\mu,i,j}^{2,\text{in}}(t) \right) \tilde{\phi}_\mu^i(t) \hat{\xi}_{\mu i}^{g^{1,\phi}}(t) - \hat{\xi}_{\mu,j}^{2,\text{in}}(t) \hat{\xi}_{\mu,i,j}^{2,\text{in}}(t) \tilde{\phi}_\mu^i(t) \right] \\ &\quad - \sum_{\mu,\nu} \int dt \int ds \left[ \hat{\Phi}_{\mu\nu,i}^{2,\text{in}}(t,s) \sigma \left( h_{\mu,i,j}^{2,\text{in}}(t) \right) \sigma \left( h_{\nu,i,j}^{2,\text{in}}(s) \right) + \hat{G}_{\mu\nu,i}^{2,\text{in}}(t,s) g_{\mu,i,j}^{2,\text{in}}(t) g_{\nu,i,j}^{2,\text{in}}(s) \right] \\ &\quad + \sum_{\mu} \int dt N \left[ (j_{\mu i}^{2,\text{in}}(t) + i \hat{\chi}_{\mu,j}^{2,\text{in}}(t)) \chi_{\mu,i,j}^{2,\text{in}}(t) + \left( \frac{1}{\kappa N} v_{\mu,i}^{2,\text{in}}(t) + i \hat{\xi}_{\mu,i,j}^{2,\text{in}}(t) \right) \xi_{\mu,i,j}^{2,\text{in}}(t) \right]. \end{aligned}$$

One way to look at it is the following: there are four conditional factorization levels in the full partition function. The first is the standard global  $N$ -site DMFT factorization over  $n$ . The remaining three form a nested local/shared sector: a shared  $N_e$ -component factorization over  $j$ , an expert factorization over  $i$  conditional on the shared  $j$ -fields, and an intra-expert  $N_e$ -component factorization over  $(i,j)$ -local fields. The important point is that these are *conditional* factorizations.

## N.12. Saddle Point Approximation

To write down the saddle point equations, we first define the single-site distributions, where  $\mathcal{H}$  in each case is the logarithm of the integrand of the corresponding  $\mathcal{Z}$ . We can then define for each  $\mathcal{Z} \in \{\mathcal{Z}_N^{\text{global}}, \mathcal{Z}_{\text{sh}}, \mathcal{Z}_{N_e j}, \mathcal{Z}_M\}$  and the corresponding  $\mathcal{H}$  the average

$$\langle \mathcal{O}(\{\chi, \xi\}) \rangle_{\mathcal{Z}} = \frac{1}{\mathcal{Z}} \int \prod_{\mu} d\mathcal{F} \exp(-\mathcal{H}[\{\chi, \xi\}, \{j, v\}]) \mathcal{O}(\{\chi, \xi\}) \quad (\text{N.58})$$

With this apparatus in place, we can take saddle equations. We treat expert-local kernels as microvariables which are implicitly defined in terms of  $\chi, \xi$ , and so do not take saddle equations of them at this point.

### N.13. Saddle-Point Equations

As in Regime II, we can use the saddle point approximation, which consist in finding the set of equations that lead to a stationary action  $S[\mathcal{K}]$ .

We note the following:

- At zero source, all single site averages  $\langle \rangle_{\mathcal{Z}_N^{global}}, \langle \rangle_{\mathcal{Z}_{sh}}$  are equivalent.
- *Conditional on the set  $\mathcal{K}_{global}, \mathcal{F}_{sh}$  of global and shared kernels,  $\mathcal{Z}^{local}$  factorises over experts, and so we can write  $\langle \rangle_{\mathcal{Z}_{Nej}}$  for the conditional average over the distribution over the experts neurons and  $\langle \rangle_{\mathcal{Z}_M|\mathcal{K}_{global}}$  for the conditional average over the distribution of the router's neurons.*
- Expert-local variables follow single-site processes, so we can drop expert indices
- By a derivation similar to that in **K**, we can prove that conjugate kernels defined as covariances between conjugate fields vanish, since they have no physical meaning and can not influence the dynamics in addition to imposing constraint when introducing kernel definitions.

The global kernels are given by their expectation values over the single-site global distribution  $\mathcal{Z}_N^{global}$ :

$$\begin{aligned}
 \Phi_{\mu\nu}^1(s, t) &= \langle \sigma(h_\mu^1(s)) \sigma(h_\nu^1(t)) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 \Phi_{\mu\nu}^3(s, t) &= \langle h_\mu^3(s) h_\nu^3(t) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 G_{\mu\nu}^1(s, t) &= \langle g_\mu^1(s) g_\nu^1(t) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 G_{\mu\nu}^3(s, t) &= \langle g_\mu^3(s) g_\nu^3(t) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 A_{\mu\nu}^{1,\phi}(s, t) &= -i \langle \hat{\xi}_\mu^{1,\phi}(s) \sigma(h_\nu^1(t)) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 \bar{A}_{\mu\nu}^1(s, t) &= -i \langle \hat{\xi}_\mu^1(t) \sigma(h_\nu^1(s)) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}} \\
 \bar{B}_{\mu\nu}^{2,in}(s, t) &= -i \langle \hat{\chi}_\mu^3(t) g_\nu^3(s) \rangle_{\mathcal{Z}_N^{global}|\mathcal{K}_{global}}
 \end{aligned} \tag{N.59}$$

The shared kernels are given by their expectation values over the single-site shared distribution  $\mathcal{Z}_{sh}$ :

$$\tilde{G}_{\mu\nu}^{2,in}(s, t) = \langle \bar{g}_{\mu,j}^{2,in}(s) \bar{g}_{\nu,j}^{2,in}(t) \rangle_{\mathcal{Z}_{sh}|\mathcal{K}_{global}}, \tag{N.60}$$

$$\tilde{\Phi}_{\mu\nu}^{2,in}(s, t) = \langle \bar{\sigma}_{\mu,j}^{2,in}(s) \bar{\sigma}_{\nu,j}^{2,in}(t) \rangle_{\mathcal{Z}_{sh}|\mathcal{K}_{global}}, \tag{N.61}$$

$$\tilde{\Phi}_{\mu\nu}^{2,in,\phi}(s, t) = \langle \bar{\sigma}_{\mu,j}^{\phi,2,in}(s) \bar{\sigma}_{\nu,j}^{2,in}(t) \rangle_{\mathcal{Z}_{sh}|\mathcal{K}_{global}}, \tag{N.62}$$

$$\hat{\Xi}_{\mu\nu}^{2,in}(s, t) = \langle \hat{\xi}_{\mu,j}^{2,in}(s) \hat{\xi}_{\nu,j}^{2,in}(t) \rangle_{\mathcal{Z}_{sh}|\mathcal{K}_{global}} \tag{N.63}$$

$$\tilde{\Phi}_{\mu\nu}^{\phi\phi,2,in}(s, t) = \langle \bar{\sigma}_{\mu,j}^{\phi,2,in}(s) \bar{\sigma}_{\nu,j}^{\phi,2,in}(t) \rangle_{\mathcal{Z}_{sh}|\mathcal{K}_{global}} \tag{N.64}$$

$$\tilde{A}_{\mu\nu}^{2,\text{in}}(s, t) = \left\langle \bar{\sigma}_{\mu,j}^{2,\text{in}}(s) \bar{\xi}_{\nu,j}^{2,\text{in}}(t) \right\rangle_{\mathcal{Z}_{sh} | \mathcal{K}_{\text{global}}} \quad (\text{N.65})$$

$$\tilde{A}_{\mu\nu}^{\phi,2,\text{in}}(s, t) = \left\langle \bar{\sigma}_{\mu,j}^{\phi,2,\text{in}}(s) \bar{\xi}_{\nu,j}^{2,\text{in}}(t) \right\rangle_{\mathcal{Z}_{sh} | \mathcal{K}_{\text{global}}} \quad (\text{N.66})$$

$$\tilde{B}_{\mu\nu}^{2,\text{in}}(s, t) = i\kappa G_{\mu\nu}^3(s, t) A_{\mu\nu}^{g^{1,\phi}}(s, t) \quad (\text{N.67})$$

$$\tilde{B}_{\mu\nu}^{\phi,2,\text{in}}(s, t) = -i\kappa G_{\mu\nu}^3(s, t) \quad (\text{N.68})$$

The router kernels require expectations over the expert ensemble distribution  $\mathcal{Z}_M$ , conditioned on the global kernels and shared fields:

$$\begin{aligned} \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) &= \left\langle \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \Phi_{\mu\nu}^{2,\text{in}}(s, t) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) &= \left\langle G_{\mu\nu}^{2,\text{in}}(s, t) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ G_{\mu\nu}^{1,\phi}(s, t) &= \left\langle g_\mu^{1,\phi}(s) g_\nu^{1,\phi}(t) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ \mathcal{S}_\mu(t) &= \left\langle e^{\psi_\mu(t)} \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ A_{\mu\nu}^{g^{1,\phi}}(s, t) &= -i \left\langle \hat{\xi}_\mu^{g^{1,\phi}}(s) \tilde{\phi}_\nu(t) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ B_{\mu\nu}^{1,\phi}(s, t) &= -i\kappa \left\langle \hat{\chi}_\mu^\phi(s) g_\nu^{1,\phi}(t) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ B_{\mu\nu}^{g^{1,\phi}}(s, t) &= +i\bar{A}_{\mu\nu}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s) + A_{\mu\nu}^{g^{1,\phi}}(t, s) G_{\mu\nu}^3(t, s) \tilde{\Phi}_{\mu\nu}^{2,\text{in}}(t, s) \\ &\quad + i\nu \tilde{\Phi}_{\mu\nu}^{2,\text{in},\phi}(t, s) G_{\mu\nu}^3(t, s) + i\bar{A}_{\mu\nu}^{2,\text{in}}(t, s) G_{\mu\nu}^3(t, s) - i\nu \tilde{\Phi}_{\mu\nu}^{2,\text{in}} \langle \hat{\chi}^3 g^3 \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \end{aligned} \quad (\text{N.69})$$

#### N.14. Expert-Local Order Parameters

The local kernels are determined by independent processes within each expert's internal dimensions:

$$\begin{aligned} \Phi_{\mu\nu}^{2,\text{in}}(s, t) &= \left\langle \sigma([h_\mu^{2,\text{in}}(t)]_j) \sigma([h_\nu^{2,\text{in}}(s)]_j) \right\rangle_{\mathcal{Z}_{Ne\ j} | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ G_{\mu\nu}^{2,\text{in}}(s, t) &= \left\langle [g_\mu^{2,\text{in}}(t)]_j [g_\nu^{2,\text{in}}(s)]_j \right\rangle_{\mathcal{Z}_{Ne\ j} | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ A_{\mu\nu}^{2,\text{in}}(s, t) &= \left\langle [\hat{\xi}_\mu^{2,\text{in}}(s)]_j \sigma([h_\nu^{2,\text{in}}]_j) \right\rangle_{\mathcal{Z}_{Ne\ j} | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \\ \bar{B}_{\mu\nu}^1(s, t) &= -\frac{i}{\kappa} \left\langle \hat{\chi}_\mu^{2,\text{in}}(s) g_\nu^{2,\text{in}}(t) \right\rangle_{\mathcal{Z}_{Ne\ j} | \mathcal{K}_{\text{global}}, \mathcal{F}_{\text{shared}}} \end{aligned} \quad (\text{N.70})$$

All non-physical conjugate kernels identically vanish at the saddle point:  $\hat{\Phi}^3 = \hat{\Phi}^1 = \hat{G}^1 = \hat{\Phi}^{2,\text{in}} = \dots = 0$ .

### N.15. Hubbard-Stratonovich transformation

The Hubbard-Stratonovich trick allows us to rewrite the quadratic terms over conjugate fields as Gaussian integrals over linear conjugate fields:

$$\exp\left(-\frac{1}{2} \hat{\mathbf{x}}^\top A \hat{\mathbf{x}}\right) = \int_{\mathbb{R}^d} \frac{d\mathbf{u}}{(2\pi)^{d/2} \sqrt{\det A}} \exp\left(-\frac{1}{2} \mathbf{u}^\top A^{-1} \mathbf{u} - i \mathbf{u} \cdot \hat{\mathbf{x}}\right) = \left\langle \exp(-i \mathbf{u} \cdot \hat{\mathbf{x}}) \right\rangle_{\mathbf{u} \sim \mathcal{N}(0, A)}. \quad (\text{N.71})$$

where  $\hat{\mathbf{x}}$  is a generic conjugate fields in the partition function above. Using Stein's lemma, we can reformulate the definitions of each of the kernels A and B as response functions:

$$\begin{aligned} A_{\mu\nu,i}^\ell(t, s) &= -i \left\langle \hat{\xi}_{\nu,i}^\ell(t) \sigma(h_{\mu,i}^{\ell-1}(s)) \right\rangle_{\beta_{\nu,i}^\ell(s)} \\ &= [G_i^\ell]^{-1} \left\langle (\xi_i^\ell - B_i^\ell \sigma(h_i^{\ell-1})) \sigma(h_i^{\ell-1}) \right\rangle_{\beta_{\nu,i}^\ell(s)} \\ &= \left\langle \frac{\partial \sigma(h_{\mu,i}^{\ell-1}(t))}{\partial \beta_{\nu,i}^{\ell\top}(s)} \right\rangle_{\beta_{\nu,i}^\ell(s)} \end{aligned} \quad (\text{N.72})$$

It easier now integrate over all the  $\hat{\mathbf{x}}$ 's, since the argument of the exponential in  $\mathcal{Z}$  has been linearised with respect to them all. Doing so yields delta functions that give us the final DMFT dynamics.

### N.16. DMFT Dynamics

Piecing together the self-consistent order parameters and the exact integrations over the Hubbard-Stratonovich fields, the DMFT description of the infinite-width MoE is summarized below.

$$\begin{aligned} h_\mu^1(t) &= \alpha_\mu^1 + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu g_\nu^1(s) \Phi_{\mu\nu}^0 \quad \alpha^1 \sim \mathcal{N}(0, K^x) \\ h_{\mu,i}^{2,\text{in}}(t) &= \alpha_\mu^{2,\text{in}}(t) + \frac{\eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \left[ \Delta_\nu \Phi_{\mu\nu}^1(s, t) + \bar{A}_{\nu\mu}^1(s, t) \right] g_{\nu,i}^{2,\text{in}}(s) \quad \alpha^{2,\text{in}} \sim \mathcal{N}(0, \Phi^1(t)) \\ h_\mu^3(t) &= \bar{\alpha}_\mu^3(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left[ \Delta_\nu \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) + \mathcal{R}_{\nu\mu}^3(s, t) \right] g_\nu^3(s) \quad \bar{\alpha}^3 \sim \mathcal{N}(0, \bar{\Phi}^{2,\text{in}}(t)) \\ &\text{with } \mathcal{R}_{\nu\mu}^3(s, t) = \kappa \bar{A}_{\nu\mu}^{2,\text{in}}(s, t) + \kappa \tilde{\Phi}_{\nu\mu}^{2,\text{in},\phi}(s, t) - \kappa \tilde{\Phi}_{\nu\mu}^{2,\text{in}}(s, t) A_{\nu\mu}^{g^1,\phi}(s, t) \\ \psi_\mu(t) &= \alpha_\mu^\phi + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_{\nu=1}^P \left[ \Delta_\nu(s) \Phi_{\mu\nu}^1(s, t) + A_{\mu\nu}^{1,\phi} \right] g_\nu^{1,\phi}(s) \quad \alpha^\phi \sim \mathcal{N}\left(0, \frac{1}{\kappa} \Phi^1(t)\right) \\ \tilde{\phi}_\mu(t) &= \frac{e^{\psi_\mu(t)}}{\mathcal{S}_\mu(t)}, \quad \hat{\psi}_\mu(t) = m_\mu(t) \psi_\mu(t), \quad m_\mu(t) = \mathbf{1}(\psi_\mu(t) - \tau_\mu(t)) \\ &\text{where } \tau_\mu(t) \text{ is chosen such that } \rho = \langle m_\mu(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{global}} \end{aligned} \quad (\text{N.73})$$

$$\begin{aligned}
 g_\mu^\phi(t) &= \frac{\eta_0 \gamma_0^l}{P \kappa} \int_0^t ds \sum_\nu G_{\mu\nu}^3(t, s) \tilde{\phi}_\mu \left[ \Phi_{\mu\nu, i}^{2, \text{in}}(t, s) \tilde{\phi}_\nu - \bar{\Phi}_{\mu\nu}^2(t, s) \right] \Delta_\nu \\
 z_\mu^1(t) &= \bar{\beta}_\mu^1(t) + \beta_\mu^{1, \phi}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left[ \bar{B}_{\mu\nu}^1 + B_{\mu\nu}^{1, \phi}(t, s) + \left[ \Delta_\nu \frac{l}{\kappa} \bar{G}_{\mu\nu}^2(s, t) \right. \right. \\
 &\quad \left. \left. + \Delta_\nu G_{\mu\nu}^{1, \phi}(s, t) \right] \right] \sigma(h_\nu^1(s)) \\
 \bar{\beta}_\mu^1(t) &\sim \mathcal{N}(0, \tilde{G}^{2, \text{in}}(t, t)) \quad \beta_\mu^{1, \phi}(t) \sim \mathcal{N}(0, G^{1, \phi}(t, t)) \\
 z_\mu^{2, \text{in}}(t) &= \beta_\mu^{2, \text{in}}(t) \tilde{\phi}_\mu^i + \frac{\eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu \sigma(h_{\nu, i}^{2, \text{in}}(s)) G_{\mu\nu}^3(s, t) \tilde{\phi}_\mu \tilde{\phi}_\nu \quad \beta_\mu^{2, \text{in}}(t) \sim \mathcal{N}(0, G^3(t, t)) \\
 z_\mu^3(t) &= \beta_\mu^3(t) + \frac{\gamma_0 \eta_0}{P} \int_0^t ds \sum_\nu \Delta_\nu h_\nu^3(s) \quad \beta_\mu^3 \sim \mathcal{N}(0, 1),
 \end{aligned} \tag{N.74}$$

**O. DMFT Analysis for Regime III ( $\mu P$ )**

In this section we are considering the case in which  $M, N, N_e \rightarrow \infty$  at fixed ratios. Specifically, we define  $\kappa = \frac{M}{N}, \iota = \frac{N_e}{N}$ , where  $\kappa, \iota$  are order one in  $N$ . The forward pass, following  $\mu P$  is defined as:

$$\begin{aligned}
 \mathbb{R}^N \ni h_\mu^1 &= \frac{1}{\sqrt{D}} W^1 x_\mu \\
 \mathbb{R}^M \ni \psi_\mu &= \frac{1}{\sqrt{N}} Q \sigma(h_\mu^1) \\
 \mathbb{R}^M \ni \phi_\mu &= \text{softmax}(\psi_\mu) \\
 \mathbb{R}^{N_e} \ni h_{\mu,i}^{2,\text{in}} &= \frac{1}{\sqrt{N}} W_i^{2,\text{in}} \sigma(h_\mu^1), \\
 \mathbb{R}^N \ni h_{\mu,i}^{2,\text{out}} &= \frac{1}{\sqrt{N}} W_i^{2,\text{out}} \sigma(h_{\mu,i}^{2,\text{in}}), \\
 \mathbb{R}^N \ni h_\mu^3 &= \sum_{i=1}^M \phi_\mu^i h_{\mu,i}^{2,\text{out}} \\
 \mathbb{R} \ni h_\mu^4 &= \frac{1}{\sqrt{N}} w^{4\text{T}} h_\mu^3 \\
 f_\mu &= \frac{1}{\gamma} h_\mu^4 \\
 \eta &= \eta_0 \gamma^2 \\
 \eta_E &= \eta_0 \gamma^2 N \\
 \eta_Q &= \eta_0 \gamma^2 \kappa \\
 \gamma &= \gamma_0 \sqrt{N} \\
 \eta_0, \gamma_0 &\sim O(1)
 \end{aligned} \tag{O.1}$$

$$w_\alpha^4(0), [W_i^{2,\text{out}}(0)]_{\alpha\beta}, [W_i^{2,\text{in}}(0)]_{\alpha\beta}, W_{\alpha\beta}^1(0), Q_{\alpha\beta}(0) \sim \mathcal{N}(0, 1)$$

