

# Fixed-Point Reasoning: Stable and Adaptive Deep Looped Models

author names withheld

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

## Abstract

Looped-in-depth architectures provide an inductive bias toward learning step-by-step procedures for tasks that require compositional reasoning. The effective depth reached by looping determines the quality of the solution. Similar to deep architectures, looped architectures are prone to signal propagation issues as the halting decision is postponed. In this paper, we address these signal propagation issues by using pre-norm layers and residual scaling. Furthermore, we propose *FPRM*: a Fixed-Point Reasoning Model that uses fixed-point convergence as an end-to-end halting mechanism in a looped architecture. We show that fixed-point halting allows FPRM to adapt its compute to the difficulty of the task. FPRM proves effective on common reasoning benchmarks, namely sudoku, maze, and state tracking.

## 1. Introduction

Reasoning in neural networks has increasingly been framed as a problem of scaling test-time compute: a model should be able to spend more computation on inputs it finds harder [23, 26]. However, doing so requires two ingredients: **(1) flexibility**: a mechanism for spending a variable amount of compute on the problem, and **(2) adaptivity**: a way to decide how much compute to spend on the problem. One way to achieve this is through Chain-of-Thought (CoT) mechanism [30]. In CoT, the model scales compute through verbalization, and makes halting decisions based on predicting a specialized halting token. However, this emerging behavior requires a special training regime [14].

An alternative approach towards an end-to-end training method for reasoning models is emerging in the form of looped architectures [4, 6, 25]. Whereas in CoT, the compute scales along the sequence dimension, in looped architectures it scales along the depth dimension:

$$\mathbf{z}_i = f_\theta(\mathbf{z}_{i-1}; \mathbf{x}), \quad (1)$$

where  $f_\theta(\mathbf{z}; \mathbf{x})$  is a neural network,  $\mathbf{z}_i$  is the latent space output after  $i$ -th iteration and  $\mathbf{x}$  is the input injection. Since in looped models, compute scales independently of parameter count, it can be increased naturally, satisfying the first requirement – flexibility. The halting decision that ensures adaptivity, however, is no longer trivial. Most approaches either fix or randomly sample the number of loops [11], which eliminates adaptivity, or use external modules trained to make halting decisions [15, 29]. The latter is associated with adaptive computation time (ACT) modules [13], which introduce non-differentiable halting objectives that can only be trained via continuous relaxation.

One of the recent examples of this paradigm are Hierarchical Reasoning Model (HRM) [29] and Tiny Recursive Model (TRM) [15], which achieve strong performance on reasoning benchmarks such as sudoku, maze, and ARC-AGI [5]. Using small networks with bounded iteration counts, these models have been able to outperform prominent LLM-based reasoning models.

At the core of these methods is the idea of the need for a hierarchical looping mechanism, wherein the compute is distributed between a fast-looping component (the L-level) and a slow-looping component (the H-level). However, paired with the ACT mechanism, they ultimately fail on the adaptivity test as the compute during test-time is fixed and pre-determined. Alternative approaches such as Deep Equilibrium Models [2] and energy-based reasoning models [8] provide a more natural criterion for halting in the form of fixed-points.

One of the striking differences between looped and non-looped transformers is the common use of a post-norm in the former [6, 11, 15, 29], which is considered to cause unstable training in deep non-looped architectures compared to pre-norm [22, 31]. However, in a looped model, post-norm blocks ensure that the magnitude of the activations stay bounded, which looped architectures rely on to ensure the hidden states do not diverge as the model iterates [18]. This begs the question: *can we switch the post-norm to pre-norm, while ensuring the activations stay bounded in a different way?* If so, this could make the training of looped models stable at large depths.

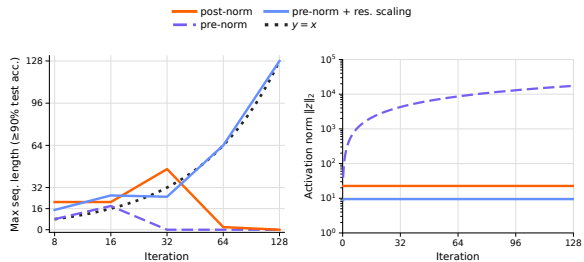
We give an affirmative answer to this question by introducing scaling parameters on the skip connections and residuals. With an additional set of modifications, we are able to train an adaptive reasoning model that loops until it converges to a fixed point. For literature review, please refer to Section A. We summarize our contributions as follows:

1. **Training the first looped Transformer with pre-norm.** We modify a transformer block to be trainable over larger depths by switching post-norm to pre-norm and adding residual scaling parameters.
2. **Stabilizing fixed-points as a halting mechanism.** To further enable stable training, we propose a theoretically grounded modification specific to fixed-point models to limit the oscillation around the fixed-point.
3. **Proposing Fixed-Point Reasoning Model (FPRM).** In our experiments, FPRM outperforms TRM and HRM on sudoku, maze and the state-tracking A5 benchmark.

An overview of the proposed architecture is available in Figure 3.

## 2. Method

As looped architectures need bounded activations across iterations, current methods mostly employ post-norm to satisfy the boundedness condition [15]. However, post-norm introduces serious signal propagation issues [7, 22, 27]. In this section, we resolve these issues separately by introducing and motivating some architecture modifications. An overview of FPRM is available in Algorithm 2.



(a) Trainability at different iteration budgets. (b) Activation norm at initialization.

Figure 1: Large effective depth, if achieved when needed, unlocks expressivity, but at the same time poses a challenge to stabilize (in a case of pre-norm) or utilize the signal (in a case of post-norm).

## 2.1. Improving signal propagation with pre-norm

We start by introducing pre-norm and post-norm in our architecture. A transformer block consists of two residual sub-layers – multi-head attention (MHA) and a feed-forward network (FFN) – interleaved with layer normalization (LN). Two canonical placements of the normalization define two variants of the block. The original *post-norm* formulation [28] applies LN after the residual addition  $\mathbf{h} = \text{LN}(\mathbf{x} + \text{MHA}(\mathbf{x}))$  and  $\mathbf{y} = \text{LN}(\mathbf{h} + \text{FFN}(\mathbf{h}))$ , while the *pre-norm* variant [31] applies LN to the input of each sub-layer and leaves the residual stream untouched  $\mathbf{h} = \mathbf{x} + \text{MHA}(\text{LN}(\mathbf{x}))$  and  $\mathbf{y} = \mathbf{h} + \text{FFN}(\text{LN}(\mathbf{h}))$ . The two are not interchangeable. Post-norm bounds activation magnitudes at every layer, but attenuates the residual signal with depth and is well known to be unstable to train in deep networks [22]. Pre-norm preserves a clean residual path, yielding well-behaved gradient flow at depth, at the cost of activations that can grow without bound.

In Figure 1(a), we observe that increasing the iteration count of a post-norm Universal Transformer [6] does not translate into improved accuracy, while a pre-norm variant continues to improve with depth. Motivated by this, our first change is to switch to pre-norm in our reasoning architecture.

## 2.2. Recovering boundedness via residual scaling

While pre-norm mitigates trainability problems in looped architectures, it removes the necessary boundedness condition that motivated the use of post-norm in them. This effect can be observed in Figure 1(b), where the activations of the pre-norm model grow with more iterations. Consequently, we propose to restore boundedness by introducing scaling parameters applied at two different scales: one within each layer of the network, and one over the looping.

