

When Do Transformers Outperform Feedforward and Recurrent Networks? A Statistical Perspective

Alireza Mousavi-Hosseini

University of Toronto and Vector Institute

MOUSAVI@CS.TORONTO.EDU

Clayton Sanford

Google Research

CHSANFORD@GOOGLE.COM

Denny Wu

New York University and Flatiron Institute

DENNYWU@NYU.EDU

Murat A. Erdogdu

University of Toronto and Vector Institute

ERDOGDU@CS.TORONTO.EDU

Abstract

Theoretical efforts to prove advantages of Transformers in comparison with classical architectures such as feedforward and recurrent neural networks have mostly focused on representational power. In this work, we take an alternative perspective and prove that even with infinite compute, feedforward and recurrent networks may suffer from larger sample complexity compared to Transformers, as the latter can adapt to a form of *dynamic sparsity*. Specifically, we consider a sequence-to-sequence data generating model on sequences of length N , where the output at each position only depends on $q \ll N$ relevant tokens, and the positions of these tokens are described in the input prompt. We prove that a single-layer Transformer can learn this model if and only if its number of attention heads is at least q , in which case it achieves a sample complexity almost independent of N , while recurrent networks require $N^{\Omega(1)}$ samples on the same problem. If we simplify this model, recurrent networks may achieve a complexity almost independent of N , while feedforward networks still require N samples. Our proposed sparse retrieval model illustrates a natural hierarchy in sample complexity across these architectures.

1. Introduction

The theoretical efforts surrounding the success of Transformers [26] have so far demonstrated various capabilities like in-context learning [3, 8, 18, 27, 29, and others] and chain-of-thought prompting along with its benefits [15, 17, 20, 21, and others] in various settings. There are fewer works that provide specific benefits of Transformers in comparison with feedforward and recurrent architectures. On the approximation side, there are tasks that Transformers can solve with size logarithmic in the input, while alternative architectures require polynomial size [24, 25]. Based on these results, [28] showed a separation between Transformers and feedforward networks by providing further optimization guarantees for gradient-based training of Transformers on a sparse token selection task.

While most prior works focused on the approximation separation between Transformers and feedforward networks (FFNs), in this work we focus on a purely statistical separation, and ask:

What function class can Transformers learn with fewer samples compared to feedforward and recurrent networks, even with infinite computational resources?

[16] approached the above problem with random features, where the query-key matrix for the attention and the first layer weights for the two-layer feedforward network were fixed at random initialization. However, this only presents a partial picture, as neural networks can learn a significantly larger class of functions once “feature learning” is allowed, i.e., parameters are trained to adapt to the structure of the underlying task [1, 6, 7, 10, 12, 13, 22].

We evaluate the statistical efficiency of Transformers and alternative architectures by characterizing how the sample complexity depends on the input sequence length. A benign length dependence (e.g., sublinear) signifies the ability to achieve low test error in longer sequences, which intuitively connects to the *length generalization* capability [5]. Our generalization bounds for bounded-norm Transformers — along with our contrasts to RNNs and feedforward neural networks — provide theoretical insights into the statistical advantages of Transformers and lay the foundation for future rigorous investigations of length generalization.

1.1. Our Contributions

We study the q -Sparse Token Regression (q STR) data generating model, a sequence-to-sequence model where the output at every position depends dynamically on a sparse subset of the input tokens. We prove that by employing the attention layer to retrieve relevant tokens at each position, single-layer Transformers can adapt to this dynamic sparsity, and learn q STR with a sample complexity almost independent of the length of input sequence N . On the other hand, we develop a new metric-entropy-based argument to derive norm and parameter-count lower bounds for RNNs that lead to a sample complexity lower bound of order $N^{\Omega(1)}$ for RNNs. Further, we show that RNNs can learn a subset of q STR where the output is a constant sequence, which we call simple- q STR, with a sample complexity polylogarithmic in N . Finally, we develop a lower bound technique for feedforward networks (FFNs) that takes advantage of the fully connected projection of the first layer to obtain a sample complexity lower bound linear in N , even when learning simple- q STR models. The following theorem summarizes our main contributions.

Theorem 1 (Informal) *We have the following hierarchy of statistical efficiency for learning q STR.*

- *A single-layer Transformer with $H \geq q$ heads can learn q STR with sample complexity almost independent of N , and cannot learn q STR when $H < q$ even with infinitely many samples.*
- *RNNs can learn simple- q STR with sample size almost independent of N , but require at least $\Omega(N^c)$ samples for some constant $c > 0$ to learn a generic q STR model, regardless of their size.*
- *Feedforward neural networks, regardless of their size, require $\Omega(Nd)$ samples to learn even simple- q STR models, where d is input token dimension.*

We empirically validate the intuitions from Theorem 1 in Figure 1.

2. Problem Setup

Statistical Model. In this paper, we will focus on the ability of different architectures for learning the following data generating model.

Definition 2 (q -Sparse Token Regression) *Suppose $\mathbf{p}, \mathbf{y} \sim \mathcal{P}$ where*

$$\mathbf{p} = \left(\left(\begin{array}{c} \mathbf{x}_1 \\ \mathbf{t}_1 \end{array} \right), \dots, \left(\begin{array}{c} \mathbf{x}_N \\ \mathbf{t}_N \end{array} \right) \right),$$

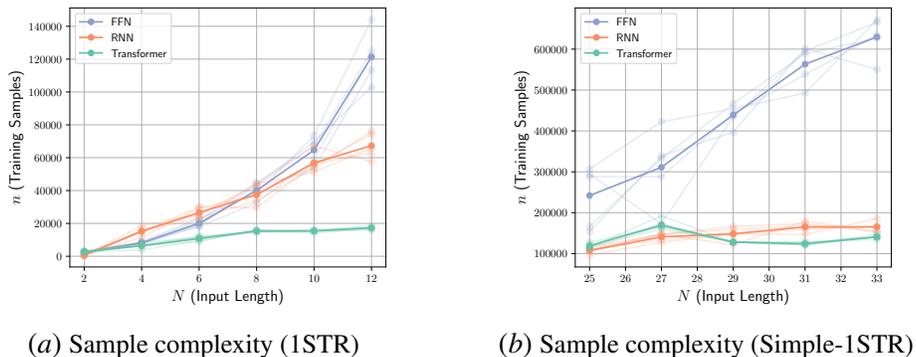


Figure 1: Number of samples to reach a certain test MSE loss threshold while training with online AdamW. We consider (a) 1STR with threshold 0.7 and (b) simple-1STR with threshold 0.02, averaged over 5 experiments. We use a linear link function, standard Gaussian input, $d = 10$ and $d_e = \lceil 5 \log(N) \rceil$. Positional encodings are sampled uniformly from the unit hypercube.

$\mathbf{t}_i \in [N]^q$ and $\mathbf{x}_i \in \mathbb{R}^d$ for $i \in [N]$. In the q -sparse token regression (q STR) data generating model, the output is given by $\mathbf{y} = (y_1, \dots, y_N)^\top \in \mathbb{R}^N$, where

$$y_i = g(\mathbf{x}_{t_{i1}}, \dots, \mathbf{x}_{t_{iq}}),$$

for some $g : \mathbb{R}^{qd} \rightarrow \mathbb{R}$. We call this model simple- q STR if the data distribution is such that $\mathbf{t}_i = \mathbf{t}$ for all $i \in [N]$ and some \mathbf{t} drawn from $[N]^q$.

The above defines a class of sequence-to-sequence functions, where the label at position i in the output sequence depends only on a subsequence of size q of the input data, determined by the set of indices \mathbf{t}_i . \mathbf{p} in the above definition denotes the prompt or context. Given the large context length of modern architectures, we are interested in a setting where $q \ll N$. In this setting, the answer at each position only depends on a few tokens, however the tokens it depends on change based on the context. Therefore, we seek architectures that are *adaptive* to this form of *dynamic sparsity* in the true data generating process, with computational and sample complexity independent of N . As a special case, choosing g as the tokens' mean recovers the *sparse averaging* model proposed in [24], where the authors separate the representational capacity of Transformers and other architectures.

Throughout the paper, we put mild assumptions on the data distribution. Specifically, we assume \mathbf{x} is sub-Gaussian and g grows at most polynomially and is approximated up to L_2 error $\varepsilon_{2\text{NN}}$ by a two-layer feedforward network with width m_g . These assumptions are formalized in Appendix A.

While Empirical Risk Minimization (ERM) is a standard abstract learning algorithm to use for generalization analysis, its standard formalizations use risk functions for scalar-valued predictions. We measure the performance of different architectures in terms of the following population risk

$$R^{\text{arc}}(\Theta) := \frac{1}{N} \mathbb{E} \left[\sum_{j=1}^N (\hat{y}_{\text{arc}}(\mathbf{p}; \Theta)_j - y_j)^2 \right] = \frac{1}{N} \mathbb{E} \left[\|\hat{\mathbf{y}}_{\text{arc}}(\mathbf{p}; \Theta) - \mathbf{y}\|_2^2 \right],$$

where arc denotes a general architecture, and $\hat{\mathbf{y}}_{\text{arc}}(\cdot; \Theta)$ denotes the output of the model parameterized by weight vector Θ . In Appendix A, we formalize several notions of ERM suitable for sequential risk formulations.

3. Transformers

A single-layer Transformer is composed of an attention and a parallel feedforward layer. We consider a standard theoretical formalization of a single-layer transformer with q heads and width m_g for the feedforward units, where m_g is defined in Assumption 2. We formally define the architecture in Appendix B.1. We consider the following parameter class $\Theta_{\text{TR}} = \{\|\text{vec}(\Theta)\|_2 \leq R\}$, and provide a learning guarantee for empirical risk minimizers over Θ_{TR} , with its proof (including the choice of R) deferred to Appendix B.2.

Theorem 3 Let $\hat{\Theta} = \arg \min_{\Theta \in \Theta_{\text{TR}}} \hat{R}_n^{\text{TR}}(\Theta)$. Under Assumptions 1 to 3, we have

$$R^{\text{TR}}(\hat{\Theta}_n) \lesssim \varepsilon_{2\text{NN}} + \tilde{O}\left(C_1 \sqrt{\frac{m_g q (d + q) + q^3 + qd^2}{n}}\right)$$

where $C_1 = R^2 q d$, with probability at least $1 - n^{-c}$ for some absolute constant $c > 0$.

Note that the sample complexity above depends on N only up to log factors. By incorporating additional structure in the ERM solution, it is possible to obtain improved sample complexities. A close study of the optimization dynamics may reveal such additional structure in the solution reached by gradient-based methods, pushing the sample complexity closer to the information-theoretic limit of $\Omega(qd)$. Figure 2 demonstrates that the attention weights achieved through standard optimization of a Transformer match our theoretical constructions – see Equation (6) – even while maintaining separate \mathbf{W}_Q and \mathbf{W}_K during training (we use the 1STR setup of Figure 1 with $N = 100$). We leave the study of optimization dynamics and the resulting sample complexity for future work.

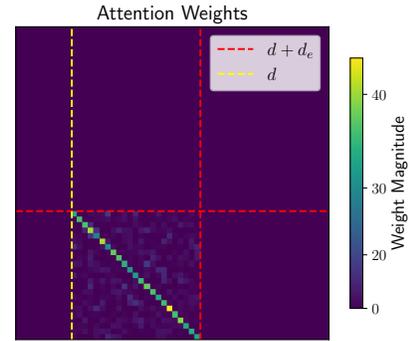


Figure 2: Trained attention weights match our theoretical construction (6).

3.1. Limitations of Transformers with Few Heads

In this section, we will demonstrate that $H \geq \Omega(q)$ is required to learn q STR models, even from a pure approximation perspective, i.e. with access to population distribution. In contrast to [4], we do not put any assumptions on the rank of the key-query projections, i.e. our lower bound applies even when the key-query projection matrix is full-rank.

Proposition 4 Consider a q STR model where $y_i = \frac{1}{\sqrt{qd}} \sum_{j=1}^q (\|\mathbf{x}_{t_{ij}}\|^2 - \mathbb{E}[\|\mathbf{x}_{t_{ij}}\|^2])$, $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma_i)$ such that $\Sigma_i = \mathbf{I}_d$ for $i < N/2$ and $\Sigma_i = 0$ for $i \geq N/2$. Then, there exists a distribution over $(\mathbf{t}_i)_{i \in [N]}$ such that for any choice of Θ_{TR} (including arbitrary $\{\mathbf{W}_{\text{QK}}^{(h)}\}_{h \in [H]}$), we have

$$\frac{1}{N} \mathbb{E} \left[\|\mathbf{y} - \hat{\mathbf{y}}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})\|_2^2 \right] \geq 1 - \frac{(q + d)H}{qd}.$$

Remark 5 We highlight the importance of the nonlinear dependence of y_i on \mathbf{x} for the above lower bound. In particular, for the sparse token averaging task introduced in [24], a single-head attention layer with a carefully constructed embedding suffices for approximation.

The above proposition implies that given sufficiently large dimensionality $d \gg q$, approximation alone necessitates at least $H = \Omega(q)$ heads. In Appendix B.3, we present the proof of Proposition 4, along with Proposition 20 which establishes an exact lower bound $H \geq q$ for all $d \geq 1$, at the expense of additional restrictions on the query-key projection matrix.

4. Recurrent Neural Networks

In this section, we first provide positive results for RNNs by proving that they can learn simple- q STR with a sample complexity only polylogarithmic in N , thus establishing a separation in their learning capability from feedforward networks. For this upper bound, we use bidirectional RNNs with deep transitions [23], formally introduced in Appendix C.1.

Theorem 6 (RNNs can learn simple- q STR) *Let $\hat{\Theta} = \arg \min_{\Theta \in \Theta_{\text{RNN}}} \hat{R}_n^{\text{RNN}}(\Theta)$ (with Θ_{RNN} defined in Equation (10)). Suppose Assumptions 1 to 3 hold with the simple- q STR model, i.e. $\mathbf{t}_i = \mathbf{t}$ for all $i \in [N]$ and some \mathbf{t} drawn from $[N]^q$. Then, with proper hyperparameters in Θ_{RNN} (see Appendix C.1), we obtain*

$$R^{\text{RNN}}(\hat{\Theta}) \lesssim \varepsilon_{2\text{NN}} + \sqrt{\frac{\text{poly}(d, q, m_g, \varepsilon_{2\text{NN}}^{-1}, \log(nN))}{n}},$$

with probability at least $1 - n^{-c}$ for some absolute constant $c > 0$.

As desired, the above sample complexity depends on N only up to polylogarithmic factors. The completed proof can be found in Appendix C.

Next, we turn to general q STR, where we provide a negative result on RNNs, proving that to learn such models their sample complexity must scale with $N^{\Omega(1)}$ regardless of model size, making them less statistically efficient than Transformers. Our lower bound covers a broad notion of bidirectional RNNs formalized in Appendix C.5, and includes the example in the upper bound.

Theorem 7 (RNNs can not learn q STR) *Consider the 1STR model where $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_{Nd})$ with a linear link function, i.e. $y_j = \langle \mathbf{u}, \mathbf{x}_{t_j} \rangle$ for some $\mathbf{u} \in \mathbb{S}^{d-1}$. Further, t_i is drawn independently from the rest of the prompt and uniformly from $[N]$ for all $i \in [N]$. Let $\hat{\Theta}_\varepsilon$ be the min-norm ε -ERM of \hat{R}_n^{RNN} , defined in (4). Then, there exist absolute constants $c_1, c_2, c_3 > 0$ such that if $n \leq \mathcal{O}(N^{c_1})$, for any $\varepsilon \geq 0$, with probability at least c_2 over the training set,*

$$\frac{1}{N} \mathbb{E} \left[\left\| \hat{\mathbf{y}}_{\text{RNN}}(\mathbf{p}; \hat{\Theta}_{n,\varepsilon}) - \mathbf{y} \right\|_2^2 \right] \geq c_3.$$

On our way to prove this theorem, we prove a novel representational lower bound for RNNs – Proposition 34 in Appendix C.5 – that captures both the number of parameters *and* the norm in the weights. This representational lower bound implies that an RNN that generalizes on the entire data distribution (hence approximates the 1STR model) requires a weight norm that scales with \sqrt{N} . On the other hand, we show that overfitting on the n training samples with zero empirical risk is possible with a $\text{poly}(n)$ weight norm. As a result, as long as $n \leq N^{c_1}$ for some small constant $c_1 > 0$, min-norm ε -ERM will choose models that overfit rather than generalize. The complete proof of Theorem 7 is presented in Appendix C.7.

5. Feedforward Neural Networks (FFNs)

In this section, we consider a general formulation of a feedforward network. Our only requirement will be that the first layer performs a fully-connected projection. The subsequent layers of the network can be arbitrarily implemented, e.g. using attention blocks or convolution filters. Specifically, the FFN implements the mapping $\mathbf{p} \mapsto f(\mathbf{T}, \mathbf{W}\mathbf{x})$ where $\mathbf{W} \in \mathbb{R}^{m_1 \times Nd}$ is the weight matrix in the first layer, $\mathbf{x} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_N^\top)^\top \in \mathbb{R}^{Nd}$, and $f : [N]^{qN} \times \mathbb{R}^{m_1} \rightarrow \mathbb{R}^N$ implements the rest of the network. Unlike Transformers, here we give the network full information of $\mathbf{T} = (\mathbf{t}_1, \dots, \mathbf{t}_N)$, and in particular, it can implement arbitrary encodings of the position variables $\mathbf{t}_1, \dots, \mathbf{t}_N$. This formulation covers usual approaches where encodings of \mathbf{t} are added to or concatenated with \mathbf{x} .

The class of algorithms we consider for training FFNs goes beyond ERM and includes stationary points of the training loss, thus covering outputs of first-order optimization algorithms. This class is formally introduced in Definition 41 in Appendix D.

For our negative result on feedforward networks, we can further restrict the class of q STR models, and only consider simple- q STR. The following minimax lower bound, with its proof deferred to Appendix C, shows that all algorithms in class \mathcal{A} fail to learn even the subset of simple- q STR models with a sample complexity sublinear in N .

Theorem 8 *Suppose $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_{Nd})$, and consider the simple-1STR model with $t_{i1} = t_1$ for all $i \in [N]$, where t_1 is drawn independently and uniformly in $[N]$, and a linear link function, i.e. $y = \langle \mathbf{u}, \mathbf{x}_{t_1} \rangle$ for some $\mathbf{u} \in \mathbb{S}^{d-1}$. Let \mathcal{A} be the class of algorithms in Definition 41. Then,*

$$\inf_{A \in \mathcal{A}} \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} R^{\text{FFN}}(f_{A(S_n)}, \mathbf{W}_{A(S_n)}) \geq 1 - \frac{n}{Nd},$$

with probability 1 over the training set S_n .

The main intuition in the proof of the above theorem follows from the stationarity property of Definition 41. With this property, the rows of the trained \mathbf{W} will always be in the span of the training data $\mathbf{x}^{(i)}$ for $i \in [n]$, and this subspace can be too small to predict y , which by randomizing \mathbf{u} , can depend on all Nd target directions.

6. Conclusion

In this paper, we established a sample complexity separation between Transformers and baseline architectures, namely feedforward and recurrent networks, for learning sequence-to-sequence models where the output at each position depends on a sparse subset of input tokens described in the input itself, coined the q STR model. We proved that Transformers can learn such a model with sample complexity almost independent of the length of the input sequence N , while feedforward and recurrent networks have sample complexity lower bounds of N and $N^{\Omega(1)}$, respectively. Further, we established a separation between FFNs and RNNs by proving that recurrent networks can learn the subset of simple- q STR models where the output at all positions is identical, whereas feedforward networks require at least N samples. An important direction for future work is to develop an understanding of the optimization dynamics of Transformers to learn q STR models, and to study sample complexity separations that highlight the role of depth in Transformers.

References

- [1] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 2552–2623. PMLR, 2023.
- [2] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 37:45614–45650, 2023.
- [3] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- [4] Noah Amsel, Gilad Yehudai, and Joan Bruna. On the benefits of rank in attention layers. *arXiv preprint arXiv:2407.16153*, 2024.
- [5] Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- [6] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation. *arXiv preprint arXiv:2205.01445*, 2022.
- [7] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017.
- [8] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2023.
- [9] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.
- [10] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks. In *Advances in Neural Information Processing Systems*, 2022.
- [11] Minshuo Chen, Xingguo Li, and Tuo Zhao. On generalization bounds of a family of recurrent neural networks. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1233–1243. PMLR, 2020.
- [12] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural Networks can Learn Representations with Gradient Descent. In *Conference on Learning Theory*, 2022.
- [13] Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023.

- [14] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [15] Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2023.
- [16] Hengyu Fu, Tianyu Guo, Yu Bai, and Song Mei. What can a single attention layer learn? a study through the random features lens. *Advances in Neural Information Processing Systems*, 36, 2023.
- [17] Juno Kim and Taiji Suzuki. Transformers provably solve parity efficiently with chain of thought. *arXiv preprint arXiv:2410.08633*, 2024.
- [18] Juno Kim, Tai Nakamaki, and Taiji Suzuki. Transformers are minimax optimal nonparametric in-context learners. In *ICML 2024 Workshop on In-Context Learning*, 2024. URL <https://openreview.net/forum?id=WjrKBQTKp>.
- [19] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pages 19565–19594. PMLR, 2023.
- [20] Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. In *The Twelfth International Conference on Learning Representations*, 2024.
- [21] William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. In *The Twelfth International Conference on Learning Representations*, 2024.
- [22] Alireza Mousavi-Hosseini, Denny Wu, and Murat A Erdogdu. Learning multi-index models with neural networks via mean-field langevin dynamics. *arXiv preprint arXiv:2408.07254*, 2024.
- [23] Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. How to construct deep recurrent neural networks. *arXiv preprint arXiv:1312.6026*, 2013.
- [24] Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations of transformers. *Advances in Neural Information Processing Systems*, 36, 2023.
- [25] Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Transformers, parallel computation, and logarithmic depth. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.

- [27] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- [28] Zixuan Wang, Stanley Wei, Daniel Hsu, and Jason D. Lee. Transformers provably learn sparse token selection while fully-connected nets cannot. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- [29] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.

Appendix A. Detailed Assumptions and Learning Procedure

Here, we state our precise assumptions on the data distribution.

Assumption 1 Suppose $\mathbb{E}[\|\mathbf{x}_i\|^r]^{1/r} \leq \sqrt{C_x d r}$ and $\mathbb{E}[|y_i|^r]^{1/r} \leq \sqrt{C_y r^s}$ for all $r \geq 1$, $i \in [N]$, and some absolute constants $s \geq 1$ and $C_x, C_y > 0$.

Learning the q STR model requires two steps: (i) extracting the relevant tokens at each position, (ii) learning the link function g . We are interested in settings where the difficulty of learning is dominated by the first step, hence we assume g can be approximated by a two-layer feedforward network.

Assumption 2 There exist $m_g \in \mathbb{N}$, $\mathbf{a}_g, \mathbf{b}_g \in \mathbb{R}^{m_g}$ and $\mathbf{W}_g \in \mathbb{R}^{m_g \times qd}$, such that $\|\mathbf{a}_g\|_2 \leq r_a/\sqrt{m_g}$, and $\|(\mathbf{W}_g, \mathbf{b}_g)\|_F \leq \sqrt{m_g} r_w$ for some constants $r_a, r_w > 0$, and

$$\left\{ \|\mathbf{x}_i\|_2 \leq \sqrt{Cd \log(nN)}, \forall i \in [q] \right\} \left| g(\mathbf{x}_1, \dots, \mathbf{x}_q) - \mathbf{a}_g^\top \sigma(\mathbf{W}_g(\mathbf{x}_1^\top, \dots, \mathbf{x}_q^\top)^\top + \mathbf{b}_g) \right|^2 \leq \varepsilon_{2\text{NN}},$$

where $C = 3C_x e$ and $\varepsilon_{2\text{NN}}$ is some absolute constant.

Ideally, $\varepsilon_{2\text{NN}}$ above is a small constant denoting the approximation error. This assumption can be verified using various universal approximation results for ReLU networks. For example, when g is an additive model of P Lipschitz functions, where each function depends only on a k -dimensional projection of the input, the above holds for every $\varepsilon_{2\text{NN}} > 0$ and $m_g = \tilde{O}((P/\sqrt{\varepsilon_{2\text{NN}}})^k)$, $r_a = \tilde{O}((P/\sqrt{\varepsilon_{2\text{NN}}})^{(k+1)/2})$, and $r_w = 1$ (we can always have $r_w = 1$ by homogeneity) [7].

Before introducing the notions of ERM that we employ, we first state several sequential risk formulations to evaluate a predictor $\hat{\mathbf{y}}_{\text{arc}}(\cdot; \Theta) \in \mathcal{F}_{\text{arc}}$ on i.i.d. training samples $\{\mathbf{p}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^n$, where arc denotes a general architecture. We define the *population risk*, *averaged empirical risk*, and *point-wise empirical risk* respectively as

$$R^{\text{arc}}(\Theta) := \frac{1}{N} \mathbb{E} \left[\sum_{j=1}^N (\hat{\mathbf{y}}_{\text{arc}}(\mathbf{p}; \Theta)_j - y_j)^2 \right] = \frac{1}{N} \mathbb{E} \left[\|\hat{\mathbf{y}}_{\text{arc}}(\mathbf{p}; \Theta) - \mathbf{y}\|_2^2 \right], \quad (1)$$

$$\hat{R}_{n,N}^{\text{arc}}(\Theta) := \frac{1}{nN} \sum_{i=1}^n \sum_{j=1}^N (\hat{\mathbf{y}}_{\text{arc}}(\mathbf{p}^{(i)}; \Theta)_j - y_j^{(i)})^2, \quad (2)$$

$$\hat{R}_n^{\text{arc}}(\Theta) := \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{y}}_{\text{arc}}(\mathbf{p}^{(i)}; \Theta)_{j^{(i)}} - y_{j^{(i)}}^{(i)})^2, \quad (3)$$

where $\{j^{(i)}\}_{i=1}^n$ are i.i.d. position indices drawn from $\text{Unif}([N])$. The goal is to minimize the population risk $R^{\text{arc}}(\Theta)$ by minimizing some empirical risk, potentially with weight regularization. We use three formalizations of learning algorithms to prove our results.

