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

Alireza Mousavi-Hosseini¹, Clayton Sanford², Denny Wu³, Murat A. Erdogdu¹

¹University of Toronto and Vector Institute, ²Google Research, ³New York University and Flatiron Institute

{mousavi,erdogdu}@cs.toronto.edu, chsanford@google.com, dennywu@nyu.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 Transformer [VSP⁺17], a neural network architecture that combines attention and feedforward blocks, forms the backbone of large language models and machine learning approaches across many domains [RNSS18, DBK⁺20, BMR⁺20]. The theoretical efforts surrounding the success of Transformers have so far demonstrated various capabilities like in-context learning [ASA⁺23, VONR⁺23, BCW⁺23, ZFB24, KNS24, and others] and chain-of-thought prompting along with its benefits [FZG⁺23, MS24, LLZM24, KS24, 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 [SHT23, SHT24]. Based on these results, [WWHL24] 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?

[FGBM23] 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

| Statistical Model | Feedforward | RNN | Transformer |
|-------------------|---------------|---------------|---------------|
| Simple- q STR | X (Theorem 9) | ✓ (Theorem 5) | ✓ (Theorem 3) |
| qSTR | ✗ (Theorem 9) | X (Theorem 7) | ✓ (Theorem 3) |

Table 1: Summary of main contributions (see Theorem 1). \checkmark indicates a sample complexity upper bound that is almost sequence length-free (up to polylogarithmic factors). \checkmark indicates a lower bound of order $N^{\Omega(1)}$.

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 [Bac17, BES+22, DLS22, BBSS22, DKL+23, AAM23, MHWE24].

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 [AWA⁺22]. While Transformers have demonstrated this ability in certain structured logical tasks, they fail in other simple settings [ZBL⁺23, LAG⁺23]. 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 (qSTR) data generating model, a sequence-to-sequence model where the output at every position depends on a sparse subset of the input tokens. Importantly, this dependence is dynamic, i.e., changes from prompt to prompt, and is described in the input itself. 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 qSTR with a sample complexity almost independent of the length of input sequence N, as long as the number of attention heads is at least q. On the other hand, we develop a new metric-entropy-based argument to derive norm and parameter-count lower bounds for RNNs approximating the qSTR model. Thanks to lower bounds on weight norm, we also obtain a sample complexity lower bound of order $N^{\Omega(1)}$ for RNNs. Further, we show that RNNs can learn a subset of qSTR where the output is a constant sequence, which we call simple-qSTR, 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-qSTR models. The following theorem and Table 1 summarize our main contributions.

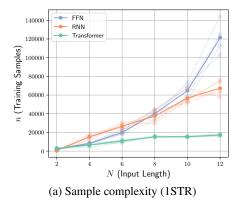
Theorem 1 (Informal). We have the following hierarchy of statistical efficiency for learning qSTR.

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

We empirically validate the intuitions from Theorem 1 in Figure 1. Observe that on a 1STR task, both FFNs and RNNs suffer from a large sample complexity for larger N. However, for a simple-1STR model, RNNs perform closer to Transformers with a much milder dependence on N than FFNs.

1.2 Related Work

While generalization is a fundamental area of study in machine learning theory, theoretical work on the generalization capabilities of Transformers remains relatively sparse. Some works analyze the inductive biases of self-attention through connections to max-margin SVM classifiers [VDT24]. Others quantify complexity in terms of the simplest programs in a formal language (such as the RASP model of [YCA23]) that solve the task and relate that to Transformer generalization [ZBL+23, CS24]. The most relevant works to our own are [EGKZ22, TT23, Tru24], which employ covering numbers to bound the sample complexity of deep Transformers with bounded weights. They demonstrate a logarithmic scaling in the sequence length, depth, and width and apply their bounds



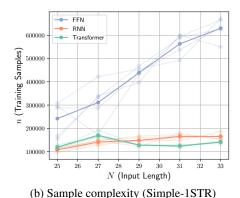


Figure 1: Number of samples required to reach a certain test MSE loss threshold while training with online AdamW. We consider (a) the 1STR model with loss threshold 0.7 and (b) the simple-1STR model with loss threshold 0.02, averaged over 5 experiments. We use a linear link function, standard Gaussian input, d=10 and $d_e=\lfloor 5\log(N)\rfloor$. Positional encodings are sampled uniformly from the unit hypercube. Experimental details and additional results on the effect of q are provided in Appendix E.

to the learnability of sparse Boolean functions. We refine these covering number bounds to better characterize generalization in sequence-to-sequence learning with dynamic sparsity [SHT23]. Our problems formalize long-context reasoning tasks, extending beyond simple retrieval to include challenges like *multi-round coreference resolution* [VOT⁺24].