Again, we let  $\theta = \text{Vec}\{W^1, W_i^{2,\text{in}}, W_i^{2,\text{out}}, w^4, Q\}$ , and define

$$\tilde{\phi}_\mu^i := \kappa N \phi_\mu^i \tag{O.2}$$

Following an analogous derivation as in Appendix L, which we omit for conciseness, we arrive at the following self-consistent set of DMFT equations.

$$\begin{aligned}
 \alpha^1(t) &\sim \mathcal{N}(0, \Phi^0) & \beta^3(t) &\sim \mathcal{N}(0, \mathbb{1}) & \beta^{1,\phi}(t) &\sim \mathcal{GP}(0, G^{1,\phi}) \\
 \alpha^\phi(t) &\sim \mathcal{GP}(0, \Phi^1) & \alpha^{2,\text{in}}(t) &\sim \mathcal{GP}(0, \Phi^1) & \beta^{2,\text{in}}(t) &\sim \mathcal{GP}(0, \tilde{\phi} \tilde{\phi} G^3) \\
 \Phi_{\mu\nu}^0 &= \frac{1}{D} x_\mu x_\nu^\top, & \Phi_{\mu\nu}^1(t, s) &= \langle \sigma(h_\mu^1(t)) \sigma(h_\nu^1(s)) \rangle_{\mathcal{Z}_N^{\text{global}}}, & \Phi_{\mu\nu}^{2,\text{in}}(s, t) &= \langle \sigma(h_\mu^{2,\text{in}}(t)) \sigma(h_\nu^{2,\text{in}}(s)) \rangle_{\mathcal{Z}_{N_e} | \mathcal{K}_{\text{global}}} \\
 \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \Phi_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} & \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) &= \langle G_{\mu\nu}^{2,\text{in}}(s, t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 \Phi_{\mu\nu}^3(s, t) &= \langle h_\mu^3(t) h_\nu^3(s) \rangle_{\mathcal{Z}_N^{\text{global}}}, & G_{\mu\nu}^1(t, s) &= \langle [\dot{\sigma}(h_\mu^1(t)) \odot z_\mu^1(t)] [\dot{\sigma}(h_\nu^1(s)) \odot z_\nu^1(s)] \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 G_{\mu\nu}^{2,\text{in}}(s, t) &= \langle [\dot{\sigma}(h_\mu^{2,\text{in}}(t)) \odot z_\mu^{2,\text{in}}(t)] [\dot{\sigma}(h_\nu^{2,\text{in}}(s)) \odot z_\nu^{2,\text{in}}(s)] \rangle_{\mathcal{Z}_{N_e} | \mathcal{K}_{\text{global}}} \\
 G_{\mu\nu}^{1,\phi}(s, t) &= \langle g_\mu^{1,\phi}(s) g_\nu^{1,\phi}(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}, & G_{\mu\nu}^3(s, t) &= \langle z_\mu^3(s) z_\nu^3(t) \rangle_{\mathcal{Z}_N^{\text{global}}} \\
 A_{\mu\nu}^{2,\text{in}}(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^{2,\text{in}}(t))}{\partial \beta_\nu^{2,\text{in}\top}(s)} \right\rangle_{\mathcal{Z}_{N_e} | \mathcal{K}_{\text{global}}}, & \bar{A}_{\mu\nu}^{2,\text{in}}(t, s) &= \frac{\kappa}{\iota} \left\langle \tilde{\phi}_\mu(t) \tilde{\phi}_\nu(s) A_{\mu\nu}^{2,\text{in}}(t, s) \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}} \\
 A_{\mu\nu}^{1,\phi}(t, s) &= \left\langle \frac{\partial \sigma(h_\mu^1(t))}{\partial \beta_\nu^{1,\phi\top}(s)} \right\rangle_{\mathcal{Z}_N^{\text{global}}}, & B_{\mu\nu}^{1,\phi}(t, s) &= \left\langle \frac{\partial g_\mu^{1,\phi}(t)}{\partial \alpha_\nu^{\phi\top}(s)} \right\rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}, & \bar{B}_{\mu\nu}^1(s, t) &= \left\langle \frac{\partial g_\mu^{2,\text{in}}(t)}{\partial \alpha_\nu^{2,\text{in}\top}(s)} \right\rangle_{\mathcal{Z}_{N_e} | \mathcal{K}_{\text{global}}} \\
 g_\mu^{1,\phi}(t) &= \frac{\iota \eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\mu\nu}^3(s, t) \left[ \Phi_{\mu\nu}^{2,\text{in}}(s, t) \tilde{\phi}_\nu(t) - \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) \right] \tilde{\phi}_\mu(s) \\
 z_\mu^1(t) &= \beta_\mu^{1,\phi}(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left\{ \bar{B}_{\nu\mu}^1(s, t) + B_{\nu\mu}^{1,\phi}(s, t) + \Delta_\nu(s) \left[ \frac{\iota}{\kappa} \bar{G}_{\mu\nu}^{2,\text{in}}(s, t) + G_{\mu\nu}^{1,\phi}(s, t) \right] \right\} \sigma(h_\nu^1(s)) \\
 z_\mu^{2,\text{in}}(t) &= \beta_\mu^{2,\text{in}}(t) + \frac{\eta_0 \gamma_0}{P \kappa} \int_0^t ds \sum_\nu \Delta_\nu(s) G_{\mu\nu}^3(s, t) \tilde{\phi}_\mu(s) \tilde{\phi}_\nu(t) \sigma(h_\nu^{2,\text{in}}(s)), \quad j \in \{1, 2, \dots, N_e\} \\
 z_\mu^3(t) &= \beta_\mu^3(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu(s) h_\nu^3(s), \\
 h_\mu^1(t) &= \alpha_\mu^1(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \Delta_\nu \Phi_{\mu\nu}^0(z_\mu^1(s) \cdot \dot{\sigma}(h_\mu^1(s))) \\
 h_\mu^{2,\text{in}}(t) &= \alpha_\mu^{2,\text{in}}(t) + \frac{\iota \eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \Delta_\nu(s) \Phi_{\mu\nu}^1(t, s) (\dot{\sigma}(h_\mu^{2,\text{in}}(t)) \odot z_\nu^{2,\text{in}}(s)), \quad j \in \{1, 2, \dots, N_e\} \\
 h_\mu^3(t) &= \frac{\iota \eta_0 \gamma_0}{\kappa P} \int_0^t ds \sum_\nu \left[ \bar{A}_{\mu\nu}^{2,\text{in}}(t, s) + \Delta_\nu(s) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(s, t) \right] z_\nu^3(s) \\
 \psi_\mu(t) &= \alpha_\mu^\phi(t) + \frac{\eta_0 \gamma_0}{P} \int_0^t ds \sum_\nu \left[ A_{\mu\nu}^{1,\phi}(s, t) + \Delta_\nu(s) \Phi_{\mu\nu}^1(s, t) \right] g_\nu^{1,\phi}(t) \\
 \tilde{\phi}_\mu(t) &= \frac{e^{\hat{\psi}_\mu(t)}}{\mathcal{S}_\mu(t)}, \quad \hat{\psi}_\mu(t) = m_\mu(t) \psi_\mu(t), \quad m_\mu(t) = \mathbf{1}(\psi_\mu(t) - \tau_\mu(t))
 \end{aligned}$$

where  $\tau_\mu(t)$  is chosen such that  $\rho = \langle m_\mu(t) \rangle_{\mathcal{Z}_M | \mathcal{K}_{\text{global}}}$

$$\begin{aligned}
 \frac{df_\mu(t)}{dt} &= \frac{\eta_0}{P} \sum_\nu \Delta_\nu \left[ G_{\mu\nu}^1(t, t) \Phi_{\mu\nu}^0(t, t) + \left[ \frac{\iota}{\kappa} \bar{G}_{\mu\nu}^{2,\text{in}}(t, t) + G_{\mu\nu}^{1,\phi}(t, t) \right] \Phi_{\mu\nu}^1(t, t) \right. \\
 &\quad \left. + \frac{\iota}{\kappa} G_{\mu\nu}^3(t, t) \bar{\Phi}_{\mu\nu}^{2,\text{in}}(t, t) + \Phi_{\mu\nu}^3(t, t) \right]
 \end{aligned}$$

*Comparison to  $\mu P$ .* Limiting dynamics under  $\mu P$  admit a simpler mean-field structure than under MSSP. Under  $\mu P$ , the expert weights are initialized independently across experts. Consequently, there is no expert-hidden disorder shared by the expert population, and the limiting theory does not require a separate shared expert-hidden single-site process. The expert averages self-average directly over independently initialized experts. The corresponding DMFT is provided in App.O.

## Part IV

# Experiments

This part covers the experimental setup for our MLP MoE as well as Transformer MoE experiments in Appendix P, provides more details concerning the main figures in Appendix P.3, followed by detailed empirical evidence for our claims about learning rate transfer and scaling properties of  $\mu$ P and MSSP for both SGD and Adam in all 3 scaling regimes in Appendix Q. We close with further empirical scaling insights that practitioners should be aware of in Sections Q.5 to Q.8.

## P. Experimental Setup

Upon acceptance, we plan to make open source code to fully reproduce our experiments publicly available on GitHub.

### P.1. MLP MoE experiments on TinyImagenet

**MoE Architecture.** We train the same embedding-MoE-readout architecture with 2-layer fully connected experts that we analyse theoretically (without scale-dependent weight multipliers) using PyTorch [45]. It consists of a linear input layer with GeLU activation function, a linear router layer, 2-layer expert MLPs with GeLU activation function, followed by  $M^{-1}$  sigmoid aggregation, followed by a linear output layer. In this architecture without normalization layers and residual connections, instabilities in the scaling procedure become apparent at more moderate scale.

**Data.** We train single pass over all 50000 available images from 100 classes of TinyImageNet [32] for 1000 update steps with batch size 50. Our initial experiments used CIFAR-10 [31], but robustly benefitting from expert specialization in MoEs requires more diverse data.

We vary the optimizer (SGD and Adam), the scaling regime (fixed  $M$ , bottleneck and all-scaling), the routing mechanism (soft routing and top- $k$  routing) and the parameterization ( $\mu$ P and MSSP).

#### P.1.1. SCALING CONFIGURATIONS

Unless otherwise specified, we initialize the last layer to zero.

**MSSP.** Depending on the optimizer and scaling regime, Table B.1 summarizes our proposed parameterization MSSP as a function of width  $N$ , expert width  $N_e$  and number of experts  $M$ . Scaling of non-MoE trainable weights remains in  $\mu$ P.

**$\mu$ P.** Table B.1 specifies  $\mu$ P for SGD and Adam in each scaling regime. The provided HP scaling rules achieve stability as well as scale-independent effective and propagating updates after sufficiently many update steps, as verified in Appendix Q.4. However their initial vanishing signal propagation induces more scale dependence in the dynamics, delayed learning and/or no monotonic improvement with scale.

For Regime I, we compare the baseline router init std  $1/N$  (standard  $\mu$ P) versus zero router initialization for improved scale independence. Both variants barely differ.

#### P.1.2. MULTIPLIER TUNING

For MLPs, we use base width 1, which uncovers scaling properties of each parameterization at more moderate scale. Before running extended evaluations or learning rate sweeps, we tune multipliers at

small model scale  $N = 128$ . The network at size  $N = 128$  has dimensions ( $N = N_e = 128$ ,  $M = 8$ ) in Regime I, ( $N = 128$ ,  $M = N/16$ ,  $N_e = 16$ ) in Regime II and ( $N = N_e = 128$ ,  $M = N/16$ ) in Regime III, from which we scale up proportionally  $N \rightarrow \infty$ .

We observe that the optimal multipliers can be as far as  $10^6$  from all ones (since the gradient RMS norm of the expert output layer decays as fast as  $N^{-2}$  in Regimes II and III (see Sections Q.4.2 and Q.4.3)), and some update terms can be negligibly small without tuning, resulting in non-balanced learning (Appendix Q.6).

For each configuration of optimizer, parameterization, routing mechanism and scaling regime, we tune the following set of multipliers:

- global initialization variance,
- global learning rate,
- layerwise learning rate multiplier (input, router, expert input, expert output, output layer)

Due to interactions between these multipliers, we sweep all six multipliers jointly at size  $N = 128$ , amounting to at least  $5^6 = 15\,625$  runs for each combination of parameterization and scaling regime. Our tuning algorithm proceeds as follows. We first tune the learning rate at width  $N = 128$ , then run a broad 6D grid of the above init and layerwise learning rate multipliers, then run a full 6D grid with 5 grid points per dimension centered at the optimum from the previous stage at multiplicative resolution 4. This is possible since several small model training runs in parallel take less than 40 seconds on one A100 GPU.

Figure P.1 shows the example of the last step for the case of SGD MSSP in Regime II. Top-5 validation accuracy optima tend to be more localized than training accuracy optima. Hence we choose the optimum based on top-5 validation accuracy. The optimal initialization variance multiplier tends to be very clearly localized and essential in all combinations of optimizer, parameterization and scaling regime.

Appendix Q.7 shows that a cheaper multiplier tuning procedure, which starts with a random search over the grid, followed by extensive 2D sweeps over all combinations of multipliers does not suffice to reliably provide a near-optimal combination of multipliers.

### P.1.3. COORDINATE CHECKS AND LEARNING RATE SWEEPS

For each scaling configuration, we repeat the same experiment for 4 independent random seeds, affecting the random weight initialization as well as data shuffling. Uncertainty bands denote  $2\sigma$ -confidence bands. If not mentioned otherwise, MLP experiments use soft routing, and the same compute budget for multiplier tuning was invested for  $\mu$ P and MSSP.

## P.2. Transformer MoE experiments

We adapt the [nanoGPT](#) and [nanoMoE](#) codebases to train transformer models with a simple but modern architecture.

**General Architecture and Training Details.** By default, models have 8 blocks with mixture-of-expert layers (no dense layers). We scale the number of attention heads proportionally with width, while the head dimension remains fixed at  $d_{\text{head}} = 64$ . We use standard  $d_{\text{head}}^{-1/2}$  attention scaling, which can be viewed as a tunable multiplier since  $d_{\text{head}}$  remains fixed. We use pre-attention and qk-RMSNorm [54]. We train for 4768 steps using AdamW with a single, tuned

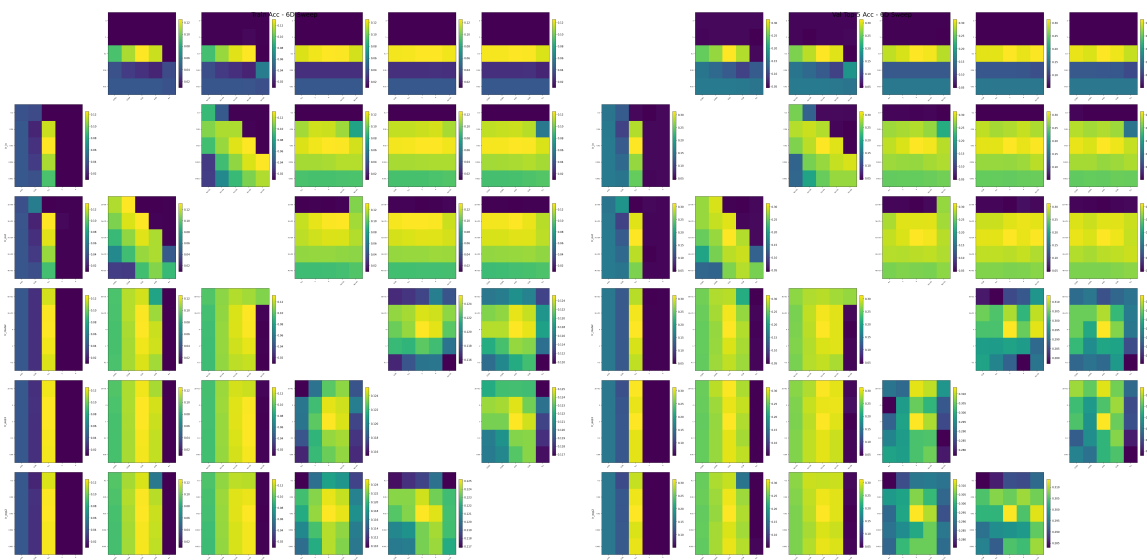


Figure P.1: **6D multiplier sweeps at small scale  $N = 128$  (MSSP, SGD, bottleneck).** 2D heatmaps showing training accuracy (left) and top-5 validation accuracy (right) of all HP pairs, while fixing all remaining HPs at the optimum. The optimum in each multiplier remains consistent across HP pairs. Train and validation optimum can slightly differ.

maximal learning rate,  $(\beta_1, \beta_2) = (0.9, 0.95)$ ,  $\epsilon = 10^{-8}$ , sequence length 1024, batch size 512, 700 steps of warmup followed by cosine learning rate decay to 10% of the maximal learning rate, weight decay 0.1, and gradient clipping.

**MoE layers.** Mixture-of-expert layers are scaled according to the respective scaling regime, using a sigmoid activation function and an auxiliary loss with weight 0.01 for load balancing [47]. We use token-choice routing and the number of active experts per token is half the total number of experts. We do not drop tokens.

**Architecture for Regime III.** In all all-scaling Regime III experiments, we use the following architecture. At base width  $N = 256$  each of the 8 MoE layers has  $n_{\text{exp}} = 8$  experts of hidden width  $N_e = N/2 = 128$ ; at  $N = 2048$  the same layer has  $n_{\text{exp}} = 64$  experts of hidden width  $N_e = 1024$ . At width  $N = 2048$ , this amounts to 2.5B total parameters.

**Architecture for Coordinate Checks in Regime II.** For the coordinate checks in bottleneck Regime II, we use 8 MoE layers with a fixed expert hidden width of  $N_e = 16$ . For keeping a total expansion ratio of 4, the expert count is  $M = 64$  at  $N = 256$  and grows to  $M = 512$  at  $N = 2048$ .

**Architecture for Learning Rate Sweeps for Regime II.** For the learning rate sweeps in bottleneck Regime II, we use a more computationally efficient bottleneck architecture due to compute constraints. We use 4 blocks alternating between dense and MoE layers (Dense, MoE, Dense, MoE). At width  $N = 256$ , each MoE block holds  $M = 32$  experts of fixed hidden width  $N_e = 32$ , scaling to  $M = 256$  experts at  $N = 2048$ . This amounts to 408M parameters at  $N = 2048$ .

**Data.** Models are trained on 2.5B tokens of [dolma3\\_mix-150B-1025](#) [44].

### P.2.1. SCALING CONFIGURATIONS

As in the MLP experiments, models are trained with the respective scaling configuration. We initialize the last layer to zero.

### P.2.2. MULTIPLIER TUNING

For every layer type, we tune initialization and learning rate multipliers. We also tune the global learning rate. Multipliers are tuned at width 256 by training on 1B tokens. We experiment both with random search and with round-robin algorithms that tune multipliers individually and across 2D grids. We evaluate multipliers across four different seeds, and find that both approaches obtain similar. An exhaustive grid search for multiplier tuning, as described for MLPs in Section P.1.2, was not computationally feasible for transformers because of the significantly longer training time.

### P.2.3. COORDINATE CHECKS AND LEARNING RATE SWEEPS

We use `torch-module-monitor` to monitor the training dynamics and perform coordinate checks. We use soft routing for the coordinate checks.