1. *Constrained ERM* minimizes an empirical risk \hat{R}_n^{arc} subject to the model parameters belonging on some (e.g., norm-constrained) set Θ . Concretely, let

$$\hat{\Theta} \in \arg \min_{\Theta \in \Theta} \hat{R}_n^{\text{arc}}(\Theta).$$

Theorem 3 considers constrained ERM algorithms for bounded-weight transformers with point-wise risk $\hat{R}_n^{\text{TR}}(\Theta)$, and Theorem 6 uses $\hat{R}_n^{\text{RNN}}(\Theta)$ for RNNs. Note that upper bounds for training with point-wise empirical risk \hat{R}_n^{arc} readily transfer to training with averaged empirical risk $\hat{R}_{n,N}^{\text{arc}}$.

2. *Min-norm ε -ERM* minimizes the norm of the parameters, subject to sufficiently small loss:

$$\hat{\Theta}_\varepsilon \in \arg \min_{\{\Theta: \hat{R}_n^{\text{arc}}(\Theta) - \min \hat{R}_n^{\text{arc}} \leq \varepsilon\}} \|\text{vec}(\Theta)\|_2. \quad (4)$$

Theorem 7 uses min-norm ε -ERM to place a sample complexity lower bound $\hat{R}_n^{\text{RNN}}(\Theta)$.

3. Beyond ERM, Theorem 8 also considers *stationary points* of the averaged or point-wise loss, with ℓ_2 regularization. This learning algorithm is presented in greater detail in Definition 41.

If Θ is defined by a norm constraint, then min-norm ε -ERM with a proper ε can be seen as an instance of constrained ERM. All three formulations are motivated by practical optimization algorithms that either minimize an explicitly regularized loss, or have an implicit bias towards min-norm solutions.

Appendix B. Details of Section 3

Here we present the omitted details and proofs of Section 3. We begin by presenting the architectural details before proving sample complexity upper bounds for Transformers.

B.1. Transformer Architectural Definition

We formally introduce the single-layer H -headed Transformer that appears in all Section 3 proofs.

Positional encoding. To break the permutation equivariance of Transformers, we append positional information to the input tokens. Given a prompt \mathbf{p} , we consider an encoding given by

$$\mathbf{Z}(\mathbf{p}) = \begin{pmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_N \\ \text{enc}(1, \mathbf{t}_1) & \dots & \text{enc}(N, \mathbf{t}_N) \end{pmatrix} \in \mathbb{R}^{D_e \times N},$$

where $\text{enc} : [N] \times [N]^q \rightarrow \mathbb{R}^{d_{\text{enc}}}$ provides the encoding of the position and of \mathbf{t}_i , and $D_e := d + d_{\text{enc}}$. We use \mathbf{z}_i to refer to the i th column above. We remark that allowing enc to take \mathbf{t}_i as input allows specific encodings of the indices \mathbf{t}_i that take advantage of the q STR structure; examples of this have been considered in prior works [28]. In practice, we expect such useful encodings to be learned automatically by previous layers in the Transformer. We remark that for a fair comparison, in our lower bounds for other architectures we allow *arbitrary processing* of \mathbf{t}_i in their encoding procedure. To specify enc , we use a set of vectors $\{\boldsymbol{\omega}_i\}_{i=1}^N$ in \mathbb{R}^{d_e} that satisfy the following property.

Assumption 3 We have $|\langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_j \rangle| \leq \frac{1}{2}$ for all $i \neq j$, and $\|\boldsymbol{\omega}_i\|^2 = 1$ for all i , with $d_e = \Theta(\log N)$.

Such a set of vectors can be obtained e.g., by sampling random Rademacher vectors from the unit cube $\{\pm 1/\sqrt{d_e}\}^{d_e}$ which will satisfy the assumption with high probability. We define

$$\text{enc}(i, \mathbf{t}_i) = \sqrt{d/q}(\boldsymbol{\omega}_i, \boldsymbol{\omega}_{t_{i1}}, \dots, \boldsymbol{\omega}_{t_{iq}})^\top \in \mathbb{R}^{(q+1)d_e},$$

hence $d_{\text{enc}} = (q+1)d_e$ and $D_e = d + (q+1)d_e$. The $\sqrt{d/q}$ prefactor ensures that \mathbf{x}_i and $\text{enc}(i, \mathbf{t}_i)$ will roughly have the same ℓ_2 norm, resulting in a balanced input to the attention layer.

Multi-head attention. Given a sequence $\{z_i\}_{i=1}^N$ where $z_i \in \mathbb{R}^{D_e}$ with D_e as the embedding dimension, a single head of attention outputs another sequence of length N in \mathbb{R}^{D_e} , given by

$$f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) = \left[\sum_{j=1}^N \mathbf{W}_V z_j \frac{e^{\langle \mathbf{W}_Q z_i, \mathbf{W}_K z_j \rangle}}{\sum_{l=1}^N e^{\langle \mathbf{W}_Q z_i, \mathbf{W}_K z_l \rangle}} \right]_{i \in [N]}.$$

Where $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V$ are the key, query, and value projection matrices respectively. We can simplify the presentation by replacing $\mathbf{W}_Q^\top \mathbf{W}_K$ with a single parameterizing matrix for query-key projections denoted by $\mathbf{W}_{\text{QK}} \in \mathbb{R}^{D_e \times D_e}$, and absorbing \mathbf{W}_V into the weights of the feedforward layer. This provides us with a simplified parameterization of attention, which we denote by $f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}})$. This simplification is standard in theoretical works (see e.g. [2, 19, 28, 29]). Our main separation results still apply when maintaining separate trainable projections.

We can concatenate the output of H attention heads with separate key-query projection matrices to obtain a multi-head attention layer with H heads. We denote the output of head $h \in [H]$ with $f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(h)})$. The output of the multi-head attention at position i is then given by

$$f_{\text{Attn}}^{(H)}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(1)}, \dots, \mathbf{W}_{\text{QK}}^{(H)})_i = (f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(1)})_i, \dots, f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(H)})_i)^\top \in \mathbb{R}^{HD_e}.$$

We will denote by $\Theta_{\text{QK}} = (\mathbf{W}_{\text{QK}}^{(1)}, \dots, \mathbf{W}_{\text{QK}}^{(H)})$ the parameters of the multi-head attention.

Finally, a two-layer neural network acts on the output of the attention to generate labels. Given input $\mathbf{h} \in \mathbb{R}^{HD_e}$, the output of the network is given by

$$f_{2\text{NN}}(\mathbf{h}; \mathbf{a}_{2\text{NN}}, \mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}}) = \mathbf{a}_{2\text{NN}}^\top \sigma(\mathbf{W}_{2\text{NN}} \mathbf{h} + \mathbf{b}_{2\text{NN}}),$$

where $\mathbf{W}_{2\text{NN}} \in \mathbb{R}^{m \times HD_e}$ are the first layer weights, $\mathbf{b}_{2\text{NN}}, \mathbf{a}_{2\text{NN}} \in \mathbb{R}^m$ are the second layer weights and biases, and m is the width. We also use the summarized notation $\Theta_{2\text{NN}} = (\mathbf{a}_{2\text{NN}}, \mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})$ to refer to the feedforward layer weights. The prediction of the transformer at position i is given by

$$\hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i = f_{2\text{NN}}(f_{\text{Attn}}^{(H)}(\mathbf{p}; \Theta_{\text{QK}})_i; \Theta_{2\text{NN}}),$$

where $\Theta_{\text{TR}} = (\Theta_{\text{QK}}, \Theta_{2\text{NN}})$ denotes the overall trainable parameters of the Transformer. We use the notation $\hat{\mathbf{y}}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}}) = (\hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_1, \dots, \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_N)^\top \in \mathbb{R}^N$ to denote the vectorized output.

B.2. Proof of Theorem 3

To prove Theorem 3, we will prove the more general theorem below.

Theorem 9 Let $\hat{\Theta} := \arg \min_{\Theta \in \Theta_{\text{TR}}} \hat{R}_n^{\text{TR}}(\Theta)$, where

$$\Theta_{\text{TR}} := \left\{ \|\mathbf{a}_{2\text{NN}}\|_2 \leq r_a / \sqrt{m}, \|(\mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})\|_{\text{F}} \leq r_w \sqrt{m}, \left\| \mathbf{W}_{\text{QK}}^{(h)} \right\|_{2,1} \leq \alpha \forall h \in [H] \right\}.$$

Suppose $H = q$, $m = m_q$, and $\alpha = \tilde{\Theta}(1)$ (given in Lemma 10). Then, under Assumptions 1 to 3, with probability at least $1 - n^{-c}$ for some absolute constant $c > 0$, we have

$$R^{\text{TR}}(\hat{\Theta}) \leq \mathcal{O}(\varepsilon_{\text{NN}}^2) + \tilde{\mathcal{O}} \left(C_1 \sqrt{\frac{(m_q q (d+q) + r_z^6 r_a^2 r_w^2 q^2 \wedge q(q^2 + d^2))}{n}} \right), \quad (5)$$

where $C_1 = q r_a^2 r_w^2 r_z^2$.

We begin with a lemma establishing the capability of Transformers in approximating q STR models.

Lemma 10 *Suppose Assumption 2 holds. Let $r_x = \sqrt{3C_x ed \log(nN)}$. Assume $H = q$ and $m_g = m$. Then, there exists Θ_{TR} such that*

$$\sup_{\{\|\mathbf{x}_j\|_2 \leq r_x, \forall j \in [N]\}} |g(\mathbf{x}_{t_{i1}}, \dots, \mathbf{x}_{t_{iq}}) - \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i| \leq 2\sqrt{\varepsilon_{2\text{NN}}},$$

and

$$\|\mathbf{a}_{2\text{NN}}\|_2 \leq \frac{r_a}{\sqrt{m}}, \quad \|(\mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})\|_{\text{F}} \leq \sqrt{m}r_w, \quad \left\| \mathbf{W}_{\text{QK}}^{(h)\top} \right\|_{2,1} \leq \frac{2d_e q}{d} \log\left(\frac{2r_a r_w r_x N \sqrt{q}}{\varepsilon_{2\text{NN}}}\right),$$

for all $h \in [H]$.

Proof In our construction, the goal of attention head h at position i will be to output $z_{t_{ih}}$. Namely, we want to achieve

$$f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(h)})_i \approx z_{t_{ih}}.$$

Note that to do so, for each key token \mathbf{z}_j , we only need to compute $\langle \boldsymbol{\omega}_{t_{ih}}, \boldsymbol{\omega}_j \rangle$. Therefore, most entries in $\mathbf{W}_{\text{QK}}^{(h)}$ can be zero. We only require a block of $d_e \times d_e$, which corresponds to comparing $\boldsymbol{\omega}_j$ and $\boldsymbol{\omega}_{t_{ih}}$ when comparing query \mathbf{z}_i and key \mathbf{z}_j . Thus, we let

$$\mathbf{W}_{\text{QK}}^{(h)} = \begin{pmatrix} \mathbf{0}_{(d+hd_e) \times d} & \mathbf{0}_{(d+hd_e) \times d_e} & \mathbf{0}_{(d+hd_e) \times qd_e} \\ \mathbf{0}_{d_e \times d} & \alpha \mathbf{I}_{d_e} & \mathbf{0}_{d_e \times qd_e} \\ \mathbf{0}_{(q-h)d_e \times d} & \mathbf{0}_{(q-h)d_e \times d_e} & \mathbf{0}_{(q-h)d_e \times qd_e} \end{pmatrix} \quad (6)$$

Then, we have $\langle \mathbf{z}_i, \mathbf{W}_{\text{QK}}^{(h)} \mathbf{z}_j \rangle = \alpha \langle \boldsymbol{\omega}_{t_{ih}}, \boldsymbol{\omega}_j \rangle d/q$. We can then verify that

$$\left\| \mathbf{A} f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(h)})_i - \mathbf{A} z_{t_{ih}} \right\|_2 \leq \sum_{j \neq t_{ih}} e^{-\alpha d/(2q)} (\|\mathbf{A} \mathbf{z}_j\| + \|\mathbf{A} z_{t_{ih}}\|_2)$$

for every matrix \mathbf{A} . We will specifically choose \mathbf{A} to be the projection onto the first d coordinates in the following. Hence, α will control the error in the softmax attention approximating a ‘‘hard-max’’ attention that would exactly choose $z_{t_{ih}}$.

To construct the weights of the feedforward layer $\mathbf{a}_{2\text{NN}}, \mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}}$, we let $\mathbf{a}_{2\text{NN}} = \mathbf{a}_g$ and $\mathbf{b}_{2\text{NN}} = \mathbf{b}_g$ from Assumption 2, and define $\mathbf{W}_{2\text{NN}}$ by extending \mathbf{W}_g with zero entries such that

$$\mathbf{W}_{2\text{NN}} \begin{pmatrix} \mathbf{z}_{t_{i1}} \\ \dots \\ \mathbf{z}_{t_{iq}} \end{pmatrix} = \mathbf{W}_g \begin{pmatrix} \mathbf{x}_{t_{i1}} \\ \dots \\ \mathbf{x}_{t_{iq}} \end{pmatrix}.$$

Then $\|\mathbf{W}_{2\text{NN}}\|_{\text{F}} = \|\mathbf{W}_g\|_{\text{F}}$. Notice that $\cdot \mapsto \mathbf{a}^\top \sigma(\mathbf{W}(\cdot) + \mathbf{b})$ is $r_a r_w$ Lipschitz. As a result, for any \mathbf{x} with $\|\mathbf{x}\| \leq r_x$ we have

$$|g(\mathbf{x}_{t_{i1}}, \dots, \mathbf{x}_{t_{iq}}) - \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i| \leq \sqrt{\varepsilon_{2\text{NN}}} + \varepsilon_{\text{Attn}},$$

where we recall

$$|g(\mathbf{x}_{t_{i1}}, \dots, \mathbf{x}_{t_{iq}}) - f_{2\text{NN}}((\mathbf{z}_{t_{i1}}, \dots, \mathbf{z}_{t_{iq}}); \mathbf{a}_{2\text{NN}}, \mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})| \leq \sqrt{\varepsilon_{2\text{NN}}},$$

and

$$\begin{aligned} \varepsilon_{\text{Attn}} &= \left| f_{2\text{NN}}(\mathbf{z}_{t_{i1}}, \dots, \mathbf{z}_{t_{iq}}; \Theta_{2\text{NN}}) - f_{2\text{NN}}(f_{\text{Attn}}^{(q)}(\mathbf{p}; \Theta_{\text{QK}}); \Theta_{2\text{NN}}) \right| \\ &\leq r_a r_w \sqrt{\sum_{h=1}^q \left\| \mathbf{A} f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(h)})_i - \mathbf{A} \mathbf{z}_{t_{ih}} \right\|_2^2} \\ &\leq 2r_a r_w r_x N \sqrt{q} e^{-\alpha d / (2q)}, \end{aligned}$$

where we recall $\mathbf{A} \mathbf{z}_j = \mathbf{x}_j$. Thus, with

$$\alpha = 2q \log(2r_a r_w r_x N \sqrt{q} / \sqrt{\varepsilon_{2\text{NN}}}) / d$$

we can guarantee the distance is at most $2\sqrt{\varepsilon_{2\text{NN}}}$. \blacksquare

Before proceeding to obtain statistical guarantees, we will show that we can consider the encodings $\mathbf{z}_j^{(i)}$ to be bounded with high probability. This will be a useful event to consider throughout the proofs of various sections.

Lemma 11 *Suppose $\{\mathbf{p}^{(i)}\}_{i=1}^n$ are n input prompts (not necessarily independent) drawn from the input distribution, with tokens denoted by $\{(\mathbf{x}_j^{(i)})_{j=1}^N\}_{i=1}^n$. Under Assumption 1, for any $r_x > 0$ we have*

$$\mathbb{P}\left(\max_{i \in [n], j \in [N]} \left\| \mathbf{x}_j^{(i)} \right\|_2 \geq r_x\right) \leq nN e^{-r_x^2 / (2C_x e d)}.$$

In particular, for $r_x = \sqrt{3C_x e d \log(nN)}$ we have

$$\mathbb{P}\left(\max_{i \in [n], j \in [N]} \left\| \mathbf{x}_j^{(i)} \right\|_2 \geq r_x\right) \leq \sqrt{\frac{1}{nN}}.$$

Proof Via Markov's inequality, for any $p > 0$ and $r_x > 0$, we have

$$\mathbb{P}\left(\max_{i,j} \left\| \mathbf{x}_j^{(i)} \right\|_2 \geq r_x\right) \leq \frac{\mathbb{E}\left[\max_{i,j} \left\| \mathbf{x}_j^{(i)} \right\|_2^p\right]}{r_x^p} \leq \frac{\mathbb{E}\left[\sum_{i,j} \left\| \mathbf{x}_j^{(i)} \right\|_2^p\right]}{r_x^p} \leq \frac{Nn(C_x p d)^{p/2}}{r_x^p}.$$

Let $p = r_x^2 / (C_x e d)$. Then,

$$\mathbb{P}\left(\max_{i,j} \left\| \mathbf{x}_j^{(i)} \right\|_2 \geq r_x\right) \leq nN e^{-r_x^2 / (2C_x e d)},$$

which proves the first statement, and the second statement follows by plugging in the specific value of r_x . \blacksquare

We are now ready to move to the generalization analysis of Transformers. First, we have to formally define the prediction function class of Transformers with a notation suitable for this section. We begin by defining the function class of attention. We have

$$\mathcal{F}_{\text{Attn}} = \{\mathbf{p}, j \mapsto f_{\text{Attn}}^{(H)}(\mathbf{p}; \Theta_{\text{QK}})_j : \Theta_{\text{QK}} \in \Theta_{\text{QK}}\},$$

where we will later specify Θ_{QK} . Additionally, we define $\mathcal{F}_{2\text{NN}}$ by

$$\mathcal{F}_{2\text{NN}} = \{\mathbf{h} \mapsto f_{2\text{NN}}(\mathbf{h}; \Theta_{2\text{NN}}) : \Theta_{2\text{NN}} \in \Theta_{2\text{NN}}\},$$

where $\Theta_{2\text{NN}} = (\mathbf{a}_{2\text{NN}}, \mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})$, and we will later specify $\Theta_{2\text{NN}}$. Then the class \mathcal{F}_{TR} can be defined as

$$\mathcal{F}_{\text{TR}} = \{\mathbf{p}, j \mapsto f_{2\text{NN}}(f_{\text{Attn}}(\mathbf{p})_j) : f_{\text{Attn}} \in \mathcal{F}_{\text{Attn}}, f_{2\text{NN}} \in \mathcal{F}_{2\text{NN}}\}.$$

Recall we use the S_n to denote the training set. To avoid extra indices, we will use the notation $\mathbf{p}, j \in S_n$ to go over $\{\mathbf{p}^{(i)}, j^{(i)}\}_{i=1}^n$. We can then define the following distances on the introduced function classes

$$\begin{aligned} d_{\infty}^{\text{TR}}(f, f') &:= \sup_{\mathbf{p}, j} |f(\mathbf{p})_j - f'(\mathbf{p})_j|, \quad \forall f, f' \in \mathcal{F}_{\text{TR}} \\ d_{\infty}^{\text{Attn}}(f, f') &:= \sup_{\mathbf{p}, j} \|f(\mathbf{p})_j - f'(\mathbf{p})_j\|_2, \quad \forall f, f' \in \mathcal{F}_{\text{Attn}} \\ d_{\infty}^{2\text{NN}}(f, f') &:= \sup_{\|\cdot\|_2 \leq \sqrt{H}r_z} |f(\cdot) - f'(\cdot)|, \quad \forall f, f' \in \mathcal{F}_{2\text{NN}}. \end{aligned}$$

We choose the radius $\sqrt{H}r_z$ for defining $d_{\infty}^{2\text{NN}}$ since on the event of Lemma 11, this will be the norm bound on the output of the attention layer at every position.

Recall that for a distance d_{∞} and a set \mathcal{F} , an ϵ -covering $\hat{\mathcal{F}}$ is a set such that for every $f \in \mathcal{F}$, there exists $\hat{f} \in \hat{\mathcal{F}}$ such that $d_{\infty}(f, \hat{f}) \leq \epsilon$. The ϵ -covering number of \mathcal{F} , denoted by $\mathcal{C}(\mathcal{F}, d_{\infty}, \epsilon)$, is the number of elements of the smallest such $\hat{\mathcal{F}}$. The following lemma relates the covering number of \mathcal{F}_{TR} to those of $\mathcal{F}_{\text{Attn}}$ and $\mathcal{F}_{2\text{NN}}$.

Lemma 12 *Suppose $f_{2\text{NN}}$ is L_f Lipschitz for every $f_{2\text{NN}} \in \mathcal{F}_{2\text{NN}}$. Then, for any $\epsilon_{2\text{NN}}, \epsilon_{\text{Attn}} > 0$, on the event of Lemma 11 we have*

$$\log \mathcal{C}(\mathcal{F}_{\text{TR}}, d_{\infty}^{\text{TR}}, \epsilon_{2\text{NN}} + L_f \epsilon_{\text{Attn}}) \leq \log \mathcal{C}(\mathcal{F}_{2\text{NN}}, d_{\infty}^{2\text{NN}}, \epsilon_{2\text{NN}}) + \log \mathcal{C}(\mathcal{F}_{\text{Attn}}, d_{\infty}^{\text{Attn}}, \epsilon_{\text{Attn}}).$$

Proof The proof simply follows from the triangle inequality, namely

$$\begin{aligned} \sup_{\mathbf{p}, j} |f_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_j - f_{\text{TR}}(\mathbf{p}; \hat{\Theta}_{\text{TR}})_j| &\leq \sup_{\|\mathbf{h}\|_2 \leq \sqrt{H}r_z} \|f_{2\text{NN}}(\mathbf{h}; \Theta_{\text{NN}}) - f_{2\text{NN}}(\mathbf{h}; \hat{\Theta}_{\text{NN}})\|_2 \\ &\quad + L_f \sup_{\mathbf{p}, j} \|f_{\text{Attn}}^{(H)}(\mathbf{p}; \Theta_{\text{QK}})_j - f_{\text{Attn}}^{(H)}(\mathbf{p}; \hat{\Theta}_{\text{QK}})_j\|_2. \end{aligned}$$

■

We have the following estimate for the covering number of $\mathcal{F}_{2\text{NN}}$.

Lemma 13 *Suppose $\|\text{vec}(\Theta_{\text{RNN}})\|_2 \leq R$ and $\|z_j^{(i)}\|_2 \leq R$ for all $i \in [n]$ and $j \in [N]$. Then,*

$$\log \mathcal{C}(\mathcal{F}_{2\text{NN}}, d_{\infty}^{2\text{NN}}, \epsilon) \lesssim m_g H D_e \log(1 + \text{poly}(R)/\epsilon).$$

This is a special case of Lemma 29, proved in Appendix C.

For the next step, define the distance

$$d_{\infty}^{\text{QK}}(\Theta_{\text{QK}}, \Theta'_{\text{QK}}) := \sup_{\mathbf{p}, j} \|\Theta_{\text{QK}}^{\top} z_j - \Theta'_{\text{QK}}{}^{\top} z_j\|_2$$

on Θ_{QK} , where we recall $\Theta_{\text{QK}} = (\mathbf{W}_{\text{QK}}^{(1)}, \dots, \mathbf{W}_{\text{QK}}^{(H)}) \in \mathbb{R}^{D_e \times HD_e}$. The following lemma relates the covering number of the multi-head attention layer to the matrix covering number of the class of attention parameters.