Expressivity of Transformers. The expressive power of Transformers has been extensively studied in prior works. Universality results establish that Transformers can approximate the output of any continuous function or Turing machine [YBR⁺19, WCM21], as well as measure-to-measure maps [GRRB24], and their memorization capacity is well-understood [MLT24]. However, complexity limitations remain for bounded-size models. Transformers with fixed model sizes are unable to solve even regular languages, such as Dyck and Parity [BAG20, Hah20]. Further work [e.g. MS23] relates Transformers to boolean circuits to establish the hardness of solving tasks like graph connectivity with even polynomial-width Transformers. Additionally, work on self-attention complexity explores how the embedding dimension and number of heads affects the ability of attention layers to approximate sparse matrices [LCW21], recover nearest-neighbor associations [AYB24], and compute sparse averages [SHT23]. The final task closely resembles our *q*STR model and has been applied to relate the capabilities of deep Transformers to parallel algorithms [SHT24]. Several works [e.g. JBKM24, BHBK24, WDL24] introduce sequential tasks where Transformers outperform RNNs or other state space models in parameter-efficient expressivity. We establish similar architectural separations with an added focus on generalization capabilities.

Statistical Separation. Our work is conceptually related to studies on feature learning and adaptivity in feedforward networks, particularly in learning models with sparsity and low-dimensional structures. Prior work has analyzed how neural networks and gradient-based optimization introduce inductive biases that facilitates the learning of low-rank and low-dimensional functions [LMZ18, WLLM19, CB20, MHPG $^+$ 23, OSSW24]. These studies often demonstrate favorable generalization properties based on certain structures of the solution such as large margin or low norm [BFT17, NLB $^+$ 18, OWSS19, WLLM19]. Our goal is to extend efficient learning of low-dimensional concepts to sequential architectures, ensuring sample complexity remains efficient in both input dimension d and context length N. Our approach, motivated by [SHT23, WWHL24], suggests that qSTR is a sequential model whose sparsity serves as a low-dimensional structure, making it the primary determinant of generalization complexity for Transformers.

Notation. For a natural number n, define $[n] \coloneqq \{1,\ldots,n\}$. We use $\|\cdot\|_p$ to denote the ℓ_p norm of vectors. For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\|\mathbf{A}\|_{p,q} \coloneqq \|(\|\mathbf{A}_{:,1}\|_p,\ldots,\|\mathbf{A}_{:,n}\|_p)\|_q$, and $\|\mathbf{A}\|_{\mathrm{op}}$ denotes the operator norm of \mathbf{A} . We use $a \lesssim b$ and $a \leq \mathcal{O}(b)$ interchangeably, which means $a \leq Cb$ for some absolute constant C. We similarly define \gtrsim and Ω . $\tilde{\mathcal{O}}$ and $\tilde{\Omega}$ hide multiplicative constants that depend polylogarithmically on problem parameters. σ denotes the ReLU activation.

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 $p, y \sim P$ where

$$oldsymbol{p} = igg(egin{pmatrix} oldsymbol{x}_1 \ oldsymbol{t}_1 \end{pmatrix}, \dots, igg(oldsymbol{x}_N \ oldsymbol{t}_N \end{pmatrix} igg),$$

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

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

for some $g: \mathbb{R}^{qd} \to \mathbb{R}$. We call this model simple-qSTR if the data distribution is such that $t_i = t$ for all $i \in [N]$ and some 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 t_i . 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 [SHT23], where the authors separate the representational capacity of Transformers and other architectures.

While our main motivation for using the qSTR model is the role of this model as a theoretical benchmark (cf. [SHT23, WWHL24]), we now present an example of how tasks similar to qSTR can arise in natural language modeling. Consider the prompt "For my vacation this summer, I'm considering either Paris or Tokyo. If I go to Paris, I want to visit their art museums, and if I end up in Tokyo, I want to try their cuisine. Can you tell me how much would my first and second option cost respectively?" In this case, t_1 is the token first and refers to the tokens Paris and art museumes, while t_2 is the token second and refers to the tokens Tokyo and cuisine. Note that for either t_1 or t_2 , the answer to the prompt only depends on two tokens out of the entire context, thus this example demonstrates the case of q=2. We refer the interested readers to the multi-round conference resolution task of [VOT+24] for more realistic examples in evaluating large models.

To obtain statistical guarantees, we will impose mild moment assumptions on the data.

Assumption 1. Suppose $\mathbb{E}[\|\mathbf{x}_i\|^r]^{1/r} \leq \sqrt{C_x dr}$ 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 $C_x, C_y > 0$.

We only require the above assumption to establish standard concentration bounds, and it is satisfied as soon as $\|x\|$ is subGaussian and y is sub-Weibull (e.g. g grows at most like a polynomial of degree s). Learning the qSTR 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 $m_g \in \mathbb{N}$, $\boldsymbol{a}_g, \boldsymbol{b}_g \in \mathbb{R}^{m_g}$ and $\boldsymbol{W}_g \in \mathbb{R}^{m_g \times qd}$, such that $\|\boldsymbol{a}_g\|_2 \leq r_a/\sqrt{m_g}$, and $\|(\boldsymbol{W}_g, \boldsymbol{b}_g)\|_{\mathrm{F}} \leq \sqrt{m_g}r_w$ for some constants $r_a, r_w > 0$, and

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

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

Ideally, ε_{2NN} 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 $\varepsilon_{\text{2NN}}>0$ and $m_g=\tilde{\mathcal{O}}\big((P/\sqrt{\varepsilon_{\text{2NN}}})^k\big)$, $r_a=\tilde{\mathcal{O}}\big((P/\sqrt{\varepsilon_{\text{2NN}}})^{(k+1)/2}\big)$, and $r_w=1$ (we can always have $r_w=1$ by homogeneity) [Bac17].