## P.3. Figure details

**Figure 1:**  $\mu$ P uses zero last-layer initialization. Under maximal stable last-layer initialization  $\sigma \simeq 1/N$ , exponents are even more unstable and learning even more delayed. More detailed evaluations are provided in Section Q.4.2. All width-scaling exponents  $\text{Expon}(v_N)$  are computed as OLS linear regression in log-log-space, hence fitting  $\alpha \in \mathbb{R}$  in a model  $v_N = C \cdot N^\alpha$  based on all available widths.

The right subplots show, at each time step the width-scaling exponents of the individual terms of Equation (MoE) and the overall updates to the post-aggregation activations  $\|\Delta h_t^l\|_{RMS}$ , which are the sum of effective and propagating updates, where the  $M^{-1}$  aggregation scaling is included in all terms.

Figure Q.26 shows width-scaling exponents of the propagating and effective updates of each linear weight matrix for the same setting as in Figure 1. The expert output weights are also evaluated pre-aggregation. Missing lines indicate that the respective term is exactly 0 (such as propagating updates in the input layer). Section Q.4.2 shows that for Adam in Regime II, while less severe, analogous width dependence in  $\mu$ P is resolved by MSSP, and monotonic improvement with scale is recovered.

**Figure 2:** We compare the top-5 training accuracy averaged over the last 50 steps of  $\mu$ P (dashed lines) and MSSP (solid lines) for the soft routing MLP MoE Adam scaling runs from each scaling regime in App. Q.4.1 for Regime I, Q.4.2 for Regime II and Q.4.3 for Regime III. Where both variants are available, we use the 0 last-layer initialization variant for the  $\mu$ P baseline for direct comparability with MSSP. These experiments use the optimal multipliers and learning rate from  $N = 128$  and transfer them to large model scales, as is common for  $\mu$ P in dense networks. Both  $\mu$ P and MSSP use the same compute budget for multiplier tuning at small size  $N = 128$ . In Regime I,  $\mu$ P and MSSP barely differ, so that the lines overlap heavily for both SGD and Adam.

## Q. Additional Experiments

### Q.1. Learning rate sweeps for Transformer MoEs

Figure Q.1 shows that learning rate transfer is achieved more cleanly in MSSP than in  $\mu$ P. The optimal learning rate in  $\mu$ P-Regime-II tends to grow (due to vanishing terms in the expert aggregation operations). Since at base width  $N = 256$  both parameterizations are equivalent, scaling differences only become apparent at large scale. Both parameterizations perform similarly well with MSSP having slightly more variance in Regime II.

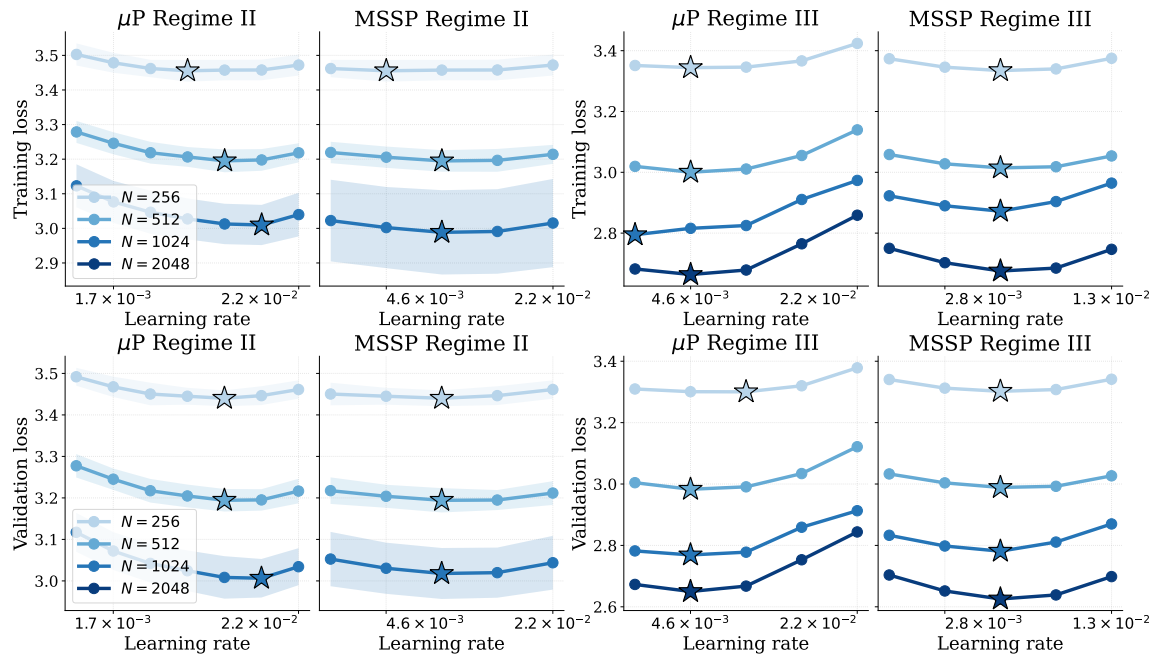


Figure Q.1: **Learning rate sweeps for  $\mu$ P and MSSP (Adam).** Training loss (top) and validation loss (bottom) for GPT MoEs trained with Adam in  $\mu$ P and MSSP for 2.5B tokens in Regime II ( $N_e \in \Theta(1)$  left) and Regime III ( $N, N_e, M, K \rightarrow \infty$  right), with  $\sigma$ -confidence bands across 4 seeds for bottleneck Regime II. Observe LR transfer and monotonic improvement with scale in MSSP.

### Q.2. Refined coordinate checks for Transformer MoEs

We use `torch-module-monitor` to measure refined coordinate checks as a sanity check for correct implementation of the MSSP scaling rules in each scaling regime. Propagating update exponents 0 verify correct initialization variance of the respective layer, given correctly scaled inputs. Effective update exponents 0 verify correct learning rate scaling of the respective layer, given correctly scaled inputs. One could isolate the effect of the layer’s learning rate scaling by normalizing the input activations  $x_{in}/\|x_{in}\|_{RMS}$ . We instead verify that  $\|x_{in}\|_{RMS} = \Theta(1)$  approximately holds in all layers (not shown).

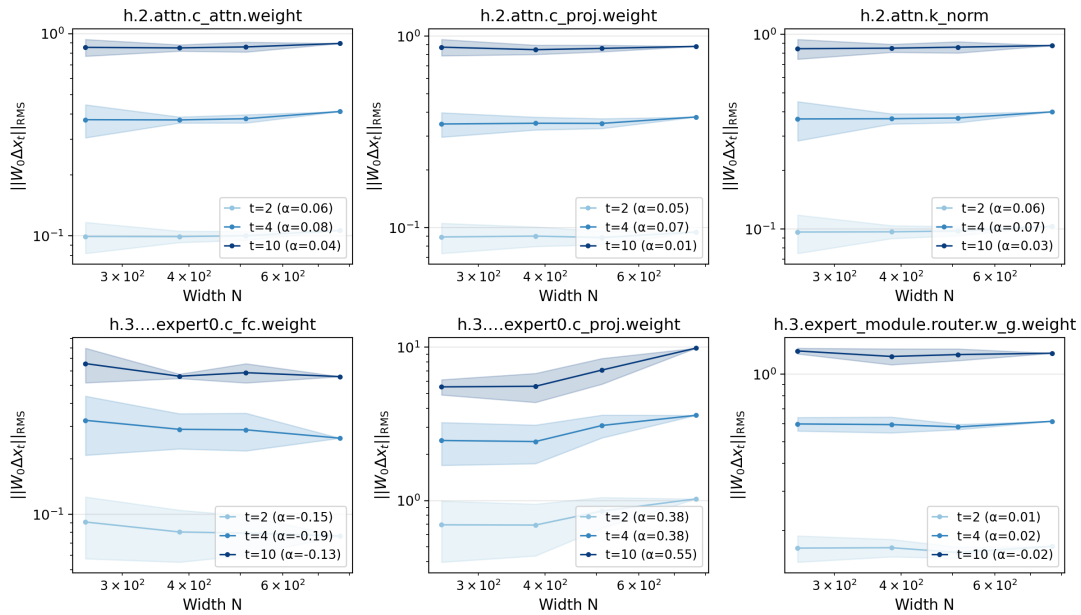


Figure Q.2: **Propagating updates in MSSP (Adam, Regime II).** Propagating updates of some example layers approximately follow the desired scaling exponents. First and second expert layer stem from a single expert. Measuring such small experts is particularly noisy, depending heavily on how many tokens are routed to them in the respective step. Recall that the desired propagating update exponent in the expert output layer `expert0.c_proj` is 0.5 in Regime II.

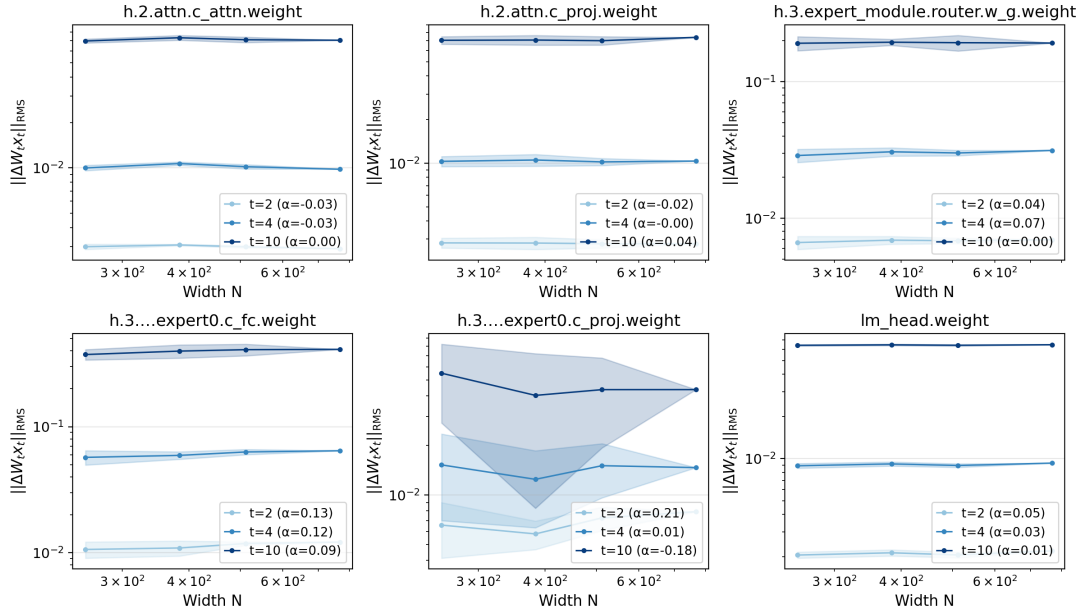


Figure Q.3: **Effective updates in MSSP (Adam, Regime II)**. Effective updates of some example layers are approximately scale-preserving at all times. First and second expert layer stem from a single expert. Measuring such small experts is particularly noisy, depending heavily on how many tokens are routed to them in the respective step.

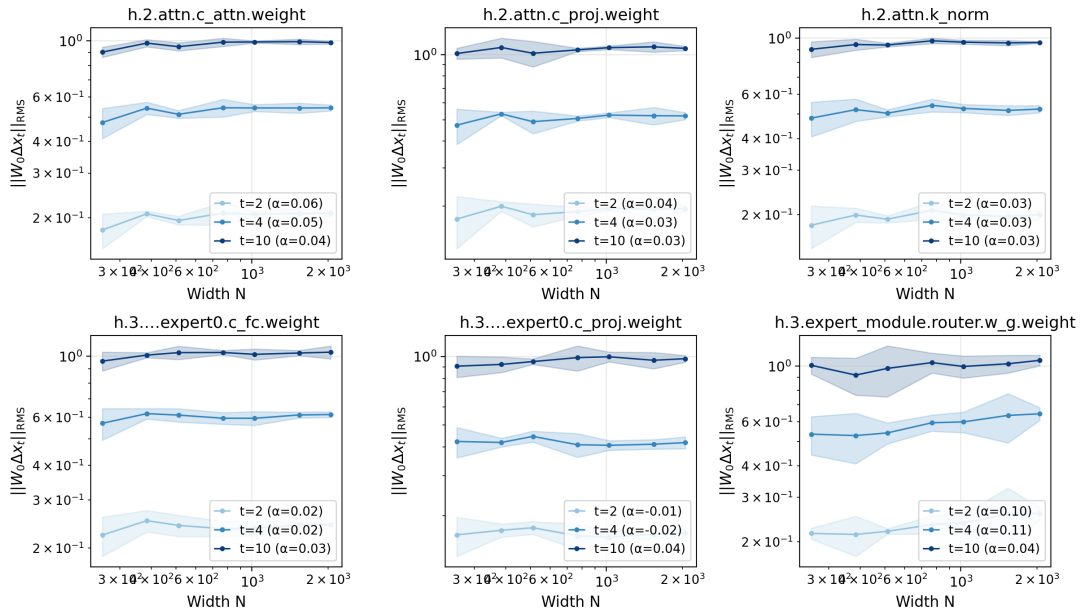


Figure Q.4: **Propagating updates in MSSP (Adam, Regime III)**. Propagating updates of some example layers are approximately scale-preserving at all times. First and second expert layer stem from a single expert.

## SCALING MOES: FROM $\mu$ P TO MSSP

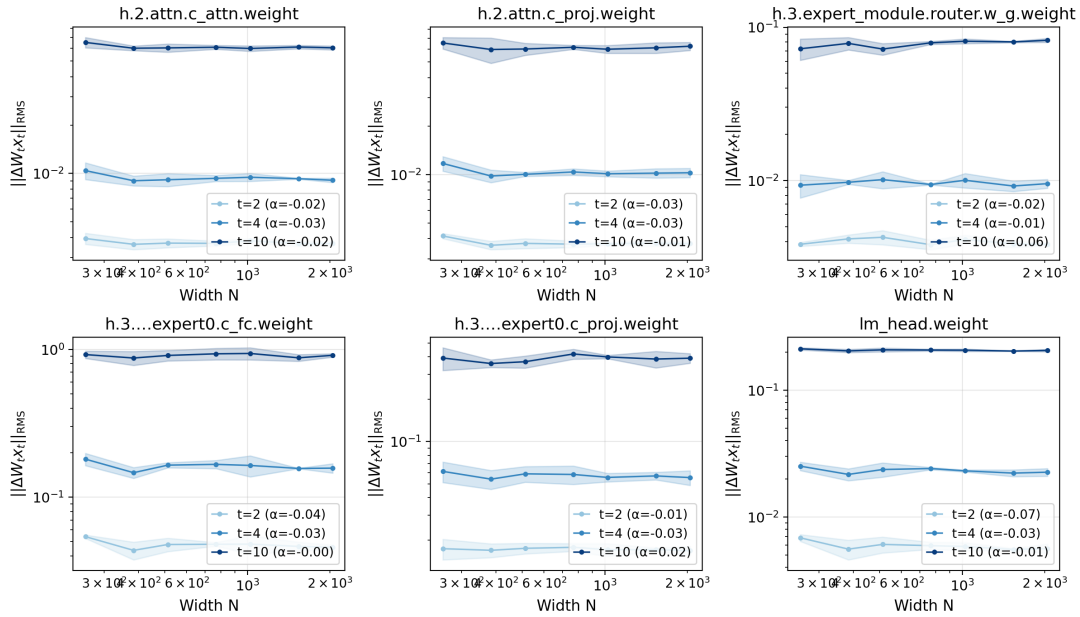


Figure Q.5: **Effective updates in MSSP (Adam, Regime III)**. Effective updates of some example layers are approximately scale-preserving at all times. First and second expert layer stem from a single expert.

### Q.3. Learning rate sweeps for MLP MoEs

Generally observe cleaner learning rate transfer across model sizes in MSSP than in  $\mu$ P. Also observe more robust monotonic improvement with scale in MSSP across optimizers and scaling regimes.

#### Q.3.1. REGIME I: FIXED NUMBER OF EXPERTS

While the optimal learning rate transfers and performance monotonically improves with scale, the performance saturates at large scales. The differing router initialization has negligible impact.

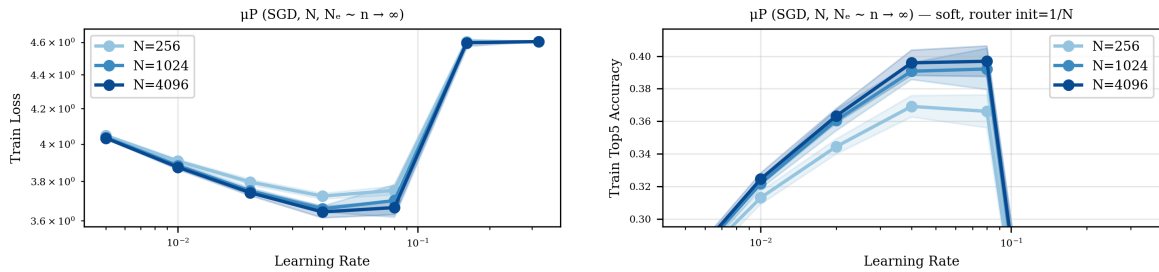


Figure Q.6: **LR sweep,  $\mu$ P with  $1/N$  router init (SGD, soft, Regime I)**.

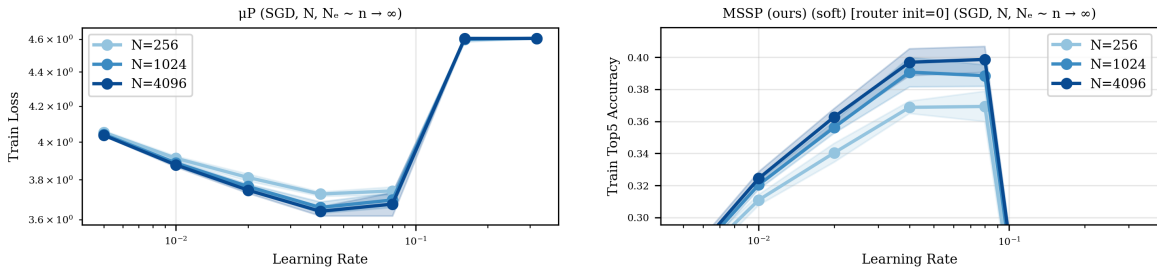


Figure Q.7: LR sweep, MSSP with zero router init (SGD, soft, Regime I).

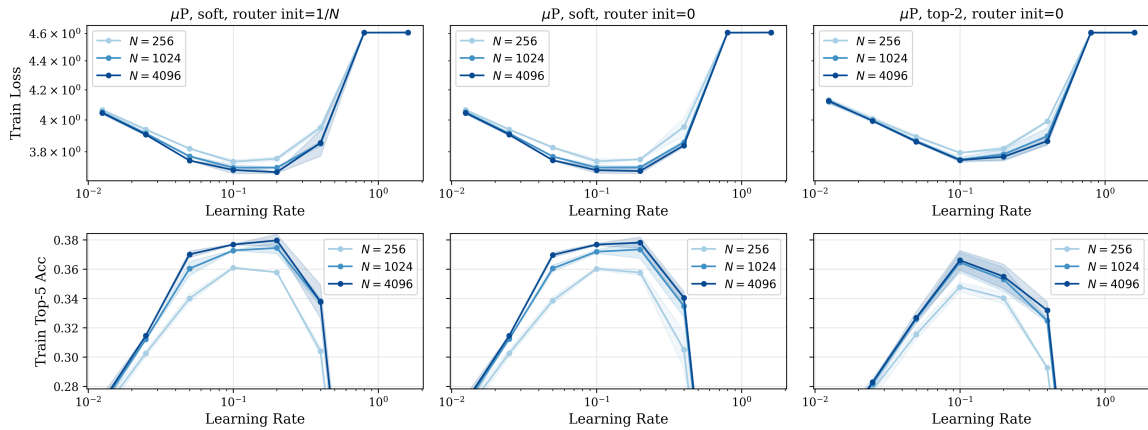


Figure Q.8: LR sweep comparison (Adam, Regime I).