**Layer-wise residual scaling.** Within a single application of  $f_\theta(\mathbf{z}, \mathbf{x})$ , the residual stream and sub-layer output  $f_{\theta_\ell}^\ell(\mathbf{z}^{\ell-1})$  (an MHA or FFN sub-block with layer normalization, as in Section 2.1) are weighted by tied scalars  $(\alpha_1, \beta_1)$  shared across all  $L$  layers:

$$\mathbf{z}^\ell = \alpha_1 \mathbf{z}^{\ell-1} + \beta_1 f_{\theta_\ell}^\ell(\mathbf{z}^{\ell-1}), \quad \ell = 1, \dots, L. \quad (2)$$

**Iteration-wise input mixing.** Between consecutive applications of  $f_\theta(\mathbf{z}, \mathbf{x})$ , we re-inject the input  $\mathbf{x}$  with tied scalars  $(\alpha_2, \beta_2)$  shared across all iterations:

$$\mathbf{z}_{i+1}^0 = \alpha_2 \mathbf{z}_i^L + \beta_2 \mathbf{x}. \quad (3)$$

The two scaling schemes are not independent. With an appropriate coupling between them, the resulting recurrence is bounded for any input, resulting in a stable looping. In the following statement, we formalize this claim:

**Theorem 1 (Boundedness of FPRM iterates)** *Consider the model defined by Equations 2 and 3, and assume each layer map satisfies  $\|f_{\theta_\ell}^\ell(\mathbf{u})\| \leq c_f$  for all  $\ell$  and  $\mathbf{u}$ . Let  $0 \leq \alpha_1, \alpha_2 < 1$ , and set  $\beta_2 = 1 - \alpha_2 \alpha_1^L$ ,  $\beta_1 = \frac{\beta_2(1-\alpha_1)}{1-\alpha_1^L}$ . Then the fixed-point iterates  $\{\mathbf{z}_i^0\}_{i \geq 0}$  are bounded, and if  $\mathbf{z}_i^0 \rightarrow \mathbf{z}_\infty^0$ , then  $\|\mathbf{z}_\infty^0\| \leq \|\mathbf{x}\| + \alpha_2 c_f$ .*

The proof is in Section B.1. Note that the boundedness condition of the sequence model in Theorem 1 is satisfied when using pre-norm [17]. However, boundedness still does not guarantee the

looping to converge to a fixed-point, which is a condition we need to satisfy the adaptivity ingredient. In Theorem 2, we show that there exists a choice of  $\alpha_2$  that satisfies this requirement. Consequently, as empirically shown by [4], the model may become locally contractive during training. We avoid limited expressivity [1, 2] by making  $\alpha_1$  and  $\alpha_2$  learnable. In practice, we find that initializing the network to be more contractive by setting  $\alpha_2$  to be small yields better performance (Table 3), in line with the common solutions to signal-propagation and rank-collapse problems [22, 27]. Together, these modifications yield a pre-norm looped transformer that maintains performance over longer looping horizons before saturating, compared to the post-norm variant (see Figure 1).

### 2.3. Oscillation around the fixed-point

So far, we have been able to guarantee that, given a small enough  $\alpha_2$ , the model introduced in Equation (3) becomes locally contractive and converges to a fixed-point. However, in practice the contraction factor introduced in Theorem 2 is only as small as the training makes it. For some inputs, we observe that the model often descends into an oscillatory behavior, causing the iteration to stay in a small region of latent space without converging. This non-convergent behavior is not in tension with Theorem 2, as the theorem gives a sufficient condition for convergence, not a complete characterization of the iteration’s behavior. In fact, oscillation around the fixed-point can happen whenever the Jacobian has complex eigenvalues with  $< 1$  real components, but with  $> 1$  magnitude.

In Theorem 3, we show that a suitable damping factor  $\eta$  eliminates the oscillations while preserving the fixed points of  $f_\theta(\mathbf{z}; \mathbf{x})$ . To choose  $\eta$  adaptively at inference time, we use a patience mechanism that decreases  $\eta$  whenever the residual stops improving. We track the smallest relative residual observed so far,  $r^* = \min_i \frac{\|\mathbf{z}_i - f_\theta(\mathbf{z}_i; \mathbf{x})\|_\infty}{\|f_\theta(\mathbf{z}_i; \mathbf{x})\|_\infty + \epsilon}$ , and when it fails to improve for  $P$  consecutive iterations, we apply a geometric decay  $\eta \leftarrow \gamma \eta$  with  $\gamma \in (0, 1)$ . The full procedure is given in Algorithm 1. In Section D, we provide further details about the optimization of a fixed-point model.

## 3. Experiments

### 3.1. Results

**Sudoku and maze navigation.** We first evaluate FPRM on the Sudoku-Extreme and Maze-Hard benchmarks introduced by HRM [29]. Sudoku-Extreme consists of challenging  $9 \times 9$  Sudoku puzzles in a small-sample learning setting, while Maze-Hard evaluates optimal path finding in hard  $30 \times 30$  mazes. These benchmarks were designed to test whether latent recurrent reasoning models can solve symbolic search problems from limited supervision. Table 1 shows that FPRM achieves the highest accuracy on both tasks, with the largest improvement on Sudoku-Extreme.

**State tracking.** We evaluate FPRM on state tracking, where a model must compose a sequence of updates to recover a latent state. This task is a proxy for stateful reasoning problems such as entity tracking, code execution, and game-state tracking [20], and longer instances require greater effective depth [19, 20]. We train on sequences of length  $k = 32$  and test up to  $k = 128$ . We compare to TRM [15], which disables its ACT head at test time and runs 16 outer iterations, and to TRM+ACT, which leaves the ACT head enabled for adaptive

Table 1: Test accuracy on Sudoku-Extreme and Maze-Hard. Architecture backbone: Transformer.

Model	Sudoku-Extreme	Maze-Hard
HRM	55.0%	74.5%
TRM	74.7%	85.3%
FPRM	<b>81.3%</b>	<b>86.4%</b>

halting. For fairness, the looping budget and number of backward passes  $K$  for BPTT are matched. As shown in Figure 2, *FPRM maintains high accuracy beyond the training regime, while both TRM variants degrade with length*. Moreover, FPRM uses few iterations on short sequences and increases its iterations smoothly with sequence length, indicating that it allocates more effective depth to harder inputs; TRM+ACT does not show the same monotonic length-sensitive halting behavior.

### 3.2. Analyzing depth-induced signal propagation issues in looped transformers

As introduced in Section 2.1, recurrent-in-depth models often achieve large effective depth, at the cost of more unstable signal propagation. To isolate and analyze this instability, we adopt the Universal Transformer framework with a fixed number of iterations, controlling for the confounding effects of dynamic halting mechanisms of looping. In Figure 1(b) we show the norm of the final activations of the residual branch of Universal Transformer at initialization. As expected, pre-norm model diverges due to unbounded activations. Conversely, having post-norm architecture or pre-norm with residual scaling ensures boundedness of the activations. In Section E, for our proposed adaptive reasoning model, we provide similar results showing that models with post-norm or pre-norm with residual scaling can achieve much bigger effective depths.

However, the boundedness-at-initialization condition from the previous section is not sufficient to ensure feature learning at large depth. We demonstrate this in Figure 1(a) using a state-tracking task, where increasing difficulty, corresponding to longer sequences, requires training at greater depths. To evaluate trainability, we plot the maximum sequence length solved with  $> 90\%$  accuracy against the number of loops. Only the pre-norm UT variant with residual scaling is capable of solving the same task length the model was trained on, in the regime where the iteration budget is matched to the task length. It is the only architecture whose maximum solvable length grows monotonically with the trained iteration count. We interpret Figure 1(a) as direct evidence that norm-controlled, non-attenuating residual propagation is a necessary condition for a looped transformer to perform sequential composition that requires larger effective depth.

## 4. Discussion and Conclusion

We present architectural modifications for looped fixed-point transformers that enable the use of pre-norm, improving the model’s ability to exploit the larger effective depth provided by looping. These modifications allow FPRM to outperform its closest baselines, HRM and TRM, on common symbolic reasoning benchmarks. We show that on state tracking it generalizes beyond its training distribution, which is to our knowledge, the first transformer-based architecture to length-generalize on this task. This capability stems from dynamically scaling depth through fixed-point iterations and improving signal propagation. We hope these architectural modifications and the accompanying insights will support further progress on latent reasoning models.

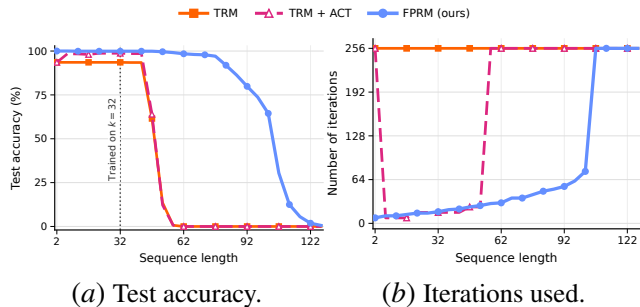


Figure 2: FPRM maintains high test accuracy substantially beyond the training regime while using far fewer iterations than the TRM baselines.