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

$$R^{\text{arc}}(\boldsymbol{\Theta}) := \frac{1}{N} \mathbb{E} \left[\sum_{j=1}^{N} (\hat{y}_{\text{arc}}(\boldsymbol{p}; \boldsymbol{\Theta})_{j} - y_{j})^{2} \right] = \frac{1}{N} \mathbb{E} \left[\|\hat{\boldsymbol{y}}_{\text{arc}}(\boldsymbol{p}; \boldsymbol{\Theta}) - \boldsymbol{y}\|_{2}^{2} \right], \quad (2.1)$$

$$\hat{R}_{n,N}^{\text{arc}}(\boldsymbol{\Theta}) \coloneqq \frac{1}{nN} \sum_{i=1}^{n} \sum_{j=1}^{N} \left(\hat{y}_{\text{arc}}(\boldsymbol{p}^{(i)}; \boldsymbol{\Theta})_{j} - y_{j}^{(i)} \right)^{2}, \tag{2.2}$$

$$\hat{R}_n^{\text{arc}}(\mathbf{\Theta}) := \frac{1}{n} \sum_{i=1}^n \left(\hat{y}_{\text{arc}}(\mathbf{p}^{(i)}; \mathbf{\Theta})_{j^{(i)}} - y_{j^{(i)}}^{(i)} \right)^2, \tag{2.3}$$

where $\{j^{(i)}\}_{i=1}^n$ are i.i.d. position indices drawn from $\mathrm{Unif}([N])$. The goal is to minimize the population risk $R^{\mathrm{arc}}(\Theta)$ 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 $\hat{R}_n^{\tt arc}$ subject to the model parameters belonging on some (e.g., norm-constrained) set Θ . Concretely, let

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

Theorem 3 considers constrained ERM algorithms for bounded-weight transformers with pointwise risk $\hat{R}_n^{\text{TR}}(\Theta)$, and Theorem 5 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^{arc}

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

$$\hat{\boldsymbol{\Theta}}_{\varepsilon} \in \underset{\{\boldsymbol{\Theta}: \hat{R}_{n}^{\operatorname{arc}}(\boldsymbol{\Theta}) - \min \hat{R}_{n}^{\operatorname{arc}} \leq \varepsilon\}}{\operatorname{arg min}} \left\| \operatorname{vec}(\boldsymbol{\Theta}) \right\|_{2}. \tag{2.4}$$

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

3. Beyond ERM, Theorem 9 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 8.

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.

3 Transformers

A single-layer Transformer is composed of an attention and a parallel feedforward layer. Given a sequence $\{z_i\}_{i=1}^N$ of input embeddings where $z_i \in \mathbb{R}^{D_e}$ with embedding dimension D_e , a single head of attention outputs another sequence of length N in \mathbb{R}^{D_e} , given by

$$f_{\mathtt{Attn}}(oldsymbol{p}; oldsymbol{W}_Q, oldsymbol{W}_K, oldsymbol{W}_V) = \left[\sum_{j=1}^N oldsymbol{W}_V oldsymbol{z}_j rac{e^{\langle oldsymbol{W}_Q oldsymbol{z}_i, oldsymbol{W}_K oldsymbol{z}_j
angle}}{\sum_{l=1}^N e^{\langle oldsymbol{W}_Q oldsymbol{z}_i, oldsymbol{W}_K oldsymbol{z}_l
angle}}
ight]_{i \in [N]}.$$

Where W_K, W_Q, W_V are the key, query, and value projection matrices respectively. The output of H units of attention can be concatenated to form multi-head attention with output $h \in \mathbb{R}^{HD_e}$. A two-layer neural network acts on h to generate the final output sequence via

$$f_{\mathtt{2NN}}(m{h};m{a}_{\mathtt{2NN}},m{W}_{\mathtt{2NN}},m{b}_{\mathtt{2NN}}) = m{a}_{\mathtt{2NN}}^ op\sigma(m{W}_{\mathtt{2NN}}m{h}+m{b}_{\mathtt{2NN}}), \quad m{W}_{\mathtt{2NN}} \in \mathbb{R}^{m imes HD_e},m{a}_{\mathtt{2NN}},m{b}_{\mathtt{2NN}} \in \mathbb{R}^m$$

Our architectural choices are standard in theoretical studies of Transformers. We provide full details, including how to obtain input embeddings by positional encoding, in Appendix A.1.

3.1 Learning Guarantees for Multi-Head Transformers

We consider the following parameter class $\Theta_{\text{TR}} = \{\|\mathrm{vec}(\mathbf{\Theta})\|_2 \leq R\}$ and provide a learning guarantee for empirical risk minimizers over Θ_{TR} , with its proof deferred to Appendix A.2.

Theorem 3. Let $\hat{\Theta} = \arg\min_{\Theta \in \Theta_{TR}} \hat{R}_n^{TR}(\Theta)$ and $m = m_g$. Suppose we set H = q and $R^2 = \tilde{\Theta}(r_a^2/m_q + m_q r_w^2 + q^2/d)$. Under Assumptions 1, 2 and 3, we have

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

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

We make the following remarks.