### Q.3.2. REGIME II: FIXED EXPERT WIDTH

Observe greatly improved learning rate transfer of MSSP over  $\mu$ P in both SGD and Adam in Regime II. Larger performance gains from small to large model width result in improved absolute performance at scale across optimizers.

The optimal learning rate in  $\mu$ P tends to grow toward the instability threshold since subcomponents of the expert aggregation dynamics are vanishing at fixed learning rate. Since the correctly scaled subcomponents of the expert aggregation would diverge when increasing the learning rate, the stability threshold does not grow with width. These conflicting objectives of stability of some terms versus effective learning in others results in worse performance at scale. By balancing all subcomponents, MSSP recovers width independence of both the maximal stable and most effective learning rate scaling for all terms in the training dynamics, so that monotonic improvement with scale is preserved. Slight saturating growth of the optimal learning rate can occur in MSSP. Analogous slight shifts have often been observed in dense architectures in  $\mu$ P [3, 17, 58].

SCALING MOES: FROM  $\mu$ P TO MSSP

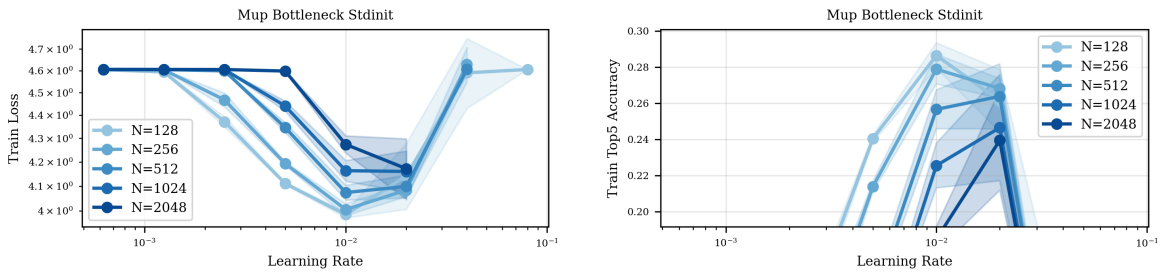


Figure Q.9: LR sweep,  $\mu$ P baseline (SGD, soft, Regime II). Ending lines denote divergence of the training run. The optimal learning rate drifts toward the maximal stable threshold, which does not increase with width. Observe delayed learning in the corresponding coordinate checks, which results in performance getting worse with model scale.

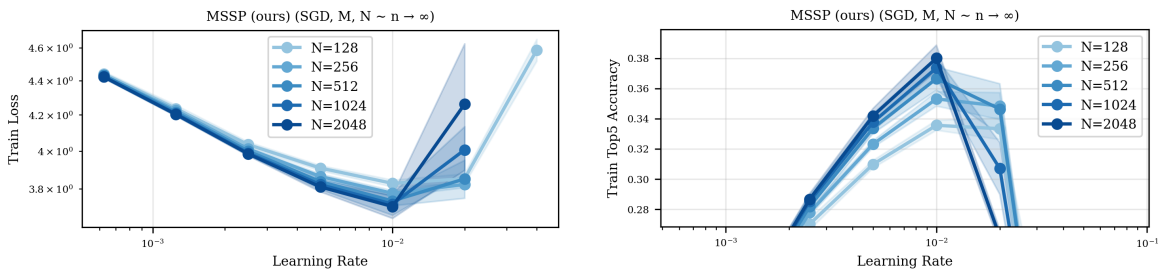


Figure Q.10: LR sweep, MSSP (SGD, soft, Regime II).

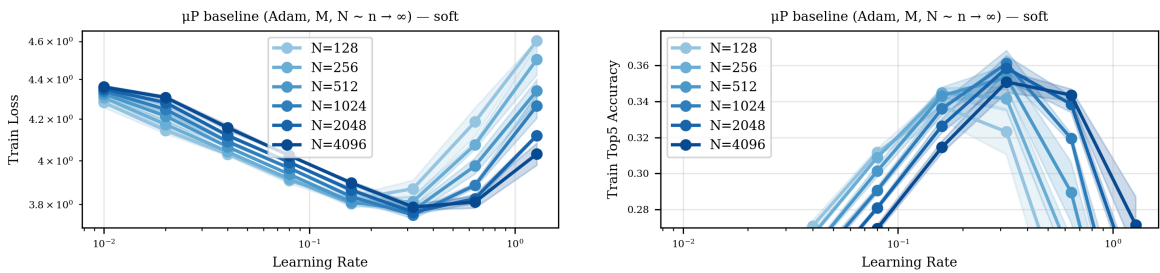


Figure Q.11: LR sweep,  $\mu$ P baseline (Adam, soft, Regime II).

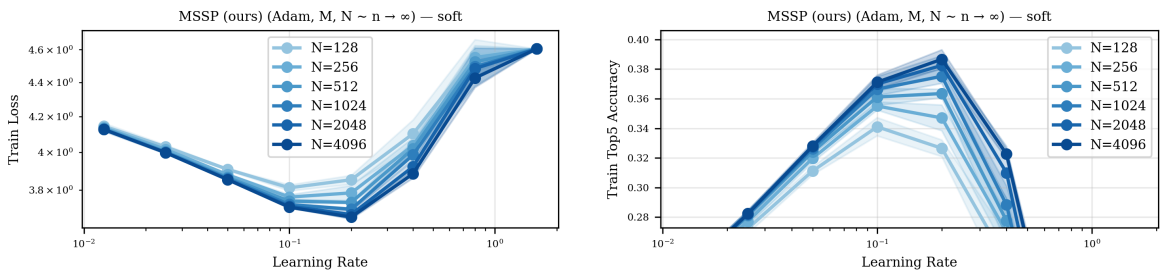


Figure Q.12: LR sweep, MSSP (Adam, soft, Regime II).

Q.3.3. REGIME III: JOINT PROPORTIONAL SCALING

Similar to Regime II, the optimal learning rate in  $\mu P$  saturates at the maximal stable learning rate. The optimal learning rate in MSSP approximately transfers from small to large scale, and achieves better performance than  $\mu P$  at large scale by consistently recovering monotonic improvement with scale.

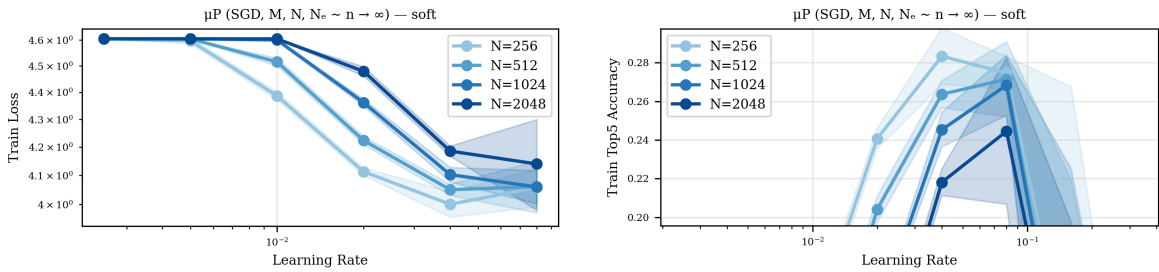


Figure Q.13: LR sweep,  $\mu P$  without shared experts (SGD, soft, Regime III).

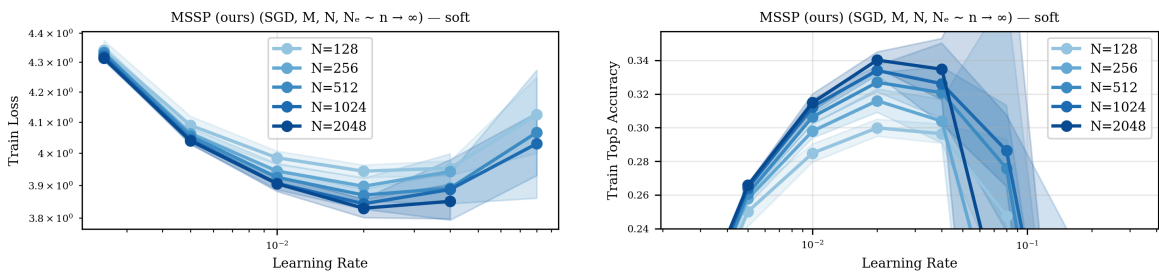


Figure Q.14: LR sweep, MSSP with shared experts (SGD, soft, Regime III).

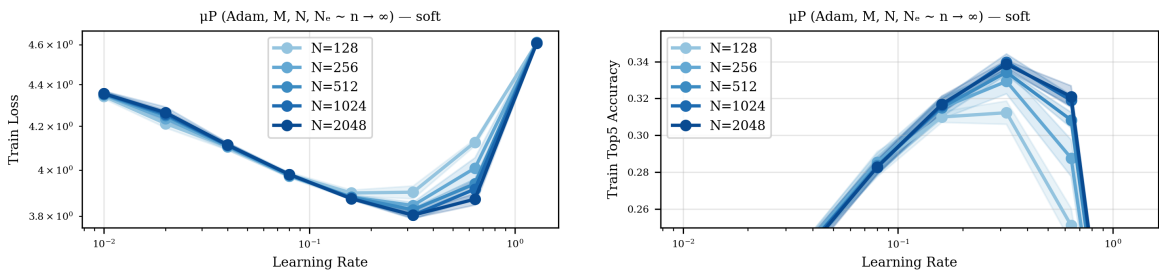


Figure Q.15: LR sweep,  $\mu P$  without shared experts (Adam, soft, Regime III).

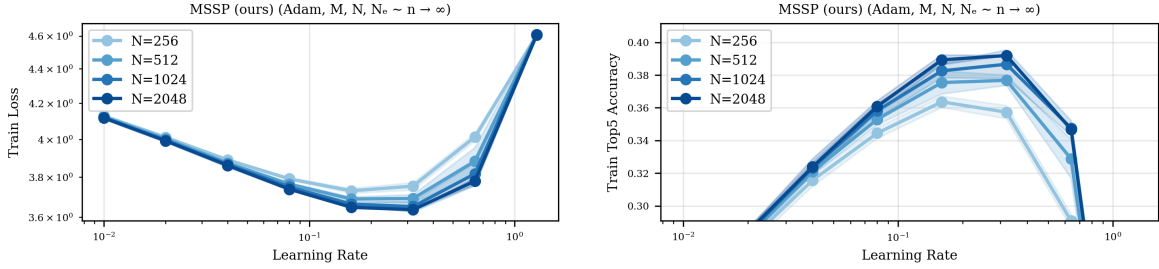


Figure Q.16: LR sweep, MSSP with shared experts (Adam, soft, Regime III).

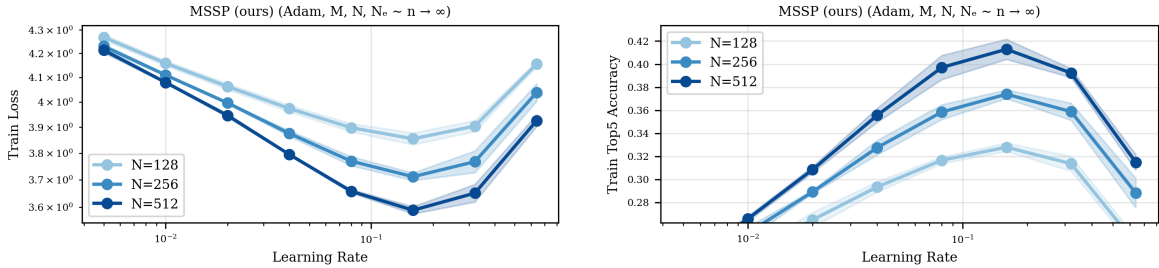


Figure Q.17: LR sweep, MSSP with shared experts (Adam, top- $k$ , Regime III).

#### Q.4. Fine-grained scaling evaluations in MLP MoEs

Here we show standardized evaluation figures for training the optimal configuration of many combination of optimizer, parameterization, routing mechanism and scaling regime. Across the board, we observe more desirable scaling properties in MSSP than in  $\mu$ P.

The first row shows training loss, training accuracy, the cumulative feature learning of the entire MoE  $\|\Delta h_t^L\|_{RMS}$ , as well as the routing logit norm  $\|\psi_t\|_{RMS}$  before the sigmoid. These allow to verify delayed learning versus monotonic improvement with scale, as well as the scale dependence in feature and router learning.

The second and third rows show refined coordinate check exponents as a function of width  $N$ . The general desired exponent for propagating and effective updates is zero, signaling scale independence of the respective component of the training dynamics (with the exception of expert output propagating update exponent 0.5 in MSSP in Regime II). The expected layerwise gradient scaling exponents vary between  $-2$  and  $0$ , depending on the scaling regime. Dashed horizontal lines indicate where expected exponents are non-zero. Decomposition exponents (third row first and second subplots) should generally be balanced, otherwise subcomponents of the dynamics strictly dominate others, which causes strong finite-width effects. Exponents should ideally also be independent of the step, otherwise non-trivial initial dynamical effects such as delayed learning are to be expected.

Missing lines in the propagating update and decomposition plots denote 0 values due to zero initialization of the relevant layer. For example under zero last-layer initialization, the initial gradient is zero except in the output layer. This helps to reduce width dependence from vanishing initial terms.

We show both raw effective updates (second row, center) and a variant with normalized incoming activation norm (second row, right) to distinguish width dependence from the combination of the current and previous layers versus isolated width dependence in the current layer.

In all regimes, the predicted exponents hold surprisingly well across training, suggesting that our theory is predictive of practical training dynamics far beyond the first few iterations. The router gradient exponents are most noisy, but the closest clean exponent of all layers still follows our prediction in MSSP remarkably well throughout training.

In  $\mu P$ , width dependence in individual subterms of the training dynamics cascades into the entire dynamics such that exponents become much more width and time dependent across layers, and often converge to intermediate values between the clean exponents  $\{-0.5, 0, 0.5\}$ .

Extensive layerwise multiplier tuning for each scaling config, optimizer and routing type is paramount for achieving stable training at all.

Q.4.1. REGIME I: FIXED NUMBER OF EXPERTS

Under maximal stable router initialization  $\sigma = 1/N$ , the propagating updates in the router are too small. Setting the router initialization to 0 removes this source of width dependence and results in cleaner scaling exponents. The effect of this intervention on the final performance is negligible.

The top- $k$  selection mechanism with  $k \asymp M$  does not change any expected scaling exponents in this paper. Indeed, Figures Q.20, Q.21, Q.24 and Q.25 verify that all exponents remain unaltered, albeit slightly more noisy.

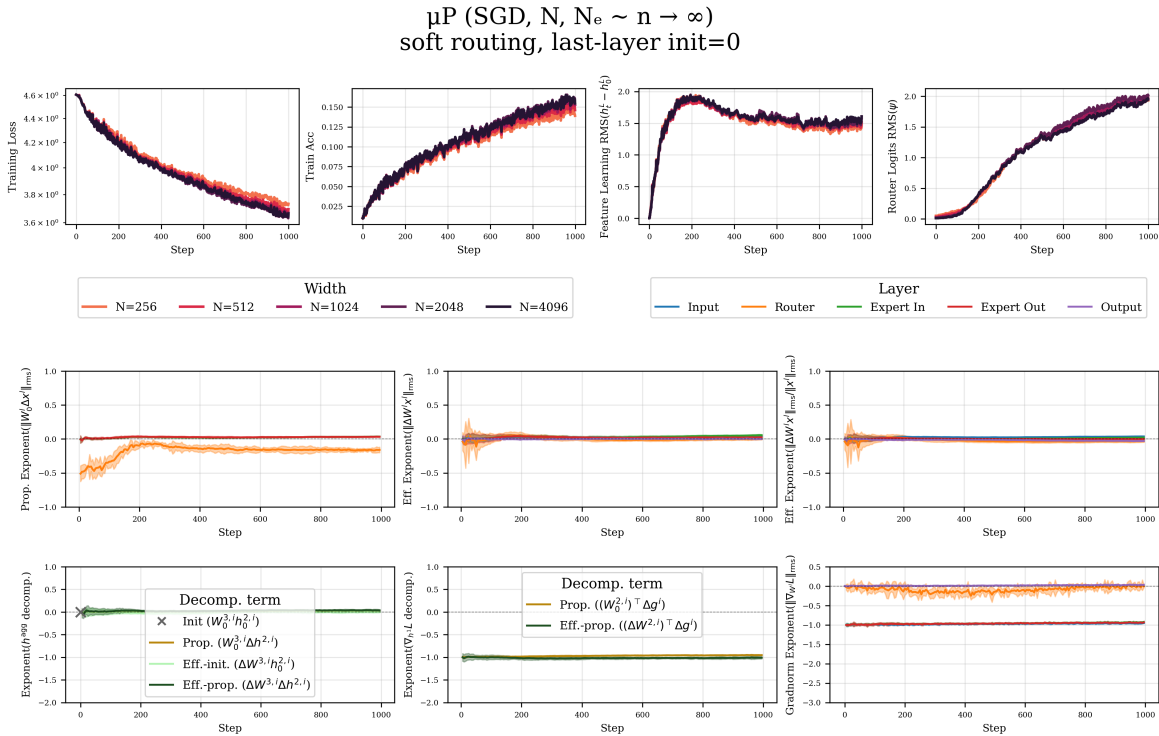


Figure Q.18:  $\mu P$  with  $1/N$  router init (SGD, Regime I).

# SCALING MOES: FROM $\mu$ P TO MSSP

MSSP (ours) (SGD,  $N, N_e \sim n \rightarrow \infty$ )  
 soft routing, router init=0, last-layer init=0

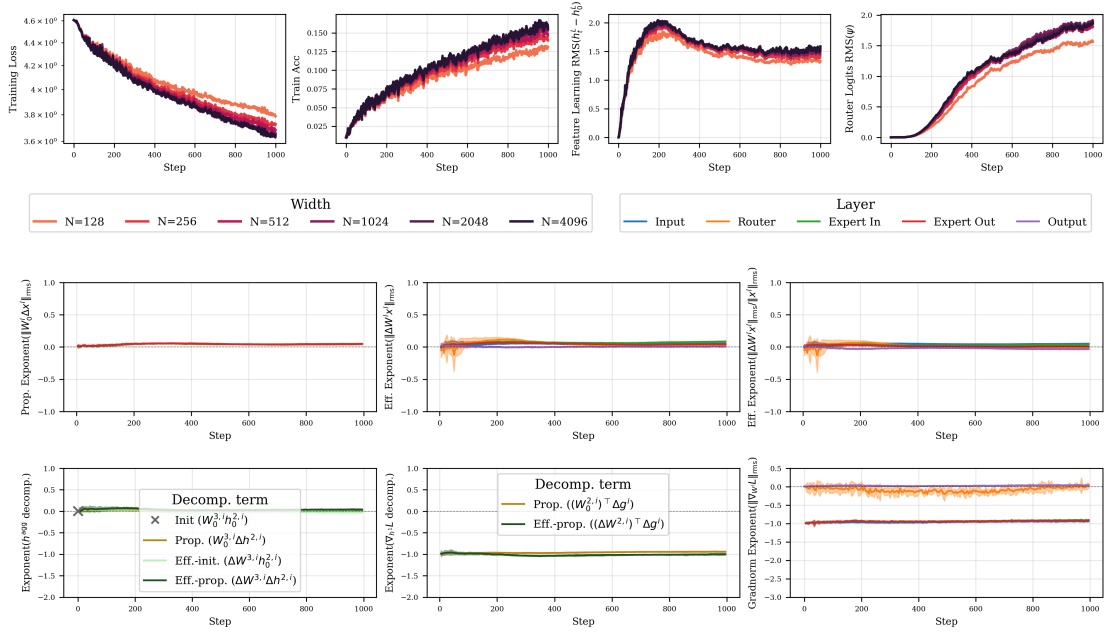


Figure Q.19:  $\mu$ P with zero router init (SGD, Regime I).