## References

- [1] Cem Anil, Ashwini Pokle, Kaiqu Liang, Johannes Treutlein, Yuhuai Wu, Shaojie Bai, J. Zico Kolter, and Roger B. Grosse. Path independent equilibrium models can better exploit test-time computation. In *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/331c41353b053683e17f7c88a797701d-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/331c41353b053683e17f7c88a797701d-Abstract-Conference.html).
- [2] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. In *NeurIPS*, pages 688–699, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/01386bd6d8e091c2ab4c7c7de644d37b-Abstract.html>.
- [3] Andrea Banino, Jan Balaguer, and Charles Blundell. Pondernet: Learning to ponder, 2021. URL <https://arxiv.org/abs/2107.05407>.
- [4] Arpit Bansal, Avi Schwarzschild, Eitan Borgnia, Zeyad Emam, Furong Huang, Micah Goldblum, and Tom Goldstein. End-to-end algorithm synthesis with recurrent networks: Extrapolation without overthinking. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/7f70331dbe58ad59d83941dfa7d975aa-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/7f70331dbe58ad59d83941dfa7d975aa-Abstract-Conference.html).
- [5] Francois Chollet, Mike Knoop, Gregory Kamradt, and Bryan Landers. Arc prize 2024: Technical report, 2025. URL <https://arxiv.org/abs/2412.04604>.
- [6] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers, 2019. URL <https://arxiv.org/abs/1807.03819>.
- [7] Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: pure attention loses rank doubly exponentially with depth. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, Proceedings of Machine Learning Research, pages 2793–2803. PMLR, 2021. URL <http://proceedings.mlr.press/v139/dong21a.html>.
- [8] Yilun Du, Shuang Li, Joshua B. Tenenbaum, and Igor Mordatch. Learning iterative reasoning through energy minimization, 2022. URL <https://arxiv.org/abs/2206.15448>.
- [9] Ying Fan, Yilun Du, Kannan Ramchandran, and Kangwook Lee. Looped transformers for length generalization, 2025. URL <https://arxiv.org/abs/2409.15647>.
- [10] Zitian Gao, Lynx Chen, Yihao Xiao, He Xing, Ran Tao, Haoming Luo, Joey Zhou, and Bryan Dai. Universal reasoning model. *arXiv preprint arXiv:2512.14693*, 2025.
- [11] Jonas Geiping, Sean McLeish, Neel Jain, John Kirchenbauer, Siddharth Singh, Brian R. Bartoldson, Bhavya Kailkhura, Abhinav Bhatele, and Tom Goldstein. Scaling up test-time compute with latent reasoning: A recurrent depth approach, 2025. URL <https://arxiv.org/abs/2502.05171>.

- [12] Zhengyang Geng, Xin-Yu Zhang, Shaojie Bai, Yisen Wang, and Zhouchen Lin. On training implicit models. In *NeurIPS*, pages 24247–24260, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/cb8da6767461f2812ae4290eac7cbc42-Abstract.html>.
- [13] Alex Graves. Adaptive computation time for recurrent neural networks. *CoRR*, abs/1603.08983, 2016. URL <http://arxiv.org/abs/1603.08983>.
- [14] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, Hao Zhang, Hanwei Xu, Honghui Ding, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jingchang Chen, Jingyang Yuan, Jinhao Tu, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaichao You, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingxu Zhou, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Tao Yun, Tian Pei, Tianyu Sun, Tao Wang, Wangding Zeng, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nat.*, 645(8081):633–638, 2025. doi: 10.1038/S41586-025-09422-Z. URL <https://doi.org/10.1038/s41586-025-09422-z>.
- [15] Alexia Jolicoeur-Martineau. Less is more: Recursive reasoning with tiny networks, 2025. URL <https://arxiv.org/abs/2510.04871>.
- [16] Ferdinand Kapl, Emmanouil Angelis, Kaitlin Maile, Johannes von Oswald, and Stefan Bauer. From growing to looping: A unified view of iterative computation in llms, 2026. URL <https://arxiv.org/abs/2602.16490>.

- [17] Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, Proceedings of Machine Learning Research*, pages 5562–5571. PMLR, 2021. URL <http://proceedings.mlr.press/v139/kim21i.html>.
- [18] Asher Labovich. Stability and generalization in looped transformers, 2026. URL <https://arxiv.org/abs/2604.15259>.
- [19] William Merrill and Ashish Sabharwal. A little depth goes a long way: The expressive power of log-depth transformers, 2025. URL <https://arxiv.org/abs/2503.03961>.
- [20] William Merrill, Jackson Petty, and Ashish Sabharwal. The illusion of state in state-space models, 2025. URL <https://arxiv.org/abs/2404.08819>.
- [21] Sajad Movahedi, Felix Sarnthein, Nicola Muca Cirone, and Antonio Orvieto. Fixed-point rnns: Interpolating from diagonal to dense, 2025. URL <https://arxiv.org/abs/2503.10799>.
- [22] Lorenzo Noci, Sotiris Anagnostidis, Luca Biggio, Antonio Orvieto, Sidak Pal Singh, and Aurélien Lucchi. Signal propagation in transformers: Theoretical perspectives and the role of rank collapse. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/ae0cba715b60c4052359b3d52a2cff7f-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ae0cba715b60c4052359b3d52a2cff7f-Abstract-Conference.html).
- [23] OpenAI. Learning to reason with LLMs. <https://openai.com/index/learning-to-reason-with-llms/>, September 2024. OpenAI blog post, accompanying the o1 release.
- [24] Zirui Ren and Ziming Liu. Are your reasoning models reasoning or guessing? a mechanistic analysis of hierarchical reasoning models, 2026. URL <https://arxiv.org/abs/2601.10679>.
- [25] Nikunj Saunshi, Nishanth Dikkala, Zhiyuan Li, Sanjiv Kumar, and Sashank J. Reddi. Reasoning with latent thoughts: On the power of looped transformers, 2025. URL <https://arxiv.org/abs/2502.17416>.
- [26] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling model parameters. *CoRR*, abs/2408.03314, 2024. doi: 10.48550/ARXIV.2408.03314. URL <https://doi.org/10.48550/arXiv.2408.03314>.
- [27] Wenfang Sun, Xinyuan Song, Pengxiang Li, Lu Yin, Yefeng Zheng, and Shiwei Liu. The curse of depth in large language models. *CoRR*, abs/2502.05795, 2025. doi: 10.48550/ARXIV.2502.05795. URL <https://doi.org/10.48550/arXiv.2502.05795>.

- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- [29] Guan Wang, Jin Li, Yuhao Sun, Xing Chen, Changling Liu, Yue Wu, Meng Lu, Sen Song, and Yasin Abbasi-Yadkori. Hierarchical reasoning model. *CoRR*, abs/2506.21734, 2025. doi: 10.48550/ARXIV.2506.21734. URL <https://doi.org/10.48550/arXiv.2506.21734>.
- [30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html).
- [31] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the transformer architecture. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, Proceedings of Machine Learning Research, pages 10524–10533. PMLR, 2020. URL <http://proceedings.mlr.press/v119/xiong20b.html>.
- [32] Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms, 2024. URL <https://arxiv.org/abs/2311.12424>.
- [33] Łukasz Kaiser and Ilya Sutskever. Neural gpus learn algorithms, 2016. URL <https://arxiv.org/abs/1511.08228>.

## Appendix A. Background and Related Works

**Looped models.** Looped architectures address the flexibility ingredient of reasoning models by decoupling effective depth from parameter count, providing variable computation per input without scaling the number of parameters. This makes looping a natural inductive bias for tasks that can be solved by repeatedly applying local or compositional subroutines. Early examples include Neural GPUs [33] and Universal Transformers [6]; more recent work shows that looped transformers can improve length generalization, learn algorithmic structure, and approximate much deeper untied models on reasoning tasks [9, 11, 16, 25, 32]. Recent recursive reasoning models such as HRM [29] and TRM [15] instantiate this principle in compact architectures. These hierarchical models were followed up by further architectural and training modifications, improving the performance [10, 24]. However, the proposed modifications from these works are orthogonal to our method, and therefore can be used jointly. However, they leave open a central design question: *how should a looped model decide how many iterations to run on each input at test-time?* FPRM addresses this through fixed-point halting.

**Adaptive computation.** Classical ACT methods answer the adaptivity question by learning an explicit halting rule, such as halting probabilities, per-token stopping decisions, or learned distributions over computation depth [3, 6]. Such mechanisms can allocate more computation to harder instances, but they introduce an additional learned decision process on top of the recurrent computation itself, with an inherently non-differentiable objective.

**Deep equilibrium and fixed-point models.** An alternative is to use the convergence of the latent dynamics as the halting criterion. In this view, computation stops when consecutive iterates become sufficiently close  $\|z_{i+1} - z_i\| \leq \epsilon$ , which corresponds to convergence toward a fixed point  $z^* = f_\theta(z^*; \mathbf{x})$ . This view is most extensively developed in Deep Equilibrium Models [2], which replace a finite stack of layers by the equilibrium point of a weight-tied transformation. It has been shown that the same fixed-point view extends to looped recurrent architectures. This provides an input-dependent notion of computation: *easy instances may reach a stable representation after few iterations, while harder instances can continue refining their state until convergence* [21]. Therefore, for looped reasoning models fixed-point convergence is not only a representation-learning device; it is also the mechanism that determines how much computation is performed. FPRM exploits this view to obtain a parameter-free halting criterion that inherits the convergence guarantees of the iteration map.