- First, the sample complexity above depends on N only up to log factors. Second, we can remove
 the C₁ factor by performing a clipping operation with a large constant on the Transformer output.
 Note that the first and second terms in the RHS above denote the approximation and estimation
 errors respectively. Extending the above guarantee to cover m ≥ mq and H ≥ q is straightforward.
- This bound provides guidance on the relative merits of scaling the parameter complexity of the feedforward versus the attention layer (which is an active research area related to Transformer scaling laws [HSSL24, JMB⁺24]), by highlighting the trade-off between the two to achieve minimal generalization error. Concretely, $m_g \gg d+q$ represents a regime where the complexity is dominated by the feedforward layer learning the downstream task g, while $m_g \ll d+q$ signifies dominance of the attention layer learning to retrieve the relevant tokens.

Finally, 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 (A.2)—even while maintaining separate \boldsymbol{W}_Q and \boldsymbol{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.

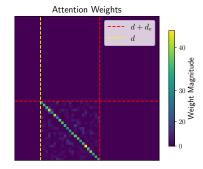


Figure 2: Trained attention weights match our theoretical construction (A.2).

3.2 Limitations of Transformers with Few Heads

We establish the necessity of the linear dependence of H on q. In contrast to [AYB24], 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 qSTR model where $y_i = \frac{1}{\sqrt{qd}} \sum_{j=1}^q (\|\mathbf{x}_{t_{ij}}\|^2 - \mathbb{E}[\|\mathbf{x}_{t_{ij}}\|^2])$, $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{\Sigma}_i)$ such that $\mathbf{\Sigma}_i = \mathbf{I}_d$ for i < N/2 and $\mathbf{\Sigma}_i = 0$ for $i \ge N/2$. Then, there exists a distribution over $(\mathbf{t}_i)_{i \in [N]}$ such that for any choice of $\mathbf{\Theta}_{TR}$ (including arbitrary $\{\mathbf{W}_{QK}^{(h)}\}_{h \in [H]}$), we have

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

Remark. We highlight the importance of the nonlinear dependence of y_i on x for the above lower bound. In particular, for the sparse token averaging task introduced in [SHT23], 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 A.3, we present the proof of Proposition 4, along with Proposition 21 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- $q{\rm STR}$ with a sample complexity only polylogarithmic in N, thus establishing a separation in their learning capability from feedforward networks. Next, we turn to general $q{\rm STR}$, 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. Throughout this section, we focus on bidirectional RNNs, since the $q{\rm STR}$ model is not necessarily causal and the output at position i may depend on future tokens.

4.1 RNNs can learn simple-qSTR

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

$$\begin{aligned} & \boldsymbol{h}_{i}^{\rightarrow} = \boldsymbol{\Pi}_{r_{h}} \big(\boldsymbol{h}_{i-1}^{\rightarrow} + f_{h}^{\rightarrow} (\boldsymbol{h}_{i-1}^{\rightarrow}, \boldsymbol{z}_{i-1}; \boldsymbol{\Theta}_{h}^{\rightarrow}) \big), \quad \forall i \in \{2, \dots, N\} \\ & \boldsymbol{h}_{i}^{\leftarrow} = \boldsymbol{\Pi}_{r_{h}} \big(\boldsymbol{h}_{i+1}^{\leftarrow} + f_{h}^{\leftarrow} (\boldsymbol{h}_{i+1}^{\leftarrow}, \boldsymbol{z}_{i+1}; \boldsymbol{\Theta}_{h}^{\leftarrow}) \big), \quad \forall i \in \{1, \dots, N-1\}, \end{aligned}$$

where $\Pi_{r_h}:\mathbb{R}^{d_h}\to\mathbb{R}^{d_h}$ 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 $\mathbf{z}_i=(\mathbf{z}_i^{\top},\operatorname{enc}(i,\mathbf{t}_i)^{\top})^{\top}$ is the encoding of \mathbf{z}_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

$$y_i = f_y(\boldsymbol{h}_i^{\rightarrow}, \boldsymbol{h}_i^{\leftarrow}, \boldsymbol{z}_i; \boldsymbol{\Theta}_y),$$

which is an L_y -layer feedforward network. Specifically, we consider an RNN with deep transitions [PGCB13] and let $f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow})$ be an L_h -layer feedforward network (see Appendix B.1 for complete definitions). We denote the complete output of the RNN via

$$\hat{m{y}}_{\mathtt{RNN}}(m{p};m{\Theta}_{\mathtt{RNN}}) = (f_y(m{h}_1^{
ightarrow},m{h}_1^{\leftarrow},m{z}_1;m{\Theta}_y),\dots,f_y(m{h}_N^{
ightarrow},m{h}_N^{\leftarrow},m{z}_N;m{\Theta}_y)) \in \mathbb{R}^N.$$

We have the following guarantee for RNNs learning simple-qSTR models.

Theorem 5. Let $\hat{\Theta} = \arg\min_{\Theta \in \Theta_{\text{RNN}}} \hat{R}_n^{\text{RNN}}(\Theta)$ (with Θ_{RNN} defined in Equation (B.2)). Suppose Assumptions 1, 2 and 3 hold with the simple-qSTR model, i.e. $t_i = t$ for all $i \in [N]$ and some t drawn from $[N]^q$. Then, with $L_h, L_y = \mathcal{O}(1)$, $r_h = \tilde{\Theta}(\sqrt{qd})$, and proper hyperparameters in Θ_{RNN} (see Appendix B.1), we obtain

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

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 dimension and norm of RNN weights, implicit in the formulation above, must have a similar polynomial scaling as evident by the proof of the above theorem in Appendix B.