$\mu$ P (SGD,  $N, N_e \sim n \rightarrow \infty$ )  
 top-2 routing, last-layer init=0

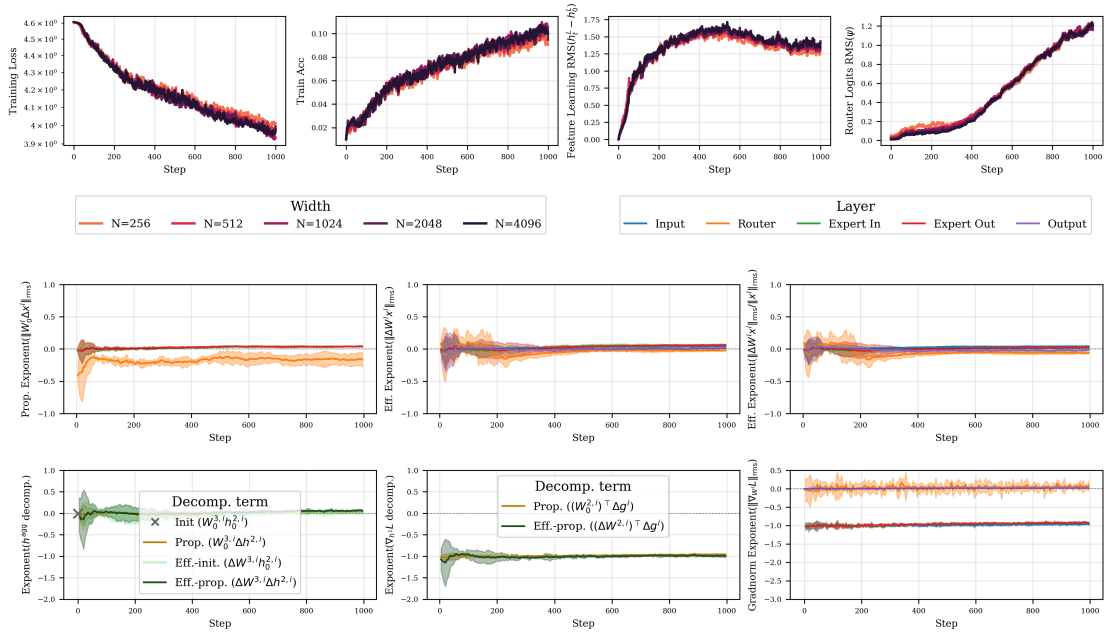


Figure Q.20:  $\mu$ P with  $1/N$  router init (SGD, Regime I, top- $k$ ).

# SCALING MOES: FROM $\mu$ P TO MSSP

MSSP (ours) (SGD,  $N, N_e \sim n \rightarrow \infty$ )  
 top-2 routing, router init=0, last-layer init=0

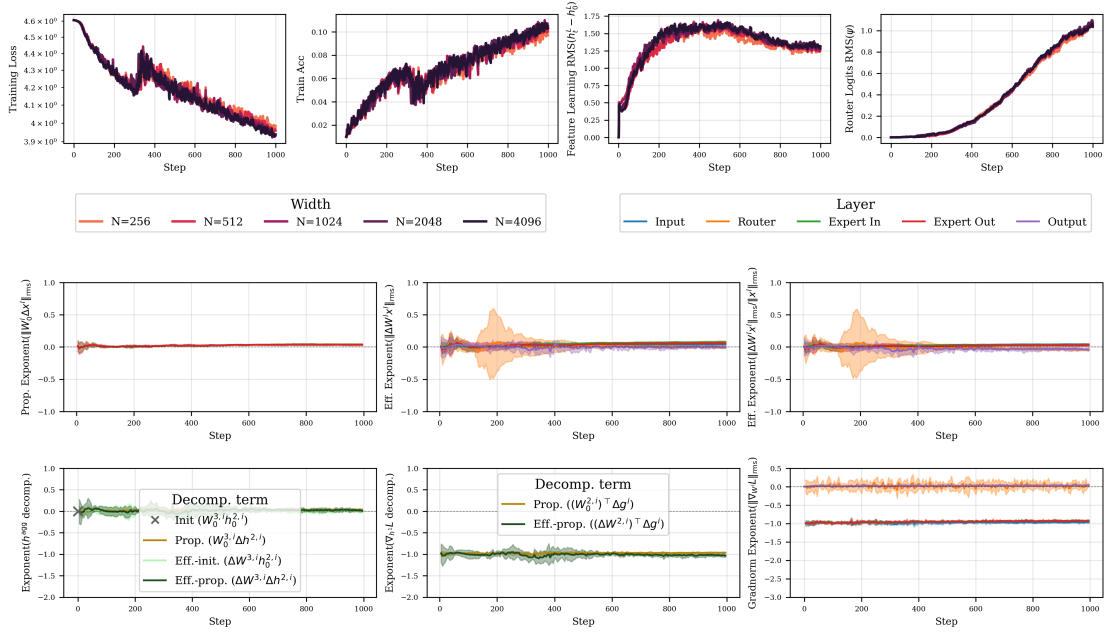


Figure Q.21: MSSP with zero router init (SGD, Regime I, top- $k$ ).

$\mu$ P (Adam,  $N, N_e \sim n \rightarrow \infty$ )  
 soft routing, last-layer init=0

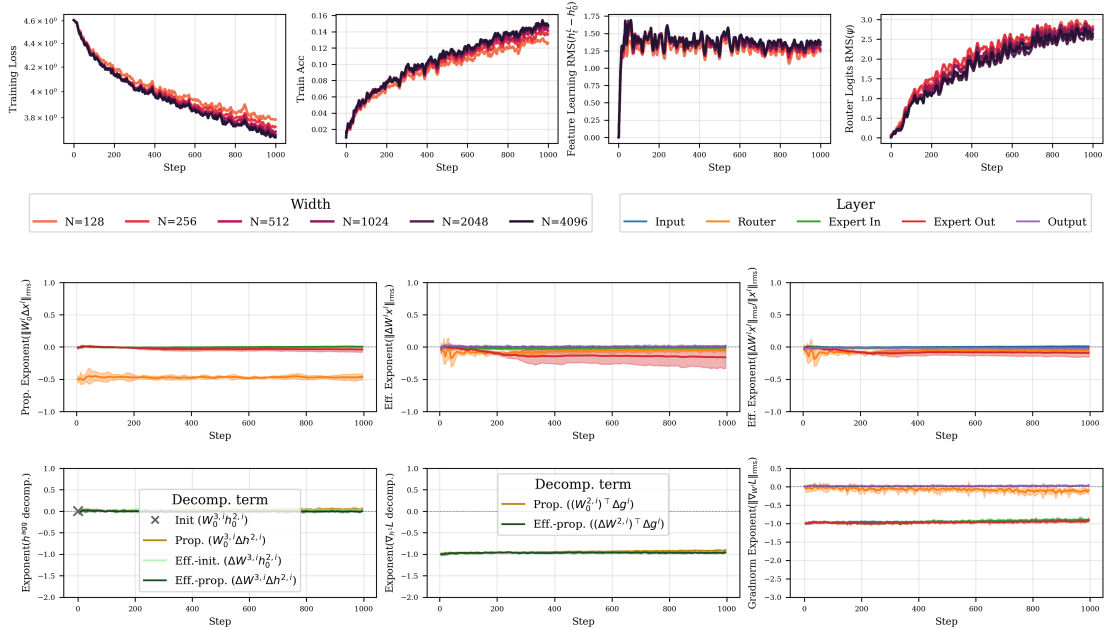


Figure Q.22:  $\mu$ P with  $1/N$  router init (Adam, Regime I).

# SCALING MOES: FROM $\mu$ P TO MSSP

MSSP (ours) (Adam,  $N, N_e \sim n \rightarrow \infty$ )  
 soft routing, router init=0, last-layer init=0

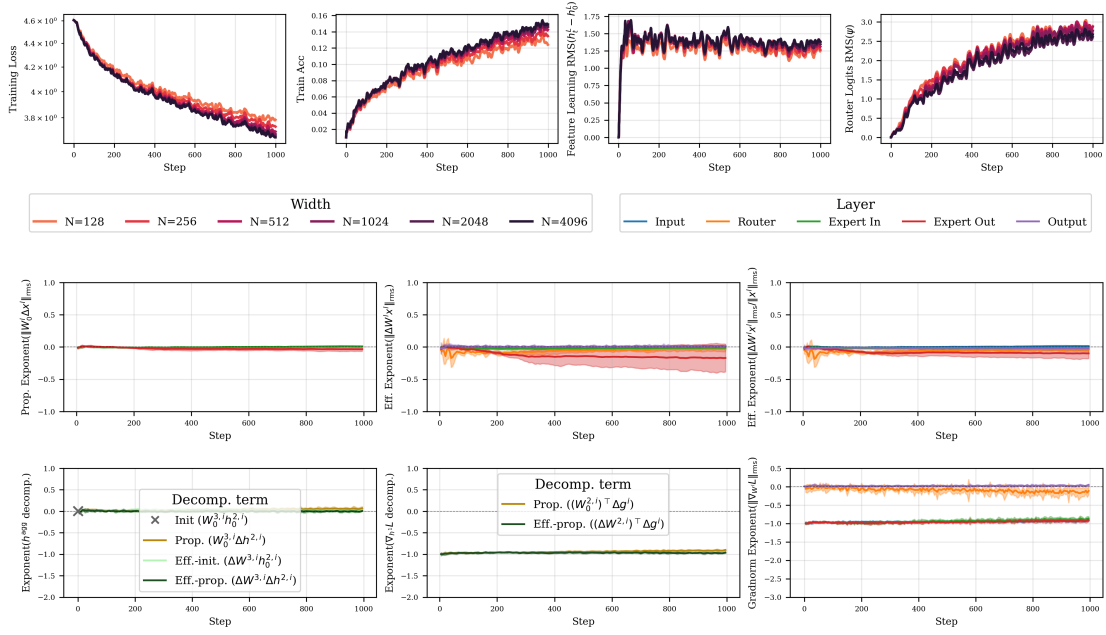


Figure Q.23: MSSP with zero router init (Adam, Regime I).

$\mu$ P (Adam,  $N, N_e \sim n \rightarrow \infty$ )  
 top-2 routing, last-layer init=0

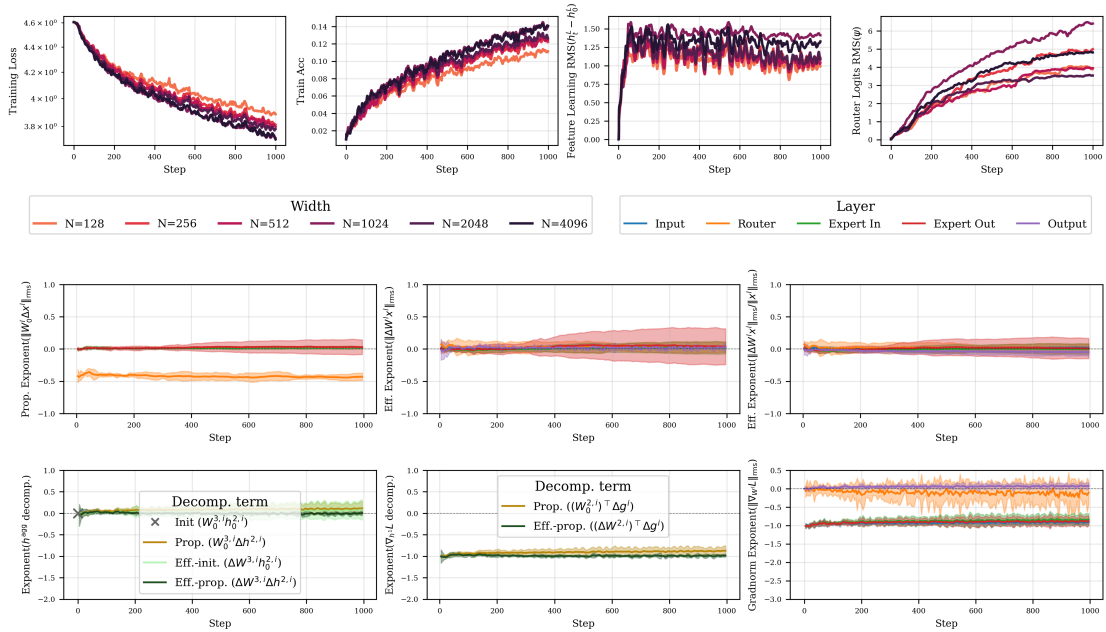


Figure Q.24:  $\mu$ P with  $1/N$  router init (Adam, Regime I, top- $k$ ).

MSSP (ours) (Adam,  $N, N_e \sim n \rightarrow \infty$ )  
 top-2 routing, router init=0, last-layer init=0

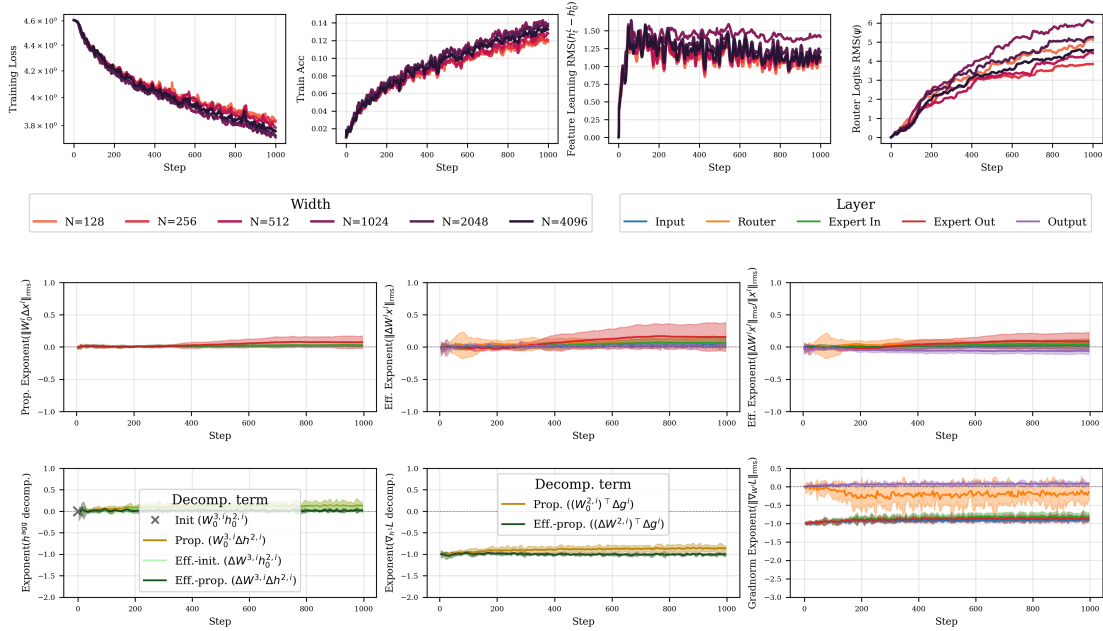


Figure Q.25: MSSP with zero router init (Adam, Regime I, top- $k$ ).

Q.4.2. REGIME II: FIXED EXPERT WIDTH

To ensure a stable baseline, we provide  $\mu$ P both with maximal stable and with 0 last-layer initialization. Both suffer from the same scaling degeneracies for SGD and Adam.

Starting with SGD, we observe strong delayed learning under maximal stable last-layer initialization  $\sigma = 1/N$ . This delay is visible not only in training loss, but also feature and router logit learning. Initial vanishing terms in the expert aggregation cascade even stronger into all layers than under last-layer 0 initialization presented in the main paper, resulting in final exponents that do not follow any clean scaling exponent from  $\{-0.5, 0, 0.5\}$ .  $\mu$ P with last-layer zero initialization starts out too small but partially self-corrects, which verifies that this parameterization indeed satisfies the  $\mu$ P desiderata, as any larger initialization or learning rate scaling would induce divergence in at least one layer after sufficiently many steps.

Adam in  $\mu$ P stabilizes much faster than SGD, and, while exponents are still not clean, they differ less from 0. Still, performance does not monotonically improve with scale and reduced feature learning at large scales is visible. Again, MSSP resolves these issues and shows monotonic improvement with scale as well as clean and balanced scaling exponents.

Propagating update scaling in  $\mu$ P remains vanishing throughout training for both SGD and Adam, as our theory predicts.

Note that, as predicted, expert output layer gradient entries decay extremely fast as  $\Theta(N^{-2})$  in both SGD and Adam. If ignored, this can cause numerical precision issues in practice at moderate model sizes. We recommend layerwise gradient scaling or equivalent Adam moment scaling to prevent numerical underflows.

SCALING MOES: FROM  $\mu$ P TO MSSP

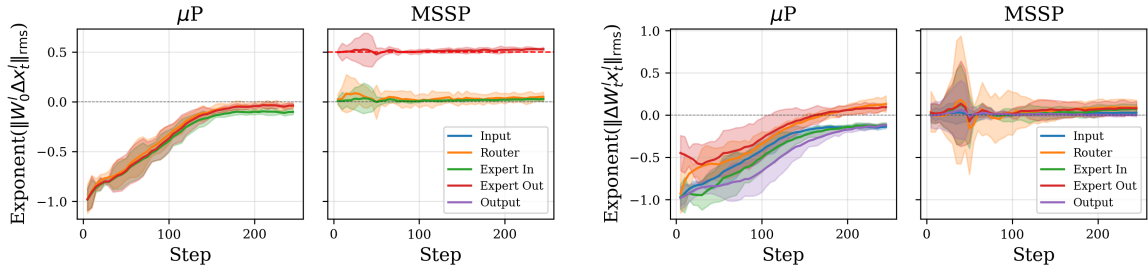


Figure Q.26: **Consistent exponents in MSSP, but not  $\mu$ P (SGD, Regime II)**. Corresponds to Figure 1. Initially vanishing post-aggregation terms in  $\mu$ P induce width-dependence to cascade into all layers. Over time, exponents partially self-correct in  $\mu$ P and almost recover width-independent propagating updates (left) as well as effective updates (right) in all layers. By allowing diverging pre-aggregation propagating updates (red dashed line), MSSP recovers approximately width-independent effective updates at all time steps.

$\mu$ P (MSSP, but smaller expert out init) (SGD,  $M, N \sim n \rightarrow \infty$ )  
soft routing, last-layer init=0

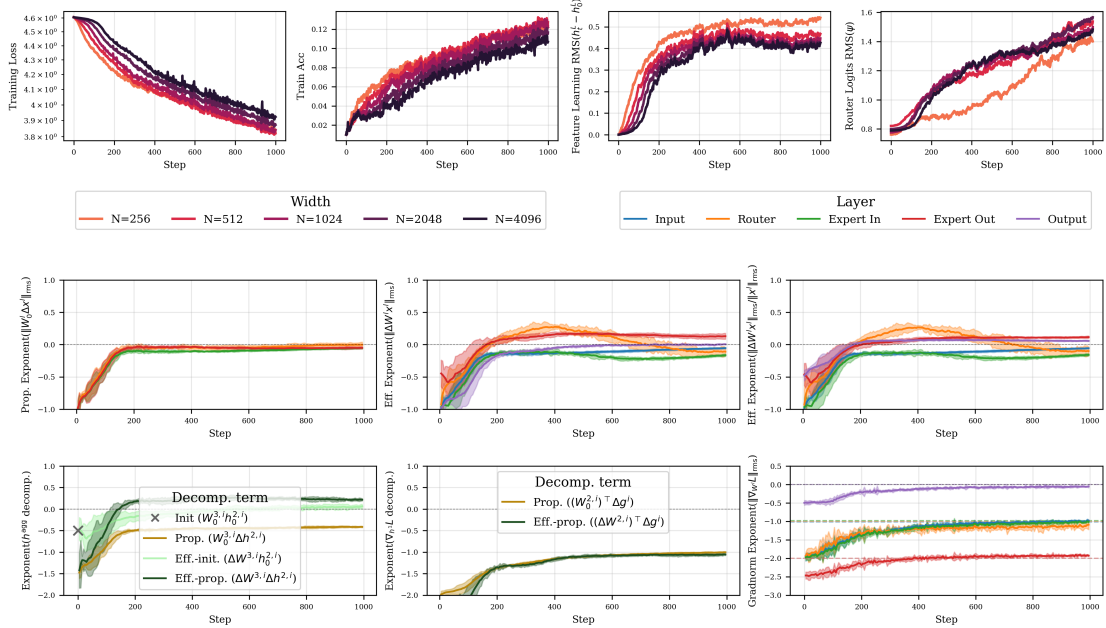


Figure Q.27:  $\mu$ P baseline (SGD, Regime II).