Lemma 14 *Suppose $\|z_j^{(i)}\|_2 \leq r_z$ for all $i \in [n]$ and $j \in [N]$. Then,*

$$\log \mathcal{C}(\mathcal{F}_{\text{Attn}}, d_\infty^{\text{Attn}}, \epsilon) \leq \log \mathcal{C}\left(\Theta_{\text{QK}}, d_\infty^{\text{QK}}, \frac{\epsilon}{2r_z^2}\right).$$

Proof We recall that $\mathbf{Z} \in \mathbb{R}^{N \times D_e}$ denotes the encoded prompt, and softmax is applied row-wise. For conciseness, Let $\Delta := \sup_{\mathbf{p}, j} \left\| f_{\text{Attn}}^{(H)}(\mathbf{p}; \Theta_{\text{QK}})_j - f_{\text{Attn}}^{(H)}(\mathbf{p}; \hat{\Theta}_{\text{QK}})_j \right\|_2^2$. Then we have

$$\begin{aligned} \Delta &= \sup_{\mathbf{p}, j \in S_n} \sum_{h \in [H]} \left\| f_{\text{Attn}}(\mathbf{p}; \mathbf{W}_{\text{QK}}^{(h)})_j - f_{\text{Attn}}(\mathbf{p}; \hat{\mathbf{W}}_{\text{QK}}^{(h)})_j \right\|_2^2 \\ &= \sup_{\mathbf{p}, j \in S_n} \sum_{h \in [H]} \left\| \text{softmax}(z_j^\top \mathbf{W}_{\text{QK}}^{(h)} \mathbf{Z}^\top) \mathbf{Z} - \text{softmax}(z_j^\top \hat{\mathbf{W}}_{\text{QK}}^{(h)} \mathbf{Z}^\top) \mathbf{Z} \right\|_2^2 \\ &\leq \sup_{\mathbf{p}, j \in S_n} \sum_{h \in [H]} \left\| \mathbf{Z}^\top \right\|_{2, \infty}^2 \left\| \text{softmax}(z_j^\top \mathbf{W}_{\text{QK}}^{(h)} \mathbf{Z}^\top)^\top - \text{softmax}(z_j^\top \hat{\mathbf{W}}_{\text{QK}}^{(h)} \mathbf{Z}^\top)^\top \right\|_1^2, \end{aligned}$$

where we used Lemma 42 for the last inequality. Moreover, by [14, Corollary A.7],

$$\begin{aligned} \left\| \text{softmax}(z_j^\top \mathbf{W}_{\text{QK}}^{(h)} \mathbf{Z}^\top)^\top - \text{softmax}(z_j^\top \hat{\mathbf{W}}_{\text{QK}}^{(h)} \mathbf{Z}^\top)^\top \right\|_1 &\leq 2 \left\| \mathbf{Z} \mathbf{W}_{\text{QK}}^{(h)\top} z_j - \mathbf{Z} \hat{\mathbf{W}}_{\text{QK}}^{(h)\top} z_j \right\|_\infty \\ &\leq 2 \left\| \mathbf{Z}^\top \right\|_{2, \infty} \left\| \mathbf{W}_{\text{QK}}^{(h)\top} z_j - \hat{\mathbf{W}}_{\text{QK}}^{(h)\top} z_j \right\|_2. \end{aligned}$$

Consequently,

$$\begin{aligned} \Delta &\leq 4r_z^4 \sup_{\mathbf{p}, j \in S_n} \sum_{h \in [H]} \left\| \mathbf{W}_{\text{QK}}^{(h)\top} z_j - \hat{\mathbf{W}}_{\text{QK}}^{(h)\top} z_j \right\|_2^2 \\ &= 4r_z^4 \sup_{\mathbf{p}, j \in S_n} \left\| \Theta_{\text{QK}}^\top z_j - \hat{\Theta}_{\text{QK}}^\top z_j \right\|_2^2, \end{aligned}$$

which completes the proof. ■

Further, we have the following covering number estimate for Θ_{QK} .

Lemma 15 *Suppose $\Theta_{\text{QK}} = \{\|\Theta_{\text{QK}}\|_{2,1} \leq R_{2,1}, \|\Theta_{\text{QK}}\|_F \leq R_F\}$ and $\|z_j^{(i)}\|_2 \leq r_z$ for all $i \in [n]$ and $j \in [N]$. Then,*

$$\log \mathcal{C}(\Theta_{\text{QK}}, d_\infty^{\text{QK}}, \epsilon) \lesssim \min\left(\frac{r_z^2 R_{2,1}^2 \log(2HD_e^2)}{\epsilon^2}, HD_e^2 \log\left(1 + \frac{2R_F r_z}{\epsilon}\right)\right).$$

Proof The first estimate comes from Maurey’s sparsification lemma [9, Lemma 3.2], while the second estimate is based on the inequality

$$\left\| \Theta_{\text{QK}}^\top \mathbf{z}_j - \hat{\Theta}_{\text{QK}}^\top \mathbf{z}_j \right\|_2 \leq r_z \left\| \Theta_{\text{QK}} - \hat{\Theta}_{\text{QK}} \right\|_{\text{F}},$$

and covering Θ_{QK} with the Frobenius norm, see e.g. Lemma 44. \blacksquare

Finally, we obtain the following covering number for \mathcal{F}_{TR} .

Proposition 16 *Suppose $\|\mathbf{a}_{2\text{NN}}\|_2 \leq r_{m,a}$, $\|(\mathbf{W}_{2\text{NN}}, \mathbf{b}_{2\text{NN}})\|_{\text{F}} \leq R_{m,w}$, and $\|\mathbf{W}_{\text{QK}}^{(h)}\|_{2,1} \leq r_{\text{QK}}$ for all $h \in [H]$. Further assume $\|\mathbf{z}_j^{(i)}\|_2 \leq r_z$ for all $i \in [n]$ and $j \in [N]$. Let $R := \max(r_{m,a}, R_{m,w}, r_z)$. Then,*

$$\log \mathcal{C}(\mathcal{F}_{\text{TR}}, d_{\mathcal{F}}, \epsilon) \lesssim m_g H D_e \log(1 + R/\epsilon) + \min \left(\frac{r_z^6 r_{m,a}^2 R_{m,w}^2 H^2 r_{\text{QK}}^2 \log(H D_e^2)}{\epsilon^2}, H D_e^2 \log \left(1 + \frac{\sqrt{H} r_{\text{QK}} r_z^3 r_{m,a} R_{m,w}}{\epsilon} \right) \right).$$

Proof The proof follows from a number of observations. First, given the parameterization in the statement of the proposition, we have $L_f = r_{m,a} R_{m,w}$ in Lemma 12. Moreover, we have $R_F \leq \sqrt{H} r_{\text{QK}}$ and $R_{2,1} \leq H r_{\text{QK}}$ in Lemma 15. The rest follows from combining the statements of the previous lemmas. \blacksquare

Next, we will use the covering number bound to provide a bound for Rademacher complexity. Recall that for a class of loss functions \mathcal{L} , the empirical and population Rademacher complexities are defined as

$$\hat{\mathfrak{R}}_n(\mathcal{L}) := \mathbb{E} \left[\sup_{\ell \in \mathcal{L}} \frac{1}{n} \sum_{i=1}^n \xi_i \ell(\mathbf{p}^{(i)}, \mathbf{y}^{(i)}, j^{(i)}) \right], \quad \mathfrak{R}_n(\mathcal{L}) := \mathbb{E}_{(\mathbf{p}, \mathbf{y}, j)} \left[\hat{\mathfrak{R}}_n(\mathcal{L}) \right]$$

respectively, where (ξ_i) are i.i.d. Rademacher random variables. Let the class of loss functions be defined by

$$\mathcal{L}_\tau := \{(\mathbf{p}, \mathbf{y}, j) \mapsto (f_{\text{TR}}(\mathbf{p})_j - y_j)^2 \wedge \tau : f_{\text{TR}} \in \mathcal{F}_{\text{TR}}\}, \quad (7)$$

for some constant $\tau > 0$ to be fixed later. We then have the following bound on Rademacher complexity.

Lemma 17 *Suppose $\max_{i \in [n], j \in [N]} \|\mathbf{z}_j^{(i)}\|_2 \leq r_z$. For the loss class \mathcal{L}_τ given by (7), we have*

$$\hat{\mathfrak{R}}_n(\mathcal{L}_\tau) \leq \tilde{O} \left(\tau \sqrt{\frac{C_1 + (C_2 \wedge C_3)}{n}} \right),$$

where $C_1 = m_g H D_e$, $C_2 = r_z^6 r_{m,a}^2 R_{m,w}^2 H^2 r_{\text{QK}}^2$, and $C_3 = H D_e^2$.

Proof Let $\mathcal{C}(\mathcal{L}, d_\infty^\mathcal{L}, \epsilon)$ denote the ϵ -covering number of \mathcal{L} , where $\ell(\mathbf{p}, \mathbf{y}, j) = (f(\mathbf{p})_j - y_j)^2 \wedge \tau$ and $\ell'(\mathbf{p}, \mathbf{y}, j) = (f'(\mathbf{p})_j - y_j)^2 \wedge \tau$. Then, for any $\alpha \geq 0$, by a standard chaining argument,

$$\begin{aligned} \hat{\mathfrak{R}}_n(\mathcal{L}_\tau) &\lesssim \alpha + \int_\alpha^\tau \sqrt{\frac{\log \mathcal{C}(\mathcal{L}, d_\infty^\mathcal{L}, \epsilon)}{n}} d\epsilon. \\ &\lesssim \alpha + \int_\alpha^\tau \sqrt{\frac{\log \mathcal{C}(\mathcal{F}, d_\infty^{\text{TR}}, \epsilon/(2\sqrt{\tau}))}{n}} \\ &\lesssim \alpha + \int_\alpha^\tau \sqrt{\frac{C_1 \log(R\sqrt{\tau}/\epsilon)}{n}} d\epsilon + \left\{ \int_\alpha^\tau \sqrt{\frac{\tau C_2 \log(HD_e^2)}{n\epsilon^2}} d\epsilon \right\} \wedge \left\{ \int_\alpha^\tau \sqrt{\frac{C_3 \log(1 + C_4\sqrt{\tau}/\epsilon)}{n}} d\epsilon \right\} \\ &\lesssim \alpha + \sqrt{\frac{\tau^2 C_1 \log(R\sqrt{\tau}/\alpha)}{n}} + \left\{ \sqrt{\frac{\tau C_2 \log(HD_e^2)}{n}} \log\left(\frac{\tau}{\alpha}\right) \right\} \wedge \left\{ \sqrt{\frac{\tau^2 C_3 \log(1 + C_4\sqrt{\tau}/\alpha)}{n}} \right\}, \end{aligned}$$

where $(C_i)_{i=1}^3$ are given in the statement of the lemma and $C_4 = \sqrt{H} r_{\text{QK}} r_z^3 r_{m,a} R_{m,w}$. Choosing $\alpha = 1/\sqrt{n}$ completes the proof. \blacksquare

Using standard symmetrization techniques, the above immediately yields a high probability upper bound for the expected truncated loss of any estimator in Θ_{TR} .

Corollary 18 Let $\hat{\Theta} = \arg \min_{\Theta \in \Theta_{\text{TR}}} \hat{R}_n^{\text{TR}}(\Theta)$, where Θ_{TR} is described in Proposition 16. Define $r_z = \sqrt{r_x^2 + d(1 + 1/q)}$ where r_x is defined in Lemma 11. Let C_1, C_2 , and C_3 be defined as in Lemma 17. Then, with probability at least $1 - \delta - (nN)^{-1/2}$ over S_n , we have

$$R_\tau^{\text{TR}}(\hat{\Theta}) - \hat{R}_n^{\text{TR}}(\hat{\Theta}) \leq \tilde{\mathcal{O}}\left(\tau \sqrt{\frac{(C_1 + C_2 \wedge C_3)}{n}}\right) + \mathcal{O}\left(\tau \sqrt{\frac{\log(1/\delta)}{n}}\right),$$

where $R_\tau^{\text{RNN}}(\hat{\Theta}) := \mathbb{E}_{\mathbf{p}, j, y} \left[(\hat{y}_{\text{TR}}(\mathbf{p}; \hat{\Theta})_j - y_j)^2 \wedge \tau \right]$

Proof The proof is a standard consequence of Rademacher-based generalization bounds, with the additional observation that

$$\frac{1}{n} \sum_{i=1}^n (\hat{y}_{\text{TR}}(\mathbf{p}^{(i)}; \hat{\Theta})_{j^{(i)}} - y_{j^{(i)}})^2 \wedge \tau \leq \hat{R}_n^{\text{TR}}(\hat{\Theta}).$$

The last step in the proof of the generalization bound is to bound $R_\tau^{\text{TR}}(\hat{\Theta})$ with $\hat{R}_n^{\text{TR}}(\hat{\Theta})$. This is achieved by the following lemma.

Lemma 19 Define $\kappa^2 := Hr_{m,a}^2 R_{m,w}^2 r_z^2$. Then, under Assumption 1, for $\tau \asymp \kappa^2 \log(\kappa^2 N \sqrt{n}) + \log(\kappa^2 \sqrt{n})^s$, we have

$$R_\tau^{\text{TR}}(\hat{\Theta}) - \hat{R}_n^{\text{TR}}(\hat{\Theta}) \leq \sqrt{\frac{1}{n}}.$$

Proof For conciseness, define $\Delta_y := \left| \hat{y}_{\text{TR}}(\mathbf{p}; \hat{\Theta})_j - y_j \right|$. By the Cuachy-Schwartz inequality, we have

$$\begin{aligned} R^{\text{TR}}(\hat{\Theta}) &= \mathbb{E}[\Delta_y^2 \mathbb{1}[\Delta_y \leq \sqrt{\tau}]] + \mathbb{E}[\Delta_y^2 \mathbb{1}[\Delta_y > \sqrt{\tau}]] \\ &\leq R_{\tau}^{\text{TR}}(\hat{\Theta}) + \mathbb{E}[\Delta_y^4]^{1/2} \mathbb{P}(\Delta_y \geq \sqrt{\tau})^{1/2}. \end{aligned}$$

Moreover,

$$\mathbb{E}[\Delta_y^4]^{1/2} \leq 2 \mathbb{E}[y_j^4]^{1/2} + 2 \mathbb{E}[\hat{y}(\mathbf{p}; \hat{\Theta})_j^4]^{1/2}.$$

By Assumption 1, we have $\mathbb{E}[y_j^4]^{1/2} \lesssim 1$. Additionally, note that

$$\begin{aligned} \left| \hat{y}(\mathbf{p}; \hat{\Theta})_j \right| &\leq \|\mathbf{a}_{2\text{NN}}\|_2 (\sqrt{H} \|\mathbf{W}_{2\text{NN}}\|_{\text{F}} \max_{l \in [N]} \|\mathbf{z}_l\|_2 + \|\mathbf{b}_{2\text{NN}}\|_2) \\ &\leq \sqrt{H} r_{m,a} R_{m,w} (1 + \max_{l \in [N]} \|\mathbf{z}_l\|_2). \end{aligned}$$

To bound $\max_{l \in [N]} \|\mathbf{z}_l\|_2$, we use the subGaussianity of $\|\mathbf{x}_l\|_2$ characterized in Assumption 1. Specifically, for all $r \geq 1$

$$\begin{aligned} \mathbb{E} \left[\max_{l \in [N]} \|\mathbf{x}_l\|_2^4 \right] &\leq \mathbb{E} \left[\max_{l \in [N]} \|\mathbf{x}_l\|_2^{4r} \right]^{1/r} \leq \mathbb{E} \left[\sum_{l=1}^N \|\mathbf{x}_l\|_2^{4r} \right]^{1/r} \\ &\leq N^{1/r} \mathbb{E} \left[\|\mathbf{x}_1\|_2^{4r} \right]^{1/r} \\ &\lesssim N^{1/r} C_x^2 d^2 r^2 \\ &\lesssim (C_x d \log(N))^2, \end{aligned}$$

where the last inequality follows from choosing $r = \log N$. As a result,

$$\mathbb{E} \left[\hat{y}(\mathbf{p}; \hat{\Theta})_j^4 \right]^{1/2} \lesssim H r_{m,a}^2 R_{m,w}^2 r_z^2 \log(N)^2 =: \kappa^2 \log(N)^2.$$

We now turn to bounding the probability. We have

$$\begin{aligned} \mathbb{P}(\Delta_y \geq \sqrt{\tau}) &\leq \mathbb{P} \left(|y_j| \geq \frac{\sqrt{\tau}}{2} \right) + \mathbb{P} \left(\left| \hat{y}(\mathbf{p}; \hat{\Theta})_j \right| \geq \frac{\sqrt{\tau}}{2} \right) \\ &\leq \exp \left(-\Omega(\tau^{1/s}) \right) + N \exp \left(-\Omega \left(\frac{\tau}{H r_{m,a}^2 R_{m,w}^2 r_z^2} \right) \right), \end{aligned}$$

where the second inequality follows from sub-Weibull concentration bounds for y and Lemma 11. Choosing $\tau = \Theta(\kappa^2 \log(\kappa^2 N \sqrt{n}) + \log(\kappa^2 \sqrt{n})^s)$ completes the proof. \blacksquare

Proof of Theorem 9. The theorem follows immediately from the approximation guarantee of Lemma 10, the generalization bound of Corollary 18, and the truncation control of Lemma 19. \blacksquare

B.3. Details on Limitations of Transformers with Few Heads

While Proposition 4 is only meaningful in the setting of $d = \Omega(q)$, the following proposition provides an exact lower bound $H \geq q$ on the number of heads for all d , at the expense of additional restrictions on the attention matrix.

Proposition 20 *Consider the q STR data model. Suppose $d = 1$ and $y_i = \frac{1}{\sqrt{q}} \sum_{j=1}^q (x_{iij}^2 - \mathbb{E}[x_{iij}^2])$. Assume $x_i \sim \mathcal{N}(0, \sigma_i^2)$ independently, such that $\sigma_i = 1$ for $i < N/2$ and $\sigma_i = 0$ for $i \geq N/2$. Further, assume the attention weights between the data and positional encoding parts of the tokens are fixed at zero, i.e. $\mathbf{W}_{\text{QK}}^{(h)} = \begin{pmatrix} \mathbf{W}_x^{(h)} & \mathbf{0}_{d \times (q+1)d_e} \\ \mathbf{0}_{(q+1)d_e \times d} & \mathbf{W}_\omega^{(h)} \end{pmatrix}$ where $\mathbf{W}_x^{(h)} \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_\omega^{(h)} \in \mathbb{R}^{(q+1)d_e \times (q+1)d_e}$ are the attention parameters, for $i \in [H]$. Then, there exists a distribution over $(\mathbf{t}_i)_{i \in [N]}$ such that for any choice of Θ_{TR} , we have*

$$\frac{1}{N} \mathbb{E} \left[\|\mathbf{y} - \hat{\mathbf{y}}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})\|_2^2 \right] \geq 1 - \frac{H}{q}.$$

Note that in our approximation constructions for learning q STR, we always fixed the attention weights between data and positional components to be zero, which is why we assume the same in Proposition 20.

Proof of Proposition 20. We will simply choose $\mathbf{t}_i = (1, \dots, q)$ deterministically for $i \geq \frac{N}{2}$ and draw \mathbf{t}_i from an arbitrary distribution for $i < N/2$. Note that we have

$$R^{\text{TR}}(\Theta_{\text{TR}}) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[(y_i - \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i)^2] \geq \frac{1}{N} \sum_{i=N/2}^N \mathbb{E}[(y_i - \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i)^2].$$

Let $\phi : \mathbb{R}^{HD_e} \rightarrow \mathbb{R}$ denote the mapping by the feedforward layer. Fix some $i \geq N/2$. Note that

$$\begin{aligned} \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i &= \phi(f_{\text{Attn}}^{(H)}(\mathbf{p}; \Theta_{\text{QK}})_i) \\ &= \phi\left(\sum_{j=1}^N \alpha_{ij}^{(1)} \mathbf{z}_j, \dots, \sum_{j=1}^N \alpha_{ij}^{(H)} \mathbf{z}_j\right) \\ &= \tilde{\phi}\left(\sum_{j=1}^q \alpha_{ij}^{(1)} x_j, \dots, \sum_{j=1}^q \alpha_{ij}^{(H)} x_j, (\mathbf{z}_l)_{l=q+1}^N\right), \end{aligned}$$

for some real-valued function $\tilde{\phi}$, where

$$\alpha_{ij}^{(h)} = \frac{e^{\langle \mathbf{z}_i, \mathbf{W}_{\text{QK}}^{(h)} \mathbf{z}_j \rangle}}{\sum_{l=1}^N e^{\langle \mathbf{z}_i, \mathbf{W}_{\text{QK}}^{(h)} \mathbf{z}_l \rangle}},$$

are the attention scores. Let $\mathbf{A}^{(i)} \in \mathbb{R}^{H \times q}$ be the matrix such that $A_{hj}^{(i)} = \alpha_{ij}^{(h)}$. Let $\mathbf{x}_{1:q} = (x_1, \dots, x_q)^\top \in \mathbb{R}^q$. Then,

$$\begin{aligned} R^{\text{TR}}(\Theta_{\text{TR}}) &\geq \frac{1}{N} \sum_{i=N/2}^N \mathbb{E} \left[\left(y_i - \tilde{\phi} \left(\mathbf{A}^{(i)} \mathbf{x}_{1:q}, (\mathbf{z}_l)_{l=q+1}^N \right) \right)^2 \right] \\ &\geq \frac{1}{Nq} \sum_{i=N/2}^N \mathbb{E} \left[\text{Var} \left(\|\mathbf{x}_{1:q}\|^2 \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q} \right) \right] \end{aligned} \quad (8)$$

where $\mathbf{V}^{(i)} \in \mathbb{R}^{H \times q}$ is a matrix whose rows form an orthonormal basis of $\text{span}(\boldsymbol{\alpha}_i^{(1)}, \dots, \boldsymbol{\alpha}_i^{(H)})$ where $\boldsymbol{\alpha}_i^{(h)} = (\alpha_{i1}^{(h)}, \dots, \alpha_{iq}^{(h)})^\top \in \mathbb{R}^q$ (note that $\mathbf{V}^{(i)}$ may have fewer than H rows, we consider the worst-case for the lower bound which is having H rows). The second inequality follows from the fact that \mathbf{z}_l is independent of $\mathbf{x}_{1:q}$ for $l \geq q+1$, and the fact that best predictor of y_i (in L_2 error) given $\mathbf{A}^{(i)} \mathbf{x}_{1:q}$ is $\mathbb{E} \left[y_i \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q} \right]$.

Next, thanks to the structural property of $\mathbf{W}_{\text{QK}}^{(h)}$ in the assumption of the proposition and the fact that $x_i = 0$ for $i \geq N/2$, $\alpha_{ij}^{(h)}$ does not depend on $(x_l)_{l \in [q]}$ for all $h \in [H]$, $i \geq N/2$, and $j \in [q]$. As a result, $\mathbf{V}^{(i)}$ is independent of $\mathbf{x}_{1:q}$. Therefore,

$$\mathbf{x}_{1:q} \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q} \sim \mathcal{N}(\mathbf{V}^{(i)\top} \mathbf{V}^{(i)} \mathbf{x}_{1:q}, \mathbf{I}_q - \mathbf{V}^{(i)\top} \mathbf{V}^{(i)}).$$

By Lemma 43, we have $\text{Var}(\|\mathbf{x}_{1:q}\|^2 \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q}) = 2(q - H)$, which combined with (8) completes the proof. \blacksquare

We now present the similarly structured proof of Proposition 4.

Proof of Proposition 4. The choice of distribution over $(\mathbf{t}_i)_{i \geq N/2}$ is similar to the one presented above, i.e. we let $\mathbf{t}_i = (1, \dots, q)$ deterministically for $i \geq \frac{N}{2}$. However, for $i < \frac{N}{2}$, we draw \mathbf{t}_i such that they are independent from \mathbf{x} . Once again, we use the fact that

$$R^{\text{TR}}(\Theta_{\text{TR}}) \geq \frac{1}{N} \sum_{i=N/2}^N \mathbb{E} [(y_i - \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i)^2].$$

Recall $\mathbf{z}_i = (\mathbf{x}_i^\top, \text{enc}(i, \mathbf{t}_i)^\top)$. Fix some $i \geq N/2$, and define

$$\tilde{\alpha}_{ij}^{(h)} = e^{\langle \text{enc}(i, \mathbf{t}_i), \mathbf{W}_{\text{QK}}^{(h,e,x)} \mathbf{x}_j \rangle + \langle \text{enc}(i, \mathbf{t}_i), \mathbf{W}_{\text{QK}}^{(h,e,e)} \text{enc}(j, \mathbf{t}_j) \rangle},$$

where we use the notation

$$\mathbf{W}_{\text{QK}}^{(h)} = \begin{pmatrix} \mathbf{W}_{\text{QK}}^{(h,x,x)} & \mathbf{W}_{\text{QK}}^{(h,x,e)} \\ \mathbf{W}_{\text{QK}}^{(h,e,x)} & \mathbf{W}_{\text{QK}}^{(h,e,e)} \end{pmatrix},$$

for the query-key matrix of each head. Recall that $\mathbf{x}_i = 0$ for $i < N/2$, thus the attention weights are given by

$$\alpha_{ij}^{(h)} = \frac{\tilde{\alpha}_{ij}^{(h)}}{\sum_{l=1}^N \tilde{\alpha}_{il}^{(h)}}.$$

Recall from the proof of Proposition 20 that we denote the feedforward layer by $\phi : \mathbb{R}^{HD_e} \rightarrow \mathbb{R}$. With this notation, we have

$$\begin{aligned} \hat{y}_{\text{TR}}(\mathbf{p}; \Theta_{\text{TR}})_i &= \phi\left(\sum_{j=1}^N \alpha_{ij}^{(1)} \mathbf{z}_j, \dots, \sum_{j=1}^N \alpha_{ij}^{(H)} \mathbf{z}_j\right) \\ &= \tilde{\phi}\left(\sum_{j=1}^q \alpha_{ij}^{(1)} \mathbf{x}_j, \dots, \sum_{j=1}^q \alpha_{ij}^{(H)} \mathbf{x}_j, (\tilde{\alpha}_{ij}^{(h)})_{h=1, j=1}^{h=H, j=N}, (\mathbf{z}_j)_{j=l+1}^N\right). \end{aligned}$$

Therefore, using the fact that \mathbf{z}_j and $\tilde{\alpha}_{ij}^{(h)}$ are independent of $\mathbf{x}_{1:q}$ for $j \geq l+1$, we have

$$\begin{aligned} R^{\text{TR}}(\Theta_{\text{TR}}) &= \frac{1}{N} \sum_{i=N/2}^N \mathbb{E} \left[\left(y_i - \tilde{\phi}\left(\sum_{j=1}^q \alpha_{ij}^{(1)} \mathbf{x}_j, \dots, \sum_{j=1}^q \alpha_{ij}^{(H)} \mathbf{x}_j, (\tilde{\alpha}_{ij}^{(h)})_{h=1, j=1}^{h=H, j=N}, (\mathbf{z}_j)_{j=l+1}^N\right) \right)^2 \right] \\ &\geq \frac{1}{Nqd} \sum_{i=N/2}^N \mathbb{E} \left[\text{Var} \left(\|\mathbf{x}_{1:q}\|^2 \mid \left(\langle \boldsymbol{\alpha}_i^{(h,r)}, \mathbf{x}_{1:q} \rangle \right)_{h=1, r=1}^{h=H, r=d}, (\tilde{\alpha}_{ij}^{(h)})_{h=1, j=1}^{h=H, j=q} \right) \right] \\ &\geq \frac{1}{Nqd} \sum_{i=N/2}^N \mathbb{E} \left[\text{Var} \left(\|\mathbf{x}_{1:q}\|^2 \mid \left(\langle \boldsymbol{\alpha}_i^{(h,r)}, \mathbf{x}_{1:q} \rangle \right)_{h=1, r=1}^{h=H, r=d}, \left(\langle \mathbf{w}_{i,j}^{(h)}, \mathbf{x}_{1:q} \rangle \right)_{h=1, j=1}^{H, q} \right) \right] \\ &= \frac{1}{Nqd} \sum_{i=N/2}^N \mathbb{E} \left[\text{var} \left(\|\mathbf{x}_{1:q}\|^2 \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q} \right) \right], \end{aligned}$$

where $\boldsymbol{\alpha}_i^{(h,r)} \in \mathbb{R}^{qd}$ such that

$$(\boldsymbol{\alpha}_i^{(h,r)})_{jl} = \begin{cases} \alpha_{ij}^{(h)}, & \text{if } l = r \\ 0, & \text{if } l \neq r, \end{cases}$$

which yields $\langle \boldsymbol{\alpha}_i^{(h,r)}, \mathbf{x}_{1:q} \rangle = \sum_{j=1}^q \alpha_{ij}^{(h)} x_{jr}$, and $\mathbf{w}_{i,j}^{(h)} \in \mathbb{R}^{qd}$ such that

$$(w_{i,j}^{(h)})_{sl} = \begin{cases} (\mathbf{W}_{\text{QK}}^{(h,e,x)})^\top \text{enc}(i, \mathbf{t}_i)_l, & \text{if } s = j \\ 0, & \text{if } s \neq j, \end{cases}$$

which yields $\langle \mathbf{w}_{i,j}^{(h)}, \mathbf{x}_{1:q} \rangle = \langle \mathbf{W}^{(h,e,x)} \text{enc}(i, \mathbf{t}_i), \mathbf{x}_j \rangle$. Finally, $\mathbf{V}^{(i)}$ is a matrix whose rows form an orthonormal basis of $\text{span}\left(\left(\boldsymbol{\alpha}_i^{(h,r)}\right)_{h=1, r=1}^{h=H, r=d}, \left(\mathbf{w}_{i,j}^{(h)}\right)_{h=1, j=1}^{h=H, j=q}\right)$. Namely, $\mathbf{V}^{(i)}$ has at most $H(d+q)$ rows. Recall that

$$\mathbf{x}_{1:q} \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q} \sim \mathcal{N}(\mathbf{V}^{(i)\top} \mathbf{V}^{(i)} \mathbf{x}_{1:q}, \mathbf{I}_{qd} - \mathbf{V}^{(i)\top} \mathbf{V}^{(i)}).$$

Once again, by Lemma 43, we conclude that $\text{var}(\|\mathbf{x}_{1:q}\|^2 \mid \mathbf{V}^{(i)} \mathbf{x}_{1:q}) \geq 2(qd - H(q+d))$, which completes the proof. \blacksquare

Appendix C. Details and Proofs of Section 4

Before presenting the proofs, we state the omitted setup and parameterization of the network in the next section.