4.2 RNNs cannot learn general qSTR

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 chracterized by

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

where $f_y: \mathbb{R}^{d_h} \times \mathbb{R}^{d_h} \times \mathbb{R}^{d} \times [N]^{q+1} \to \mathbb{R}$, f_h^{\to} , $f_h^{\leftarrow}: \mathbb{R}^{d_h} \times \mathbb{R}^{d} \times [N]^{q+1} \to \mathbb{R}^{d_h}$, U^{\to} , $U^{\leftarrow} \in \mathbb{R}^{d_h \times d_h}$, d_h is the width of the model, and $r_h > 0$ is some constant. Moreover, $\operatorname{proj}_{r_h}: \mathbb{R}^{d_h} \to \mathbb{R}^{d_h}$

is any mapping that guarantees $\|\operatorname{proj}_{r_h}(\cdot)\|_2 \leq r_h$. As mentioned before, this operation mirrors the layer normalization to ensure that h_i remains stable. Further, we assume $f_y(\cdot, \boldsymbol{x}, \boldsymbol{t})$ is \mathfrak{L}/r_h -Lipschitz for all $\boldsymbol{x} \in \mathbb{R}^d$ and $\boldsymbol{t} \in [N]^q$. This formulation covers different variants of (bidirectional) RNNs used in practice such as LSTM and GRU, and includes the RNN formulation of Section 4.1 as a special case. Define $\boldsymbol{U} \coloneqq (\boldsymbol{U}^{\rightarrow}, \boldsymbol{U}^{\leftarrow}) \in \mathbb{R}^{d_h \times 2d_h}$ 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 \boldsymbol{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 qSTR 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 Appendix B.5.

Proposition 6. 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

$$rac{1}{N} \mathbb{E} \Big[\| oldsymbol{y} - \hat{oldsymbol{y}}_{ exttt{RNN}}(oldsymbol{p}) \|^2 \Big] \leq c,$$

implies

$$d_h \geq \Omega\Big(\frac{N}{\log(1+\mathfrak{L}^2\|\boldsymbol{U}\|_{op}^2)}\Big), \quad \textit{and} \quad \|\boldsymbol{U}\|_{op}^2 \geq \Omega\Big(\frac{N}{\mathfrak{L}^2\log(1+d_h)}\Big).$$

Remark. Note that the unboundedness of Gaussian random variables is not an issue for approximation here, since $(g(\boldsymbol{x}_1),\ldots,g(\boldsymbol{x}_N))$ is highly concentrated around $\mathbb{S}^{N-1}(\sqrt{N})$. In fact, one can directly assume $(g(\boldsymbol{x}_1),\ldots,g(\boldsymbol{x}_N))\sim \mathrm{Unif}(\mathbb{S}^{N-1}(\sqrt{N}))$ 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 qSTR 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 [SHT23] using different tools. More importantly, the norm lower bound $\|U\|_F \geq \tilde{\Omega}(\sqrt{N})$ 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 $f_h^{\rightarrow}(h,x,t;\Theta_h^{\rightarrow})$ to denote a general parameterized function (similarly with f^{\leftarrow}). We will assume f_y is given by a feedforward network,

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

where $\boldsymbol{h}=(\boldsymbol{h}^{\rightarrow},\boldsymbol{h}^{\leftarrow})\in\mathbb{R}^{2d_h},~\boldsymbol{z}=(\boldsymbol{x}_i,f_E(\boldsymbol{t}_i,i))\in\mathbb{R}^{d+d_E}.$ Here, $f_E(\boldsymbol{t}_i,i)$ is an arbitrary encoding function with arbitrary dimension d_E . Then $\boldsymbol{\Theta}_y=(\boldsymbol{U},\boldsymbol{W}_y,\boldsymbol{b}_y,\boldsymbol{W}_2,\boldsymbol{b}_2,\ldots,\boldsymbol{W}_{L_y}),$ and $\boldsymbol{\Theta}_{\text{RNN}}=(\boldsymbol{U},\boldsymbol{\Theta}_y,\boldsymbol{\Theta}_h^{\rightarrow},\boldsymbol{\Theta}_h^{\leftarrow}).$ Note that thanks to the homogeneity of ReLU, we can always reparameterize the network by taking $\bar{\boldsymbol{h}}=\boldsymbol{h}/r_h,\bar{\boldsymbol{W}}_y=\boldsymbol{W}_y/r_h,\bar{\boldsymbol{b}}_y=\boldsymbol{b}_y/r_h,$ and $\bar{\boldsymbol{W}}_2=\boldsymbol{W}_2/r_h$ without changing the prediction function. Thus, in the following, we take $r_h=1$ without losing the expressive power of the network. We then have the following sample complexity lower bound.