**Signal propagation in looped models.** As unrolled looped architectures can be viewed as deep networks, they are exposed to the same signal-propagation difficulties that arise in very deep transformers. In deep sequence models, increasing depth can make optimization harder and can prevent later layers from being effectively used, a phenomenon often discussed as the curse of depth [7, 22,

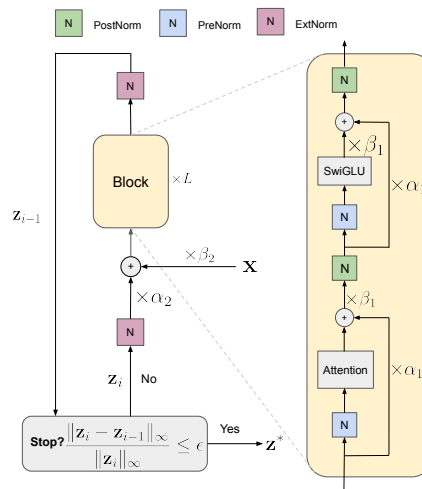


Figure 3: **FPRM**: our fixed-point looped transformer with better signal propagation due to pre-norm and residual scalings.

27]. Pre-norm is the standard remedy for these issues, but in a looped setting it removes a property the architecture relies on: **bounded activations**. Bounded activations, important for stable training, are typically enforced through post-norm [11, 15, 29]. FPRM combines pre-norm with residual scaling to recover bounded and stable dynamics while still allowing gradients and representations to propagate through many iterations.

## Appendix B. Theorems and Proofs

### B.1. Fixed-point iterations bound

**Proof** We start by unrolling the computation across the  $L$  layers within a single fixed-point iteration. By recursively applying the layer update, we obtain

$$\begin{aligned}
 \mathbf{z}_i^L &= \alpha_1 \cdot \mathbf{z}_i^{L-1} + \beta_1 \cdot f_{\theta_L}^L(\mathbf{z}_i^{L-1}) \\
 &= \alpha_1^2 \cdot \mathbf{z}_i^{L-2} + \alpha_1 \beta_1 \cdot f_{\theta_{L-1}}^{L-1}(\mathbf{z}_i^{L-2}) + \beta_1 \cdot f_{\theta_L}^L(\mathbf{z}_i^{L-1}) \\
 &= \alpha_1^3 \cdot \mathbf{z}_i^{L-3} + \alpha_1^2 \beta_1 \cdot f_{\theta_{L-2}}^{L-2}(\mathbf{z}_i^{L-3}) + \alpha_1 \beta_1 \cdot f_{\theta_{L-1}}^{L-1}(\mathbf{z}_i^{L-2}) + \beta_1 \cdot f_{\theta_L}^L(\mathbf{z}_i^{L-1}) \\
 &\quad \vdots \\
 &= \alpha_1^L \cdot \mathbf{z}_i^0 + \beta_1 \cdot \left( \sum_{j=0}^{L-1} \alpha_1^j \cdot f_{\theta_{L-j}}^{L-j}(\mathbf{z}_i^{L-j-1}) \right).
 \end{aligned}$$

Now, substituting this expression into the fixed-point update gives

$$\begin{aligned}
 \mathbf{z}_{i+1}^0 &= \alpha_2 \cdot \mathbf{z}_i^L + \beta_2 \cdot \mathbf{x} \\
 &= \alpha_2 \cdot \left( \alpha_1^L \cdot \mathbf{z}_i^0 + \beta_1 \cdot \sum_{j=0}^{L-1} \alpha_1^j \cdot f_{\theta_{L-j}}^{L-j}(\mathbf{z}_i^{L-j-1}) \right) + \beta_2 \cdot \mathbf{x} \\
 &= \alpha_2 \alpha_1^L \cdot \mathbf{z}_i^0 + \alpha_2 \beta_1 \cdot \left( \sum_{j=0}^{L-1} \alpha_1^j \cdot f_{\theta_{L-j}}^{L-j}(\mathbf{z}_i^{L-j-1}) \right) + \beta_2 \cdot \mathbf{x}.
 \end{aligned}$$

For compactness, define  $\rho = \alpha_2 \alpha_1^L$  and  $\mathbf{s}_i = \sum_{j=0}^{L-1} \alpha_1^j \cdot f_{\theta_{L-j}}^{L-j}(\mathbf{z}_i^{L-j-1})$ . Then the fixed-point iteration can be written as  $\mathbf{z}_{i+1}^0 = \rho \cdot \mathbf{z}_i^0 + \alpha_2 \beta_1 \cdot \mathbf{s}_i + \beta_2 \cdot \mathbf{x}$ . Unrolling this recursion over fixed-point iterations gives

$$\mathbf{z}_{i+1}^0 = \rho^{i+1} \cdot \mathbf{z}_0^0 + \beta_2 \left( \sum_{k=0}^i \rho^k \right) \cdot \mathbf{x} + \alpha_2 \beta_1 \left( \sum_{k=0}^i \rho^k \cdot \mathbf{s}_{i-k} \right). \quad (4)$$

Since  $0 \leq \alpha_1, \alpha_2 < 1$ , we have  $0 \leq \rho < 1$ . Therefore, the geometric series is convergent. Taking norms and using the boundedness of each layer map, we get

$$\|\mathbf{z}_{i+1}^0\| \leq \rho^{i+1} \|\mathbf{z}_0^0\| + \beta_2 \left( \sum_{k=0}^i \rho^k \right) \|\mathbf{x}\| + \alpha_2 \beta_1 \left( \sum_{k=0}^i \rho^k \right) \left( \sum_{j=0}^{L-1} \alpha_1^j \right) c_f. \quad (5)$$

Letting  $i \rightarrow \infty$ , the first term vanishes and the two geometric sums converge, which gives

$$\limsup_{i \rightarrow \infty} \|\mathbf{z}_i^0\| \leq \frac{\beta_2}{1 - \rho} \|\mathbf{x}\| + \frac{\alpha_2 \beta_1}{1 - \rho} \left( \frac{1 - \alpha_1^L}{1 - \alpha_1} \right) c_f. \quad (6)$$

Substituting back  $\rho = \alpha_2 \alpha_1^L$ , we obtain

$$\limsup_{i \rightarrow \infty} \|\mathbf{z}_i^0\| \leq \frac{\beta_2}{1 - \alpha_2 \alpha_1^L} \|\mathbf{x}\| + \frac{\alpha_2 \beta_1 (1 - \alpha_1^L)}{(1 - \alpha_2 \alpha_1^L) (1 - \alpha_1)} c_f. \quad (7)$$

We now set  $\beta_2 = 1 - \alpha_2 \alpha_1^L$ . This makes the coefficient of  $\|\mathbf{x}\|$  equal to 1. Furthermore, setting  $\beta_1 = \frac{\beta_2 (1 - \alpha_1)}{(1 - \alpha_1^L)}$  makes the coefficient of  $c_f$  equal to  $\alpha_2$ . Therefore,

$$\limsup_{i \rightarrow \infty} \|\mathbf{z}_i^0\| \leq \|\mathbf{x}\| + \alpha_2 \cdot c_f. \quad (8)$$

In particular, if the fixed-point iteration converges to  $\mathbf{z}_\infty^0$ , then

$$\|\mathbf{z}_\infty^0\| \leq \|\mathbf{x}\| + \alpha_2 \cdot c_f. \quad (9)$$

This completes the proof. ■

## B.2. Contractive mapping

**Theorem 2 (Small  $\alpha_2$  implies convergence)** *Let  $\tilde{f}_\theta(\mathbf{u})$  denote one full pass through the  $L$  layers defined by Equation 2, and define the iteration map*

$$T_{\alpha_2}(\mathbf{z}; \mathbf{x}) := \tilde{f}_\theta(\alpha_2 \mathbf{z} + \beta_2 \mathbf{x}).$$

*If  $\tilde{f}_\theta$  is Lipschitz with constant  $\lambda_f$  and  $\alpha_2 \lambda_f < 1$ , then  $T_{\alpha_2}(\cdot; \mathbf{x})$  is a contraction, and the fixed-point iteration  $\mathbf{z}_{i+1} = T_{\alpha_2}(\mathbf{z}_i; \mathbf{x})$  converges to a unique  $\mathbf{z}^*$  at a linear rate:*

$$\|T_{\alpha_2}(\mathbf{z}_i; \mathbf{x}) - \mathbf{z}_i\| \leq (\alpha_2 \lambda_f)^i \|T_{\alpha_2}(\mathbf{z}_0; \mathbf{x}) - \mathbf{z}_0\|.$$