C.1. Complete Setup of RNNs in the Upper Bound

A bidirectional RNN maintains, for each position in the sequence, a forward and a reverse hidden state, denoted by $(\mathbf{h}_i^{\rightarrow})_{i=1}^N$ and $(\mathbf{h}_i^{\leftarrow})_{i=1}^N$, where $\mathbf{h}_i^{\rightarrow}, \mathbf{h}_i^{\leftarrow} \in \mathbb{R}^{d_h}$. These hidden states are obtained by initializing $\mathbf{h}_1^{\rightarrow} = \mathbf{h}_N^{\leftarrow} = \mathbf{0}_{d_h}$ and recursively applying

$$\begin{aligned} \mathbf{h}_i^{\rightarrow} &= \Pi_{r_h}(\mathbf{h}_{i-1}^{\rightarrow} + f_h^{\rightarrow}(\mathbf{h}_{i-1}^{\rightarrow}, \mathbf{z}_{i-1}; \Theta_h^{\rightarrow})), \quad \forall i \in \{2, \dots, N\} \\ \mathbf{h}_i^{\leftarrow} &= \Pi_{r_h}(\mathbf{h}_{i+1}^{\leftarrow} + f_h^{\leftarrow}(\mathbf{h}_{i+1}^{\leftarrow}, \mathbf{z}_{i+1}; \Theta_h^{\leftarrow})), \quad \forall i \in \{1, \dots, N-1\}, \end{aligned}$$

where $\Pi_{r_h} : \mathbb{R}^{d_h} \rightarrow \mathbb{R}^{d_h}$ is the projection $\Pi_{r_h} \mathbf{h} = (1 \wedge r_h / \|\mathbf{h}\|_2) \mathbf{h}$, and f_h^{\rightarrow} and f_h^{\leftarrow} are implemented by feedforward networks, parameterized by Θ_h^{\rightarrow} and Θ_h^{\leftarrow} respectively. Recall $\mathbf{z}_i = (\mathbf{x}_i^{\top}, \text{enc}(i, \mathbf{t}_i)^{\top})^{\top}$ is the encoding of \mathbf{x}_i . We remark that while we add Π_{r_h} for technical reasons, it resembles layer normalization which ensures stability of the state transitions on very long inputs; a more involved analysis can replace Π_{r_h} with standard formulations of layer normalization. Additionally, directly adding $\mathbf{h}_{i-1}^{\rightarrow}$ and $\mathbf{h}_{i+1}^{\leftarrow}$ to the output of transition functions represents residual or skip connections. The output at position i is generated by

$$y_i = f_y(\mathbf{h}_i^{\rightarrow}, \mathbf{h}_i^{\leftarrow}, \mathbf{z}_i; \Theta_y),$$

which is an L_y -layer feedforward network. Specifically, we consider an RNN with deep transitions [23] and let $f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow})$ be an L_h -layer feedforward network, given by

$$f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}) = \mathbf{W}_{L_h}^{\rightarrow} \sigma(\mathbf{W}_{L_h-1}^{\rightarrow} \dots \sigma(\mathbf{W}_2^{\rightarrow} \sigma(\mathbf{W}_1^{\rightarrow}(\cdot) + \mathbf{b}_1^{\rightarrow}) + \mathbf{b}_2^{\rightarrow}) \dots + \mathbf{b}_{L_h-1}^{\rightarrow}), \quad (9)$$

with $\Theta_h^{\rightarrow} = (\mathbf{W}_1^{\rightarrow}, \mathbf{b}_1^{\rightarrow}, \dots, \mathbf{W}_{L_h-1}^{\rightarrow}, \mathbf{b}_{L_h-1}^{\rightarrow}, \mathbf{W}_{L_h}^{\rightarrow})$ and a similar equation for $f_h^{\leftarrow}(\cdot; \Theta_h^{\leftarrow})$. As will be evident from its proof, Theorem 6 only requires $L_h, L_y = \mathcal{O}(1)$.

We denote the complete output of the RNN via

$$\hat{\mathbf{y}}_{\text{RNN}}(\mathbf{p}; \Theta_{\text{RNN}}) = (f_y(\mathbf{h}_1^{\rightarrow}, \mathbf{h}_1^{\leftarrow}, \mathbf{z}_1; \Theta_y), \dots, f_y(\mathbf{h}_N^{\rightarrow}, \mathbf{h}_N^{\leftarrow}, \mathbf{z}_N; \Theta_y)) \in \mathbb{R}^N.$$

We now define the constraint set of this architecture. Let

$$\Theta_{\text{RNN}} = \left\{ \Theta : \|\text{vec}(\Theta)\|_2 \leq R, \|\mathbf{W}_{L_h}^{\rightarrow}\|_{\text{op}} \dots \|\mathbf{W}_{1,h}^{\rightarrow}\|_{\text{op}} \leq \alpha_N, \|\mathbf{W}_{L_h}^{\leftarrow}\|_{\text{op}} \dots \|\mathbf{W}_{1,h}^{\leftarrow}\|_{\text{op}} \leq \alpha_N \right\}, \quad (10)$$

where $\mathbf{W}_{1,h}^{\rightarrow}$ contains the first d_h columns of $\mathbf{W}_1^{\rightarrow}$, and the conditions above are introduced to ensure f_h^{\rightarrow} and f_h^{\leftarrow} are at most α_N -Lipschitz with respect to the hidden state input. One way to meet this requirement is to multiply $\mathbf{W}_{1,h}^{\rightarrow}$ by a factor of $\alpha_N / \prod_{l=2}^{L_h} \|\mathbf{W}_l^{\rightarrow}\|_{\text{op}}$ in the forward pass. Without this Lipschitzness constraint, current techniques for proving uniform RNN generalization bounds will suffer from a sample complexity linear in N , see e.g. [11].

For Theorem 6 we only require $\alpha_N \leq N^{-1}$. In particular, we can choose $\alpha_N = 0$ and fix $\mathbf{W}_{1,h}^{\rightarrow} = \mathbf{W}_{1,h}^{\leftarrow} = \mathbf{0}$, which would simplify the parameterization of the network. Namely, in our construction f^{\rightarrow} and f^{\leftarrow} do not need to depend on \mathbf{h}^{\rightarrow} and \mathbf{h}^{\leftarrow} respectively.

C.2. Overview of the Proof of Theorem 6

The following is the roadmap we will take for the proof of Theorem 6. The goal here is to implement a bi-directional RNN in such a way that

$$\mathbf{h}_i^{\rightarrow} \approx (\mathbf{x}_{t_1} \mathbb{1}[t_1 < i], \dots, \mathbf{x}_{t_q} \mathbb{1}[t_q < i]),$$

and

$$\mathbf{h}_i^{\leftarrow} \approx (\mathbf{x}_{t_1} \mathbb{1}[t_1 > i], \dots, \mathbf{x}_{t_q} \mathbb{1}[t_q > i]).$$

Throughout this section, we will use the notation

$$\Psi(\mathbf{x}, \mathbf{t}, i) = (\mathbf{x}^{\top} \mathbb{1}[t_1 = i], \dots, \mathbf{x}^{\top} \mathbb{1}[t_q = i])^{\top}.$$

We can obtain the hidden states above through the following updates

$$\mathbf{h}_{i+1}^{\rightarrow} = \mathbf{h}_i^{\rightarrow} + \Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i),$$

and

$$\mathbf{h}_{i-1}^{\leftarrow} = \mathbf{h}_i^{\leftarrow} + \Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i).$$

where

$$\Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i)_l = \frac{\mathbf{x}_i \sigma(\langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_{t_l} \rangle - \delta)}{1 - \delta} = \mathbf{x}_i \mathbb{1}[t_l = i], \quad \forall l \in [q]$$

where we recall $\boldsymbol{\omega}_t = (\boldsymbol{\omega}_{t_1}, \dots, \boldsymbol{\omega}_{t_q})$, and σ is ReLU. As a result, our network must approximate

$$f_h^{\rightarrow}(\mathbf{h}_i^{\rightarrow}, \mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i; \boldsymbol{\Theta}_h^{\rightarrow}) = f_h^{\leftarrow}(\mathbf{h}_i^{\leftarrow}, \mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i; \boldsymbol{\Theta}_h^{\leftarrow}) \approx \Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i).$$

A core challenge in this approximation is that if we simply control

$$\|f_h^{\rightarrow}(\mathbf{h}_i^{\rightarrow}, \mathbf{z}_i; \boldsymbol{\Theta}_h^{\rightarrow}) - \Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i)\|_2 \leq \varepsilon, \quad (11)$$

this error will propagate through the forward pass, and we will have

$$\left\| \mathbf{h}_i^{\rightarrow} - \sum_{j=1}^{i-1} \Psi(\mathbf{x}_j, \boldsymbol{\omega}_t, \boldsymbol{\omega}_j) \right\|_2 \lesssim N\varepsilon.$$

As a result, we would like an implementation that satisfies the following

$$\|f_h^{\rightarrow}(\mathbf{h}_i^{\rightarrow}, \mathbf{z}_i; \boldsymbol{\Theta}_h^{\rightarrow})_l - \Psi(\mathbf{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i)_l\|_2 \leq \begin{cases} 0 & t_l \neq i \\ \varepsilon & t_l = i. \end{cases} \quad (12)$$

Note that

$$\mathbf{h}_i^{\rightarrow} = \sum_{j=1}^{i-1} f_h^{\rightarrow}(\mathbf{h}_j^{\rightarrow}, \mathbf{z}_j; \boldsymbol{\Theta}_h^{\rightarrow}).$$

Since for each $l \in [q]$, $t_l = j$ is possible for at most one $j \in [N]$, (12) implies

$$\left\| \mathbf{h}_i^{\rightarrow} - \sum_{j=1}^{i-1} \Psi(\mathbf{x}_j, \boldsymbol{\omega}_t, \boldsymbol{\omega}_j) \right\|_2 \leq \sqrt{q}\varepsilon,$$

for all $i \in [N]$, hence, we can avoid dependence on N .

We can implement f_h^{\rightarrow} to satisfy (11) with a depth three network, where the first two layers implements $\langle \omega_i, \omega_{t_j} \rangle$ (as a sum of Lipschitz 2-dimensional functions, an example of their approximation is given by [7, Proposition 6]), and the third performs coordinate-wise product between \mathbf{x}_i and $\sigma(\langle \omega_i, \omega_{t_j} \rangle - 1/2)$ (which for each coordinate is a Lipschitz two-dimensional function). To ensure f_h^{\rightarrow} satisfies (12), we can pass the outputs to a fourth layer which rectifies its input near zero to be exactly zero using ReLU activations.

To generate y_i from $\mathbf{h}_i^{\rightarrow}$ and $\mathbf{h}_i^{\leftarrow}$, we first calculate

$$\begin{aligned} \mathbf{h}_i &= f_{hh}(\mathbf{h}_i^{\rightarrow}, \mathbf{h}_i^{\leftarrow}, \mathbf{x}_i, \omega_i, \omega_t) \\ &\approx \mathbf{h}_i^{\rightarrow} + \mathbf{h}_i^{\leftarrow} + \Psi(\mathbf{x}_i, \omega_t, \omega_i) \\ &\approx (\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_q}). \end{aligned}$$

Finally, y_i can be generated from \mathbf{h}_i by applying the two-layer neural network from Assumption 2 that approximates $y_i = g(\mathbf{x}_t)$.

Note that the construction above has a complexity $\text{poly}(d, q, \log(nN))$ (both in terms of number and weight of parameters), only depending on N up to log factors. As a result, by a simple parameter-counting approach, the sample complexity of regularized ERM would also be (almost) independent of N . We also simply use the encoding

$$\mathbf{z}_i = (\mathbf{x}_i, \omega_i, \omega_{t_1}, \dots, \omega_{t_q})^\top,$$

for the RNN positive result. The scaling difference with the encoding for Transformers is only made to simplify the exposition, as we no longer keep explicit dependence on d and q .

C.3. Approximation Upper Bounds for RNNs

As explained above, to implement f_h^{\rightarrow} we first construct a depth three neural network (with two layers of non-linearity) which approximately performs the following mapping

$$\begin{pmatrix} \mathbf{h} \\ \mathbf{x} \\ \omega_i \\ \omega_{t_1} \\ \vdots \\ \omega_{t_q} \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{x} \\ \langle \omega_i, \omega_{t_1} \rangle \\ \vdots \\ \langle \omega_i, \omega_{t_q} \rangle \end{pmatrix} \mapsto \begin{pmatrix} 2\mathbf{x}\sigma(\langle \omega_i, \omega_{t_1} \rangle - 1/2) \\ \vdots \\ 2\mathbf{x}\sigma(\langle \omega_i, \omega_{t_q} \rangle - 1/2) \end{pmatrix}.$$

The first mapping will be provided by

$$\chi_1 = \mathbf{A}_1 \sigma(\mathbf{W}_1 \chi_0 + \mathbf{b}_1),$$

where $\chi_0 = (\mathbf{h}^\top, \mathbf{x}^\top, \omega_i^\top, \omega_{t_1}^\top, \dots, \omega_{t_q}^\top)^\top \in \mathbb{R}^{d_h + d + (q+1)d_e}$, $\mathbf{W}_1 \in \mathbb{R}^{m_1 \times (d_h + d + (q+1)d_e)}$, $\mathbf{b}_1 \in \mathbb{R}^{m_1}$, and $\mathbf{A}_1 \in \mathbb{R}^{(d+q) \times m_1}$, with m_1 as the width of the first layer. We will use the notation

$$\chi_1 = (\chi_1^{\mathbf{x}}, \chi_1^\omega(1), \dots, \chi_1^\omega(q))$$

to refer for the first d coordinates and the rest of the q coordinates of χ_1 respectively, thus ideally $\chi_1^{\mathbf{x}} = \mathbf{x}$ and $\chi_1^\omega(l) = \langle \omega_i, \omega_{t_l} \rangle$. The second mapping is provided by

$$\chi_2 = \mathbf{A}_2 \sigma(\mathbf{W}_2 \chi_1 + \mathbf{b}_2),$$

where $\mathbf{W}_2 \in \mathbb{R}^{m_2 \times (d+q)}$, $\mathbf{b}_2 \in \mathbb{R}^{m_2}$, and $\mathbf{A}_2 \in \mathbb{R}^{dq \times m_2}$. We will similarly use the notation $\chi_2 = (\chi_2(1), \dots, \chi_2(q))$, where our goal is to have $\chi_2(l) \approx 2\mathbf{x}\sigma(\langle \omega_i, \omega_{t_l} \rangle - 1/2)$. To implement the first mapping, we rely on the following lemma.

Lemma 21 *Let σ be the ReLU activation. For any $\varepsilon > 0$ and positive integer d_e , there exists $m = \mathcal{O}(d_e^3(\log(d_e/\varepsilon)/\varepsilon)^2)$, $\mathbf{a} \in \mathbb{R}^m$, $\mathbf{W} \in \mathbb{R}^{m \times 2d_e}$, and $\mathbf{b} \in \mathbb{R}^m$, such that*

$$\sup_{\omega_1, \omega_2 \in \mathbb{S}^{d_e-1}} \left| \langle \omega_1, \omega_2 \rangle - \mathbf{a}^\top \sigma \left(\mathbf{W} \begin{pmatrix} \omega_1 \\ \omega_2 \end{pmatrix} + \mathbf{b} \right) \right| \leq \varepsilon,$$

and

$$\|\mathbf{a}\|_2 \leq \mathcal{O}\left(d_e^{5/2}(\log(d_e/\varepsilon)/\varepsilon)^{3/2}/\sqrt{m}\right), \quad \|\mathbf{W}^\top\|_{1,\infty} \leq 1, \quad \|\mathbf{b}\|_\infty \leq 1.$$

Proof Consider the mapping $e_{1j}, e_{2j} \mapsto e_{1j}e_{2j}$. Note that when $|e_{1j}| \leq 1$ and $|e_{2j}| \leq 1$, this mapping is $\sqrt{2}$ -Lipschitz, and the output is bounded between $[-1, 1]$. Then, by Lemma 45, for every $\varepsilon_j > 0$, there exists $m_j \leq \mathcal{O}((1/\varepsilon_j \log(1/\varepsilon_j))^2)$, $\mathbf{a}_j \in \mathbb{R}^{m_j}$, $\mathbf{W}_j \in \mathbb{R}^{m_j \times 2d_e}$, and $\mathbf{b}_j \in \mathbb{R}^{m_j}$, such that

$$\sup_{|e_{1j}| \leq 1, |e_{2j}| \leq 1} \left| e_{1j}e_{2j} - \sum_{l=1}^m a_{jl} \sigma \left(\langle \mathbf{w}_{jl}, (\omega_1^\top, \omega_2^\top)^\top \rangle + b_{jl} \right) \right| \leq \varepsilon_j,$$

$\|\mathbf{a}_j\|_2 \leq \mathcal{O}((\log(1/\varepsilon_j)/\varepsilon_j)^{3/2}/\sqrt{m_j})$, $\|\mathbf{b}_j\|_\infty \leq 1$, and $\|\mathbf{w}_{jl}\|_1 \leq 1$. Specifically, the only non-zero coordinates of \mathbf{w}_{jl} are the j th and $d_e + j$ th coordinates.

Let $\varepsilon_j = \varepsilon/d_e$ and $m = \sum_{j=1}^{d_e} m_j = \mathcal{O}(d_e^3(\log(d_e/\varepsilon)/\varepsilon)^2)$. Construct $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$ and $\mathbf{W} \in \mathbb{R}^{m \times 2d_e}$ by concatenating (\mathbf{a}_j) , (\mathbf{b}_j) , and (\mathbf{W}_j) respectively. The resulting network satisfies

$$\sup_{\omega_1, \omega_2 \in \mathbb{S}^{d_e-1}} \left| \langle \omega_1, \omega_2 \rangle - \mathbf{a}^\top \sigma \left(\mathbf{W} \begin{pmatrix} \omega_1 \\ \omega_2 \end{pmatrix} + \mathbf{b} \right) \right| \leq \varepsilon,$$

while $\|\mathbf{a}\|_2 \leq \mathcal{O}(d_e^{5/2}(\log(d_e/\varepsilon)/\varepsilon)^{3/2}/\sqrt{m})$, $\|\mathbf{b}\|_\infty \leq 1$, and $\|\mathbf{W}^\top\|_{1,\infty} \leq 1$, completing the proof. \blacksquare

We can now specify $\mathbf{A}_1, \mathbf{W}_1$, and \mathbf{b}_1 in our construction.

Lemma 22 *For any $\varepsilon > 0$, let $\bar{m}_1 = \mathcal{O}(d_e^3(\log(d_e/\varepsilon)/\varepsilon)^2)$ and $m_1 = 2d + q\bar{m}_1$. Then, there exist $\mathbf{A}_1 \in \mathbb{R}^{(d+q) \times m_1}$, $\mathbf{W}_1 \in \mathbb{R}^{m_1 \times (d_n + d + (q+1)d_e)}$, and $\mathbf{b}_1 \in \mathbb{R}^{m_1}$, given by Equations (13) to (17), such that*

$$\chi_1^{\mathbf{x}} = \mathbf{x}, \quad |\chi_1^\omega(l) - \langle \omega_i, \omega_{t_l} \rangle| \leq \varepsilon,$$

for all $\mathbf{h} \in \mathbb{R}^{d_n}$, $\mathbf{x} \in \mathbb{R}^d$, $\omega_i, (\omega_{t_j})_{j \in [q]} \in \mathbb{S}^{d_e-1}$, and $l \in [q]$. Furthermore, we have the following guarantees

$$\|\mathbf{W}_1^\top\|_{1,\infty} \leq \mathcal{O}(1), \quad \|\mathbf{b}_1\|_\infty \leq \mathcal{O}(1), \quad \|\mathbf{A}_1^\top\|_{1,\infty} \leq \mathcal{O}(d_e^{5/2}(\log(d_e/\varepsilon)/\varepsilon)^{3/2}).$$

Proof We define the decompositions

$$\mathbf{W}_1 = \begin{pmatrix} \mathbf{W}_{11} \\ \mathbf{W}_{12} \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} \mathbf{b}_{11} \\ \mathbf{b}_{12} \end{pmatrix}, \quad \mathbf{A}_1 = \begin{pmatrix} \mathbf{A}_{11} \\ \mathbf{A}_{12} \end{pmatrix}, \quad (13)$$

where $\mathbf{W}_{11} \in \mathbb{R}^{2d \times (d_h + d + d_e)}$, $\mathbf{W}_{12} \in \mathbb{R}^{q\bar{m}_1 \times (d_h + d + d_e)}$, $\mathbf{b}_{11} \in \mathbb{R}^{2d}$, $\mathbf{b}_{12} \in \mathbb{R}^{q\bar{m}_1}$, $\mathbf{A}_{11} \in \mathbb{R}^{d \times m_1}$, and $\mathbf{A}_{12} \in \mathbb{R}^{q \times m_1}$. Let $\mathbf{v}_1, \dots, \mathbf{v}_d$ denote the standard basis of \mathbb{R}^d , and notice that $\sigma(z) - \sigma(-z) = z$. Therefore, we can implement the identity part of the mapping by letting

$$\mathbf{W}_{11} = \begin{pmatrix} \mathbf{0}_{d_h} & \mathbf{v}_1^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \mathbf{0}_{d_h} & -\mathbf{v}_1^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \vdots & \vdots & \\ \mathbf{0}_{d_h} & \mathbf{v}_d^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \mathbf{0}_{d_h} & -\mathbf{v}_d^\top & \mathbf{0}_{(q+1)d_e}^\top \end{pmatrix}, \quad (14)$$

as well as

$$\mathbf{b}_1 = \mathbf{0}_{2d}, \quad \text{and} \quad \mathbf{A}_{11} = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & \mathbf{0}_{q\bar{m}_1}^\top \\ 0 & 0 & 1 & -1 & \dots & 0 & \mathbf{0}_{q\bar{m}_1}^\top \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 & 1 & -1 & \mathbf{0}_{q\bar{m}_1}^\top \end{pmatrix} \quad (15)$$

Notice that $\|\mathbf{W}_{11}^\top\|_{1,\infty} = 1$ and $\|\mathbf{A}_{11}^\top\|_{1,\infty} = 2$. To implement the inner product part of the mapping, we take the construction of weights, biases, and second layer weights from Lemma 21, and rename them as $\tilde{\mathbf{W}}_1 \in \mathbb{R}^{\bar{m}_1 \times 2d_e}$, $\tilde{\mathbf{b}}_1 \in \mathbb{R}^{\bar{m}_1}$, and $\tilde{\mathbf{a}}_1 \in \mathbb{R}^{\bar{m}_1}$. Let us introduce the decomposition $\tilde{\mathbf{W}}_1 = (\tilde{\mathbf{W}}_{11} \quad \tilde{\mathbf{W}}_{12})$, where $\tilde{\mathbf{W}}_{11}, \tilde{\mathbf{W}}_{12} \in \mathbb{R}^{\bar{m}_1 \times d_e}$. With this decomposition, we can separate the projections applied to the first and second vectors in Lemma 21. We can then define

$$\mathbf{W}_{12} = \begin{pmatrix} \mathbf{0}_{\bar{m}_1 \times (d_h + d)} & \tilde{\mathbf{W}}_{11} & \tilde{\mathbf{W}}_{12} & \mathbf{0}_{\bar{m}_1 \times d_e} & \dots & \mathbf{0}_{\bar{m}_1 \times d_e} \\ \mathbf{0}_{\bar{m}_1 \times (d_h + d)} & \tilde{\mathbf{W}}_{11} & \mathbf{0}_{\bar{m}_1 \times d_e} & \tilde{\mathbf{W}}_{12} & \dots & \mathbf{0}_{\bar{m}_1 \times d_e} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{\bar{m}_1 \times (d_h + d)} & \tilde{\mathbf{W}}_{11} & \mathbf{0}_{\bar{m}_1 \times d_e} & \mathbf{0}_{\bar{m}_1 \times d_e} & \dots & \tilde{\mathbf{W}}_{12} \end{pmatrix}, \quad (16)$$

as well as

$$\mathbf{b}_{12} = \begin{pmatrix} \tilde{\mathbf{b}}_1 \\ \vdots \\ \tilde{\mathbf{b}}_1 \end{pmatrix}, \quad \text{and} \quad \mathbf{A}_{12} = \begin{pmatrix} \mathbf{0}_{2d}^\top & \tilde{\mathbf{a}}_1^\top & \mathbf{0}_{\bar{m}_1}^\top & \dots & \mathbf{0}_{\bar{m}_1}^\top \\ \mathbf{0}_{2d}^\top & \mathbf{0}_{\bar{m}_1}^\top & \tilde{\mathbf{a}}_1^\top & \dots & \mathbf{0}_{\bar{m}_1}^\top \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{2d}^\top & \mathbf{0}_{\bar{m}_1}^\top & \dots & \mathbf{0}_{\bar{m}_1}^\top & \tilde{\mathbf{a}}_1^\top \end{pmatrix}. \quad (17)$$

From Lemma 21, we have $\|\mathbf{W}_{12}^\top\|_{1,\infty} \leq 1$, $\|\mathbf{b}_{12}\|_\infty \leq 1$, and

$$\|\mathbf{A}_{12}^\top\|_{1,\infty} = \|\tilde{\mathbf{a}}_1\|_1 \leq \mathcal{O}(d_e^{5/2}(\log(d_e/\varepsilon)/\varepsilon)^{3/2}),$$

which completes the proof. \blacksquare

To introduce the construction of the next layer, we rely on the following lemma which establishes the desired approximation for a single coordinate, the proof of which is similar to that of Lemma 21.