Theorem 7. Consider the 1STR model of Proposition 6. 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 $\|\operatorname{vec}(\Theta_{RNN})\|_2 \leq e^{N^c/L_y}$ 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 (2.4). 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,

$$rac{1}{N}\,\mathbb{E}igg[\left\|\hat{oldsymbol{y}}_{\mathtt{RNN}}(oldsymbol{p};\hat{oldsymbol{\Theta}}_{n,arepsilon})-oldsymbol{y}
ight\|_2^2igg]\geq c_3.$$

Remark. It is possible to remove the subexponential bound on $\|\text{vec}(\Theta_{\text{RNN}})\|$ by allowing the learner to search over families of RNNs with arbitrary $d_h \leq e^{N^c}$ rather than fixing a single d_h . Additionally, one would avoid solutions that violate this norm constraint in practice due to numerical instability.

To prove the above theorem, we use the fact that an RNN that generalizes on the entire data distribution (hence approximates the 1STR model) requires a weight norm that scales with \sqrt{N} , while overfitting on the n samples in the training set with zero empirical risk is possible with a $\operatorname{poly}(n)$ 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. A similar approach was taken in [POW⁺24] to prove sample complexity separations between two and three-layer feedforward networks. The complete proof is presented in Appendix B.6.

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 $\boldsymbol{p} \mapsto f(\boldsymbol{T}, \boldsymbol{W}\boldsymbol{x})$ where $\boldsymbol{W} \in \mathbb{R}^{m_1 \times Nd}$ is the weight matrix in the first layer, $\boldsymbol{x} = (\boldsymbol{x}_1^\top, \dots, \boldsymbol{x}_N^\top)^\top \in \mathbb{R}^{Nd}$, and $f: [N]^{qN} \times \mathbb{R}^{m_1} \to \mathbb{R}^N$ implements the rest of the network. Unlike the Transformer architecture, here we give the network full information of $\boldsymbol{T} = (\boldsymbol{t}_1, \dots, \boldsymbol{t}_N)$, and in particular the network can implement arbitrary encodings of the position variables $\boldsymbol{t}_1, \dots, \boldsymbol{t}_N$. This formulation covers usual approaches where encodings of \boldsymbol{t} are added to or concatenated with \boldsymbol{x} .

For our negative result on feedforward networks, we can further restrict the class of qSTR models, and only look at simple-qSTR where \hat{R}_n of (2.3) and $\hat{R}_{n,N}$ of (2.2) will be equivalent. Additionally, the lower bound of this section holds regardless of the loss function used for training; for some arbitrary loss $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$, we define the empirical risk of the FFN as

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

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

Definition 8. Let A_{SP} denote the set of algorithms that return a stationary point of the regularized empirical risk. Specifically, for every $A \in A_{SP}$, $A(S_n)$ returns $f_{A(S_n)}$, $W_{A(S_n)}$, such that

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

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 $A \in \mathcal{A}_{ERM}$ returns

$$A(S_n) = \mathop{rg\min}_{\{f, oldsymbol{W}: \hat{\mathcal{L}}^{\mathtt{FFN}}(f, oldsymbol{W}) \leq arepsilon\}} \|oldsymbol{W}\|_F,$$

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

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. The following minimax lower bound shows that all algorithms in class \mathcal{A} fail to learn even the subset of simple-qSTR models with a sample complexity sublinear in N.

Theorem 9. 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 8. Then,

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

with probability 1 over the training set S_n .

Remark. The above lower bound implies that learning the simple 1STR model with FFNs requires at least Nd samples. Note that here we do not have any assumption on m_1 , i.e. the network can have infinite width. This is a crucial difference with the lower bounds in [SHT23, WWHL24] which are computational, i.e., a similar model cannot be learned unless $m_1 \ge Nd$.

The main intuition is that from the stationarity property of Definition 8, the rows of the trained W will always be in the span of the training data $x^{(i)}$ for $i \in [n]$. This is an n-dimensional subspace, and the best predictor that only depends on this subspace still has a loss determined by the variance of y conditioned on this subspace. By randomizing the target direction u, the label y can depend on all Nd target directions. As a result, as long as n < Nd, this variance will be bounded away from zero, leading to the failure of FFNs, even with infinite compute/width. See Appendix D for detailed proof.

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 $q{\rm STR}$ 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- $q{\rm STR}$ 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 $q{\rm STR}$ models, and to study sample complexity separations that highlight the role of depth in Transformers.

Acknowledgments and Disclosure of Funding

The authors thank Alberto Bietti and Song Mei for useful discussions. MAE was partially supported by the NSERC Grant [2019-06167], the CIFAR AI Chairs program, the CIFAR Catalyst grant, and the Ontario Early Researcher Award.

References

- [AAM23] Emmanuel Abbe, Enric Boix Adsera, and Theodor Misiakiewicz, *Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics*, The Thirty Sixth Annual Conference on Learning Theory, PMLR, 2023, pp. 2552–2623.
- [ACDS23] 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** (2023), 45614–45650.
- [ASA⁺23] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou, *What learning algorithm is in-context learning? investigations with linear models*, The Eleventh International Conference on Learning Representations, 2023.
- [AWA⁺22] 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** (2022), 38546–38556.
 - [AYB24] Noah Amsel, Gilad Yehudai, and Joan Bruna, On the benefits of rank in attention layers, arXiv preprint arXiv:2407.16153 (2024).
 - [Bac17] Francis Bach, *Breaking the curse of dimensionality with convex neural networks*, Journal of Machine Learning Research **18** (2017), no. 19, 1–53.
 - [BAG20] S. Bhattamishra, Kabir Ahuja, and Navin Goyal, *On the ability and limitations of transformers to recognize formal languages*, Conference on Empirical Methods in Natural Language Processing, 2020.
- [BBSS22] Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song, *Learning single-index models with shallow neural networks*, Advances in Neural Information Processing Systems, 2022.