*For  $\mathbf{z}_0 = 0$ , the right-hand side specialises to  $(\alpha_2 \lambda_f)^i \|\tilde{f}_\theta(\beta_2 \mathbf{x})\|$ . The proof is in Section B.2*

**Proof** We first show that  $T_{\alpha_2}(\cdot; \mathbf{x})$  is a contraction with respect to  $\mathbf{z}$ . For any  $\mathbf{z}, \mathbf{z}'$ , using the Lipschitzness of  $\tilde{f}_\theta(\cdot)$ , we have

$$\begin{aligned} \|T_{\alpha_2}(\mathbf{z}; \mathbf{x}) - T_{\alpha_2}(\mathbf{z}'; \mathbf{x})\| &= \left\| \tilde{f}_\theta(\alpha_2 \cdot \mathbf{z} + \beta_2 \cdot \mathbf{x}) - \tilde{f}_\theta(\alpha_2 \cdot \mathbf{z}' + \beta_2 \cdot \mathbf{x}) \right\| \\ &\leq \lambda_f \|\alpha_2 \cdot \mathbf{z} + \beta_2 \cdot \mathbf{x} - \alpha_2 \cdot \mathbf{z}' - \beta_2 \cdot \mathbf{x}\| \\ &= \alpha_2 \lambda_f \|\mathbf{z} - \mathbf{z}'\|. \end{aligned}$$

Therefore, if  $0 \leq \alpha_2 \lambda_f < 1$ , the map  $T_{\alpha_2}(\cdot; \mathbf{x})$  is strictly contractive. By the Banach fixed-point theorem, it has a unique fixed point  $\mathbf{z}^*$ , and the iteration  $\mathbf{z}_{i+1} = T_{\alpha_2}(\mathbf{z}_i; \mathbf{x})$  converges to  $\mathbf{z}^*$ .

We now prove the residual bound. Since  $\mathbf{z}_{i+1} = T_{\alpha_2}(\mathbf{z}_i; \mathbf{x})$ , the residual at iteration  $i$  can be written as

$$\|T_{\alpha_2}(\mathbf{z}_i; \mathbf{x}) - \mathbf{z}_i\| = \|\mathbf{z}_{i+1} - \mathbf{z}_i\|.$$

Using the contraction property of  $T_{\alpha_2}(\cdot; \mathbf{x})$ , we get

$$\begin{aligned} \|\mathbf{z}_{i+1} - \mathbf{z}_i\| &= \|T_{\alpha_2}(\mathbf{z}_i; \mathbf{x}) - T_{\alpha_2}(\mathbf{z}_{i-1}; \mathbf{x})\| \\ &\leq \alpha_2 \lambda_f \|\mathbf{z}_i - \mathbf{z}_{i-1}\|. \end{aligned}$$

Applying this inequality recursively gives

$$\|\mathbf{z}_{i+1} - \mathbf{z}_i\| \leq (\alpha_2 \lambda_f)^i \|\mathbf{z}_1 - \mathbf{z}_0\|.$$

Since  $\mathbf{z}_1 = T_{\alpha_2}(\mathbf{z}_0; \mathbf{x})$ , we obtain

$$\|T_{\alpha_2}(\mathbf{z}_i; \mathbf{x}) - \mathbf{z}_i\| \leq (\alpha_2 \lambda_f)^i \|T_{\alpha_2}(\mathbf{z}_0; \mathbf{x}) - \mathbf{z}_0\|.$$

Finally, if  $\mathbf{z}_0 = 0$ , then  $T_{\alpha_2}(\mathbf{z}_0; \mathbf{x}) = \tilde{f}_\theta(\beta_2 \cdot \mathbf{x})$ , and therefore

$$\|T_{\alpha_2}(\mathbf{z}_i; \mathbf{x}) - \mathbf{z}_i\| \leq (\alpha_2 \lambda_f)^i \left\| \tilde{f}_\theta(\beta_2 \cdot \mathbf{x}) \right\|.$$

This completes the proof. ■

### Appendix C. Proof of Theorem 4

**Proof** Since  $\|\mathbf{J}\|_2 = \sigma < 1$ , the Neumann series is convergent, and we have  $(\mathbf{I} - \mathbf{J})^{-1} = \sum_{j=0}^{\infty} \mathbf{J}^j$ . Therefore, the error of the  $k$ -term truncated approximation is

$$\begin{aligned} \left\| (\mathbf{I} - \mathbf{J})^{-1} - \sum_{j=0}^{k-1} \mathbf{J}^j \right\|_F &= \left\| \sum_{j=k}^{\infty} \mathbf{J}^j \right\|_F \\ &\leq \sum_{j=k}^{\infty} \|\mathbf{J}^j\|_F. \end{aligned}$$

Using the relation  $\|\mathbf{A}\|_F \leq \sqrt{D} \|\mathbf{A}\|_2$  for  $\mathbf{A} \in \mathbb{R}^{D \times D}$ , together with submultiplicativity of the spectral norm, we get

$$\begin{aligned} \|\mathbf{J}^j\|_F &\leq \sqrt{D} \|\mathbf{J}^j\|_2 \\ &\leq \sqrt{D} \|\mathbf{J}\|_2^j \\ &= \sqrt{D} \cdot \sigma^j. \end{aligned}$$

Substituting this into the previous inequality gives

$$\begin{aligned} \left\| (\mathbf{I} - \mathbf{J})^{-1} - \sum_{j=0}^{k-1} \mathbf{J}^j \right\|_F &\leq \sqrt{D} \sum_{j=k}^{\infty} \sigma^j \\ &= \sqrt{D} \cdot \frac{\sigma^k}{1 - \sigma}. \end{aligned}$$

Thus, the approximation error decays as  $\mathcal{O}(\sigma^k)$ .

To make the corresponding gradient statement explicit, let  $\mathbf{P} = \frac{\partial f_\theta}{\partial \theta}(\mathbf{z}^*; \mathbf{x})$  and  $\boldsymbol{\delta} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^*}$ . The exact implicit gradient is  $\nabla_\theta \mathcal{L} = \mathbf{P}^\top (\mathbf{I} - \mathbf{J})^{-\top} \boldsymbol{\delta}$ , whereas the  $k$ -step truncated BPTT gradient is  $\widehat{\nabla}_\theta^{(k)} \mathcal{L} = \mathbf{P}^\top \left( \sum_{j=0}^{k-1} (\mathbf{J}^\top)^j \right) \boldsymbol{\delta}$ . Therefore,

$$\begin{aligned} \left\| \nabla_\theta \mathcal{L} - \widehat{\nabla}_\theta^{(k)} \mathcal{L} \right\|_2 &\leq \|\mathbf{P}\|_2 \left\| (\mathbf{I} - \mathbf{J})^{-\top} - \sum_{j=0}^{k-1} (\mathbf{J}^\top)^j \right\|_2 \|\boldsymbol{\delta}\|_2 \\ &\leq \|\mathbf{P}\|_2 \left( \sum_{j=k}^{\infty} \|\mathbf{J}\|_2^j \right) \|\boldsymbol{\delta}\|_2 \\ &= \|\mathbf{P}\|_2 \frac{\sigma^k}{1 - \sigma} \|\boldsymbol{\delta}\|_2. \end{aligned}$$

Hence, the truncated BPTT gradient error also decays exponentially with the number of backward passes  $k$ . This completes the proof.  $\blacksquare$

### C.1. Damping of fixed-point iterations

**Theorem 3 (Damping stabilizes oscillatory fixed-point dynamics)** *Suppose  $f_\theta(\cdot; \mathbf{x})$  is continuously differentiable in a neighborhood of a fixed point  $\mathbf{z}^*$ , and that every eigenvalue  $\lambda_i$  of the Jacobian  $\mathbf{J}$  at  $\mathbf{z}^*$  satisfies  $\Re(\lambda_i) < 1$ . Define the damped iteration map  $g_{\eta, \theta}(\mathbf{z}; \mathbf{x}) := \eta f_\theta(\mathbf{z}; \mathbf{x}) + (1 - \eta) \mathbf{z}$ . Then there exists  $\eta_0 \in (0, 1)$  such that, for every  $\eta \in (0, \eta_0)$ , the iteration  $\mathbf{z}_{i+1} = g_{\eta, \theta}(\mathbf{z}_i; \mathbf{x})$  converges locally to  $\mathbf{z}^*$ . Moreover,  $g_{\eta, \theta}(\cdot; \mathbf{x})$  and  $f_\theta(\cdot; \mathbf{x})$  have the same fixed points.*

The proof is in Section C.1.