Lemma 23 *Let σ be the ReLU activation. Suppose $|h| \leq r_\infty^h$, $|x| \leq r_\infty^x$ and $|z| \leq 1$. Let $R := \sqrt{1 + r_\infty^x{}^2 + r_\infty^h{}^2}$. For any $\varepsilon > 0$, there exists $m = \mathcal{O}(R^6(\log(R/\varepsilon)/\varepsilon)^3)$, $\mathbf{a} \in \mathbb{R}^m$,*

$\mathbf{W} \in \mathbb{R}^{m \times 2}$, and $\mathbf{b} \in \mathbb{R}^m$, such that

$$\sup_{|h| \leq r_\infty^h, |x| \leq r_\infty^x, |z| \leq 1} \left| h + 2x\sigma(z - 1/2) - \mathbf{a}^\top \sigma(\mathbf{W}(h, x, z)^\top + \mathbf{b}) \right| \leq \varepsilon$$

and

$$\|\mathbf{a}\|_2 \leq \mathcal{O}(R^6(\log(R/\varepsilon)/\varepsilon)^2/\sqrt{m}), \quad \|\mathbf{W}^\top\|_{1,\infty} \leq R^{-1}, \quad \|\mathbf{b}\|_\infty \leq 1.$$

Additionally, if $r_\infty^h = 0$, we have the improved bounds

$$m = \mathcal{O}(R^4(\log(R/\varepsilon)/\varepsilon)^2), \quad \|\mathbf{a}\|_2 \leq \mathcal{O}(R^5(\log(R/\varepsilon)/\varepsilon)^{3/2}/\sqrt{m})$$

Proof Note that $(h, x, z) \mapsto h + 2x\sigma(z - 1/2)$ is $2R$ -Lipschitz, and $|h + 2x\sigma(z - 1/2)| \leq R$. The proof follows from Lemma 45 with dimension 3 when $r_\infty^h \neq 0$ and dimension 2 otherwise. \blacksquare

With that, we can now construct the weights for the second mapping in the network.

Lemma 24 Suppose $\|\chi_1^x\|_\infty \leq r_x$ and $\max_l |\chi^\omega(l)| \leq 1$. Let $R := \sqrt{1 + r_x^2}$. Then, for every $\varepsilon > 0$ and absolute constant $\delta \in (0, 1)$, there exists $\bar{m}_2 \leq \mathcal{O}(R^4(\log(R/\varepsilon)/\varepsilon)^{3/2})$, $m_2 := q\delta\bar{m}_2$, and $\mathbf{A}_2 \in \mathbb{R}^{d_h \times m_2}$, $\mathbf{W}_2 \in \mathbb{R}^{m_2 \times (d+q)}$, and $\mathbf{b}_2 \in \mathbb{R}^{m_2}$ given by Equations (18) and (19) such that

$$\|\chi_2(l) - 2\chi_1^x \sigma(\chi_1^\omega(l) - 1/2)\|_\infty \leq \varepsilon,$$

for all such χ_1 and $l \in [q]$, where we recall $\chi_2 = \mathbf{A}_2 \sigma(\mathbf{W}_2 \chi_1 + \mathbf{b}_2)$. Moreover, we have

$$\|\mathbf{A}_2^\top\|_{1,\infty} \leq \mathcal{O}(R^4(\log(R/\varepsilon)/\varepsilon)^{3/2}), \quad \|\mathbf{W}_2^\top\|_{1,\infty} \leq R^{-1}, \quad \|\mathbf{b}_2\|_\infty \leq 1.$$

Proof Let $\tilde{\mathbf{W}} = (\tilde{w}_{21} \quad \tilde{w}_{22})$, $\tilde{\mathbf{b}}$, and $\tilde{\mathbf{a}}$ be the weights obtained from Lemma 23, where $\tilde{w}_{21}, \tilde{w}_{22}, \tilde{\mathbf{b}}, \tilde{\mathbf{a}} \in \mathbb{R}^{\bar{m}_2}$. To construct \mathbf{W}_2 and \mathbf{b}_2 , we let

$$\mathbf{W}_2 = \begin{pmatrix} \mathbf{W}_2(1, 1) \\ \vdots \\ \mathbf{W}_2(1, d) \\ \vdots \\ \mathbf{W}_2(q, 1) \\ \vdots \\ \mathbf{W}_2(q, d) \end{pmatrix}, \quad \mathbf{b}_2 = \begin{pmatrix} \mathbf{b}_2(1, 1) \\ \vdots \\ \mathbf{b}_2(1, d) \\ \vdots \\ \mathbf{b}_2(q, 1) \\ \vdots \\ \mathbf{b}_2(q, d) \end{pmatrix}. \quad (18)$$

where $\mathbf{W}_2(l, j) \in \mathbb{R}^{\bar{m}_2 \times (d+q)}$ is given by

$$\mathbf{W}_2(l, j) = (\mathbf{0}_{\bar{m}_2 \times (j-1)} \quad \tilde{w}_{21} \quad \mathbf{0}_{\bar{m}_2 \times (d-j)} \quad \mathbf{0}_{\bar{m}_2 \times (l-1)} \quad \tilde{w}_{22} \quad \mathbf{0}_{\bar{m}_2 \times (q-l)}),$$

and $\mathbf{b}_2(l, j) = \tilde{\mathbf{b}}_2$. Consequently, $\|\mathbf{W}_2^\top\|_{1,\infty} \leq 1$ and $\|\mathbf{b}_2\|_\infty \leq 1$. Finally, we have

$$\mathbf{A}_2 = \begin{pmatrix} \tilde{\mathbf{a}}_2^\top & \mathbf{0}_{\bar{m}_2}^\top & \cdots & \mathbf{0}_{\bar{m}_2}^\top \\ \mathbf{0}_{\bar{m}_2}^\top & \tilde{\mathbf{a}}_2^\top & \cdots & \mathbf{0}_{\bar{m}_2}^\top \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{\bar{m}_2}^\top & \cdots & \mathbf{0}_{\bar{m}_2}^\top & \tilde{\mathbf{a}}_2^\top \end{pmatrix}. \quad (19)$$

Consequently, we obtain $\|\mathbf{A}_2^\top\|_{1,\infty} \leq \mathcal{O}(R^4(\log(R/\varepsilon)/\varepsilon)^{3/2})$, completing the proof. \blacksquare

We are now ready to provide the four-layer feedforward construction of $f^\rightarrow(\mathbf{h}, \mathbf{x}, \mathbf{t}; \Theta_h^\rightarrow)$.

Proposition 25 *Let $\mathbf{z} = (\mathbf{x}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_{t_1}, \dots, \boldsymbol{\omega}_{t_q})$. Then, for every $\varepsilon > 0$, there exists a feedforward network with $L_h = 4$ layers given by*

$$f^{\rightarrow}(\mathbf{h}, \mathbf{z}; \boldsymbol{\Theta}_h^{\rightarrow}) = \mathbf{W}_{L_h} \sigma \left(\dots \sigma \left(\mathbf{W}_2 \sigma \left(\mathbf{W}_1 (\mathbf{h}^{\top}, \mathbf{z}^{\top})^{\top} + \mathbf{b}_1 \right) + \mathbf{b}_2 \right) \dots \right)$$

where $\mathbf{W}_i \in \mathbb{R}^{m_i \times m_{i-1}}$, $\mathbf{b}_i \in \mathbb{R}_i^{m_i}$ for $i \in \{2, \dots, L_h - 1\}$, $\mathbf{W}_1 \in \mathbb{R}^{m_1 \times d_h + d + (q+1)d_e}$, $\mathbf{b}_1 \in \mathbb{R}^{m_1}$, and $\mathbf{W}_{L_h} \in \mathbb{R}^{d_h \times m_{L_h-1}}$ that satisfies the following:

1. If $t_l = i$, then

$$\left\| f^{\rightarrow}(\mathbf{h}, \mathbf{z}; \hat{\boldsymbol{\Theta}}_h^{\rightarrow})_l - \mathbf{x} \right\|_2 \leq \varepsilon$$

2. Else $f^{\rightarrow}(\mathbf{h}, \mathbf{z}; \hat{\boldsymbol{\Theta}}_h^{\rightarrow})_l = \mathbf{0}_d$,

for all $l \in [q]$, $\mathbf{h} \in \mathbb{R}^{d_h}$ and $\|\mathbf{x}\|_2 \leq r_x$. Additionally $\|\mathbf{W}_i\|_F \leq \text{poly}(r_x, D_e, \varepsilon^{-1})$ for all $i \in [L_h]$ and $m_i, \|\mathbf{b}_i\|_2 \leq \text{poly}(r_x, D_e, \varepsilon^{-1})$ for all $i \in [L_h - 1]$, where we recall $D_e = d + (q+1)d_e$.

Proof Let $\tilde{\mathbf{A}}_1 \in \mathbb{R}^{(d+q) \times m_1}$, $\tilde{\mathbf{W}}_1 \in \mathbb{R}^{m_1 \times (d_h + d + (q+1)d_e)}$, $\tilde{\mathbf{b}}_1 \in \mathbb{R}^{m_1}$ be given by Lemma 22 with error parameter ε_1 and $\tilde{\mathbf{A}}_2 \in \mathbb{R}^{d_h \times m_2}$, $\tilde{\mathbf{W}}_2 \in \mathbb{R}^{m_2 \times (d+q)}$, $\tilde{\mathbf{b}}_2 \in \mathbb{R}^{m_2}$ be given by Lemma 24 with error parameter ε_2 . Recall that

$$\chi_1 = \tilde{\mathbf{A}}_1 \sigma(\tilde{\mathbf{W}}_1 \chi_0 + \tilde{\mathbf{b}}_1), \quad \chi_2 = \tilde{\mathbf{A}}_2 \sigma(\tilde{\mathbf{W}}_2 \chi_1 + \tilde{\mathbf{b}}_2).$$

By the triangle inequality,

$$\begin{aligned} \left\| \Psi(\mathbf{x}, t, i) - \tilde{\mathbf{A}}_2 \sigma(\tilde{\mathbf{W}}_2 \chi_1 + \tilde{\mathbf{b}}_2) \right\|_{\infty} &\leq \left\| \Psi(\mathbf{x}, t, i) - \tilde{\mathbf{A}}_2 \sigma(\tilde{\mathbf{W}}_2 \bar{\chi}_1 + \tilde{\mathbf{b}}_2) \right\|_{\infty} \\ &\quad + \left\| \tilde{\mathbf{A}}_2 \sigma(\tilde{\mathbf{W}}_2 \bar{\chi}_1 + \tilde{\mathbf{b}}_2) - \tilde{\mathbf{A}}_2 \sigma(\tilde{\mathbf{W}}_2 \chi_1 + \tilde{\mathbf{b}}_2) \right\|_{\infty} \\ &\leq \varepsilon_2 + \left\| \tilde{\mathbf{A}}_2^{\top} \right\|_{1, \infty} \left\| \tilde{\mathbf{W}}_2 \right\|_{1, \infty} \|\chi_1 - \bar{\chi}_1\|_{\infty} \\ &\leq \varepsilon_2 + \left\| \tilde{\mathbf{A}}_2 \right\|_{1, \infty} \left\| \tilde{\mathbf{W}}_2 \right\|_{1, \infty} \varepsilon_1, \end{aligned}$$

where $\bar{\chi}_1 = (\mathbf{x}^{\top}, \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_{t_1} \rangle, \dots, \langle \boldsymbol{\omega}_i, \boldsymbol{\omega}_{t_q} \rangle)^{\top}$. By letting $\varepsilon_2 = \varepsilon/4$, we obtain

$$m_2, \left\| \tilde{\mathbf{A}}_2 \right\|_F, \left\| \tilde{\mathbf{W}}_2 \right\|_F, \left\| \tilde{\mathbf{b}}_2 \right\|_2 \leq \text{poly}(r_x, D_e, \varepsilon^{-1}).$$

Similarly, we can let $\varepsilon_1 = \varepsilon/(4 \left\| \tilde{\mathbf{A}}_2 \right\|_{1, \infty} \left\| \tilde{\mathbf{W}}_2 \right\|_{1, \infty})$, which yields

$$m_1, \left\| \tilde{\mathbf{A}}_2 \right\|_F, \left\| \tilde{\mathbf{W}}_2 \right\|_F, \left\| \tilde{\mathbf{b}}_2 \right\|_2 \leq \text{poly}(r_x, D_e, \varepsilon^{-1}).$$

Let

$$\mathbf{W}_2 = \tilde{\mathbf{W}}_2 \tilde{\mathbf{A}}_1, \quad \mathbf{W}_1 = \tilde{\mathbf{W}}_1, \quad \mathbf{b}_1 = \tilde{\mathbf{b}}_1, \quad \mathbf{b}_2 = \tilde{\mathbf{b}}_2.$$

Then,

$$\chi_2 = \tilde{\mathbf{A}}_2 \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 (\mathbf{h}^{\top}, \mathbf{z}^{\top})^{\top} + \mathbf{b}_1) + \mathbf{b}_2),$$

satisfies $\|\chi_2 - \Psi(\mathbf{x}, \mathbf{t}, i)\|_\infty \leq \varepsilon/2$ for all $\|\mathbf{x}\|_2 \leq r_x$.

Recall that when $t_l \neq i$ for some $l \in [q]$, we would like to guarantee the output of the network to be equal to $\Psi(\mathbf{x}, \mathbf{t}, i)_l = \mathbf{0}_d$. To do so, we rely on the fact that $z \mapsto \sigma(z - b) - \sigma(-z - b)$ is zero for $|z| \leq b$, and has an L_∞ distance of b from the identity, i.e. $|z - \sigma(z - b) + \sigma(-z - b)| \leq b$. This mapping needs to be applied element-wise to χ_2 . Let $\tilde{\mathbf{W}}_3 \in \mathbb{R}^{2d_h \times d_h}$, $\mathbf{b}_3 \in \mathbb{R}^{2d_h}$, and $\mathbf{W}_4 \in \mathbb{R}^{d_h \times 2d_h}$ via

$$\tilde{\mathbf{W}}_3 = \begin{pmatrix} \mathbf{v}_1^\top \\ -\mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_d^\top \\ -\mathbf{v}_d^\top \end{pmatrix}, \quad \mathbf{b}_3 = -\frac{\varepsilon}{2} \mathbf{1}_{2d_h}, \quad \mathbf{W}_4 = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & -1 & \dots & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 1 & -1 \end{pmatrix}.$$

As a result, $\chi_3 = \mathbf{W}_4 \sigma(\tilde{\mathbf{W}}_3 \chi_2 + \mathbf{b}_3)$ satisfies

$$|(\chi_3)_j - (\chi_2)_j| \leq \begin{cases} 0 & |(\chi_2)_j| \leq \varepsilon/2 \\ \varepsilon/2 & |(\chi_2)_j| > \varepsilon/2 \end{cases}, \quad \forall j \in [d_h]. \quad (20)$$

We thus make two observations. First, $\|\chi_3 - \chi_2\|_\infty \leq \varepsilon/2$, and consequently $\|\chi_3(l) - \Psi(\mathbf{x}, \mathbf{t}, i)_l\|_\infty \leq \varepsilon$ for all $l \in [q]$. Second, when $t_l \neq i$, we have $\Psi(\mathbf{x}, \mathbf{t}, i)_l = \mathbf{0}_d$ and $|(\chi_2)_j| \leq \varepsilon/2$ for all $j \in [d]$ since $\|\chi_2(l) - \Psi(\mathbf{x}, \mathbf{t}, i)_l\|_\infty \leq \varepsilon/2$. Consequently, by the first case in (20), we have $\chi_3(l)_j = 0$ for all $j \in [d]$. We can summarize these two observations as follows

$$\|\chi_3(l) - \Psi(\mathbf{x}, \mathbf{t}, i)_l\|_\infty \leq \begin{cases} 0 & t_l \neq i \\ \varepsilon & t_l = i \end{cases},$$

which completes the proof. \blacksquare

With the above implementation of $f^\rightarrow(\mathbf{h}, \mathbf{z}; \Theta_h^\rightarrow)$, we have the following guarantee on \mathbf{h}_i^\rightarrow for all $i \in [N]$.

Corollary 26 *Let f_h^\rightarrow be given by the construction in Proposition 25, and suppose $r_h \geq \sqrt{q}(r_x + \sqrt{d}\varepsilon)$. Then, \mathbf{h}_i^\rightarrow satisfies the following guarantees for all $i \in [N]$ and $l \in [q]$:*

1. If $t_l \geq i$, then $\mathbf{h}_i^\rightarrow(l) = \mathbf{0}_d$
2. If $t_l < i$, then $\|\mathbf{h}_i^\rightarrow(l) - \mathbf{x}_{t_l}\|_\infty \leq \varepsilon$.

Proof We can prove the statement by induction. Note that it holds for $i = 1$ since $\mathbf{h}_1^\rightarrow = \mathbf{0}_d$. For the induction step, suppose it holds up to some i , and recall

$$\mathbf{h}_{i+1}^\rightarrow = \mathbf{h}_i^\rightarrow + f_h^\rightarrow(\mathbf{h}_i^\rightarrow, \mathbf{z}_i; \Theta_h^\rightarrow).$$

- If $t_l \geq i + 1$, then $\mathbf{h}_i^\rightarrow(l) = \mathbf{0}_d$ and $f_h^\rightarrow(\mathbf{h}_i^\rightarrow, \mathbf{z}_i; \Theta_h^\rightarrow) = \mathbf{0}_d$ by Proposition 25.
- If $t_l < i < i + 1$, then $\|\mathbf{h}_i^\rightarrow(l) - \mathbf{x}_{t_l}\|_\infty \leq \varepsilon$ by induction hypothesis, and $f_h^\rightarrow(\mathbf{h}_i^\rightarrow, \mathbf{z}_i; \Theta_h^\rightarrow) = \mathbf{0}_d$.

- Finally, if $t_l = i < i + 1$, then $\mathbf{h}_i^{\rightarrow}(l) = 0$ and $\|f_h^{\rightarrow}(\mathbf{h}_i^{\rightarrow}, \mathbf{z}_i; \Theta_h^{\rightarrow}) - \mathbf{x}_{t_l}\|_{\infty} \leq \varepsilon$.

Note that since $\|\mathbf{h}_j^{\rightarrow}\|_2 \leq r_h$ for all $j \in [N]$, the projection Π_{r_h} will always be identity through the forward pass, concluding the proof. \blacksquare

By symmetry, the same construction for f_h^{\leftarrow} would yield a similar guarantee on $\mathbf{h}_j^{\leftarrow}$. The last step is to design $f_y(\mathbf{h}^{\rightarrow}, \mathbf{h}^{\leftarrow}, \mathbf{z}; \Theta_y)$ such that

$$f_y(\mathbf{h}^{\rightarrow}, \mathbf{h}^{\leftarrow}, \mathbf{z}; \Theta_y) \approx g(\mathbf{h}^{\rightarrow} + \mathbf{h}^{\leftarrow} + (\mathbf{x}_i^{\top} \mathbb{1}[t_1 = i], \dots, \mathbf{x}_i^{\top} \mathbb{1}[t_q = i])^{\top}).$$

The following proposition provides the end-to-end RNN guarantee for approximating simple q STR models.

Proposition 27 *Suppose g satisfies Assumption 2. Then there exist RNN weights Θ_{RNN} with $\text{vec}(\Theta_{\text{RNN}}) \in \mathbb{R}^p$ (i.e. with p parameters) and $r_h \geq \sqrt{q}r_x + \sqrt{\varepsilon_{2\text{NN}}}/(r_a r_w)$, such that*

$$\sup_{i \in [N]} |g(\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_q}) - \hat{y}(\mathbf{p}; \Theta_{\text{RNN}})_i|^2 \leq 4\varepsilon_{2\text{NN}} \quad (21)$$

for all $\mathbf{t} \in [N]^q$ and $\|\mathbf{x}_j\|_2 \leq r_x$ for all $j \in [N]$. Additionally, we have

$$\|\text{vec}(\Theta_{\text{RNN}})\|_2 \leq \text{poly}(r_x, D_e, r_w, r_a, \varepsilon_{2\text{NN}}^{-1}), \quad p \leq \text{poly}(r_x, D_e, m_g, r_w, r_a, \varepsilon_{2\text{NN}}^{-1}), \quad (22)$$

and $f_h^{\rightarrow}, f_h^{\leftarrow}$ do not depend on \mathbf{h}^{\rightarrow} and \mathbf{h}^{\leftarrow} , namely the first d_h columns of $\mathbf{W}_1^{\rightarrow}$ and $\mathbf{W}_1^{\leftarrow}$ that are multiplied by \mathbf{h}^{\rightarrow} and \mathbf{h}^{\leftarrow} respectively are zero.

Proof As the proof of this proposition mostly follows from the previous proofs in this section, we only state the procedure for obtaining the desired weights.