- [BCW⁺23] 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).
- [BES⁺22] 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).
- [BFT17] 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).
- [BHBK24] Satwik Bhattamishra, Michael Hahn, Phil Blunsom, and Varun Kanade, Separations in the representational capabilities of transformers and recurrent architectures, 2024.
- [BMR⁺20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al., *Language models are few-shot learners*, Advances in neural information processing systems **33** (2020), 1877–1901.
 - [CB20] Lenaic Chizat and Francis Bach, *Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss*, 2020.
 - [CLZ20] Minshuo Chen, Xingguo Li, and Tuo Zhao, *On generalization bounds of a family of recurrent neural networks*, Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics, Proceedings of Machine Learning Research, vol. 108, PMLR, 2020, pp. 1233–1243.
 - [CS24] Sourav Chatterjee and Timothy Sudijono, Neural networks generalize on low complexity data, ArXiv abs/2409.12446 (2024).
- [DBK⁺20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al., *An image is worth 16x16 words: Transformers for image recognition at scale*, arXiv preprint arXiv:2010.11929 (2020).
- [DKL⁺23] 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).
 - [DLS22] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi, *Neural Networks can Learn Representations with Gradient Descent*, Conference on Learning Theory, 2022.
- [EGKZ22] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang, *Inductive biases and variable creation in self-attention mechanisms*, International Conference on Machine Learning, PMLR, 2022, pp. 5793–5831.
- [FGBM23] 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).
- [FZG⁺23] 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).
- [GRRB24] Borjan Geshkovski, Philippe Rigollet, and Domènec Ruiz-Balet, *Measure-to-measure interpolation using transformers*, arXiv preprint arXiv:2411.04551 (2024).
 - [Hah20] Michael Hahn, *Theoretical limitations of self-attention in neural sequence models*, Transactions of the Association for Computational Linguistics **8** (2020), 156–171.
- [HSSL24] Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li, What matters in transformers? not all attention is needed, 2024.

- [JBKM24] Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and Eran Malach, *Repeat after me: Transformers are better than state space models at copying*, ArXiv **abs/2402.01032** (2024).
- [JMB⁺24] Samy Jelassi, Clara Mohri, David Brandfonbrener, Alex Gu, Nikhil Vyas, Nikhil Anand, David Alvarez-Melis, Yuanzhi Li, Sham M. Kakade, and Eran Malach, *Mixture of parrots: Experts improve memorization more than reasoning*, 2024.
- [KNS24] Juno Kim, Tai Nakamaki, and Taiji Suzuki, *Transformers are minimax optimal non-parametric in-context learners*, ICML 2024 Workshop on In-Context Learning, 2024.
 - [KS24] Juno Kim and Taiji Suzuki, *Transformers provably solve parity efficiently with chain of thought*, arXiv preprint arXiv:2410.08633 (2024).
- [LAG⁺23] Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang, *Exposing attention glitches with flip-flop language modeling*, 2023.
- [LCW21] Valerii Likhosherstov, Krzysztof Choromanski, and Adrian Weller, *On the expressive power of self-attention matrices*, ArXiv **abs/2106.03764** (2021).
- [LIPO23] Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oy-mak, *Transformers as algorithms: Generalization and stability in in-context learning*, International Conference on Machine Learning, PMLR, 2023, pp. 19565–19594.
- [LLZM24] Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma, *Chain of thought empowers transformers to solve inherently serial problems*, The Twelfth International Conference on Learning Representations, 2024.
- [LMZ18] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang, Algorithmic regularization in overparameterized matrix sensing and neural networks with quadratic activations, Conference On Learning Theory, PMLR, 2018, pp. 2–47.
- [MHPG⁺23] Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A Erdogdu, *Neural networks efficiently learn low-dimensional representations with sgd*, The Eleventh International Conference on Learning Representations, 2023.
- [MHWE24] 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).
 - [MLT24] Sadegh Mahdavi, Renjie Liao, and Christos Thrampoulidis, *Memorization capacity of multi-head attention in transformers*, The Twelfth International Conference on Learning Representations, 2024.
 - [MS23] William Merrill and Ashish Sabharwal, *The expressive power of transformers with chain of thought*, 2023.
 - [MS24] ______, *The expressive power of transformers with chain of thought*, The Twelfth International Conference on Learning Representations, 2024.
- [NLB⁺18] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro, *Towards understanding the role of over-parametrization in generalization of neural networks*, arXiv preprint arXiv:1805.12076 (2018).
- [OSSW24] Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu, *Pretrained transformer efficiently learns low-dimensional target functions in-context*, arXiv preprint arXiv:2411.02544 (2024).
- [OWSS19] Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro, *A function space view of bounded norm infinite width relu nets: The multivariate case*, 2019.
- [PGCB13] Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio, *How to construct deep recurrent neural networks*, arXiv preprint arXiv:1312.6026 (2013).