**Proof** We first show that the fixed points of  $g_{\eta, \theta}(\cdot; \mathbf{x})$  and  $f_\theta(\cdot; \mathbf{x})$  coincide. Since

$$g_{\eta, \theta}(\mathbf{z}; \mathbf{x}) - \mathbf{z} = \eta (f_\theta(\mathbf{z}; \mathbf{x}) - \mathbf{z}),$$

and  $\eta > 0$ , we have  $g_{\eta, \theta}(\mathbf{z}; \mathbf{x}) = \mathbf{z}$  if and only if  $f_\theta(\mathbf{z}; \mathbf{x}) = \mathbf{z}$ .

We now study the local stability of the damped iteration around  $\mathbf{z}^*$ . The Jacobian of  $g_{\eta, \theta}(\cdot; \mathbf{x})$  at  $\mathbf{z}^*$  is

$$\frac{\partial g_{\eta, \theta}}{\partial \mathbf{z}}(\mathbf{z}^*; \mathbf{x}) = (1 - \eta) \mathbf{I} + \eta \mathbf{J},$$

so for every eigenvalue  $\lambda_i$  of  $\mathbf{J}$ , the corresponding eigenvalue of the damped Jacobian is

$$\mu_i(\eta) = 1 - \eta + \eta \lambda_i = 1 + \eta(\lambda_i - 1).$$

Local asymptotic stability is implied by  $|\mu_i(\eta)| < 1$  for all  $i$ . Writing  $\lambda_i = a_i + i b_i$  with  $a_i = \Re(\lambda_i)$ ,

$$|\mu_i(\eta)|^2 = 1 + 2\eta(a_i - 1) + \eta^2 |\lambda_i - 1|^2.$$

Hence  $|\mu_i(\eta)| < 1$  if and only if  $\eta |\lambda_i - 1|^2 < 2(1 - a_i)$ , i.e.,

$$0 < \eta < \frac{2(1 - \Re(\lambda_i))}{|\lambda_i - 1|^2}.$$

By assumption  $\Re(\lambda_i) < 1$  for every  $i$ , so each upper bound is strictly positive. Setting

$$\eta_0 = \min \left\{ 1, \min_i \frac{2(1 - \Re(\lambda_i))}{|\lambda_i - 1|^2} \right\} > 0,$$

we obtain  $|\mu_i(\eta)| < 1$  for every  $i$  and every  $\eta \in (0, \eta_0)$ . Therefore  $\mathbf{z}^*$  is locally asymptotically stable under the damped iteration  $\mathbf{z}_{t+1} = g_{\eta, \theta}(\mathbf{z}_t; \mathbf{x})$ , and the iterates converge to  $\mathbf{z}^*$  from any sufficiently close initialization.  $\blacksquare$

## C.2. A toy failure mode for recurrent PostLN

Figure 4 gives a minimal version of the normalization tradeoff discussed in the main text. We take random  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{100 \times 2}$ , set  $\mathbf{z}^0 = \mathbf{x}$ , and tie the same rank-one map for  $L = 20$  iterations,

$$W(\mathbf{w}) = \mathbf{w}\mathbf{1}^\top, \quad \mathbf{w} = (w_1, w_2)^\top.$$

This choice allows 2d visualization while making the normalization operation non-trivial. Here  $\mathbf{N}$  is row-wise  $\ell_2$  normalization,

$$\mathbf{N}(\mathbf{z})_{i,:} = \frac{\mathbf{z}_{i,:}}{\|\mathbf{z}_{i,:}\|_2 + \varepsilon}.$$

The post-norm recurrence is

$$\mathbf{z}_{\text{post}}^{\ell+1} = \mathbf{N} \left( \mathbf{z}_{\text{post}}^\ell + \mathbf{z}_{\text{post}}^\ell W(\mathbf{w}) \right),$$

while the scaled pre-norm recurrence is

$$\mathbf{z}_{\text{pre}}^{\ell+1} = (1 - \beta)\mathbf{z}_{\text{pre}}^\ell + \beta \mathbf{N} \left( \mathbf{z}_{\text{pre}}^\ell W(\mathbf{w}) \right), \quad \beta = \frac{1}{2}.$$

For both models we sweep a  $200 \times 200$  grid over  $[-5, 5]^2$  and plot

$$\Delta \mathcal{L}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) - \min_{\mathbf{v} \in \mathcal{G}} \mathcal{L}(\mathbf{v}), \quad \mathcal{L}(\mathbf{w}) = \frac{1}{nd} \|\mathbf{z}^L(\mathbf{w}) - \mathbf{y}\|_F^2.$$

The figure illustrates that **boundedness does not imply trainability**. Post-norm keeps the recurrent state bounded by construction: after every step, each row is projected back to the unit sphere. However, the same projection also removes radial information at every iteration. In this two-parameter slice, the resulting loss is organized into thin angular sectors with sharp ridges and narrow low-loss regions. Thus a random initialization of  $\mathbf{w}$  is likely to start in a bounded but poorly conditioned part of the landscape, where the gradient does not point into a useful basin. This is the toy analogue of the optimization difficulty of recurrent post-norm blocks.

The right panel is not bare pre-norm; it is pre-norm with residual scaling. This matters because naive pre-norm removes the projection that controls the recurrent state and can lead to activation growth, as discussed above. With the scaled update, each row satisfies

$$\|\mathbf{z}_{\text{pre},i}^{\ell+1}\|_2 \leq (1 - \beta)\|\mathbf{z}_{\text{pre},i}^\ell\|_2 + \beta,$$

so the toy dynamics remain bounded while preserving a live residual stream. The broader low-loss region in Figure 4 is therefore consistent with the architecture we use in Theorem 1: pre-normalization improves signal propagation, while residual scaling replaces the boundedness mechanism that post-normalization provided.