Let $(\mathbf{v}_j)_{j=1}^{d_h}$ denote the standard basis of \mathbb{R}^{d_h} . Since $\sigma(z) - \sigma(-z) = z$, we can implement the identity mapping in \mathbb{R}^{d_h} via a two-layer feedforward network with the following weights

$$\mathbf{W}_{\text{id}} = \begin{pmatrix} \mathbf{v}_1^{\top} \\ -\mathbf{v}_1^{\top} \\ \vdots \\ \mathbf{v}_{d_h}^{\top} \\ -\mathbf{v}_{d_h}^{\top} \end{pmatrix}, \quad \mathbf{b}_{\text{id}} = \mathbf{0}_{2d_h}, \quad \mathbf{A}_{\text{id}} = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & -1 & \dots & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 1 & -1 \end{pmatrix},$$

where $\mathbf{W}_{\text{id}} \in \mathbb{R}^{2d_h \times d_h}$, $\mathbf{b}_{\text{id}} \in \mathbb{R}^{2d_h}$, and $\mathbf{A}_{\text{id}} \in \mathbb{R}^{d_h \times 2d_h}$. Let $\mathbf{W}_1, \mathbf{b}_1, \tilde{\mathbf{A}}_1, \tilde{\mathbf{W}}_2, \mathbf{b}_2, \tilde{\mathbf{A}}_2$ be given as in the proof of Proposition 25, for achieving an L_{∞} error of $\tilde{\varepsilon}$, to be fixed later. Recall $\mathbf{z}_i = (\mathbf{x}_i^{\top}, \boldsymbol{\omega}_i^{\top}, \boldsymbol{\omega}_{t_1}^{\top}, \dots, \boldsymbol{\omega}_{t_q}^{\top})^{\top}$. In the following, we remove the zero columns of \mathbf{W}_1 corresponding to the \mathbf{h} part of the input (see Lemma 22), which does not change the resulting function. Our construction can then be denoted by

$$\begin{array}{llll} \mathbf{h}_i^{\rightarrow} & \xrightarrow{\mathbf{A}_{\text{id}}\sigma(\mathbf{W}_{\text{id}}\cdot)} & \mathbf{h}_i^{\rightarrow} & \xrightarrow{\mathbf{A}_{\text{id}}\sigma(\mathbf{W}_{\text{id}}\cdot)} & \mathbf{h}_i^{\rightarrow} & \searrow \\ \mathbf{h}_i^{\leftarrow} & \xrightarrow{\mathbf{A}_{\text{id}}\sigma(\mathbf{W}_{\text{id}}\cdot)} & \mathbf{h}_i^{\leftarrow} & \xrightarrow{\mathbf{A}_{\text{id}}\sigma(\mathbf{W}_{\text{id}}\cdot)} & \mathbf{h}_i^{\leftarrow} & \rightarrow \\ \mathbf{z}_i & \xrightarrow{\tilde{\mathbf{A}}_1\sigma(\mathbf{W}_1+\mathbf{b}_1)} & \boldsymbol{\chi}_1 & \xrightarrow{\tilde{\mathbf{A}}_2\sigma(\tilde{\mathbf{W}}_2+\mathbf{b}_2)} & \boldsymbol{\chi}_2 & \nearrow \end{array} \quad \mathbf{h}_i^{\rightarrow} + \mathbf{h}_i^{\leftarrow} + \boldsymbol{\chi}_2 \xrightarrow{\mathbf{a}_g^{\top}\sigma(\mathbf{W}_g+\mathbf{b}_g)} \hat{y}_{\text{RNN}}(\mathbf{p}; \Theta_{\text{RNN}})_i$$

Note that the addition above can be implemented exactly by using the fact that $\sigma(z_1 + z_2 + z_3) - \sigma(-z_1 - z_2 - z_3) = z_1 + z_2 + z_3$. Specifically, the weights of this layer are given by

$$\mathbf{W}_{\text{add}} = \begin{pmatrix} \mathbf{v}_1^\top & \mathbf{v}_1^\top & \mathbf{v}_1^\top \\ -\mathbf{v}_1^\top & -\mathbf{v}_1^\top & -\mathbf{v}_1^\top \\ \vdots & \vdots & \vdots \\ \mathbf{v}_{d_h}^\top & \mathbf{v}_{d_h}^\top & \mathbf{v}_{d_h}^\top \\ -\mathbf{v}_{d_h}^\top & -\mathbf{v}_{d_h}^\top & -\mathbf{v}_{d_h}^\top \end{pmatrix}, \quad \mathbf{b}_{\text{add}} = \mathbf{0}_{2d_h}, \quad \mathbf{A}_{\text{add}} = \mathbf{A}_{\text{id}},$$

where $\mathbf{W}_{\text{add}} \in \mathbb{R}^{2d_h \times 3d_h}$, $\mathbf{b}_{\text{add}} \in \mathbb{R}^{2d_h}$, $\mathbf{A}_{\text{add}} \in \mathbb{R}^{d_h \times 2d_h}$.

Let Θ_h^\rightarrow (and similarly Θ_h^\leftarrow) be given by Proposition 25 with corresponding error ε_h . Using the shorthand notation $\mathbf{x}_t = (\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_q}) \in \mathbb{R}^{dq}$ and $\hat{\mathbf{x}}_t = \mathbf{h}_i^\rightarrow + \mathbf{h}_i^\leftarrow + \chi_2$, we have

$$\begin{aligned} \|\mathbf{h}_i^\rightarrow + \mathbf{h}_i^\leftarrow + \chi_2 - \hat{\mathbf{x}}_t\|_2 &\leq \left\| \mathbf{h}_i^\rightarrow - \sum_{j=1}^{i-1} \Psi(\mathbf{x}_j, \mathbf{t}, j) \right\|_2 + \left\| \mathbf{h}_i^\leftarrow - \sum_{j=N}^{i+1} \Psi(\mathbf{x}_j, \mathbf{t}, j) \right\|_2 + \|\chi_2 - \Psi(\mathbf{x}_i, \mathbf{t}, i)\|_2 \\ &\leq \sqrt{qd}(2\varepsilon_h + \tilde{\varepsilon}), \end{aligned}$$

which holds for all input prompts \mathbf{p} with $\|\mathbf{x}_j\|_2 \leq r_x$ for all $j \in [N]$. Finally, we have

$$\begin{aligned} \sup_{\|\mathbf{x}_j\|_2 \leq r_x, \forall j \in [N]} \left| g(\mathbf{x}_t) - \mathbf{a}_g^\top \sigma(\mathbf{W}_g \hat{\mathbf{x}}_t + \mathbf{b}_g) \right| &\leq \sup_{\|\mathbf{x}_j\|_2 \leq r_x, \forall j \in [N]} \left| g(\mathbf{x}_t) - \mathbf{a}_g^\top \sigma(\mathbf{W}_g \mathbf{x}_t + \mathbf{b}_g) \right| \\ &\quad + \sup_{\|\mathbf{x}_j\|_2 \leq r_x, \forall j \in [N]} \left| \mathbf{a}_g^\top \sigma(\mathbf{W}_g \mathbf{x}_t + \mathbf{b}_g) - \mathbf{a}_g^\top \sigma(\mathbf{W}_g \hat{\mathbf{x}}_t + \mathbf{b}_g) \right| \\ &\leq \sqrt{\varepsilon_{2\text{NN}}} + r_a r_w \sqrt{qd}(2\varepsilon_h + \tilde{\varepsilon}). \end{aligned}$$

Choosing $\varepsilon_h = \sqrt{\varepsilon_{2\text{NN}}}/(4\sqrt{qd}r_a r_w)$ and $\tilde{\varepsilon} = \sqrt{\varepsilon_{2\text{NN}}}/(2\sqrt{qd}r_a r_w)$, we obtain RNN weights that satisfy $\|\text{vec}(\Theta_{\text{RNN}})\|_2 \leq \text{poly}(r_x, D_e, r_a, r_w, \varepsilon_{2\text{NN}}^{-1})$, completing the proof. \blacksquare

C.4. Generalization Upper Bounds for RNNs

Recall the state transitions

$$\begin{aligned} \mathbf{h}_{j+1}^\rightarrow &= \Pi_{r_h}(\mathbf{h}_j^\rightarrow + f_h^\rightarrow(\mathbf{h}_j^\rightarrow, \mathbf{z}_j; \Theta_h^\rightarrow)) \\ \mathbf{h}_{j-1}^\leftarrow &= \Pi_{r_h}(\mathbf{h}_j^\leftarrow + f_h^\leftarrow(\mathbf{h}_j^\leftarrow, \mathbf{z}_j; \Theta_h^\leftarrow)). \end{aligned}$$

We will use the notation $\mathbf{h}_j^\rightarrow(\mathbf{p}; \Theta_h^\rightarrow)$ and $\mathbf{h}_j^\leftarrow(\mathbf{p}; \Theta_h^\leftarrow)$ to highlight the dependence of the hidden states on the prompt \mathbf{p} and parameters Θ_h^\rightarrow and Θ_h^\leftarrow . We then define the prediction function as $F(\mathbf{p}; \Theta_h^\rightarrow, \Theta_h^\leftarrow, \Theta_y)$ where

$$F(\mathbf{p}; \Theta_h^\rightarrow, \Theta_h^\leftarrow, \Theta_y)_j = f_y(\mathbf{h}_j^\rightarrow(\mathbf{p}; \Theta_h^\rightarrow), \mathbf{h}_j^\leftarrow(\mathbf{p}; \Theta_h^\leftarrow), \mathbf{z}_j; \Theta_y).$$

We can now define the function class

$$\mathcal{F}_{\text{RNN}} = \{\mathbf{p}, j \mapsto F(\mathbf{p}; \Theta_h^\rightarrow, \Theta_h^\leftarrow, \Theta_y)_j : \Theta_h^\rightarrow, \Theta_h^\leftarrow, \Theta_y \in \Theta_{\text{RNN}}\}.$$

We can then define our distance function by going over $\{\mathbf{p}, j \in S_n\}$,

$$d_\infty(F, \hat{F}) = \sup_{\mathbf{p}, j \in S_n} \left| F(\mathbf{p}; \boldsymbol{\Theta}_h^\rightarrow, \boldsymbol{\Theta}_h^\leftarrow, \boldsymbol{\Theta}_y)_j - F(\mathbf{p}; \hat{\boldsymbol{\Theta}}_h^\rightarrow, \hat{\boldsymbol{\Theta}}_h^\leftarrow, \boldsymbol{\Theta}_y)_j \right|.$$

We will further use the notation

$$f_y(\cdot; \boldsymbol{\Theta}_y) = \mathbf{W}_{L_y}^y \sigma(\mathbf{W}_{L_y-1}^y \dots \sigma(\mathbf{W}_{L_1}^1(\cdot) + \mathbf{b}_1^y) \dots + \mathbf{b}_{L_y-1}^y) \in \mathcal{F}_{\text{NN}, L_y}^y,$$

and

$$f_h^\rightarrow(\cdot; \boldsymbol{\Theta}_h^\rightarrow) = \mathbf{W}_{L_h}^\rightarrow \sigma(\mathbf{W}_{L_h-1}^\rightarrow \dots \sigma(\mathbf{W}_1^\rightarrow(\cdot) + \mathbf{b}_1^\rightarrow) \dots + \mathbf{b}_{L_h-1}^\rightarrow) \in \mathcal{F}_{\text{NN}, L_h}^\rightarrow.$$

We similarly define $\mathcal{F}_{\text{NN}, L_h}^\leftarrow$. The covering number of \mathcal{F}_{RNN} can be related to that of $\mathcal{F}_{\text{NN}, L_y}^y$, $\mathcal{F}_{\text{NN}, L_h}^\rightarrow$, and $\mathcal{F}_{\text{NN}, L_h}^\leftarrow$, through the following lemma.

Lemma 28 *Suppose for every $\boldsymbol{\Theta}_h^\rightarrow, \boldsymbol{\Theta}_h^\leftarrow, \boldsymbol{\Theta}_y \in \Theta_{\text{RNN}}$ we have*

$$\left\| \mathbf{W}_{L_y}^y \dots \mathbf{W}_1^y \right\|_{\text{op}} \leq C_W^y, \quad \left\| \mathbf{W}_{L_h}^\rightarrow \right\|_{\text{op}} \dots \left\| \mathbf{W}_{1,h}^\rightarrow \right\|_{\text{op}} \leq \alpha_N, \quad \left\| \mathbf{W}_{L_h}^\leftarrow \right\|_{\text{op}} \dots \left\| \mathbf{W}_{1,h}^\leftarrow \right\|_{\text{op}} \leq \alpha_N,$$

where $\alpha_N \leq N^{-1}$. Then,

$$\begin{aligned} \log \mathcal{C}(\mathcal{F}_{\text{RNN}}, d_\infty, \epsilon) &\leq \log \mathcal{C}(\mathcal{F}_{\text{NN}, L_y}^y, d_\infty, \epsilon/2) + \log \mathcal{C}\left(\mathcal{F}_{\text{NN}, L_h}^\rightarrow, d_\infty, \frac{\epsilon}{4eC_W^y N}\right) \\ &\quad + \log \mathcal{C}\left(\mathcal{F}_{\text{NN}, L_h}^\leftarrow, d_\infty, \frac{\epsilon}{4eC_W^y N}\right) \end{aligned}$$

Proof Throughout the proof, we will use the shorthand notation $\mathbf{h}_j^\rightarrow = \mathbf{h}_j^\rightarrow(\mathbf{p}; \boldsymbol{\Theta}_h^\rightarrow)$ and $\hat{\mathbf{h}}_j^\rightarrow = \mathbf{h}_j^\rightarrow(\mathbf{p}; \hat{\boldsymbol{\Theta}}_h^\rightarrow)$, with similarly define \mathbf{h}_j^\leftarrow and $\hat{\mathbf{h}}_j^\leftarrow$. We begin by observing

$$\sup_{\mathbf{p}, j \in S_n} \left| f_y(\mathbf{h}_j^\rightarrow, \mathbf{h}_j^\leftarrow, \mathbf{z}_j; \boldsymbol{\Theta}_y) - f_y(\hat{\mathbf{h}}_j^\rightarrow, \hat{\mathbf{h}}_j^\leftarrow, \mathbf{z}_j; \hat{\boldsymbol{\Theta}}_y) \right| \leq \mathcal{E}_1 + \mathcal{E}_2$$

where

$$\begin{aligned} \mathcal{E}_1 &:= \sup_{\mathbf{p}, j \in S_n} \left| f_y(\mathbf{h}_j^\rightarrow, \mathbf{h}_j^\leftarrow, \mathbf{z}_j; \boldsymbol{\Theta}_y) - f_y(\mathbf{h}_j^\rightarrow, \mathbf{h}_j^\leftarrow, \mathbf{z}_j; \hat{\boldsymbol{\Theta}}_y) \right| \\ \mathcal{E}_2 &:= \sup_{\mathbf{p}, j \in S_n} \left| f_y(\mathbf{h}_j^\rightarrow, \mathbf{h}_j^\leftarrow, \mathbf{z}_j; \hat{\boldsymbol{\Theta}}_y) - f_y(\hat{\mathbf{h}}_j^\rightarrow, \hat{\mathbf{h}}_j^\leftarrow, \mathbf{z}_j; \hat{\boldsymbol{\Theta}}_y) \right|. \end{aligned}$$

Then, we observe that $\mathcal{E}_1 = d_\infty(f_y(\cdot; \boldsymbol{\Theta}_y), f_y(\cdot; \hat{\boldsymbol{\Theta}}_y))$. Thus, we can ensure $\mathcal{E}_1 \leq \epsilon/2$ with a covering $\{\hat{\boldsymbol{\Theta}}_y\}$ of size $\mathcal{C}(\mathcal{F}_{\text{NN}, L_y}^y, d_\infty, \epsilon/2)$. Hence, we move to \mathcal{E}_2 .

Using the Lipschitzness of f_y , we obtain

$$\begin{aligned} \mathcal{E}_2 &\leq \left\| \mathbf{W}_{L_y}^y \dots \mathbf{W}_1^y \right\|_{\text{op}} \left(\sup_{\mathbf{p}, j} \left\| \mathbf{h}_j^\rightarrow - \hat{\mathbf{h}}_j^\rightarrow \right\|_2 + \sup_{\mathbf{p}, j} \left\| \mathbf{h}_j^\leftarrow - \hat{\mathbf{h}}_j^\leftarrow \right\|_2 \right) \\ &\leq C_W^y \left(\sup_{\mathbf{p}, j} \left\| \mathbf{h}_j^\rightarrow - \hat{\mathbf{h}}_j^\rightarrow \right\|_2 + \sup_{\mathbf{p}, j} \left\| \mathbf{h}_j^\leftarrow - \hat{\mathbf{h}}_j^\leftarrow \right\|_2 \right). \end{aligned}$$

Further, by Lipschitzness of Π_{r_h} , we have

$$\begin{aligned} \sup_{\mathbf{p},j} \left\| \mathbf{h}_j^{\rightarrow} - \hat{\mathbf{h}}_j^{\rightarrow} \right\|_2 &\leq \sup_{\mathbf{p},j} \left\| \mathbf{h}_{j-1}^{\rightarrow} - \hat{\mathbf{h}}_{j-1}^{\rightarrow} \right\|_2 + \underbrace{\sup_{\mathbf{p},j} \left\| f_h^{\rightarrow}(\mathbf{h}_{j-1}^{\rightarrow}, \mathbf{z}_{j-1}; \hat{\Theta}_h^{\rightarrow}) - f_h^{\rightarrow}(\hat{\mathbf{h}}_{j-1}^{\rightarrow}, \mathbf{z}_{j-1}; \hat{\Theta}_h^{\rightarrow}) \right\|_2}_{=:\mathcal{E}_1^h} \\ &\quad + \underbrace{\sup_{\mathbf{p},j} \left\| f_h^{\rightarrow}(\mathbf{h}_{j-1}^{\rightarrow}, \mathbf{z}_{j-1}; \Theta_h^{\rightarrow}) - f_h^{\rightarrow}(\mathbf{h}_{j-1}^{\rightarrow}, \mathbf{z}_{j-1}; \hat{\Theta}_h^{\rightarrow}) \right\|_2}_{=:\mathcal{E}_2^h}. \end{aligned}$$

By the Lipschitzness of f_h^{\rightarrow} , for the second term we have

$$\mathcal{E}_1^h \leq \left\| \hat{\mathbf{W}}_{L_h}^{\rightarrow} \dots \hat{\mathbf{W}}_{1,h}^{\rightarrow} \right\|_{\text{op}} \left\| \mathbf{h}_{j-1}^{\rightarrow} - \hat{\mathbf{h}}_{j-1}^{\rightarrow} \right\|_2 \leq \alpha_N \left\| \mathbf{h}_{j-1}^{\rightarrow} - \hat{\mathbf{h}}_{j-1}^{\rightarrow} \right\|_2.$$

Moreover, we have $\mathcal{E}_2^h \leq d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow}))$. Consequently, we obtain

$$\begin{aligned} \sup_{\mathbf{p},j} \left\| \mathbf{h}_j^{\rightarrow} - \hat{\mathbf{h}}_j^{\rightarrow} \right\|_2 &\leq (1 + \alpha_N) \sup_{\mathbf{p},j} \left\| \mathbf{h}_{j-1}^{\rightarrow} - \hat{\mathbf{h}}_{j-1}^{\rightarrow} \right\|_2 + d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow})) \\ &\leq \sum_{l=0}^{j-2} (1 + \alpha_N)^l d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow})) \\ &\leq \frac{(1 + \alpha_N)^{j-1} - 1}{\alpha_N} d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow})) \\ &\leq eN d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow})). \end{aligned}$$

We can similarly obtain an upper bound on $\sup_{\mathbf{p},j} \left\| \mathbf{h}_j^{\leftarrow} - \hat{\mathbf{h}}_j^{\leftarrow} \right\|_2$. Hence, we have

$$\mathcal{E}_2 \leq eC_w^y N \left\{ d_{\infty}(f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow}), f_h^{\rightarrow}(\cdot; \hat{\Theta}_h^{\rightarrow})) + d_{\infty}(f_h^{\leftarrow}(\cdot; \Theta_h^{\leftarrow}), f_h^{\leftarrow}(\cdot; \hat{\Theta}_h^{\leftarrow})) \right\}.$$

Therefore, by constructing $\epsilon/(2eC_w^y N)$ coverings $\{\hat{\Theta}_h^{\rightarrow}\}$ and $\{\hat{\Theta}_h^{\leftarrow}\}$ which have sizes

$$\mathcal{C}(\mathcal{F}_{\text{NN},L_h}^{\rightarrow}, \epsilon/(4eC_w^y N)), \quad \text{and,} \quad \mathcal{C}(\mathcal{F}_{\text{NN},L_h}^{\leftarrow}, \epsilon/(4eC_w^y N))$$

respectively, we complete the covering of \mathcal{F}_{RNN} . ■

The next step is to bound the covering number of the class of feedforward networks, as performed by the following lemma.

Lemma 29 *Let*

$$\mathcal{F}_{\text{NN},L} = \{\mathbf{x} \mapsto \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\dots \mathbf{W}_2(\sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)) \dots + \mathbf{b}_{L-1})) : \Theta_{\text{NN}} \in \Theta_{\text{NN}}\},$$

where $\Theta_{\text{NN}} = (\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_{L-1}, \mathbf{b}_{L-1}, \mathbf{W}_L)$ and $\text{vec}(\Theta_{\text{NN}}) \in \mathbb{R}^p$. Further, define the distance function

$$d_{\infty}(f, f') = \sup_{\|\mathbf{x}\| \leq R} |f(\mathbf{x}) - f'(\mathbf{x})|, \quad \forall f, f' \in \mathcal{F}_{\text{NN},L}.$$

Suppose $\|\mathbf{W}_l\|_F, \|\mathbf{b}_l\|_2 \leq R$ for all l . Then, for any absolute constant depth $L = \mathcal{O}(1)$, we have

$$\log \mathcal{C}(\mathcal{F}_{\text{NN},L}, d_{\infty}, \epsilon) \leq p \log(1 + \text{poly}(R)/\epsilon).$$

Proof Let $\mathbf{x}_0 = \mathbf{x}$, $\mathbf{x}_l = \sigma(\mathbf{W}_l \mathbf{x}_{l-1} + \mathbf{b}_l)$ for $l \in [L-1]$, and $\mathbf{x}_L = \mathbf{W}_L \mathbf{x}_{L-1}$. Also let $(\hat{\mathbf{x}}_l)$ be the corresponding definitions under weights and biases $(\hat{\mathbf{W}}_l)$ and $(\hat{\mathbf{b}}_l)$. First, we remark that for $l \in [L-1]$,

$$\|\mathbf{x}_l\|_2 \leq \|\mathbf{W}_l\|_{\text{op}} \|\mathbf{x}_{l-1}\|_2 + \|\mathbf{b}_l\|_2 \quad (23)$$

$$\begin{aligned} &\leq \prod_{i=1}^l \|\mathbf{W}_i\|_{\text{op}} \|\mathbf{x}_0\|_2 + \sum_{i=0}^{l-1} \|\mathbf{b}_{l-i-1}\|_2 \prod_{j=0}^i \|\mathbf{W}_{l-j}\|_{\text{op}} + \|\mathbf{b}_l\|_2 \\ &\leq \text{poly}(R), \end{aligned} \quad (24)$$

where we used the fact that L is an absolute constant. Next, for $l \in [L-1]$, we have

$$\begin{aligned} \|\mathbf{x}_l - \hat{\mathbf{x}}_l\|_2 &\leq \left\| \mathbf{W}_l \mathbf{x}_{l-1} - \hat{\mathbf{W}}_l \hat{\mathbf{x}}_{l-1} \right\|_2 + \left\| \mathbf{b}_l - \hat{\mathbf{b}}_l \right\|_2 \\ &\leq \|\mathbf{W}_l\|_{\text{op}} \|\mathbf{x}_{l-1} - \hat{\mathbf{x}}_{l-1}\|_2 + \|\hat{\mathbf{x}}_{l-1}\|_2 \|\mathbf{W}_l - \hat{\mathbf{W}}_l\|_{\text{op}} + \|\mathbf{b}_l - \hat{\mathbf{b}}_l\|_2 \\ &\leq \text{poly}(R) \left\{ \|\mathbf{x}_{l-1} - \hat{\mathbf{x}}_{l-1}\|_2 + \|\mathbf{W}_l - \hat{\mathbf{W}}_l\|_{\text{F}} + \|\mathbf{b}_l - \hat{\mathbf{b}}_l\|_2 \right\}. \end{aligned}$$

Once again, using the fact that L is an absolute constant and by expanding the above inequality, we obtain

$$\|\mathbf{x}_l - \hat{\mathbf{x}}_l\|_2 \leq \text{poly}(R) \left\{ \sum_{i=1}^l \|\mathbf{W}_i - \hat{\mathbf{W}}_i\|_{\text{F}} + \|\mathbf{b}_i - \hat{\mathbf{b}}_i\|_2 \right\}.$$

Finally, we have the bound

$$\begin{aligned} \|\mathbf{x}_L - \hat{\mathbf{x}}_L\|_2 &\leq \|\mathbf{W}_L\|_{\text{op}} \|\mathbf{x}_{L-1} - \hat{\mathbf{x}}_{L-1}\|_2 + \|\hat{\mathbf{x}}_{L-1}\|_2 \|\mathbf{W}_L - \hat{\mathbf{W}}_L\|_{\text{op}} \\ &\leq \text{poly}(R) \left\| \text{vec}(\Theta_{\text{NN}}) - \text{vec}(\hat{\Theta}_{\text{NN}}) \right\|_2. \end{aligned}$$

Consequently, we have

$$\begin{aligned} \log \mathcal{C}(\mathcal{F}_{\text{NN},L}, d_\infty, \epsilon) &\leq \log \mathcal{C}(\{\Theta \in \mathbb{R}^p : \|\Theta\|_2 \leq \text{poly}(d, q)\}, \|\cdot\|_2, \epsilon / \text{poly}(R)) \\ &\leq p \log(1 + \text{poly}(R)/\epsilon), \end{aligned}$$

where the last inequality follows from Lemma 44. ■

Therefore, we immediately obtain the following bound on the covering number of \mathcal{F}_{RNN} .

Corollary 30 *Suppose $\Theta_{\text{RNN}} \subseteq \{\Theta \in \mathbb{R}^p : \|\text{vec}(\Theta)\|_2 \leq R\}$ and $\|\mathbf{z}_j^{(i)}\|_2 \leq R$ for all $i \in [n]$ and $j \in [N]$. Then,*

$$\log \mathcal{C}(\mathcal{F}_{\text{RNN}}, d_\infty, \epsilon) \leq p \log(1 + \text{poly}(R)N/\epsilon).$$

We can now proceed with standard Rademacher complexity based arguments. Similar to the argument in Appendix B.2, we define a truncated version of the loss by considering the loss class

$$\mathcal{L}_\tau^{\text{RNN}} = \{(\mathbf{p}, \mathbf{y}, j) \mapsto (f_{\text{RNN}}(\mathbf{p})_j - y_j)^2 \wedge \tau : f_{\text{RNN}} \in \mathcal{F}_{\text{RNN}}\},$$

where the constant $\tau > 0$ will be chosen later. We then have the following bound on the empirical Rademacher complexity of $\mathcal{L}_\tau^{\text{RNN}}$.

Lemma 31 *In the same setting as Corollary 30 and with $\tau \geq 1$, we have*

$$\hat{\mathfrak{R}}_n(\mathcal{L}_\tau^{\text{RNN}}) \leq \mathcal{O}\left(\tau \sqrt{\frac{p \log(RNn\tau)}{n}}\right).$$

Proof By a standard discretization bound for Rademacher complexity, for all $\epsilon > 0$ we have

$$\begin{aligned} \hat{\mathfrak{R}}_n(\mathcal{L}_\tau^{\text{RNN}}) &\leq \epsilon + \tau \sqrt{\frac{2 \log \mathcal{C}(\mathcal{L}_\tau^{\text{RNN}}, d_\infty, \epsilon)}{n}} \\ &\leq \epsilon + \tau \sqrt{\frac{2 \log \mathcal{C}(\mathcal{F}_{\text{RNN}}, d_\infty, \epsilon/(2\sqrt{\tau}))}{n}} \\ &\leq \epsilon + \tau \sqrt{\frac{2p \log(1 + \text{poly}(R)N\sqrt{\tau}/\epsilon)}{n}}, \end{aligned}$$

where the second inequality follows from Lipschitzness of $(\cdot)^2 \wedge \tau$. We conclude the proof by choosing $\epsilon = 1/\sqrt{n}$. \blacksquare

We can directly turn the above bound on the empirical Rademacher complexity into a bound on generalization gap.