# SCALING MOES: FROM $\mu$ P TO MSSP

$\mu$ P (MSSP, but smaller expert out init) (SGD, M, N  $\sim n \rightarrow \infty$ )  
 soft routing, last-layer init=1/N

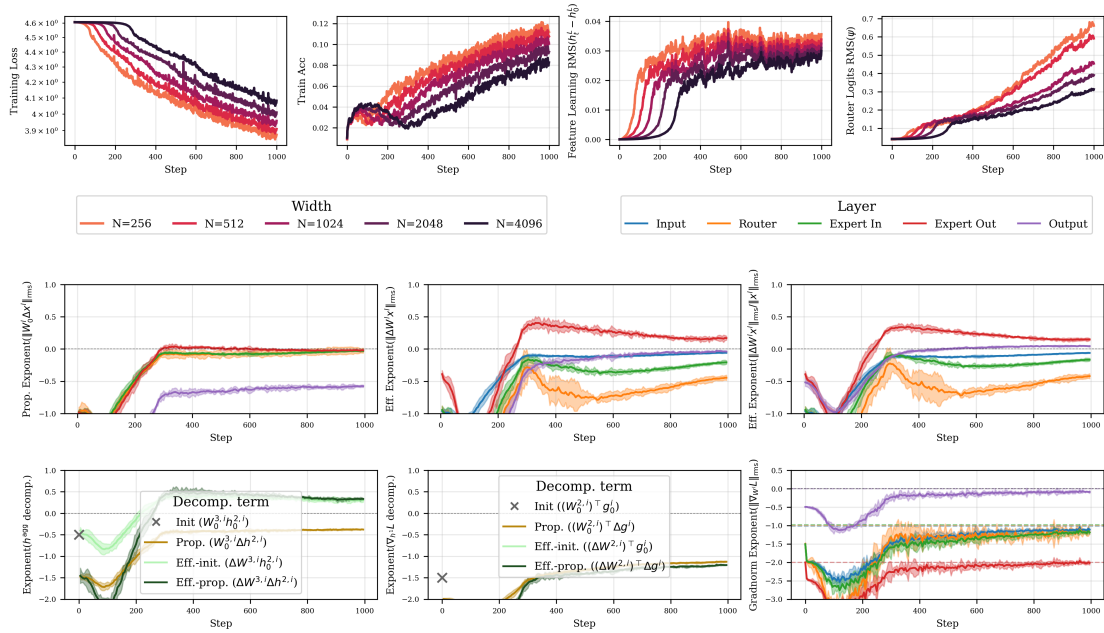


Figure Q.28:  $\mu$ P baseline with 1/N last-layer init (SGD, Regime II). Interesting initial growth in training accuracy even though training loss and features are barely moving. Learning is still delayed with increasing width.

SCALING MOES: FROM  $\mu P$  TO MSSP

MSSP (ours) (SGD, M, N  $\sim n \rightarrow \infty$ )  
 soft routing, last-layer init=0

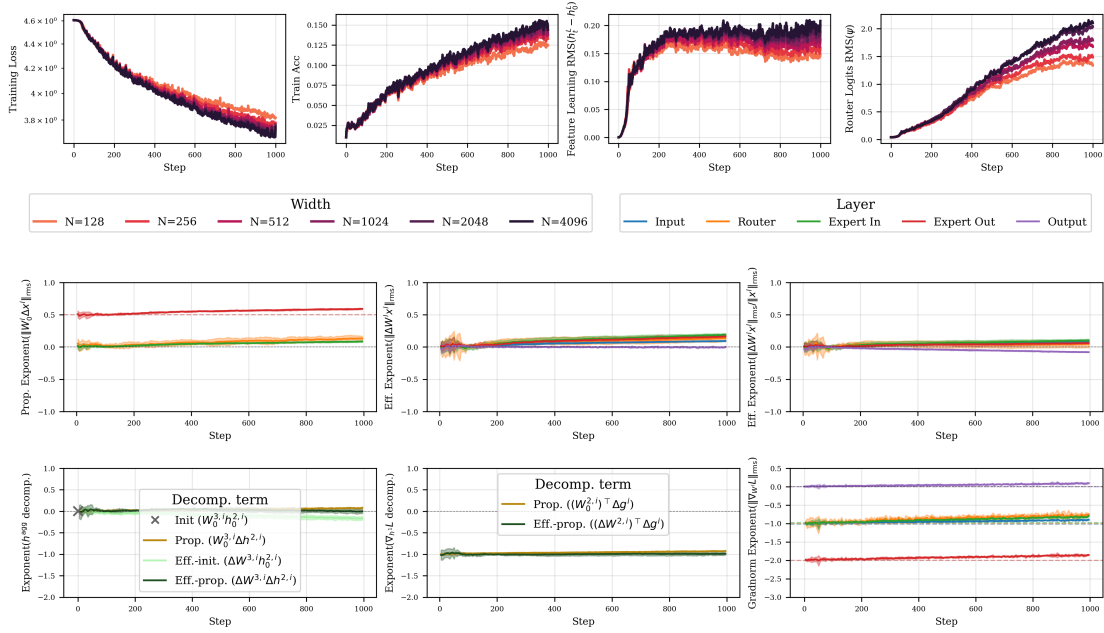


Figure Q.29: **MSSP (SGD, Regime II)**. Note that the propagating update exponent 0.5 of the expert output layer is desired here, since it is cancelled out by a clean CLT effect in the expert aggregation. The width independence of all other quantities throughout training verifies that the divergence of the individual expert outputs  $h^{2, \text{out}}$  through this propagating update term results in approximately width-independent training dynamics.

SCALING MOES: FROM  $\mu P$  TO MSSP

$\mu P$  (MSSP, but smaller expert out init) (Adam,  $M, N \sim n \rightarrow \infty$ )

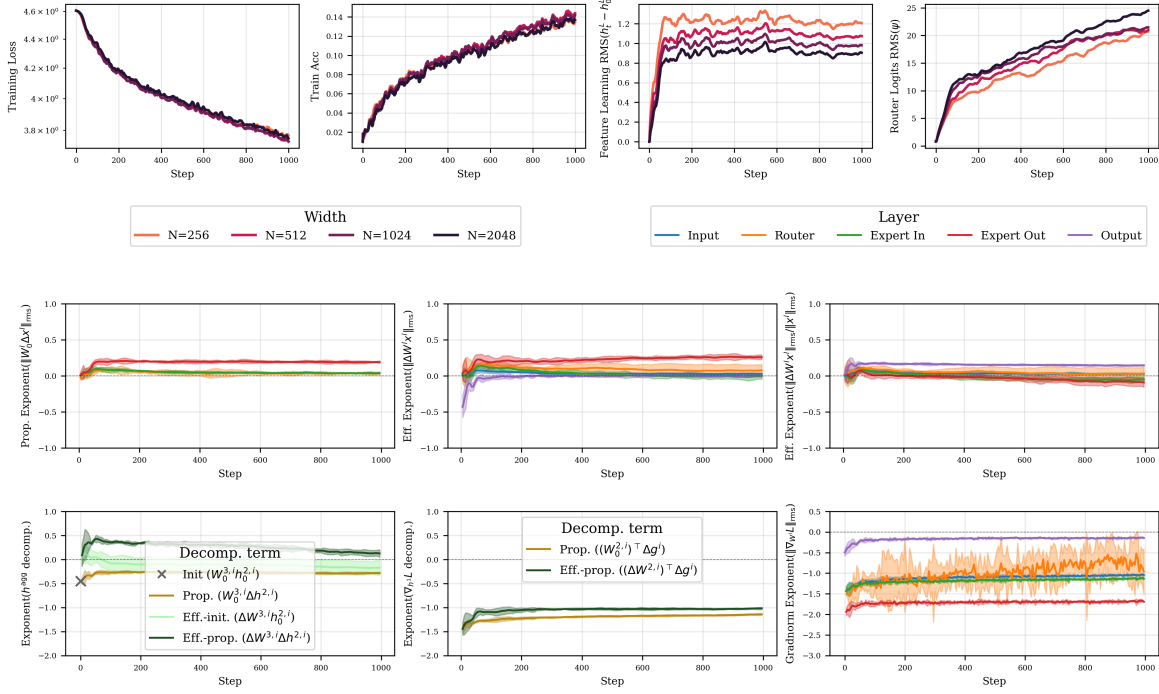


Figure Q.30:  $\mu P$  baseline (Adam, Regime II).

$\mu P$  (MSSP, but smaller expert out init) (Adam,  $M, N \sim n \rightarrow \infty$ )  
soft routing, last-layer init=1/N

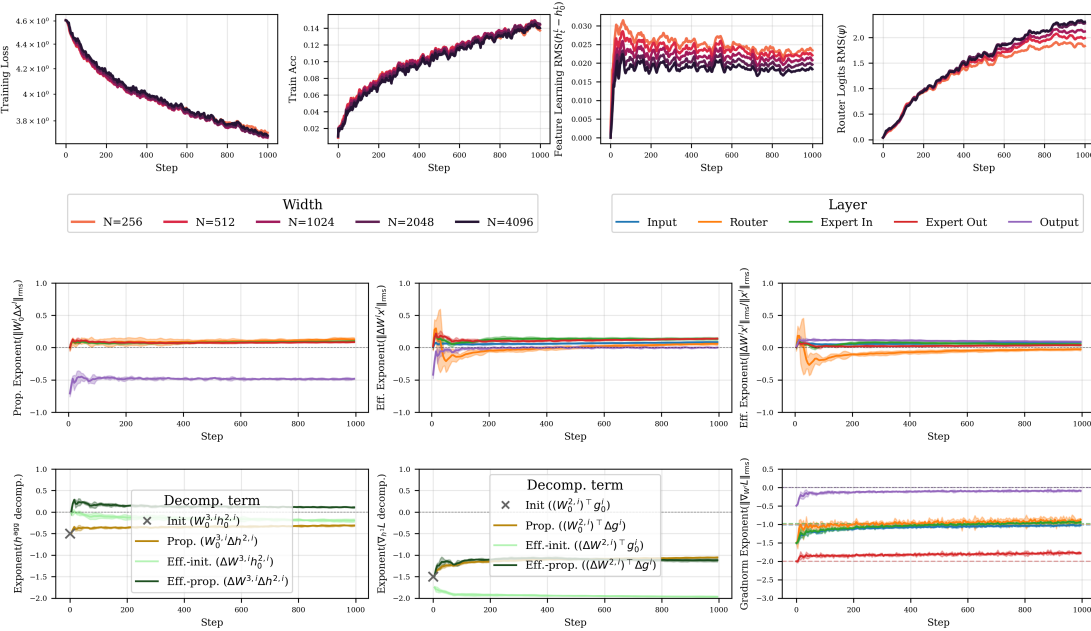


Figure Q.31:  $\mu P$  baseline with 1/N last-layer init (Adam, Regime II).

SCALING MOES: FROM  $\mu P$  TO MSSP

MSSP (ours) (Adam,  $M, N \sim n \rightarrow \infty$ )  
 soft routing, last-layer init=0

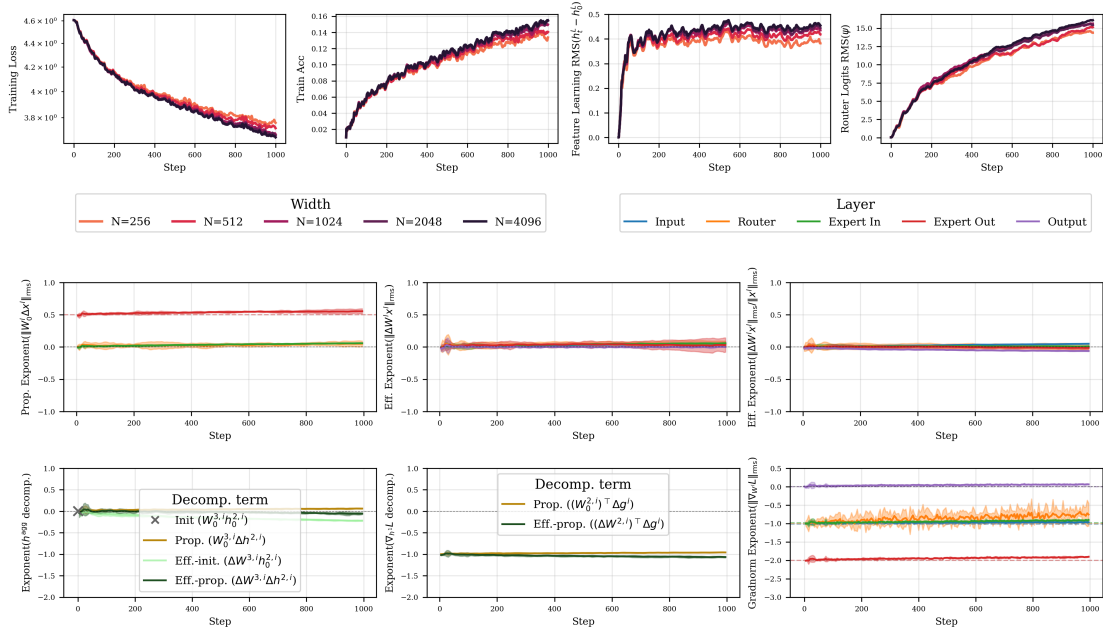


Figure Q.32: **MSSP (Adam, Regime II)**. Note that the propagating update exponent 0.5 of the expert output layer is desired here, since it is cancelled out by a clean CLT effect in the expert aggregation. The width independence of all other quantities throughout training verifies that the divergence of the individual expert outputs  $h^{2,\text{out}}$  through this propagating update term results in approximately width-independent training dynamics.

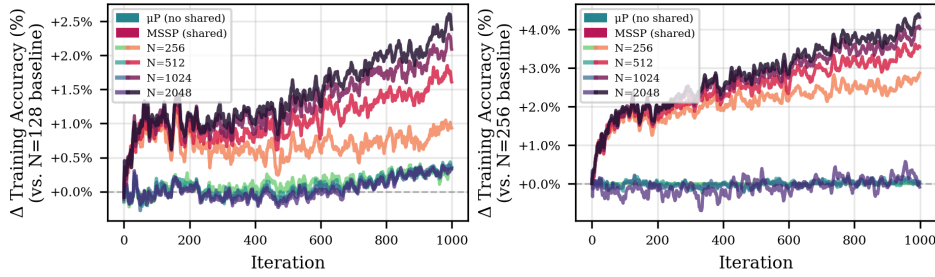


Figure Q.34: **Monotonic improvement with scale only in MSSP, not in  $\mu\text{P}$  (Adam, Regime III).** Training accuracy difference compared to  $\mu\text{P}$  with  $1/N$  (left) and  $0$  (right) last-layer initialization at  $N = 128$ , with separately tuned multipliers. MSSP outperforms both versions of  $\mu\text{P}$ .

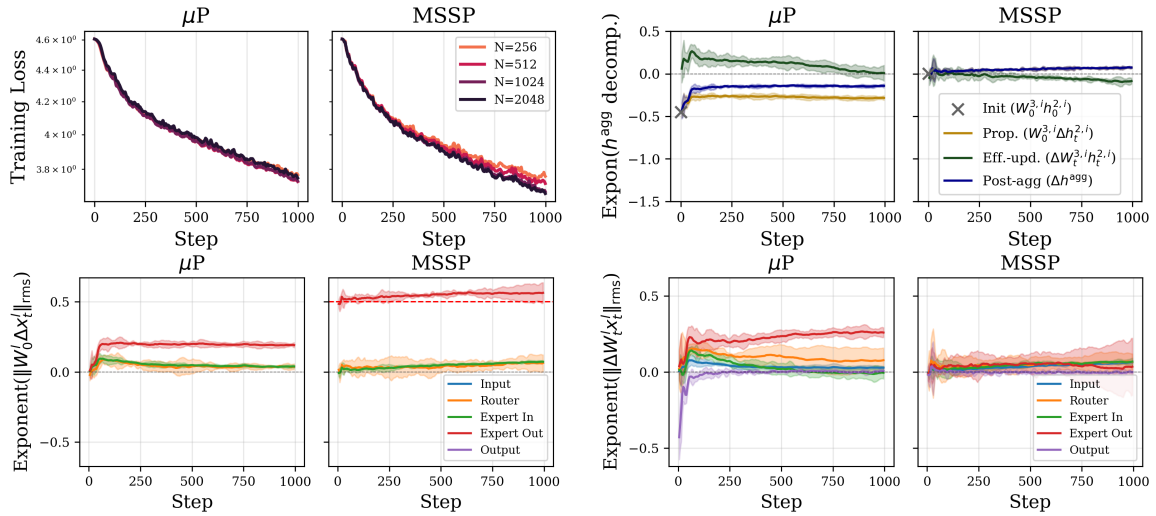


Figure Q.33: **Monotonic improvement with scale only in MSSP, not in  $\mu\text{P}$  (Adam, Regime II).** Same as Figure 1, but for Adam. While Adam’s stability is less severely impacted at moderate model scale, expert output exponents are not clean and performance does not improve with model scale in  $\mu\text{P}$ .

Q.4.3. REGIME III: JOINT PROPORTIONAL SCALING

For SGD, observe delayed learning in  $\mu\text{P}$  similar to Regime II. In both  $\mu\text{P}$  and MSSP, the router gradient norm is very noisy, but this does not affect the width-independent signal propagation in MSSP.

Adam in  $\mu\text{P}$  is surprisingly stable in the all-scaling Regime III. Still, Adam’s performance only improves robustly and significantly with scale in MSSP, both under soft and top-k routing.

top-k selection does not affect qualitative scaling properties as theoretically predicted.

As in Regime II, note that expert input and output layer gradients are decaying extremely fast with exponent  $-2$  as predicted in both SGD and Adam, requiring layerwise gradient or Adam moment scaling to prevent numerical underflows at moderate scale.

## SCALING MOES: FROM $\mu P$ TO MSSP

$\mu P$  (SGD,  $M, N, N_e \sim n \rightarrow \infty$ )  
no shared experts, soft routing, last-layer init=0

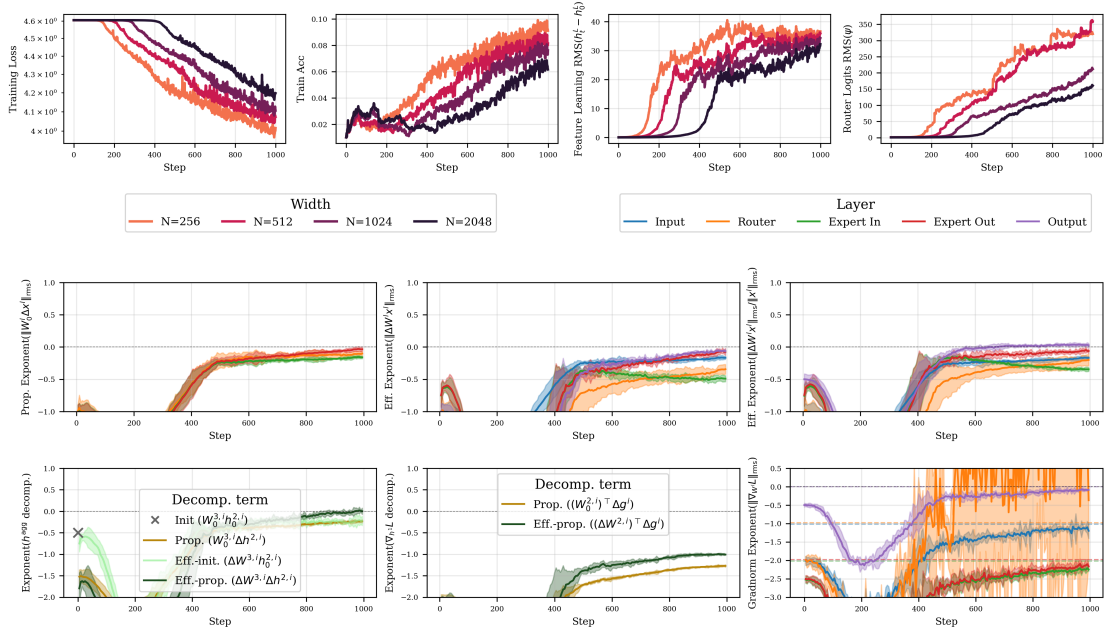


Figure Q.35:  $\mu P$  without shared experts (SGD, Regime III).