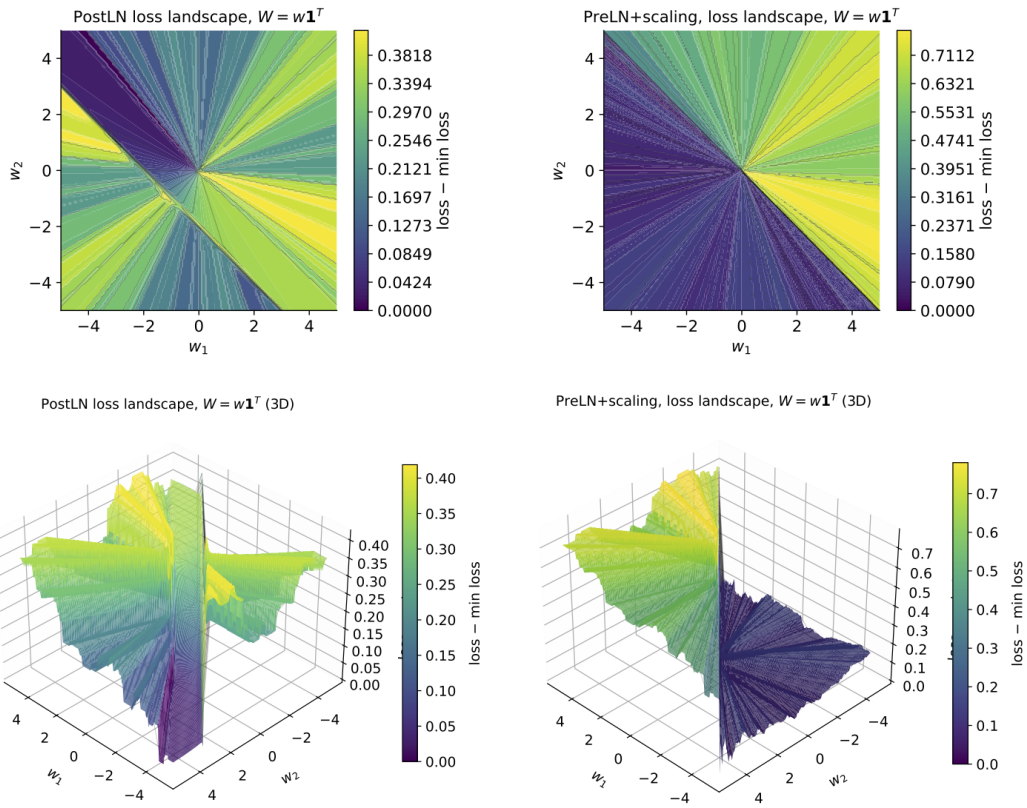


Figure 4: Landscape visualization for the setup proposed in Section C.2

---

**Algorithm 1** Fixed-point optimizer FPOPT: one damped step with patience-based decay
 

---

**Input:** initial damping  $\eta_0$ ; decay  $\gamma \in (0, 1)$ ; patience  $P$ 
**Internal state:**  $\eta \leftarrow \eta_0$ ,  $p \leftarrow P$ ,  $r^* \leftarrow \infty$ 
**Procedure** STEP( $\mathbf{z}, \tilde{\mathbf{z}}$ ):

```

     $r \leftarrow \|\mathbf{z} - \tilde{\mathbf{z}}\|_\infty / (\|\tilde{\mathbf{z}}\|_\infty + \epsilon)$ ; // residual
     $\mathbf{z} \leftarrow \eta \tilde{\mathbf{z}} + (1 - \eta) \mathbf{z}$ ; // damped update
    if  $r < r^*$  then
         $r^* \leftarrow r$ ,  $p \leftarrow P$ ; // progress: reset
    else
         $p \leftarrow p - 1$  if  $p = 0$  then
             $\eta \leftarrow \gamma \eta$ ,  $p \leftarrow P$ ; // decay  $\eta$ 
        end
    end
    return  $\mathbf{z}, r$ 

```

---

**Algorithm 2** FPRM training loop with truncated BPTT and deep supervision
 

---

**Input:** Model  $f_\theta$ ; prediction head  $h_\phi$ ; fixed-point optimiser FPOPT; model optimiser MODELOPT; input  $\mathbf{x}$ ; target  $\mathbf{y}$ ; BPTT depth  $K$ ; initial state  $\mathbf{z}_0$ 

```

 $\mathbf{z} \leftarrow \mathbf{z}_0$  while FPOPT.CONT() do // outer loop
    for  $k = 1, \dots, K$  do // BPTT window
         $\mathbf{y}_k \leftarrow f_\theta(\mathbf{z}; \mathbf{x})$   $\mathbf{z} \leftarrow \text{FPOPT.STEP}(\mathbf{z}, \mathbf{y}_k)$ 
    end
     $\hat{\mathbf{y}} \leftarrow h_\phi(\mathbf{z})$ ; // deep supervision
     $\mathcal{L} \leftarrow \text{CROSSENTROPY}(\hat{\mathbf{y}}, \mathbf{y})$   $\text{MODELOPT.BACKWARD}(\mathcal{L})$   $\mathbf{z} \leftarrow \text{detach}(\mathbf{z})$ 
end

```

---

## Appendix D. Further Details About the Architecture

**Optimization of fixed-point models.** One advantage of contractive fixed-point models is that they can be trained using truncated back-propagation through time (BPTT) [12]. Specifically, following the implicit function theorem we can write the gradient w.r.t. the parameters of the model as [2]:

$$\frac{\partial \mathbf{z}^*}{\partial \theta} = \left( \mathbf{I} - \frac{\partial f_\theta}{\partial \mathbf{z}}(\mathbf{z}^*; \mathbf{x}) \right)^{-1} \frac{\partial f_\theta}{\partial \theta}(\mathbf{z}^*; \mathbf{x}). \quad (10)$$

As the contractiveness entails a contractive Jacobian, the Neumann series corresponding to the first term on the RHS of Equation (10) is converging, allowing us to estimate the gradient as:

$$\frac{d\mathbf{z}^*}{d\theta} \approx \sum_{j=0}^{k-1} \mathbf{J}^j \frac{\partial f_\theta}{\partial \theta}(\mathbf{z}^*; \mathbf{x}), \quad (11)$$

which in our setting translates into doing  $k$  backward passes after we reach the fixed-point. The following proposition shows that truncated BPTT becomes exponentially accurate as an estimation of the gradient of fixed-point models:

**Proposition 4 (Exponential decay of truncated-BPTT error)**

Let  $\mathbf{z}^* = f_\theta(\mathbf{z}^*; \mathbf{x})$  and  $\mathbf{J} = \frac{\partial f_\theta}{\partial \mathbf{z}}(\mathbf{z}^*; \mathbf{x}) \in \mathbb{R}^{D \times D}$ . If  $\mathbf{J}$  is contractive in spectral norm,  $\|\mathbf{J}\|_2 = \sigma < 1$ , then for every  $k \geq 0$ ,

$$\left\| (\mathbf{I} - \mathbf{J})^{-1} - \sum_{j=0}^{k-1} \mathbf{J}^j \right\|_F \leq \sqrt{D} \frac{\sigma^k}{1 - \sigma}.$$

The proof is in Section C. An essentially equivalent result appears in the proof of Theorem 2 of [12].

Theorem 4 allows for a fixed memory footprint during training, essentially decoupling the number of loops from the memory complexity of the model. In the same spirit, HRM [29] and TRM [15] approximate the gradient with a small number of backward passes at the fixed-point, although TRM argues that the fixed-point condition is unnecessary in practice. However, we note that reaching fixed-points in this paper is not strictly done for the purpose of optimization, but rather as a halting mechanism to replace ACT. In the following, we provide a full description of the proposed method.

**Fixed-point solver.** Let  $\mathbf{z}^{(k)} \in \mathbb{R}^{B \times T \times d}$  denote the iterate at step  $k$  of the solver (with  $B$  batch,  $T$  sequence,  $d$  hidden), and  $\mathbf{z}^{(k+1)}$  the proposed update. Convergence is measured per sample  $\mathbf{z}_b^{(k)}$  by the relative  $\ell_\infty$  residual,

$$\mathbf{r}_b^{(k)} = \frac{\|\mathbf{z}_b^{(k+1)} - \mathbf{z}_b^{(k)}\|_\infty}{\|\mathbf{z}_b^{(k+1)}\|_\infty + \epsilon}, \mathbf{r}_b^{(k)} \in \mathbb{R}.$$

A sample is declared converged when  $r_b^{(k)} < \tau$ . In practice, we set  $\tau$  to 0.1. During training, the solver terminates as soon as the fraction of converged samples falls below a tolerated quantile  $p$ ,

$$\sum_{b=1}^B \mathbf{1}[r_b^{(k)} \geq \tau] < p \cdot B.$$

Two safeguards bound the loop: a hard cap  $k \leq K_{\max}$  on iterations, and early termination if the adaptive step size collapses below a minimum. At evaluation we set  $p = 0$ , requiring every sample to satisfy  $r_b^{(k)} < \tau$ .

**Deep supervision.** We adopt a similar deep supervision mechanism to HRM [29] and TRM [15]. Let  $T_{\text{sup}}$  denote the deep supervision interval. Every  $T_{\text{sup}}$  iterations the intermediate activations of the model are decoded through the output head, the task loss is computed, and backpropagation is performed through the most recent step only. Then, the activations and the solver state are detached, implementing truncated backpropagation through time (TBPTT). After the outer solver converges, a final iteration is performed with gradients enabled and its output is supervised as well. The number of backward passes per forward pass is therefore  $\lceil K/T_{\text{sup}} \rceil + 1$ , where  $K$  is the number of iterations actually taken by the solver on that batch. Setting  $T_{\text{sup}} = 1$  supervises every intermediate-step, which is the typical configuration.

## Appendix E. Additional Experiments

**Introducing residual scaling to TRM.** Since TRM could potentially benefit from FPRM’s architectural changes, we add them to TRM individually and in combination, measuring test sequence accuracy on Sudoku Extreme. Results are shown in Table 2, reported as the change relative to the TRM baseline (74.32%). No configuration matches the baseline, and the effects are not additive. The optimal configuration for FPRM turned out to be the most detrimental for TRM.

Table 2: Effect of adding FPRM’s architectural modifications to TRM: pre-norm, residual scaling ( $\alpha_1, \alpha_2$ ), and a final post-activation output norm, individually and in combination. Values are the change in test sequence accuracy on Sudoku Extreme (%) relative to the TRM baseline of 74.32%.

Configuration	$\Delta$ Seq. Acc. (%)
Original TRM (post-norm, no scale)	—
+ output norm	−0.47
+ residual scaling ( $\alpha_1$ only)	−13.96
+ residual scaling ( $\alpha_1, \alpha_2$ )	−6.16
+ pre-norm & residual scaling ( $\alpha_1, \alpha_2$ )	− <b>54.55</b>
+ pre-norm, residual scaling, final norm (full FPRM stack)	−9.45

Table 3: Sensitivity of FPRM to residual scaling initialization on Sudoku Extreme dataset. Each cell reports best test sequence accuracy (%) for a given pair of initial values.