Corollary 32 *Let $\hat{\Theta} = \arg \min_{\Theta \in \Theta_{\text{RNN}}} \hat{R}_n^{\text{RNN}}(\Theta)$. Suppose $\Theta_{\text{RNN}} \subseteq \{\Theta \in \mathbb{R}^p : \|\text{vec}(\Theta)\|_2 \leq R\}$, and additionally $\sqrt{3C_x \text{ed} \log(nN) + q + 1} \leq R$. Then, for every $\delta > 0$, with probability at least $1 - \delta - (nN)^{-1/2}$ over the training set, we have*

$$R_\tau^{\text{RNN}}(\hat{\Theta}) - \hat{R}_\tau^{\text{RNN}}(\hat{\Theta}) \leq \mathcal{O}\left(\tau \sqrt{\frac{p \log(RNn\tau)}{n}} + \tau \sqrt{\frac{\log(1/\delta)}{n}}\right).$$

Proof We highlight that for the specified R , Lemma 11 guarantees $\|z_j^{(i)}\|_2 \leq R$ for all $i \in [n]$ and $j \in [N]$ with probability at least $1 - (nN)^{-1/2}$. Standard Rademacher complexity generalization arguments applied to Lemma 31 complete the proof. \blacksquare

Note that $\hat{R}_\tau^{\text{RNN}}(\hat{\Theta}) \leq \hat{R}_n^{\text{RNN}}(\hat{\Theta})$ which is further controlled in the approximation section by Proposition 27. Therefore, the last step is to demonstrate that choosing $\tau = \text{poly}(d, q, \log n)$ suffices to achieve a desirable bound on $R^{\text{RNN}}(\hat{\Theta})$ through $R_\tau^{\text{RNN}}(\hat{\Theta})$.

Lemma 33 *Consider the setting of Corollary 32, and additionally assume $R \geq r_h$. Then, for some $\tau = \text{poly}(R, \log n)$, we have*

$$R^{\text{RNN}}(\hat{\Theta}) - R_\tau^{\text{RNN}}(\hat{\Theta}) \leq \sqrt{\frac{1}{n}}.$$

Proof The proof of this lemma proceeds similarly to the proof of Lemma 19. By defining

$$\Delta_y := \left| \hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j - y_j \right|$$

and following the same steps (where we recall $j \sim \text{Unif}([N])$), we obtain

$$\begin{aligned} R^{\text{RNN}}(\hat{\Theta}) &= \mathbb{E}[\Delta_y^2 \mathbb{1}[\Delta_y \leq \sqrt{\tau}]] + \mathbb{E}[\Delta_y^2 \mathbb{1}[\Delta_y > \sqrt{\tau}]] \\ &\leq R_\tau^{\text{RNN}}(\hat{\Theta}) + \mathbb{E}[\Delta_y^4]^{1/2} \mathbb{P}(\Delta_y \geq \sqrt{\tau})^{1/2}, \end{aligned}$$

where

$$\mathbb{E}[\Delta_y^4]^{1/2} \leq 2 \mathbb{E}[y_j^4]^{1/2} + 2 \mathbb{E}[\hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j^4]^{1/2}$$

and

$$\mathbb{P}(\Delta_y > \sqrt{\tau}) \leq \mathbb{P}\left(|y_j| \geq \frac{\sqrt{\tau}}{2}\right) + \mathbb{P}\left(|\hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j| \geq \frac{\sqrt{\tau}}{2}\right)$$

From Assumption 1, we have $\mathbb{E}[y_j^4]^{1/2} \lesssim 1$ and $\mathbb{P}(|y_j| \geq \sqrt{\tau}/2) \leq e^{-\Omega(\tau^{1/s})}$. For the prediction of the RNN, we have the following bound (see (24) for the derivation)

$$\left| \hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j \right| \leq \prod_{l=1}^{L_y} \|\mathbf{W}_l^y\|_{\text{op}} \|\mathbf{h}_j^{\rightarrow}, \mathbf{h}_j^{\leftarrow}, \mathbf{z}_j\|_2 + \sum_{i=0}^{L_y-1} \|\mathbf{b}_{L_y-i-1}^y\|_2 \prod_{l=0}^i \|\mathbf{W}_{L_y-l}^y\|_{\text{op}}.$$

As a result,

$$\left| \hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j \right| \leq \text{poly}(R)(1 + r_h + \|\mathbf{z}_j\|).$$

As a result, by the fact that $r_h \leq R$ and Assumption 1, after taking an expectation, we immediately have

$$\mathbb{E}[\hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j^4]^{1/2} \leq \text{poly}(R).$$

On the other hand, from Lemma 11 (with $n = N = 1$), we obtain

$$\mathbb{P}\left(\left| \hat{y}_{\text{RNN}}(\mathbf{p}; \hat{\Theta})_j \right| \geq \frac{\sqrt{\tau}}{2}\right) \leq e^{-\Omega(\tau/\text{poly}(R))}$$

Therefore, for some $\tau = \text{poly}(R, \log n)$ we can obtain the bound stated in the lemma. \blacksquare

We can summarize the above facts into the proof of Theorem 6.

Proof of Theorem 6. From the approximation bound of Proposition 27, we know that for some $R = \text{poly}(d, q, r_a, r_w, \varepsilon_{2\text{NN}}^{-1}, \log(nN))$ and the constraint set

$$\Theta_{\text{RNN}} = \left\{ \Theta : \|\text{vec}(\Theta)\|_2 \leq R, \|\mathbf{W}_{L_h}^{\rightarrow}\|_{\text{op}} \dots \|\mathbf{W}_{1,h}^{\rightarrow}\|_{\text{op}} \leq \alpha_N, \|\mathbf{W}_{L_h}^{\leftarrow}\|_{\text{op}} \dots \|\mathbf{W}_{1,h}^{\leftarrow}\|_{\text{op}} \leq \alpha_N \right\}$$

with any $\alpha_N \leq N^{-1}$, we have $\hat{R}^{\text{RNN}}(\hat{\Theta}) \lesssim \varepsilon_{2\text{NN}}$. The proof is then completed by letting $r_h = \sqrt{q}r_x + \sqrt{\varepsilon_{2\text{NN}}}/(r_a r_w)$, invoking the generalization bound of Corollary 32, and the bound on truncation error given in Lemma 33, with $R = \text{poly}(d, q, r_a, r_w, \varepsilon_{2\text{NN}}^{-1}, \log(nN))$. \blacksquare

C.5. RNN Lower Bound Formulation and Details

For our lower bound, we will consider a broad class of recurrent networks, without restricting to a specific form of parametrization. Specifically, we consider bidirectional RNNs characterized by

$$\begin{aligned} \mathbf{h}_{i+1}^{\rightarrow} &= \text{proj}_{r_h}(f_h^{\rightarrow}(\mathbf{h}_i^{\rightarrow}, \mathbf{x}_i, \mathbf{t}_i, i)), \quad \forall i \in \{1, \dots, N-1\} \\ \mathbf{h}_{i-1}^{\leftarrow} &= \text{proj}_{r_h}(f_h^{\leftarrow}(\mathbf{h}_i^{\leftarrow}, \mathbf{x}_i, \mathbf{t}_i, i)), \quad \forall i \in \{2, \dots, N\} \\ y_i &= f_y(\mathbf{U}^{\rightarrow} \mathbf{h}_i^{\rightarrow}, \mathbf{U}^{\leftarrow} \mathbf{h}_i^{\leftarrow}, \mathbf{x}_i, \mathbf{t}_i, i), \quad \forall i \in [N] \end{aligned}$$

where $f_y : \mathbb{R}^{d_h} \times \mathbb{R}^{d_h} \times \mathbb{R}^d \times [N]^{q+1} \rightarrow \mathbb{R}$, $f_h^{\rightarrow}, f_h^{\leftarrow} : \mathbb{R}^{d_h} \times \mathbb{R}^d \times [N]^{q+1} \rightarrow \mathbb{R}^{d_h}$, $\mathbf{U}^{\rightarrow}, \mathbf{U}^{\leftarrow} \in \mathbb{R}^{d_h \times d_h}$, d_h is the width of the model, and $r_h > 0$ is some constant. Moreover, $\text{proj}_{r_h} : \mathbb{R}^{d_h} \rightarrow \mathbb{R}^{d_h}$ is any mapping that guarantees $\|\text{proj}_{r_h}(\cdot)\|_2 \leq r_h$. As mentioned before, this operation mirrors the layer normalization to ensure that \mathbf{h}_i remains stable. Further, we assume $f_y(\cdot, \mathbf{x}, \mathbf{t})$ is \mathfrak{L}/r_h -Lipschitz for all $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{t} \in [N]^q$. This formulation covers different variants of (bidirectional) RNNs used in practice such as LSTM and GRU, and includes the RNN formulation of Theorem 6 as a special case. Define $\mathbf{U} := (\mathbf{U}^{\rightarrow}, \mathbf{U}^{\leftarrow}) \in \mathbb{R}^{d_h \times 2d_h}$ for conciseness. Note that in practice $f_y, f_h^{\rightarrow}, f_h^{\leftarrow}$ are determined by additional parameters. However, the only weight that we explicitly denote in this formulation is \mathbf{U} , since our lower bound will directly involve this projection, and we keep the rest of the parameters implicit for our representational lower bound.

Our technique for proving the RNN lower bound differs significantly from that of FFNs. In particular, we will control the representation cost of the q STR model, i.e., a lower bound on the norm of Θ_{RNN} .

We will now present the RNN lower bound, with its proof deferred to Section C.6.

Proposition 34 *Consider the 1STR model where $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_{Nd})$ with a linear link function, i.e. $y_j = \langle \mathbf{u}, \mathbf{x}_{t_j} \rangle$ for some $\mathbf{u} \in \mathbb{S}^{d-1}$. Further, t_i is drawn independently from the rest of the prompt and uniformly from $[N]$ for all $i \in [N]$. Then, there exists an absolute constant $c > 0$, such that*

$$\frac{1}{N} \mathbb{E} \left[\|\mathbf{y} - \hat{\mathbf{y}}_{\text{RNN}}(\mathbf{p})\|^2 \right] \leq c,$$

implies

$$d_h \geq \Omega \left(\frac{N}{\log(1 + \mathfrak{L}^2 \|\mathbf{U}\|_{\text{op}}^2)} \right), \quad \text{and} \quad \|\mathbf{U}\|_{\text{op}}^2 \geq \Omega \left(\frac{N}{\mathfrak{L}^2 \log(1 + d_h)} \right).$$

Remark 35 *Note that the unboundedness of Gaussian random variables is not an issue for approximation here, since $(g(\mathbf{x}_1), \dots, g(\mathbf{x}_N))$ is highly concentrated around $\mathbb{S}^{N-1}(\sqrt{N})$. In fact, one can directly assume $(g(\mathbf{x}_1), \dots, g(\mathbf{x}_N)) \sim \text{Unif}(\mathbb{S}^{N-1}(\sqrt{N}))$ and derive a similar lower bound. The choice of Gaussian above is only made to simplify the presentation of the proof.*

The above proposition has two implications. First, it has a *computational* consequence, implying that any RNN representing the q STR models requires a width that grows at least linearly with the context-length N . A similar lower bound in terms of bit complexity was derived in [24] using different tools. More importantly, the norm lower bound $\|\mathbf{U}\|_{\text{F}} \geq \tilde{\Omega}(\sqrt{N})$ has a *generalization* consequence, which we discuss below.

To translate the above representational cost result to a sample complexity lower bound, we now introduce the parametrization of the output function f_y . The exact parametrization of the transition

functions will be unimportant, and we will use the notation $f_h^{\rightarrow}(\mathbf{h}, \mathbf{x}, \mathbf{t}; \Theta_h^{\rightarrow})$ to denote a general parameterized function (similarly with f^{\leftarrow}). We will assume f_y is given by a feedforward network,

$$f_y(\mathbf{U}^{\rightarrow} \mathbf{h}^{\rightarrow}, \mathbf{U}^{\leftarrow} \mathbf{h}^{\leftarrow}, \mathbf{x}, \mathbf{t}; \Theta_y) = \mathbf{W}_{L_y} \sigma(\dots \sigma(\mathbf{W}_2 \sigma(\mathbf{U} \mathbf{h} + \mathbf{W}_y \mathbf{z} + \mathbf{b}_y) + \mathbf{b}_2) \dots),$$

where $\mathbf{h} = (\mathbf{h}^{\rightarrow}, \mathbf{h}^{\leftarrow}) \in \mathbb{R}^{2d_h}$, $\mathbf{z} = (\mathbf{x}_i, f_E(\mathbf{t}_i, i)) \in \mathbb{R}^{d+d_E}$. Here, $f_E(\mathbf{t}_i, i)$ is an arbitrary encoding function with arbitrary dimension d_E . Then $\Theta_y = (\mathbf{U}, \mathbf{W}_y, \mathbf{b}_y, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_{L_y})$, and $\Theta_{\text{RNN}} = (\mathbf{U}, \Theta_y, \Theta_h^{\rightarrow}, \Theta_h^{\leftarrow})$. Note that thanks to the homogeneity of ReLU, we can always reparameterize the network by taking $\bar{\mathbf{h}} = \mathbf{h}/r_h$, $\bar{\mathbf{W}}_y = \mathbf{W}_y/r_h$, $\bar{\mathbf{b}}_y = \mathbf{b}_y/r_h$, and $\bar{\mathbf{W}}_2 = \mathbf{W}_2/r_h$ without changing the prediction function. Thus, in Theorem 7, we take $r_h = 1$ without losing the expressive power of the network. We now state the rigorous version of Theorem 7.

Theorem 36 *Consider the 1STR model of Proposition 34. Suppose the size of the hidden state, the depth of the prediction function, and the weight norm respectively satisfy $d_h \leq e^{N^c}$, $2 \leq L_y \leq C$, and $\|\text{vec}(\Theta_{\text{RNN}})\|_2 \leq e^{N^c/L_y}$ for some absolute constants $c < 1$ and $C \geq 2$, and recall we set $r_h = 1$ due to homogeneity of the network. Let $\hat{\Theta}_\varepsilon$ be the min-norm ε -ERM of \hat{R}_n^{RNN} , defined in (4). Then, there exist absolute constants $c_1, c_2, c_3 > 0$ such that if $n \leq \mathcal{O}(N^{c_1})$, for any $\varepsilon \geq 0$, with probability at least c_2 over the training set,*

$$\frac{1}{N} \mathbb{E} \left[\left\| \hat{\mathbf{y}}_{\text{RNN}}(\mathbf{p}; \hat{\Theta}_{n,\varepsilon}) - \mathbf{y} \right\|_2^2 \right] \geq c_3.$$

C.6. Proof of Proposition 34

The crux of the proof of Proposition 34 is to show the following position, which provides a lower bound on the prediction error at any fixed position in the prompt.

Proposition 37 *Consider the same setting as in Proposition 34. There exists an absolute constant $c > 0$, such that for any fixed $j \in [N]$, if*

$$\mathbb{E}[(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \leq c,$$

then

$$d_h \geq \Omega\left(\frac{N}{\log(1 + \mathfrak{L}^2 \|\mathbf{U}\|_{\text{op}}^2)}\right), \quad \text{and} \quad \|\mathbf{U}\|_{\text{op}}^2 \geq \Omega\left(\frac{N}{\mathfrak{L}^2 \log(1 + d_h)}\right).$$

We shortly remark that the statement of Proposition 34 directly follows from that of Proposition 37.

Proof of Proposition 34. Let c be the constant given by Proposition 37. Suppose that

$$\frac{1}{N} \mathbb{E} \left[\|\hat{\mathbf{y}}_{\text{RNN}}(\mathbf{p}) - \mathbf{y}\|_2^2 \right] \leq c.$$

Then,

$$\min_{j \in [N]} \mathbb{E}[(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \leq \frac{1}{N} \sum_{j=1}^N \mathbb{E}[(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \leq c.$$

As a result, there exists some $j \in [N]$ such that $\mathbb{E}[(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \leq c$. We can then invoke Proposition 37 to obtain lower bounds on d_h and $\|\mathbf{U}\|_{\text{op}}$, completing the proof of Proposition 34. ■

We now present the proof of Proposition 37.

Proof of Proposition 37. Let $\mathbf{h}_j = (\mathbf{U}^{\rightarrow} \mathbf{h}_j^{\rightarrow}, \mathbf{U}^{\leftarrow} \mathbf{h}_j^{\leftarrow}) \in \mathbb{R}^{2d_h}$, and define

$$\Phi(\mathbf{h}_j) := \left(f_y(\mathbf{h}_j, \mathbf{x}_j, (1), j), \dots, f_y(\mathbf{h}_j, \mathbf{x}_j, (j-1), j), f_y(\mathbf{h}_j, \mathbf{x}_j, (j+1), j), \dots, f_y(\mathbf{h}_j, \mathbf{x}_j, (N), j) \right)^{\top} \in \mathbb{R}^{N-1}.$$

In other words, $\Phi : \mathbb{R}^{2d_h} \rightarrow \mathbb{R}^{N-1}$ captures all possible outcomes of $\hat{y}_{\text{RNN}}(\mathbf{p})_j$ depending on the value of t_j (excluding the case where $t_j = j$). Ideally, we must have $f_y(\mathbf{h}_j, \mathbf{x}_j, (k), j) \approx g(\mathbf{x}_k)$.

Let $\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(P)}$ be an i.i.d. sequence of prompts, then modify them to share the j th input token, i.e. $\mathbf{x}_j^{(i)} = \mathbf{x}_j^{(1)}$ for all $i \in [P]$, with P to be determined later. Note that by our assumption on prompt distribution, this operation does not change the marginal distribution of each $\mathbf{p}^{(i)}$. Similarly, define

$$\mathbf{g}^{(i)} := (g(\mathbf{x}_1^{(i)}), \dots, g(\mathbf{x}_{j-1}^{(i)}), g(\mathbf{x}_{j+1}^{(i)}), \dots, g(\mathbf{x}_N^{(i)}))^{\top} \in \mathbb{R}^{N-1}$$

for each prompt. We also let $\mathbf{h}_j^{(i)\rightarrow}, \mathbf{h}_j^{(i)\leftarrow}$ be the corresponding hidden states obtained from passing these prompts through the RNN, and define $\mathbf{h}_j^{(i)}$ using them. Note that $\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(P)}$ is an i.i.d. sequence of vectors drawn from $\mathcal{N}(0, \mathbf{I}_{N-1})$.

We now define two events E_1 and E_2 , where

$$E_1 = \left\{ \forall i \neq k, \quad \left\| \mathbf{g}^{(i)} - \mathbf{g}^{(k)} \right\|_2 \geq \varepsilon_g \sqrt{N-1} \right\},$$

and

$$E_2 = \left\{ \sum_{i=1}^P \mathbb{1} \left[\left\| \Phi(\mathbf{h}_j^{(i)}) - \mathbf{g}^{(i)} \right\|_2 \geq \frac{\varepsilon \sqrt{N}}{\delta} \right] \leq 2\delta^2 P \right\},$$

where $\delta \in (0, 1)$ will be chosen later. In other words, E_1 is the event in which $\mathbf{g}^{(i)}$ are ‘‘packed’’ in the space, while E_2 is the event where the RNN will be ‘‘wrong’’ at position j on at most $2\delta^2$ fraction of the prompts. We will now attempt to lower bound $\mathbb{P}(E_1 \cap E_2)$.

Note that $\mathbf{g}^{(i)} - \mathbf{g}^{(k)} \stackrel{(d)}{=} \sqrt{2}\mathbf{g}$ where $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_{N-1})$. By a union bound we have

$$\begin{aligned} \mathbb{P}(E_1^C) &\leq \sum_{i \neq k} \mathbb{P} \left(\left\| \mathbf{g}^{(i)} - \mathbf{g}^{(k)} \right\|_2 \leq \varepsilon_g \sqrt{N-1} \right) \\ &\leq P^2 \mathbb{P} \left(\sqrt{2} \|\mathbf{g}\|_2 \leq \varepsilon_g \sqrt{N-1} \right) \\ &\leq P^2 \mathbb{P} \left(\|\mathbf{g}\|_2 - \mathbb{E}[\|\mathbf{g}\|_2] \leq \left(\frac{\varepsilon_g}{\sqrt{2}} - c \right) \sqrt{N-1} \right) \\ &\leq P^2 e^{-(c - \varepsilon_g/\sqrt{2})^2 (N-1)/2}, \end{aligned}$$

for all $\varepsilon_g \leq c\sqrt{2}$, where $c > 0$ is an absolute constant such that $c\sqrt{N-1} \leq \mathbb{E}[\|\mathbf{g}\|]$, and the last inequality holds by subGaussianity of the norm of a standard Gaussian random vector. From here on, we will choose $\varepsilon_g = c/\sqrt{2}$ (and simply denote $\varepsilon_g \asymp 1$), which implies $\mathbb{P}(E_1^C) \leq P^2 e^{-c^2(N-1)/8}$.

To lower bound $\mathbb{P}(E_2)$, consider a random prompt-label pair \mathbf{p}, \mathbf{y} and the corresponding \mathbf{g} . Note that in the prompt \mathbf{p} , the index t_j is drawn independently of the rest of \mathbf{p} , and has a uniform distribution in $[N]$. Let $\mathbf{p}[t_j \mapsto k]$ denote a modification of \mathbf{p} where we set t_j equal to k , and let

$\mathbf{y}[t_j \mapsto k]$ be the labels corresponding to this modified prompt. We then have

$$\begin{aligned} \frac{1}{N} \|\Phi(\mathbf{h}_j) - \mathbf{g}\|_2^2 &= \frac{1}{N} \sum_{k \neq j} (\hat{y}_{\text{RNN}}(\mathbf{p}[t_j \mapsto k])_j - g(\mathbf{x}_k))^2 \\ &\leq \frac{1}{N} \sum_{k=1}^N (\hat{y}_{\text{RNN}}(\mathbf{p}[t_j \mapsto k])_j - y(\mathbf{p}[t_j \mapsto k])_j)^2 \\ &= \mathbb{E}_{t_j} [(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \end{aligned}$$

As a result, via a Markov inequality, we obtain

$$\begin{aligned} \mathbb{P}\left(\frac{1}{N} \|\Phi(\mathbf{h}_j) - \mathbf{g}\|_2^2 \geq \frac{\varepsilon^2}{\delta^2}\right) &= \mathbb{P}\left(\mathbb{E}_{t_j} [(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2] \geq \frac{\varepsilon^2}{\delta^2}\right) \\ &\leq \frac{\delta^2 \mathbb{E}[(\hat{y}_{\text{RNN}}(\mathbf{p})_j - y_j)^2]}{\varepsilon^2} \\ &\leq \delta^2. \end{aligned}$$

Going back to our lower bound on $\mathbb{P}(E_2)$, define the Bernoulli random variable

$$z^{(i)} = \mathbb{1}\left[\left\|\Phi(\mathbf{h}_j^{(i)}) - \mathbf{g}^{(i)}\right\|_2 \geq \frac{\varepsilon\sqrt{N}}{\delta}\right].$$

Note that $(z^{(i)})$ are i.i.d. since $\mathbf{h}_j^{(i)}$ and $\mathbf{g}^{(i)}$ do not depend on \mathbf{x}_j . Then, by Hoeffding's inequality,

$$\mathbb{P}(E_2^C) = \mathbb{P}\left(\sum_{j=1}^P z^{(i)} \geq 2\delta^2 P\right) \leq e^{-2P\delta^4}.$$

We now have our desired lower bound on $\mathbb{P}(E_1 \cap E_2)$, given by

$$\mathbb{P}(E_1 \cap E_2) \geq 1 - \mathbb{P}(E_1^C) - \mathbb{P}(E_2^C) \geq 1 - e^{-2P\delta^4} - P^2 e^{-c^2(N-1)/8}.$$

Suppose $\delta \geq e^{-c'N}$ for some absolute constant $c' > 0$. Then, choosing $P = \lfloor e^{c''N} \rfloor$ for some absolute constant $c'' > 0$ would ensure $\mathbb{P}(E_1 \cap E_2) > 0$, and allows us to look at this intersection.

Let $\mathcal{I} = \{i : z^{(i)} = 0\}$. On E_1 , and for $i, k \in \mathcal{I}$ with $i \neq k$ we have

$$\begin{aligned} \left\|\Phi(\mathbf{h}_j^{(i)}) - \Phi(\mathbf{h}_j^{(k)})\right\|_2 &\geq \left\|\mathbf{g}^{(i)} - \mathbf{g}^{(k)}\right\|_2 - \left\|\Phi(\mathbf{h}_j^{(i)}) - \mathbf{g}^{(i)}\right\|_2 - \left\|\Phi(\mathbf{h}_j^{(k)}) - \mathbf{g}^{(k)}\right\|_2 \\ &\geq \varepsilon_g \sqrt{N-1} - \frac{2\varepsilon\sqrt{N}}{\delta} =: \mathfrak{L}\sqrt{N}\varepsilon_h. \end{aligned}$$

Note that from the Lipschitzness of f_y , we have $\left\|\Phi(\mathbf{h}_j^{(i)}) - \Phi(\mathbf{h}_j^{(k)})\right\|_2 \leq \frac{\mathfrak{L}\sqrt{N}}{r_h} \left\|\mathbf{h}_j^{(i)} - \mathbf{h}_j^{(k)}\right\|_2$.