MSSP (ours) (SGD,  $M, N, N_e \sim n \rightarrow \infty$ )  
shared experts, soft routing, last-layer init=0

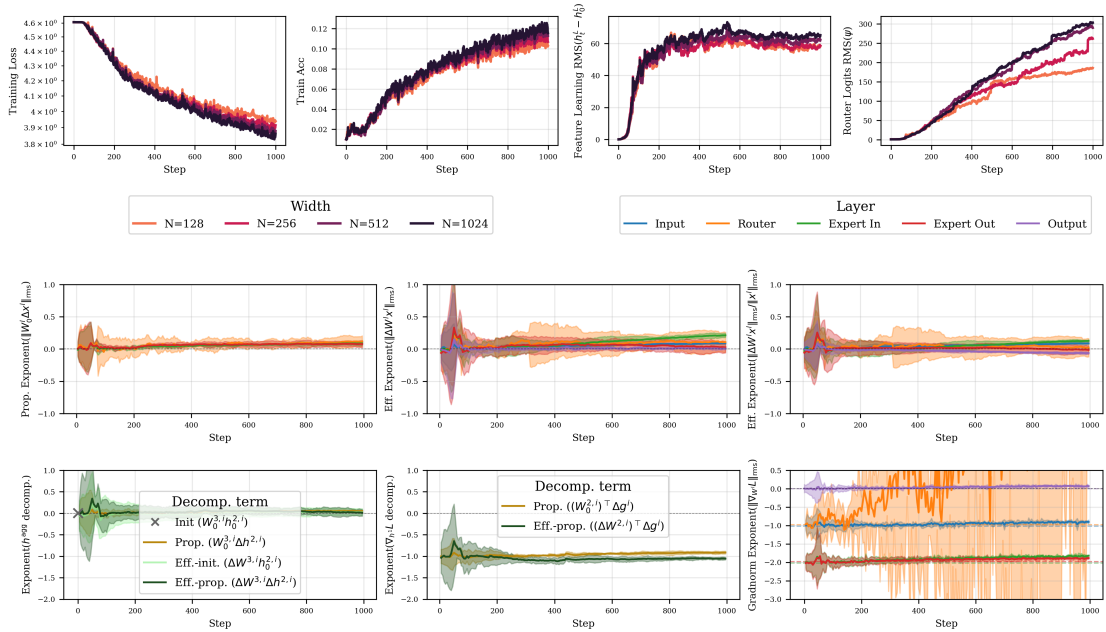


Figure Q.36: MSSP with shared experts (SGD, Regime III).

# SCALING MOES: FROM $\mu P$ TO MSSP

$\mu P$  (SGD,  $M, N, N_e \sim n \rightarrow \infty$ )  
no shared experts, top-2 routing, last-layer init=1/N

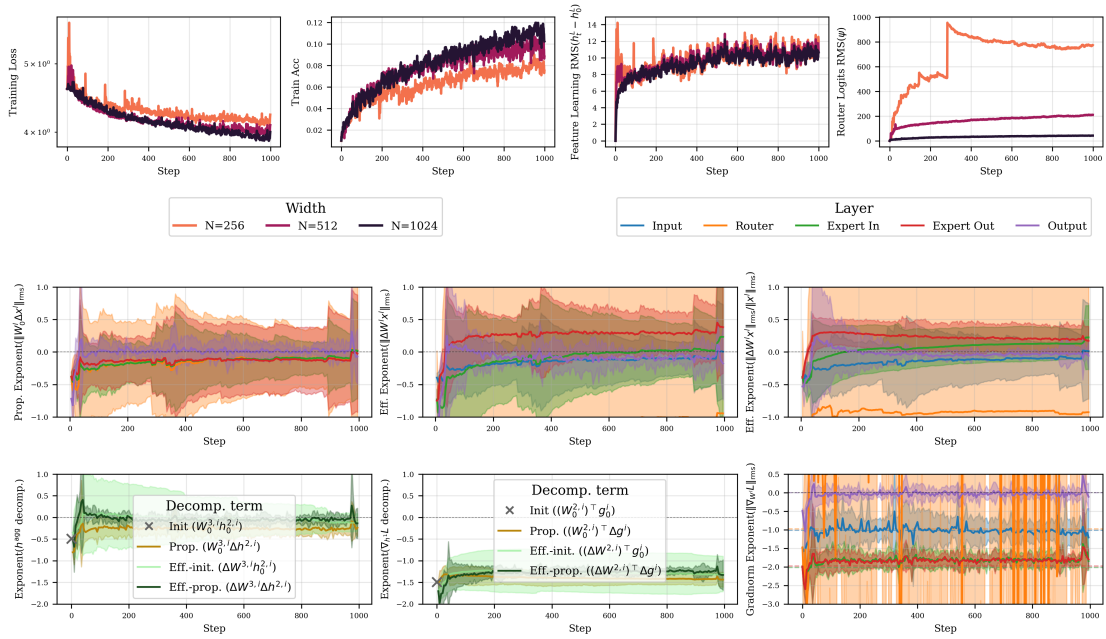


Figure Q.37:  $\mu P$  without shared experts with 1/N last-layer init (SGD, Regime III, top- $k$ ).

MSSP (ours) (SGD,  $M, N, N_e \sim n \rightarrow \infty$ )  
shared experts, top-2 routing, last-layer init=0

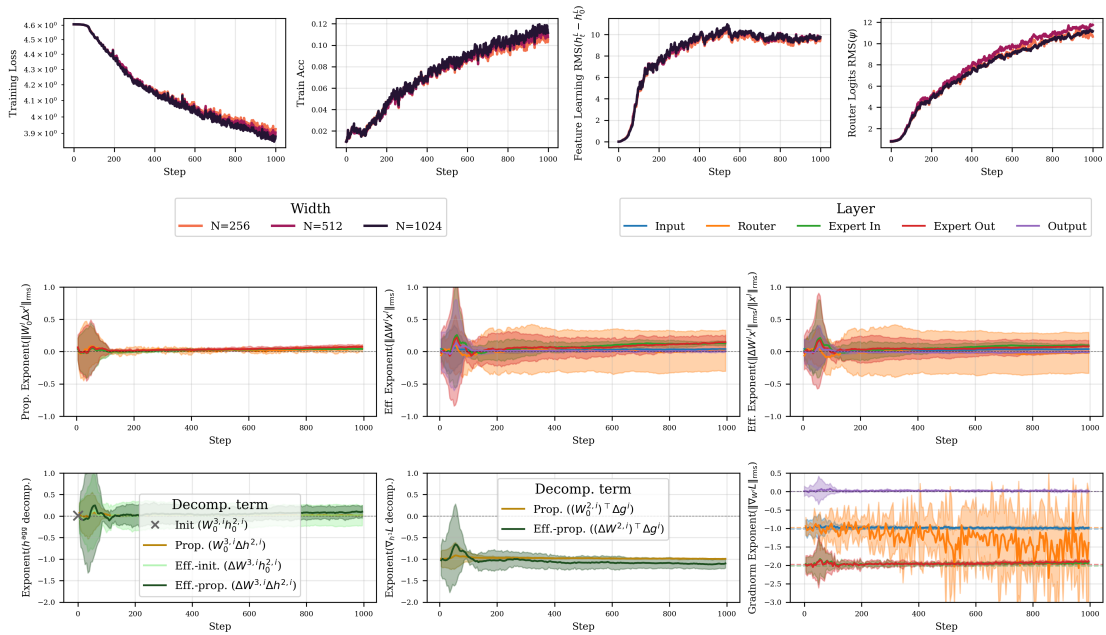


Figure Q.38: MSSP with shared experts (SGD, Regime III, top- $k$ ).

## SCALING MOES: FROM $\mu$ P TO MSSP

$\mu$ P (Adam,  $M, N, N_e \sim n \rightarrow \infty$ )  
no shared experts, soft routing, last-layer init=0

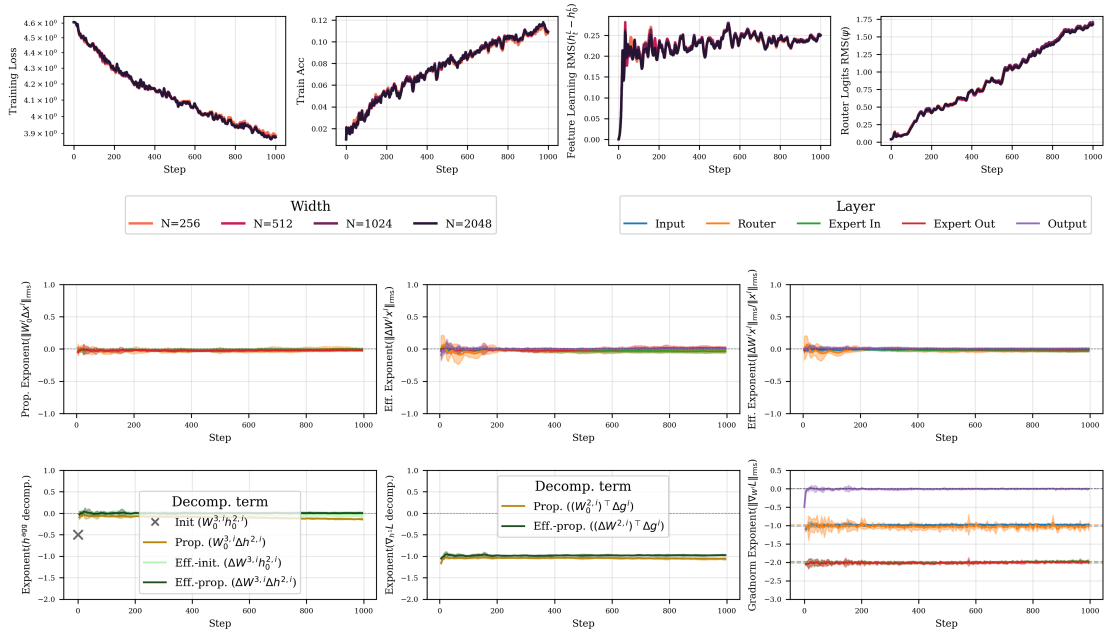


Figure Q.39:  $\mu$ P without shared experts (Adam, Regime III).

$\mu$ P (Adam,  $M, N, N_e \sim n \rightarrow \infty$ )  
no shared experts, soft routing, last-layer init=1/N

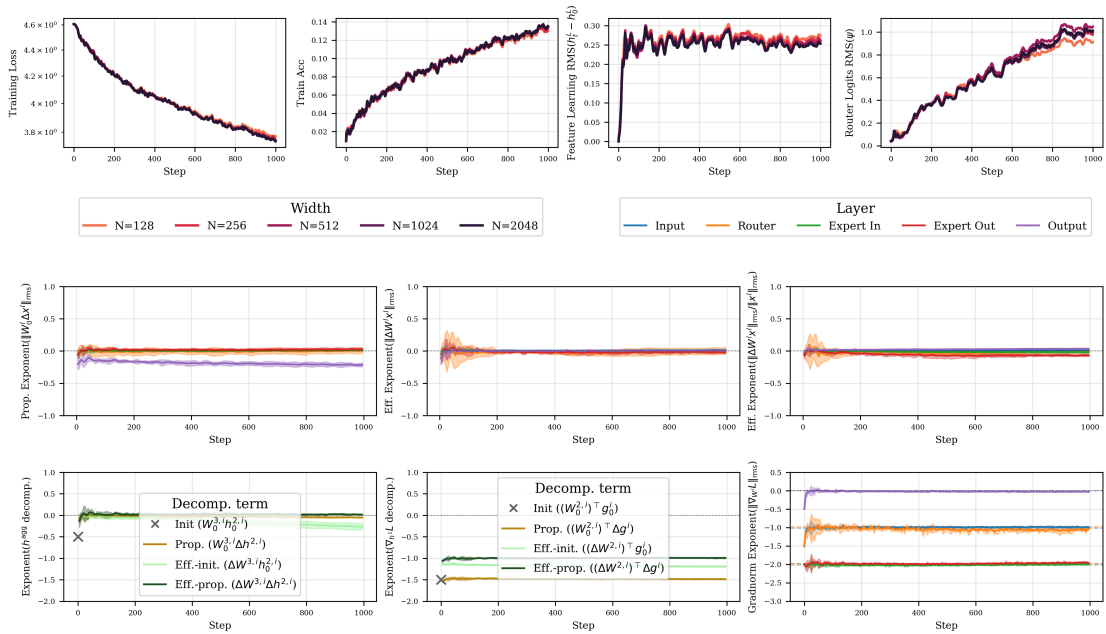


Figure Q.40:  $\mu$ P with  $1/N$  last-layer init, no shared experts (Adam, Regime III).

# SCALING MOES: FROM $\mu$ P TO MSSP

MSSP (ours) (Adam, M, N,  $N_e \sim n \rightarrow \infty$ )  
 shared experts, soft routing, last-layer init=0

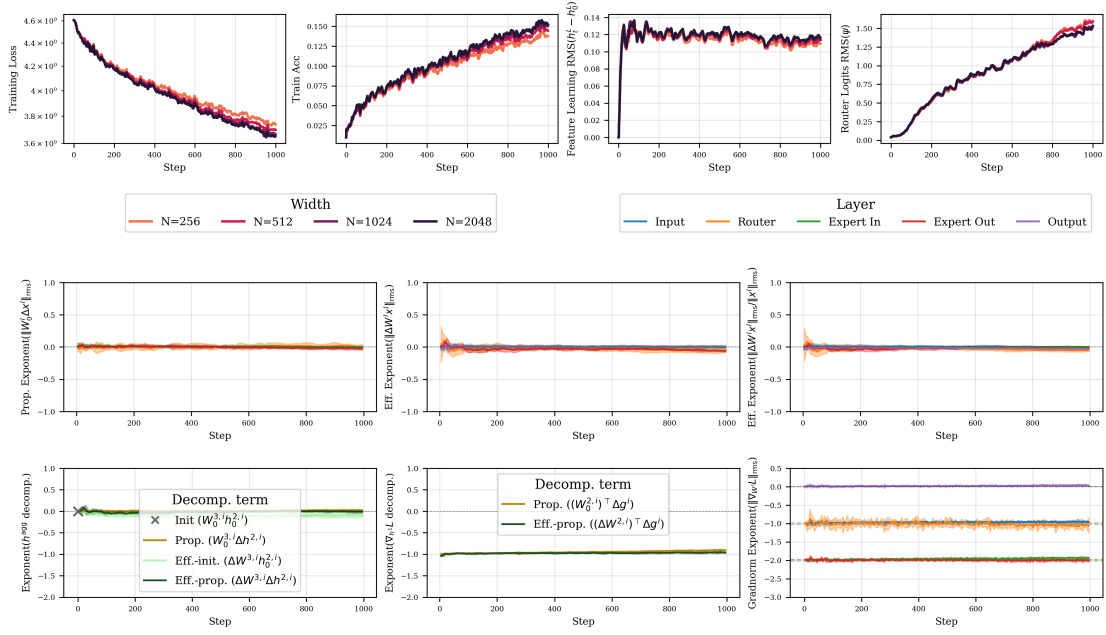


Figure Q.41: MSSP with shared experts (Adam, Regime III).

$\mu$ P (Adam, M, N,  $N_e \sim n \rightarrow \infty$ )  
 no shared experts, last-layer init=1/N

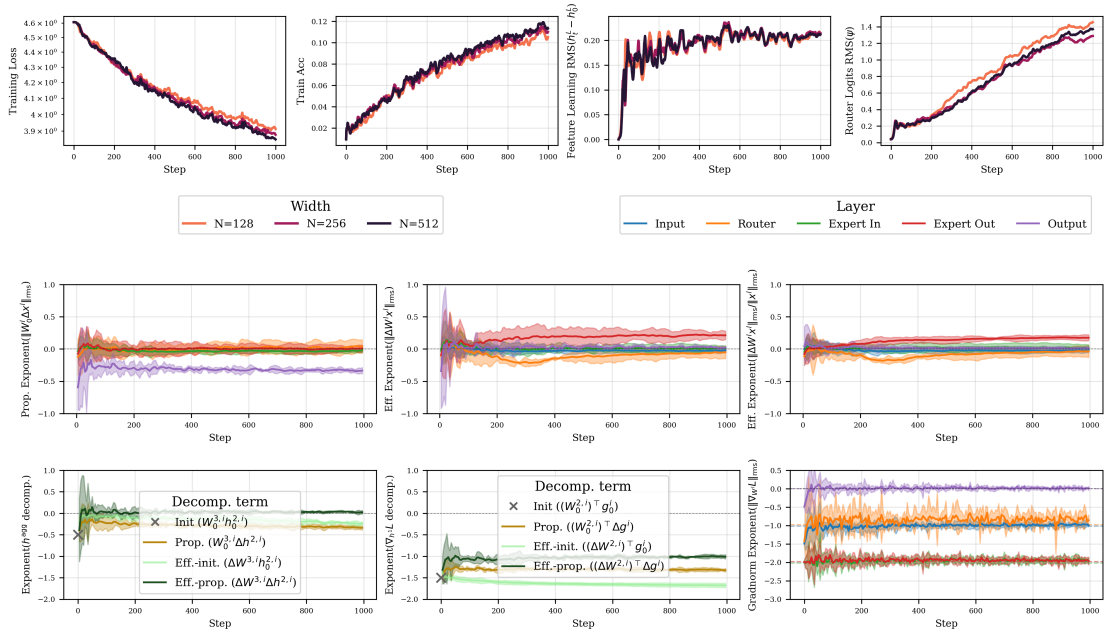


Figure Q.42:  $\mu$ P without shared experts with  $1/N$  last-layer init (Adam, Regime III, top- $k$ ).