$\alpha_1$ init	$\alpha_2$ init		
	0.25	0.50	0.75
0.25	67.04	67.41	67.72
0.50	73.55	70.99	70.88
0.75	<b>74.91</b>	72.51	71.60

**Loop-utilization.** In Section 3.2, we showed that post-norm and pre-norm transformer architectures suffer from training instability as the model reaches greater effective depth, which translates into an inability to solve tasks that require scaling effective depth. In this section, we switch to a class of models that can adapt their compute budget at inference time. Specifically, we train FPRM with pre-norm and residual scaling, and FPRM with post-norm, both on the Sudoku task until their fixed-point halting mechanisms stop them. In Figure 5(a)subfigure, we demonstrate that test-time compute scaling improves performance for both types of normalization. However, once the fixed point is reached, additional compute no longer improves performance. We further show that the gap between the pre-norm model with residual scaling and the post-norm model in Figure 5(a)subfigure could arise from underutilization of effective depth in the post-norm FPRM. This manifests as faster performance saturation at a lower value. The results in Figure 5 further support our hypothesis about signal-propagation problems, justifying the modifications introduced in this paper.

## Appendix F. Additional Experimental Details

**Weight initialization.** It seems like initializing the weights using a truncated normal distribution (LeCun initialization) is common practice in looped architectures. In our experiments, it accelerates the convergence but there’s very little material difference in sequence accuracy after convergence.

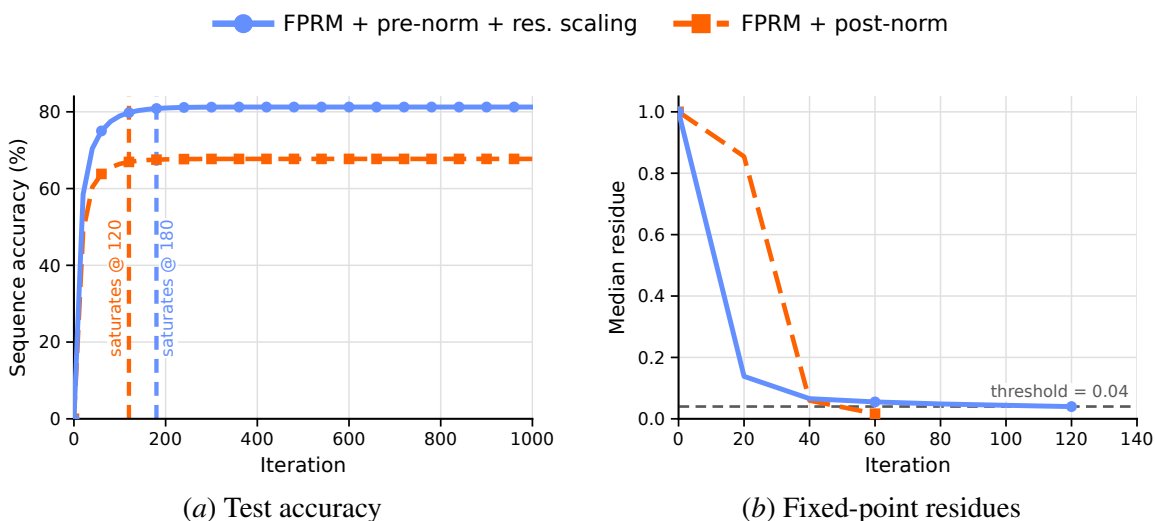


Figure 5: **Loop-utilization of FPRM on Sudoku.** (a) test accuracy of FPRM with pre-norm with residual scales vs. post-norm. (b) median residue. The pre-norm model is better at loop utilization, while both have similar residues. This indicates similar latent-space convergence, with more meaningful updates in the pre-norm variant, resulting in improved performance.

**Grokking.** There is some evidence for grokking in looped architectures, but in maze we observe convergence on the train data. And training the models for a longer period (up to 7 days) did not yield better performance.

**Hyperparameters, device specification and estimated GPU hours** Based on our approximation, around 1k GPU-hours produced the numbers in the paper, out of 14k GPU-hours burned overall in the project. Next, we provide the values for some of the most important hyperparameters in the paper, per each model and dataset.

Table 4: Hyperparameters for Sudoku-Extreme experiments (Table 1 of the paper). Shared across all models:  $1 \times A100$ -40GB, batch 768, 60 000 epochs, ADAM\_ATAN2 optimizer, lr=  $10^{-4}$  (constant after 2 000-step warm-up), weight-decay 1.0, EMA enabled, puzzle-embedding length 16.

	HRM	TRM	FPRM
<i>Looping structure</i>			
<i>H</i> -cycles	2	3	3
<i>L</i> -cycles	2	6	6
<i>H</i> -layers	4	0	0
<i>L</i> -layers	4	2	2
$n_{\text{back},L}$	–	= <i>L</i> -cycles	4
<i>Halting</i>			
mechanism	ACT	ACT	fixed-point
halt_max_steps	16	16	–
max_iter	–	–	1000
iter. distribution	–	–	exponential ( $s=16$ )
<i>Block / signal-prop modifications</i>			
norm type	post-norm	post-norm	pre-norm
norm placement	–	–	output
residual scaling	–	–	input-independent
$\alpha_1, \alpha_2$ init	–	–	0.5, 0.5
conv branch	–	none	conv2d ( $k=3$ )

Table 5: Hyperparameters for Maze-Hard experiments (Table 1 of the paper). HRM uses the published configuration; TRM and FPRM both train on maze-30x30-hard-1k with 4x A100-80GB, 60 000 epochs, lr=  $10^{-4}$  constant, weight-decay 1.0, EMA enabled, puzzle-embedding length 16.

	HRM	TRM	FPRM
optimizer	ADAM_ATAN2	ADAM_ATAN2	ADAM_ATAN2
batch size	768	512	768
lr warm-up steps	2000	0	2000
<i>Looping structure</i>			
<i>H</i> -cycles	2	3	3
<i>L</i> -cycles	2	4	6
<i>H</i> -layers	4	0	0
<i>L</i> -layers	4	2	2
$n_{\text{back},L}$	–	= <i>L</i> -cycles	4
<i>Halting</i>			
mechanism	ACT	ACT	fixed-point
halt_max_steps	16	16	–
max_iter	–	–	1000
iter. distribution	–	–	exponential ( $s=16$ )
<i>Block / signal-prop modifications</i>			
norm type	post-norm	post-norm	pre-norm
norm placement	–	–	output
residual scaling	–	–	input-independent
$\alpha_1, \alpha_2$ init	–	–	0.75, 0.5
conv branch	–	none	conv2d ( $k=3$ )

Table 6: Hyperparameters for state-tracking experiments on  $A_5$  (Figure 3 of the paper) and for the Universal Transformer signal-propagation analysis (Figure 2). All runs use  $1 \times A100$ -80GB, global batch 1024, ADAM\_ATAN2, lr-warm-up 0, weight-decay  $10^{-2}$ , EMA disabled, no puzzle embedding (`puzzle_emb_len=0`). Trained at  $k_{\text{train}}=32$ , evaluated for  $k \in [2, 128]$ .

	TRM	TRM+ACT	FPRM	UT (Fig. 2)
epochs	50	50	20	30
learning rate	$10^{-4}$	$10^{-4}$	$10^{-3}$	$10^{-4}$
<i>Looping structure</i>				
$H$ -cycles	2	2	3	3
$L$ -cycles	4	4	6	6
$L$ -layers	4	4	2	2
$n_{\text{back},L}$	= $L$ -cycles	= $L$ -cycles	4	4
<i>Halting</i>				
mechanism	fixed iters	ACT	fixed-point	fixed iters
halt_max_steps	16	16	–	–
max_iter	–	–	128	= $k_{\text{train}}$
iter. distribution	–	–	deterministic	deterministic
<i>Block / signal-prop modifications</i>				
norm type	post-norm	post-norm	pre-norm	sweep <sup>†</sup>
norm placement	–	–	none	none
residual scaling	–	–	input-indep.	sweep <sup>†</sup>
$\alpha_1, \alpha_2$ init	–	–	0.5, 0.5	sweep <sup>†</sup>
spec.-norm linear	–	–	yes	no
conv branch	none	none	conv1d ( $k=4$ )	none

<sup>†</sup> UT row sweeps the {post-norm, pre-norm, pre-norm + residual-scaling} variants from Figure Figure 1(a)subfigure; the residual-scaling variant uses  $\alpha_1=0.75, \alpha_2=0.5$  and  $k_{\text{train}} \in \{8, 16, 32, 64\}$ .