As a result, the set $\{\mathbf{h}_j^{(i)} : i \in \mathcal{I}\}$ is an $r_h \varepsilon_h$ -packing for $\{\mathbf{h} : \|\mathbf{h}\|_2 \leq \sqrt{2}\|\mathbf{U}\|_{\text{op}} r_h\}$. Using Lemma 44, the log packing number can be bounded by

$$\log \mathcal{I} \leq \left\{d_h \log\left(1 + \frac{2\sqrt{2}\|\mathbf{U}\|_{\text{op}}}{\varepsilon_h}\right)\right\} \wedge \left\{\frac{2\|\mathbf{U}\|_{\text{op}}^2}{\varepsilon_h^2} \left(1 + \log\left(1 + \frac{M\varepsilon_h^2}{2\|\mathbf{U}\|_{\text{op}}^2}\right)\right)\right\}.$$

On $E_1 \cap E_2$, we have $\mathcal{I} \geq (1 - 2\delta^2)P \geq (1 - 2\delta^2)e^{cN}$ for some absolute constant $c > 0$. Therefore,

$$\frac{\log(1 - 2\delta^2) + cN}{\log(1 + 2\sqrt{2}\|\mathbf{U}\|_{\text{op}}/\varepsilon_h)} \leq d_h,$$

and

$$\frac{\varepsilon_h^2(\log(1 - 2\delta^2) + cN)}{2 + 2\log(1 + d_h\varepsilon_h^2/(2\|\mathbf{U}\|_{\text{op}}^2))} \leq \|\mathbf{U}\|_{\text{op}}^2.$$

Choosing $\delta = 1/2$ and recalling $\varepsilon_g \asymp 1$, we obtain $\varepsilon_h \gtrsim (1 - C\varepsilon)/\mathfrak{L}$ for some absolute constant $C > 0$, which concludes the proof. \blacksquare

C.7. Proof of Theorem 7

We first provide an estimate for the capacity of two-layer feedforward networks to interpolate n samples.

Lemma 38 *Suppose $\{\mathbf{x}^{(i)}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ and let $y^{(i)} = \langle \mathbf{u}, \mathbf{x}_{t_i} \rangle$ for arbitrary $t_i \in [N]$ and $\mathbf{u} \in \mathbb{S}^{d-1}$. Then, there exists an absolute constant $c > 0$ such that for all $m \geq n$ and with probability at least c , there exist data dependent weights $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$ and $\mathbf{W} \in \mathbb{R}^{m \times d}$, such that*

$$\mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}) = y^{(i)}, \quad \forall i \in [n]$$

and

$$\|\mathbf{a}\|_2^2 + \|\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{b}\|_2^2 \leq \mathcal{O}(n^3).$$

Proof The proof of Lemma 38 is an immediate consequence of two lemmas.

1. Lemma 39 shows that the inputs $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ can be projected to sufficiently separated scalar values with a unit vector \mathbf{v} .
2. Lemma 40 perfectly fits n univariate samples using a two-layer ReLU neural network. When invoking this lemma, we use $\|\mathbf{z}\|_2 = \mathcal{O}(\sqrt{n})$ and $\varepsilon = \Omega(1/n^2)$ as given by Lemma 39.

The only missing piece is to upper bound $\|\mathbf{y}\|_2$ appearing in the final bound of Lemma 40. To that end, we apply the following Markov inequality,

$$\mathbb{P}\left(\|\mathbf{y}\|_2^2 \geq 6n\right) \leq \frac{\mathbb{E}\left[\|\mathbf{y}\|_2^2\right]}{6n} \leq \frac{1}{6}.$$

As the statement of Lemma 39 holds with probability at least $\frac{1}{3}$, this suggests that the statement of Lemma 38 holds with probability at least $\frac{1}{6}$, concluding the proof. \blacksquare

Lemma 39 *Suppose $\{\mathbf{x}^{(i)}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$. Then, with probability at least $1/3$, there exists some $\mathbf{v} \in \mathbb{S}^{d-1}$ (dependent on $\{\mathbf{x}^{(i)}\}$) such that for all $i \neq j$,*

$$\left| \mathbf{v}^\top \mathbf{x}^{(i)} - \mathbf{v}^\top \mathbf{x}^{(j)} \right| = \Omega\left(\frac{1}{n^2}\right). \quad (25)$$

and $\sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}^{(i)})^2 = \mathcal{O}(n)$.

Proof The proof follows the probabilistic method. Sample $\mathbf{v} \sim \text{Unif}(\mathbb{S}^{d-1})$ independent of $\{\mathbf{x}^{(i)}\}$. For each $i \neq j$, let

$$a_{i,j} = \mathbf{u}^\top (\mathbf{x}^{(i)} - \mathbf{x}^{(j)})$$

and note that $a_{i,j} \mid \mathbf{v} \sim \mathcal{N}(0, 2)$. We apply basic Gaussian anti-concentration to place a lower bound on the probability of any $a_{i,j}$ being close to zero,

$$\mathbb{P}(\exists i, j \text{ s.t. } |a_{i,j}| \leq \epsilon) \leq \sum_{i \neq j} \mathbb{P}(|a_{i,j}| \leq \epsilon) = \sum_{i \neq j} \mathbb{E}[\mathbb{P}(|a_{i,j}| \leq \epsilon \mid \mathbf{v})] \leq \frac{n^2 \epsilon}{\sqrt{\pi}} \leq \frac{1}{3},$$

where the last inequality follows by taking $\epsilon = \sqrt{\pi}/(3n^2)$. Furthermore,

$$\mathbb{P}\left(\sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}^{(i)})^2 \geq 3n\right) \leq \frac{\sum_{i=1}^n \mathbb{E}[(\mathbf{v}^\top \mathbf{x}^{(i)})^2]}{3n} = \frac{1}{3},$$

by Markov's inequality. Combining the two events completes the proof. \blacksquare

Lemma 40 Consider some $\mathbf{z} = (z^{(1)}, \dots, z^{(n)})^\top \in \mathbb{R}^n$ and $\mathbf{y} = (y^{(1)}, \dots, y^{(n)})^\top \in \mathbb{R}^n$, such that $|z^{(i)} - z^{(j)}| \geq \epsilon$ for all $i \neq j$. For simplicity, assume $\epsilon \leq 1$. Then, there exists a two-layer ReLU neural network

$$g(t) = \sum_{j=1}^m a_j \sigma(w_j t + b_j)$$

that satisfies $g(z^{(i)}) = y^{(i)}$ for all $i \in [n]$, $m = n$, and

$$\|\mathbf{a}\|_2^2 + \|\mathbf{w}\|_2^2 + \|\mathbf{b}\|_2^2 = \mathcal{O}\left(\frac{\|\mathbf{y}\|_2 \sqrt{n + \|\mathbf{z}\|_2^2}}{\epsilon}\right). \quad (26)$$

Proof Without loss of generality, we assume that $z^{(1)} \leq \dots \leq z^{(n)}$. Then, we define the neural network g as follows:

$$\begin{aligned} g(t) &= \sum_{i=1}^n a'_i \sigma(w'_i t - b'_i) = y^{(1)} \sigma(t - z^{(1)} + 1) + \left(\frac{y^{(2)} - y^{(1)}}{z^{(2)} - z^{(1)}} - y^{(1)}\right) \sigma(t - z^{(1)}) \\ &\quad + \sum_{i=3}^n \left(\frac{y^{(i)} - y^{(i-1)}}{z^{(i)} - z^{(i-1)}} - \frac{y^{(i-1)} - y^{(i-2)}}{z^{(i-1)} - z^{(i-2)}}\right) \sigma(t - z^{(i-1)}). \end{aligned}$$

One can verify by induction that $g(z^{(i)}) = y^{(i)}$ for every i by noting that the slope of g is

$$(y^{(i)} - y^{(i-1)}) / (z^{(i)} - z^{(i-1)})$$

between $(z^{(i-1)}, y^{(i-1)})$ and $(z^{(i)}, y^{(i)})$. From the above, we have $w'_i = 1$, $\|\mathbf{b}'\|_2^2 \lesssim \|\mathbf{z}\|_2^2 + 1$, and $\|\mathbf{a}'\|_2^2 \lesssim \|\mathbf{y}\|_2^2 / \epsilon^2$. For $\alpha = ((\|\mathbf{z}\|_2^2 + n)\epsilon^2 / \|\mathbf{y}\|_2^2)^{1/4}$, let $\mathbf{u} = \alpha \mathbf{u}'$, $\mathbf{w} = \mathbf{w}' / \alpha$, and $\mathbf{b} = \mathbf{b}' / \alpha$. By homogeneity, the neural network with weights $(\mathbf{u}, \mathbf{w}, \mathbf{b})$ has identical outputs to that of $(\mathbf{u}', \mathbf{w}', \mathbf{b}')$ and satisfies (26), completing the proof. \blacksquare

We are now ready to present the proof of the sample complexity lower bound for RNNs.

Proof of Theorem 7. First, consider the case where $d_h < n$. Note that as a function of $U\mathbf{h} = (U^{\rightarrow}\mathbf{h}^{\rightarrow}, U^{\leftarrow}\mathbf{h}^{\leftarrow})$, f_y is \mathfrak{L} -Lipschitz with

$$\mathfrak{L} = \|\mathbf{W}_{L_y}\|_{\text{op}} \|\mathbf{W}_{L_y-1}\|_{\text{op}} \cdots \|\mathbf{W}_2\|_{\text{op}}.$$

Using the AM-GM inequality,

$$\left(\mathfrak{L}^2\|\mathbf{U}\|_{\text{op}}^2\right)^{1/L_y} \leq \frac{1}{L_y} \|\text{vec}(\Theta)\|_2^2 \leq e^{N^c/L_y}.$$

As a result, we have $\mathfrak{L}\|\mathbf{U}\|_{\text{op}} \leq e^{N^c/2}$. By invoking Proposition 25, to obtain population risk less than some absolute constant $c_3 > 0$, we need

$$d_h \geq \Omega\left(\frac{N}{\log(1 + \mathfrak{L}^2\|\mathbf{U}\|_{\text{op}}^2)}\right) \geq \Omega(N^{1-c}).$$

This implies $n \geq d_h \geq \Omega(N^{1-c})$. By taking c_1 in the theorem statement to be less than $1 - c$, we obtain a contradiction. Therefore, we must have either a population risk at least c_3 or $d_h \geq n$.

Suppose now that $d_h \geq n$. We show that with constant probability, we can construct an RNN that interpolates the n training samples with norm independent of n . We simply let $\Theta_h^{\rightarrow} = \mathbf{0}$, $\Theta_h^{\leftarrow} = \mathbf{0}$, $\mathbf{U} = \mathbf{0}$, and describe the construction of $\mathbf{W}_{L_y}, \dots, \mathbf{W}_2, \mathbf{W}_y$, and (\mathbf{b}_l) in the following. Using the construction of Lemma 38, we can let

$$\mathbf{W}_y = \begin{pmatrix} \mathbf{W} & \mathbf{0}_{n \times d_E} \\ \mathbf{0}_{(m-n) \times d} & \mathbf{0}_{(m-n) \times d_E} \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} \mathbf{b} \\ \mathbf{0}_{m-n} \end{pmatrix}, \quad \mathbf{W}_2 = \begin{pmatrix} \mathbf{a}^\top & \mathbf{0}_{m-n}^\top \\ -\mathbf{a}^\top & \mathbf{0}_{m-n}^\top \\ \mathbf{0}_{(m-2) \times n} & \mathbf{0}_{(m-2) \times (m-n)} \end{pmatrix},$$

where $\mathbf{W} \in \mathbb{R}^{n \times d}$, and $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are given by Lemma 38. Then,

$$\mathbf{W}_2^\top \sigma(\mathbf{W}_y \mathbf{x}_{j^{(i)}}^{(i)} + \mathbf{b}_y) = (y_{j^{(i)}}^{(i)}, -y_{j^{(i)}}^{(i)}, 0, \dots, 0)^\top.$$

For $(\mathbf{W}_l)_{l=3}^{L_y-1}$, we let $(W_l)_{11} = (W_l)_{22} = 1$, and choose the rest of the coordinates of \mathbf{W}_l to be zero. Therefore, the output of the l th layer is given by

$$(\sigma(y_{j^{(i)}}^{(i)}), \sigma(-y_{j^{(i)}}^{(i)}), 0, \dots, 0)^\top.$$

For the final layer, we let $\mathbf{W}_{L_y} = (1, -1, 0, \dots, 0)$. Using the fact that $\sigma(z) - \sigma(-z) = z$, we obtain

$$f_y(U^{\rightarrow}\mathbf{h}_j^{\rightarrow}, U^{\leftarrow}\mathbf{h}_j^{\leftarrow}, \mathbf{z}_{j^{(i)}}^{(i)}; \Theta_y) = y_{j^{(i)}}^{(i)}$$

We have found Θ such that $\hat{R}_n^{\text{RNN}}(\Theta) = 0$ and $\|\text{vec}(\Theta)\|_2^2 \leq \mathcal{O}(n^3)$ (recall that $L_y \leq \mathcal{O}(1)$). As a result, $\hat{\Theta}_\varepsilon$ must also satisfy $\|\text{vec}(\hat{\Theta}_\varepsilon)\|_2^2 \leq \mathcal{O}(n^3)$.

On the other hand, notice that as a function of $U\mathbf{h} = (U^{\rightarrow}\mathbf{h}^{\rightarrow}, U^{\leftarrow}\mathbf{h}^{\leftarrow})$, f_y is \mathfrak{L} -Lipschitz with

$$\mathfrak{L} = \|\mathbf{W}_{L_y}\|_{\text{op}} \|\mathbf{W}_{L_y-1}\|_{\text{op}} \cdots \|\mathbf{W}_2\|_{\text{op}}.$$

From Proposition 34, using the fact that $\|\cdot\|_{\text{op}} \leq \|\cdot\|_{\text{F}}$ and the AM-GM inequality, we obtain

$$\frac{1}{L_y} \|\text{vec}(\Theta)\|_2^2 \geq \left(\mathcal{L}^2 \|\mathbf{U}\|_{\text{op}}^2 \right)^{1/L_y} \geq \Omega \left(\left(\frac{N}{\log d_h} \right)^{1/L_y} \right)$$

to achieve population risk less than some absolute constant $c_3 > 0$. Recall that $\log d_h \leq N^c$ for some $c < 1$. The proof is completed by noticing that unless $n \geq \Omega(N^{c_1})$ for some absolute constant $c_1 > 0$, $\|\text{vec}(\hat{\Theta}_\varepsilon)\|_2$ will always be less than the lower bound above, with some absolute constant probability $c_2 > 0$ over the training set. \blacksquare

Appendix D. Details of Section 5

We first define the class of algorithms considered in Theorem 8. Importantly, this lower bound holds regardless of the loss function used for training; for some arbitrary loss $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, we define the empirical risk of the FFN as

$$\hat{\mathcal{L}}^{\text{FFN}}(f, \mathbf{W}) := \frac{1}{nN} \sum_{i=1}^n \sum_{j=1}^N \ell(y_j^{(i)}, f(\mathbf{T}^{(i)}, \mathbf{W} \mathbf{x}^{(i)}_j)),$$

where $\mathbf{T}^{(i)} = (t_1^{(i)}, \dots, t_N^{(i)})$. We still use $R^{\text{FFN}}(f, \mathbf{W})$ for expected squared loss. Our lower bound covers a broad set of algorithms, characterized by the following definition.

Definition 41 Let \mathcal{A}_{SP} denote the set of algorithms that return a stationary point of the regularized empirical risk. Specifically, for every $A \in \mathcal{A}_{\text{SP}}$, $A(S_n)$ returns $f_{A(S_n)}$, $\mathbf{W}_{A(S_n)}$, such that

$$\nabla_{\mathbf{W}} \hat{\mathcal{L}}^{\text{FFN}}(f_{A(S_n)}, \mathbf{W}_{A(S_n)}) + \lambda \mathbf{W}_{A(S_n)} = 0,$$

for some $\lambda > 0$ depending on A . S_n above denotes the training set. Let \mathcal{A}_{ERM} denote the set of algorithms that return the min-norm approximate ERM. Specifically, every $A \in \mathcal{A}_{\text{ERM}}$ returns

$$A(S_n) = \arg \min_{\{f, \mathbf{W} : \hat{\mathcal{L}}^{\text{FFN}}(f, \mathbf{W}) \leq \varepsilon\}} \|\mathbf{W}\|_{\text{F}},$$

for some $\varepsilon \geq 0$. Define $\mathcal{A} := \mathcal{A}_{\text{SP}} \cup \mathcal{A}_{\text{ERM}}$.

In particular, \mathcal{A} goes beyond constrained ERM in that it also includes the (ideal) output of first-order optimization algorithms with weight decay, or ERM with additional ℓ_2 penalty on the weights.

D.1. Proof of Theorem 8

Let \mathbf{u} be sampled uniformly from \mathbb{S}^{d-1} independently from $\mathbf{p} = (t_1, \mathbf{x})$, and note that we have

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \mathbb{E}[(y_j - f_{A(S_n)}(t_1, \mathbf{W}_{A(S_n)} \mathbf{x})_j)^2] \geq \mathbb{E}_{\mathbf{u} \sim \text{Unif}(\mathbb{S}^{d-1}), j, y, \mathbf{p} \sim \mathcal{P}} [(y_j - f_{A(S_n)}(t_1, \mathbf{W}_{A(S_n)} \mathbf{x})_j)^2],$$

for all $A \in \mathcal{A}$. From this point, we will simply use f for $f_{A(S_n)}$ and \mathbf{W} for $\mathbf{W}_{A(S_n)}$. Next, we argue that the output weights of any algorithm in \mathcal{A} satisfy

$$\mathbf{w}_k = \sum_{i=1}^n \alpha_k^{(i)} \mathbf{x}^{(i)}, \quad \forall k \in [m_1],$$

for some coefficients $(\alpha_k^{(i)})_{i \in [n], k \in [m_1]}$. This is straightforward to verify for $A \in \mathcal{A}_{\text{SP}}$, as

$$\nabla_{\mathbf{w}_k} \hat{\mathcal{L}}^{\text{FFN}}(f, \mathbf{W}) \in \text{span}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}).$$

For $A \in \mathcal{A}_{\text{ERM}}$, note that $\hat{\mathcal{L}}^{\text{FFN}}$ only depends on \mathbf{w}_k through its projection on $\text{span}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$. As a result, any minimum-norm ε -ERM would satisfy $\mathbf{w}_k \in \text{span}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$.

Note that for $n \leq Nd$, the span of $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ is n -dimensional with probability 1 over S_n . Let $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}$ denote an orthonormal basis of $\text{span}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$, and let $\mathbf{V} = (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)})^\top \in \mathbb{R}^{n \times Nd}$. Recall that for the simple-1STR model considered here, $y_j = y = \langle \mathbf{u}, \mathbf{x}_{t_j} \rangle$ for $j \in [N]$. Then,

$$\mathbb{E}_{\mathbf{u}, y, j, \mathbf{p}} [(y_j - f(t_1, \mathbf{W}\mathbf{x}_j))^2] \geq \mathbb{E}_{\mathbf{u}, t_1, \mathbf{V}\mathbf{x}} [\text{Var}(y | \mathbf{u}, t_1, \mathbf{V}\mathbf{x})] = \mathbb{E}_{\mathbf{u}, t_1, \mathbf{V}\mathbf{x}} [\text{Var}(\langle \mathbf{P}_{t_1} \mathbf{u}, \mathbf{x} \rangle | \mathbf{u}, t_1, \mathbf{V}\mathbf{x})],$$

where $\mathbf{P}_{t_1} \in \mathbb{R}^{Nd \times d}$ has the form $(\underbrace{\mathbf{0}_d, \dots, \mathbf{I}_d, \dots, \mathbf{0}_d}_{t_1})^\top$. The conditioning above comes from the

fact that via training, f and \mathbf{W} can depend on \mathbf{u} , but the prediction depends on \mathbf{x} only through $\mathbf{V}\mathbf{x}$. Consequently, we replace the prediction of the FFN by the best predictor having access to \mathbf{u} , t_1 , and $\mathbf{V}\mathbf{x}$. Note that t_1 , \mathbf{u} , and $\mathbf{V}\mathbf{x}$ are jointly independent, and the joint distribution $(\langle \mathbf{P}_{t_1} \mathbf{u}, \mathbf{x} \rangle, \mathbf{V}\mathbf{x})$ is given by $\mathcal{N}\left(0, \begin{pmatrix} 1 & \mathbf{V}\mathbf{P}_{t_1}\mathbf{u} \\ \mathbf{u}^\top \mathbf{P}_{t_1}^\top \mathbf{V}^\top & \mathbf{I}_n \end{pmatrix}\right)$, thus we have

$$\text{Var}(\langle \mathbf{P}_{t_1} \mathbf{u}, \mathbf{x} \rangle | \mathbf{u}, t_1, \mathbf{V}\mathbf{x}) = 1 - \|\mathbf{V}\mathbf{P}_{t_1}\mathbf{u}\|^2.$$

In particular,

$$\mathbb{E}_{\mathbf{u}} [\text{Var}(\langle \mathbf{P}_{t_1} \mathbf{u}, \mathbf{x} \rangle | \mathbf{u}, t_1, \mathbf{V}\mathbf{x})] = 1 - \frac{1}{d} \sum_{i=1}^n \|\mathbf{P}_{t_1}^\top \mathbf{v}^{(i)}\|^2,$$

and

$$\begin{aligned} \mathbb{E}_{\mathbf{u}, t_1} [\text{Var}(\langle \mathbf{P}_{t_1} \mathbf{u}, \mathbf{x} \rangle | \mathbf{u}, t_1, \mathbf{V}\mathbf{x})] &= 1 - \frac{1}{Nd} \sum_{t_1=1}^N \sum_{i=1}^n \|\mathbf{P}_{t_1}^\top \mathbf{v}^{(i)}\|^2 \\ &= 1 - \frac{1}{Nd} \sum_{i=1}^n \|\mathbf{v}^{(i)}\|^2 = 1 - \frac{n}{Nd}. \end{aligned}$$

■

Appendix E. Auxiliary Lemmas

Lemma 42 Suppose $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ and $\mathbf{B} \in \mathbb{R}^{d_2 \times d_3}$. Then, for all $r, s \geq 1$ and $p, q \geq 1$ such that $1/p + 1/q = 1$, we have

$$\|\mathbf{A}\mathbf{B}\|_{r,s} \leq \|\mathbf{A}\|_{r,p} \|\mathbf{B}\|_{q,s}.$$

Proof First, we note that for any vector $\mathbf{b} \in \mathbb{R}^{d_2}$ we have

$$\|\mathbf{A}\mathbf{b}\|_r = \left\| \sum_{j=1}^{d_2} b_j \mathbf{A}_{:,j} \right\|_r \leq \sum_{j=1}^{d_2} |b_j| \|\mathbf{A}_{:,j}\|_r \leq \|\mathbf{A}\|_{r,p} \|\mathbf{b}\|_q,$$

where the last inequality holds for all conjugate indices p, q and follows from Hölder's inequality. We now have

$$\|\mathbf{AB}\|_{r,s}^s = \sum_{j=1}^{d_3} \|\mathbf{AB}_{:,j}\|_r^s \leq \sum_{j=1}^{d_3} \|\mathbf{A}\|_{r,p}^s \|\mathbf{B}_{:,j}\|_q^s = \|\mathbf{A}\|_{r,p} \|\mathbf{B}\|_{q,s}.$$

■

The next lemma follows from standard Gaussian integration.

Lemma 43 *Suppose $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then $\text{Var}(\|\mathbf{x}\|^2) = 2 \text{tr}(\boldsymbol{\Sigma}^\top \boldsymbol{\Sigma}) + 4\boldsymbol{\mu}^\top \boldsymbol{\Sigma} \boldsymbol{\mu}$.*

The following lemma combines two different techniques for establishing a packing number over the unit ball, the first construction uses volume comparison, whereas the second construction uses Maurey's sparsification lemma, both of which are well-established in the literature.

Lemma 44 *Let \mathcal{P} denote the ϵ -packing number of the unit ball in \mathbb{R}^d . We have*

$$\log \mathcal{P} \leq \left\{ d \log \left(1 + \frac{2}{\epsilon} \right) \right\} \wedge \left\{ \frac{1}{\epsilon^2} (1 + \log(1 + 2d\epsilon^2)) \right\}.$$

Finally, the lemma below allows us to approximate arbitrary Lipschitz functions with two-layer feedforward networks.

Lemma 45 ([7, Propositions 1 and 6]) *Suppose $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies $|f(\mathbf{x})| \leq LR$ and $|f(\mathbf{x}) - f(\mathbf{x}')| \leq L\|\mathbf{x} - \mathbf{x}'\|_2$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ with $\|\mathbf{x}\|_2 \leq R$ and $\|\mathbf{x}'\|_2 \leq R$ and some constants $L, R > 0$. Then, for every $\epsilon > 0$, there exists a positive integer m and $\mathbf{W} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{a} \in \mathbb{R}^m$, such that*

$$\sup_{\|\mathbf{x}\|_2 \leq R} \left| f(\mathbf{x}) - \mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \right| \leq \epsilon.$$

Additionally, we have

$$m \leq C_d \left(\frac{LR(1 + \log(LR/\epsilon))}{\epsilon} \right)^d, \quad \|\mathbf{W}^\top\|_{2,\infty} \leq \frac{1}{R}, \quad \|\mathbf{b}\|_\infty \leq 1, \quad \|\mathbf{a}\|_2 \leq \frac{C_d LR}{\sqrt{m}} \cdot \left(\frac{LR(1 + \log(LR/\epsilon))}{\epsilon} \right)^{\frac{d+1}{2}}.$$