MSSP (ours) (Adam,  $M, N, N_e \sim n \rightarrow \infty$ )  
 shared experts, last-layer init=0

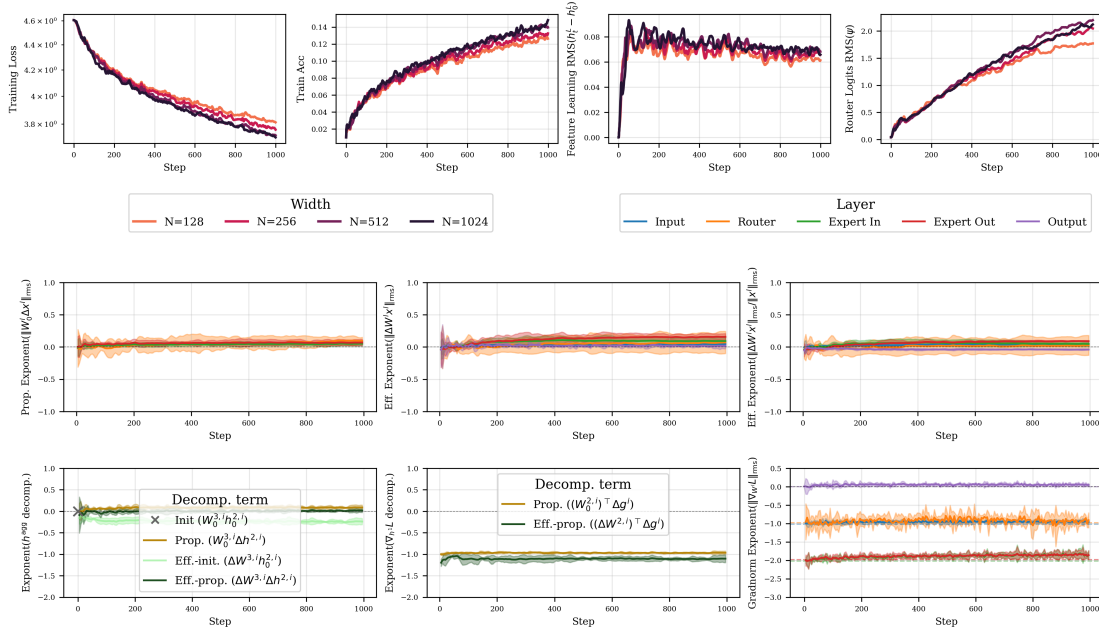


Figure Q.43: MSSP with shared experts (Adam, Regime III, top- $k$ ).

**Q.5. Soft softmax routing collapses to uniform for  $\mu P$  in Regime I**

Here we train with single pass SGD with the optimal learning rate under MSE loss using soft softmax routing over binary classification from CIFAR-10 using all data points from class 'airplane' +1 and class 'automobile' -1 with batch size 64.

Figure Q.44 shows learning rate sweeps scaling with  $\mu P$  versus standard parameterization (SP) under MSE loss. Lines ending on the right indicates divergence with NaNs under larger learning rates. Validation performance gets worse with scale in SP. Due to router collapse, performance does not improve with scale under soft routing in  $\mu P$  either. The optimal and maximal stable learning rate shrinks with width in SP, but transfers in  $\mu P$ .

As theoretically predicted, the router gradient vanishes with increased width (Figure Q.45), which results in vanishing router updates inducing vanishing router logits inducing uniform routing. Figure Q.46 shows that this collapse prevails over the entire course of training.

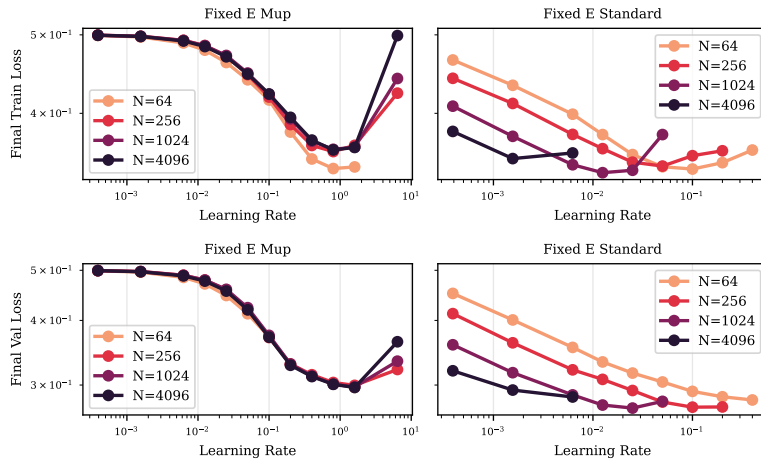


Figure Q.44: **Learning rate transfer in Regime I (soft softmax routing, MSE loss, CIFAR-10).** Top: training loss. Bottom: validation loss. Left:  $\mu P$ . Right: SP. Ending lines denote divergence. Observe the maximal stable and optimal learning rate shrinking in SP, but staying consistent in  $\mu P$ . Performance does not monotonically improve due to router and expert collapse.

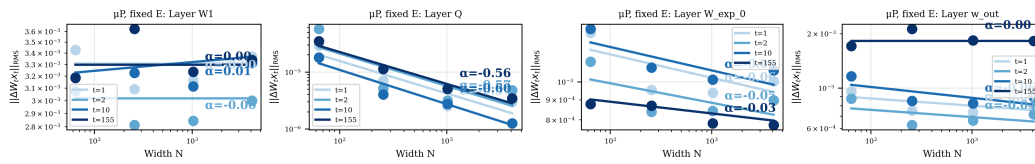


Figure Q.45: **Effective updates across layers ( $\mu P$ , soft, Regime I).** Columns (left to right): input layer (W1), router layer (Q), first expert layer (W\_exp\_0), output layer (w\_out). Approximate width-independence in  $\mu P$ , except in the router, which collapses under soft softmax routing.

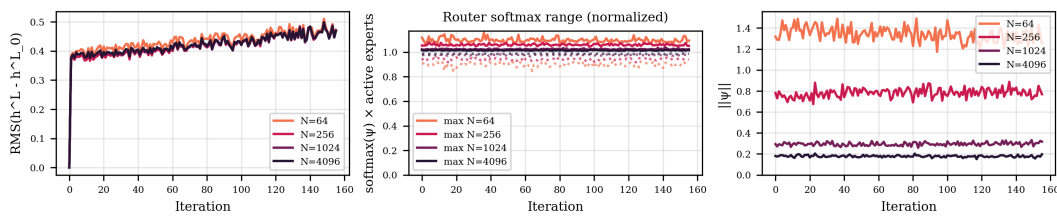


Figure Q.46: **Soft softmax routing collapses to uniform in  $\mu P$  in Regime I.** Here, we train single pass SGD in  $\mu P$  over the CIFAR-10 subset described in this section. Columns (left to right): feature learning  $\|\Delta h^L\|_{RMS}$ , normalized post-softmax routing weights, router logit norm  $\|\psi\|_{RMS}$ . Routing collapses as routing logits converge to 0 with scale over the entire course of training. Feature learning is still approximately width-independent.

Across several settings, we found that expert specialization does not necessarily improve performance on CIFAR-10, as the dataset is too small and not diverse enough. Hence we use TinyImageNet for all other MLP experiments.

**Q.6. MoEs require layerwise learning rate tuning**

Without layerwise learning rate multiplier tuning, effective updates in router and expert input layer can initially vanish even in MSSP (Figure Q.47).

Figure Q.48 shows that, while all exponents are approximately 0 as theoretically predicted, the propagating updates dominate the effective updates in absolute scale in the first steps by a factor of more than  $10^6$  in both the forward and backward pass aggregation operations.

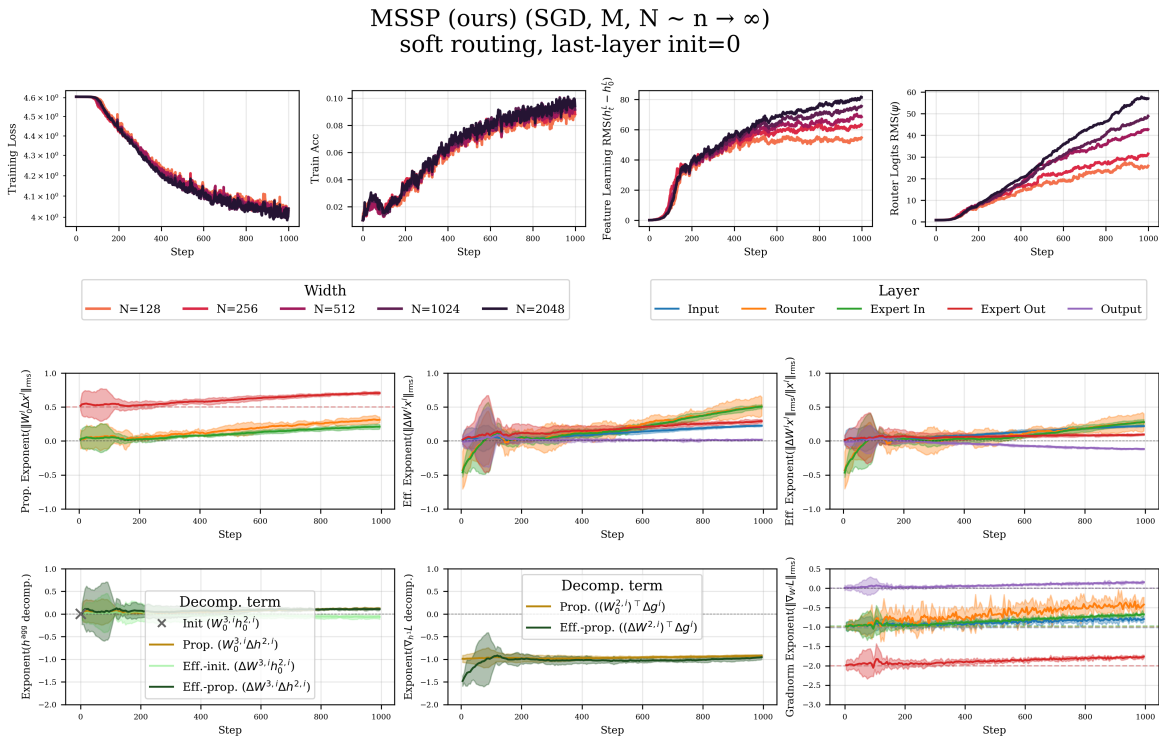


Figure Q.47: **MSSP without tuned multipliers (SGD, Regime II)**. Here we set all initialization variance and layerwise learning rate multipliers to 1.0 and only tune the global learning rate.

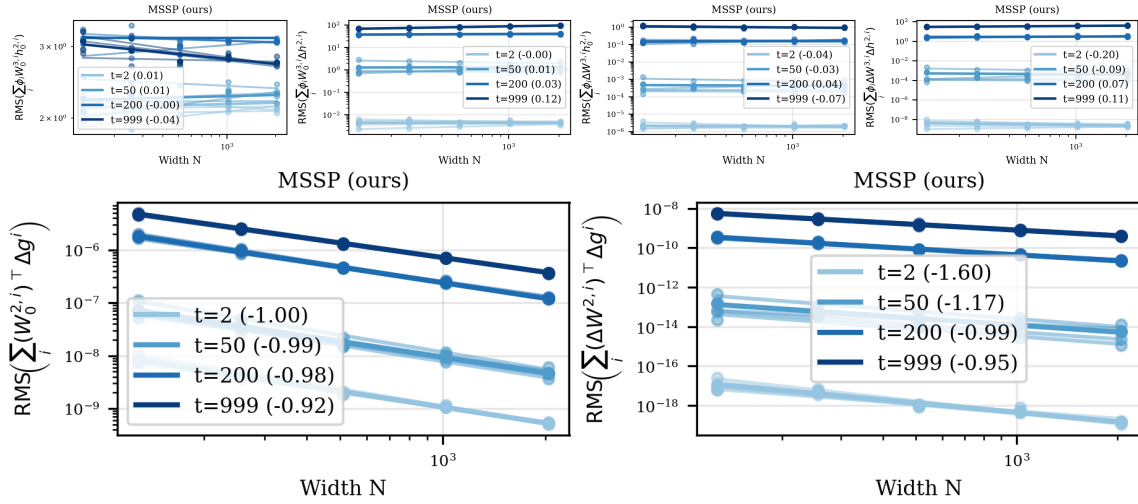


Figure Q.48: **Without tuned multipliers, sub-term contributions differ by orders of magnitude (MSSP, SGD, Regime II).** The forward pass expert aggregation components (top) and the backward pass components to  $\nabla_{h^1} \mathcal{L}_t$ . Propagating updates dominate the effective updates in absolute scale for several hundred steps. In the forward pass, propagating updates start out at around  $10^0$  and effective updates around  $10^{-6}$ . In the backward pass, the propagating update term starts around  $10^{-9}$  at  $N = 1024$  and  $t = 2$ , whereas the effective update term starts around much smaller  $10^{-18}$ . Hence the weight updates contribute vanishingly to the overall first-layer gradient  $\nabla_{h^1} \mathcal{L}_t$ .

Hence, in parameterizations where correctly scaled effective updates should dominate vanishing propagating updates, the propagating updates can still dominate in absolute value for several hundred steps at realistic scales. Hence it can seem that the empirical exponents do not follow the predicted ones, when not measured in sufficient granularity.

After layerwise learning rate tuning, which can consequently require grids with ranges beyond  $10^6$ , we observe that propagating and effective updates generally end up having a similar order of magnitude in the first 1000 steps, and hence we accurately measure the predicted limiting exponents at moderate model scales.

Overall this highlights that extensively tuning layerwise learning rate multipliers is even more essential in MoEs than in dense models.

### Q.7. Random search and 2D multiplier tuning do not suffice

Instead of full 6D sweeps, one could also tune multipliers more efficiently by running a random search, and, starting from the HPs found in the first stage, doing 2D sweeps of all HP pairs around the optimum. This approach finds near-optimal HPs if the suboptimality of the random stage is predominantly restricted to a subspace allowing 2 interacting multipliers.

However, Figure Q.49 shows that 2D sweeps after a random sweep do not suffice. For example in SGD  $\mu P$  Regime III (left), all 2D sweeps containing the expert input or output layer lr multiplier suggest that smaller values would perform better, but the joint change of (lr router, lr out) dominates and does not allow a third change. Similarly in SGD  $\mu P$  Regime II (center) the pair (lr router, lr expert2) marginally dominates other improvements such as even smaller expert output learning rate. Hence, after such a 2D sweep, it remains unclear whether the new expert output learning rate, and

all other HPs, are indeed robustly near-optimal, or whether higher-order interactions continue to induce an indefinite 2D update without ever converging to a local optimum. In Adam MSSP Regime I (right), (lr router, lr in) are updated and it remains unclear whether the global optimum would also include a reduced expert output lr multiplier.

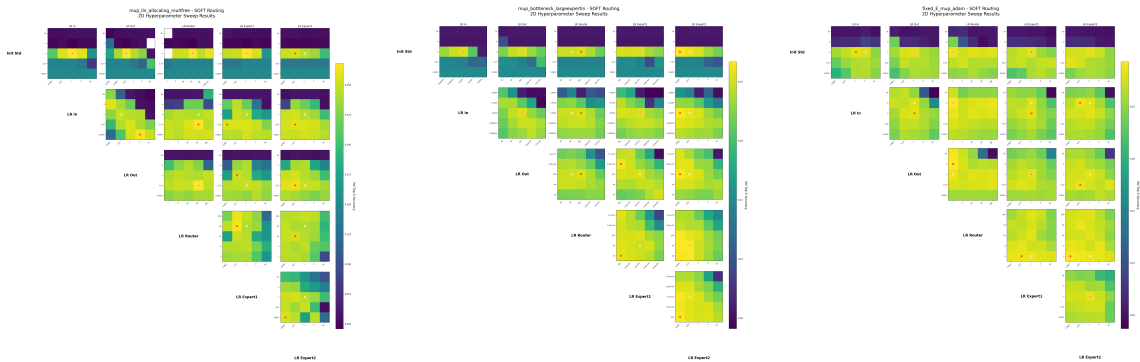


Figure Q.49: **2D multiplier sweeps at small scale  $N = 128$  from random search stage.** 2D heatmaps showing top-5 validation accuracy of all HP pairs, while fixing all other multipliers at the optimum found from a random sweep over at least 1024 runs. From left to right: SGD  $\mu$ P Regime III, SGD  $\mu$ P Regime II, Adam MSSP Regime I. White circles denote the optimum after this 2D stage, red 'x' denote the optimum within the respective 2D sweep. Often these do not align, suggesting that higher dimensional interactions make this hyperparameter tuning strategy insufficient.

### Q.8. Global Adam $\epsilon$ induces width dependence at sufficient scale

Here, we mimick the effects of further scaling – which we cannot run due to compute constraints – by increasing Adam  $\epsilon$ , since gradient RMS norms decay with model scale. We compare naive constant Adam  $\epsilon$  versus our layerwise Adam  $\epsilon$  scaling under allscaling and layerwise LR multipliers tuned at width 128. Figure Q.50 shows more time-dependent and less clean effective and propagating update exponents under global  $\epsilon$ , especially of expert effective updates. Figure Q.51 shows feature learning and router learning is reduced with scale under global epsilon, as updates are beginning to vanish.

## SCALING MOES: FROM $\mu P$ TO MSSP

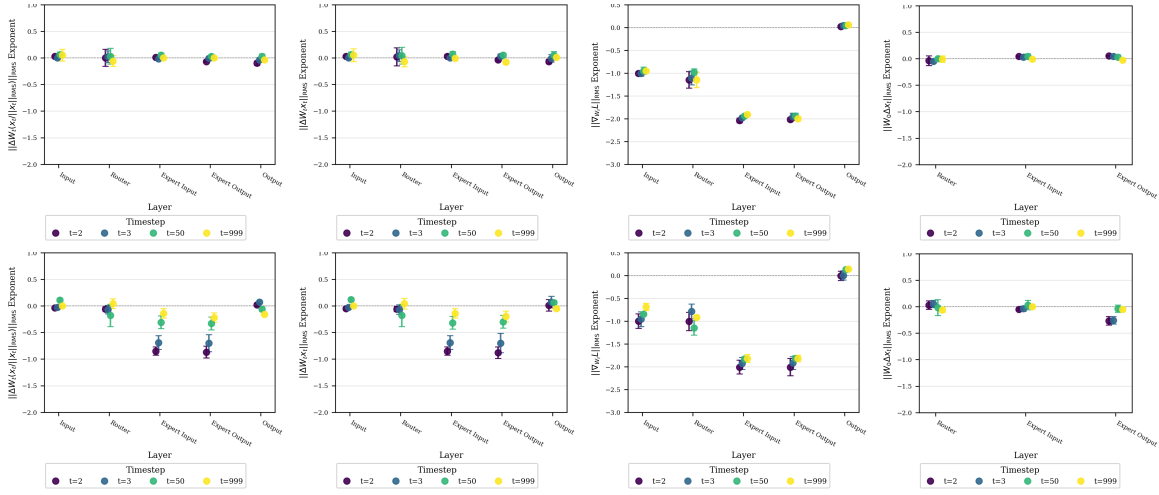


Figure Q.50: **Global Adam  $\epsilon$  induces entangled exponents.** RCC exponents of Adam MSSP in Regime III with layerwise  $\epsilon$  scaling (Ours, top), the same scaling rule with constant global  $\epsilon$  (bottom). Columns (left to right): Effective Updates (Normalized), Effective Updates (Raw), Gradient Norms (Raw), Propagating Updates (Raw).

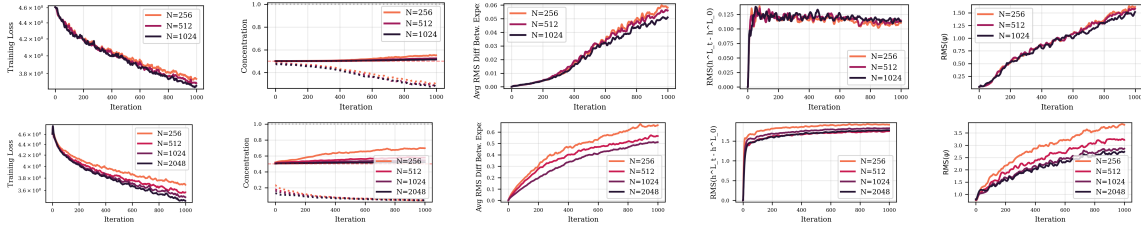


Figure Q.51: **Reduced feature learning and router learning under global Adam  $\epsilon$ .** Adam MSSP in Regime III with layerwise  $\epsilon$  scaling (top) versus the same scaling rule with constant global Adam  $\epsilon$  (bottom). Columns (left to right): Training loss, router concentration (0.5 indicates uniform routing), solid lines denote the maximal and dotted lines the minimal routing weight in relation to uniform routing), average RMS norm difference between experts, accumulated feature learning ( $\|\Delta h^L\|_{RMS}$ ), router logit norm  $\|\psi\|_{RMS}$ .