- [POW⁺24] Suzanna Parkinson, Greg Ongie, Rebecca Willett, Ohad Shamir, and Nathan Srebro, *Depth separation in norm-bounded infinite-width neural networks*, The Thirty Seventh Annual Conference on Learning Theory, PMLR, 2024, pp. 4082–4114.
- [RNSS18] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever, *Improving language understanding by generative pre-training*, OpenAI Blog (2018).
- [SHT23] Clayton Sanford, Daniel J Hsu, and Matus Telgarsky, Representational strengths and limitations of transformers, Advances in Neural Information Processing Systems 36 (2023).
- [SHT24] Clayton Sanford, Daniel Hsu, and Matus Telgarsky, *Transformers, parallel computation, and logarithmic depth*, Proceedings of the 41st International Conference on Machine Learning, 2024.
- [Tru24] Lan V Truong, On rank-dependent generalisation error bounds for transformers, arXiv preprint arXiv:2410.11500 (2024).
- [TT23] Jacob Trauger and Ambuj Tewari, Sequence length independent norm-based generalization bounds for transformers, 2023.
- [VDT24] Bhavya Vasudeva, Puneesh Deora, and Christos Thrampoulidis, *Implicit bias and fast convergence rates for self-attention*, ArXiv **abs/2402.05738** (2024).
- [VONR+23] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov, *Transformers learn in-context by gradient descent*, International Conference on Machine Learning, PMLR, 2023, pp. 35151–35174.
 - [VOT+24] Kiran Vodrahalli, Santiago Ontanon, Nilesh Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna, Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, Rohan Anil, Ethan Dyer, Siamak Shakeri, Roopali Vij, Harsh Mehta, Vinay Ramasesh, Quoc Le, Ed Chi, Yifeng Lu, Orhan Firat, Angeliki Lazaridou, Jean-Baptiste Lespiau, Nithya Attaluri, and Kate Olszewska, *Michelangelo: Long context evaluations beyond haystacks via latent structure queries*, 2024.
 - [VSP+17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, Attention is all you need, Advances in Neural Information Processing Systems, vol. 30, 2017.
 - [WCM21] Colin Wei, Yining Chen, and Tengyu Ma, Statistically meaningful approximation: a case study on approximating turing machines with transformers, 2021.
 - [WDL24] Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu, Rnns are not transformers (yet): The key bottleneck on in-context retrieval, 2024.
- [WLLM19] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma, Regularization matters: Generalization and optimization of neural nets vs their induced kernel, Advances in Neural Information Processing Systems **32** (2019).
- [WWHL24] Zixuan Wang, Stanley Wei, Daniel Hsu, and Jason D. Lee, *Transformers provably learn sparse token selection while fully-connected nets cannot*, Proceedings of the 41st International Conference on Machine Learning, 2024.
- [YBR⁺19] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar, *Are transformers universal approximators of sequence-to-sequence functions?*, 2019.
- [YCA23] Andy Yang, David Chiang, and Dana Angluin, *Masked hard-attention transformers* recognize exactly the star-free languages, 2023.
- [ZBL⁺23] Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio, and Preetum Nakkiran, *What algorithms can transformers learn? a study in length generalization*, ArXiv abs/2310.16028 (2023).
- [ZFB24] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett, *Trained transformers learn linear models in-context*, Journal of Machine Learning Research **25** (2024), no. 49, 1–55.

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.

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 equivaraince of Transformers, we append positional information to the input tokens. Given a prompt p, we consider an encoding given by

$$oldsymbol{Z}(oldsymbol{p}) = egin{pmatrix} oldsymbol{x}_1 & \dots & oldsymbol{x}_N \ \operatorname{enc}(1,oldsymbol{t}_1) & \dots & \operatorname{enc}(N,oldsymbol{t}_N) \end{pmatrix} \in \mathbb{R}^{D_e imes N},$$

where $\mathrm{enc}:[N]\times[N]^q\to\mathbb{R}^{d_{\mathrm{enc}}}$ provides the encoding of the position and of t_i , and $D_e\coloneqq d+d_{\mathrm{enc}}$. We use z_i to refer to the *i*th column above. We remark that allowing enc to take t_i as input allows specific encodings of the indices t_i that take advantage of the qSTR structure; examples of this have been considered in prior works [WWHL24]. 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 t_i in their encoding procedure. To specify enc, we use a set of vectors $\{\omega_i\}_{i=1}^N$ 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

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

hence $d_{\rm enc}=(q+1)d_e$ and $D_e=d+(q+1)d_e$. The $\sqrt{d/q}$ prefactor ensures that \boldsymbol{x}_i and ${\rm enc}(i,\boldsymbol{t}_i)$ 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

$$f_{\mathtt{Attn}}(oldsymbol{p}; oldsymbol{W}_Q, oldsymbol{W}_K, oldsymbol{W}_V) = \left[\sum_{j=1}^N oldsymbol{W}_V oldsymbol{z}_j rac{e^{\langle oldsymbol{W}_Q oldsymbol{z}_i, oldsymbol{W}_K oldsymbol{z}_j
angle}}{\sum_{l=1}^N e^{\langle oldsymbol{W}_Q oldsymbol{z}_i, oldsymbol{W}_K oldsymbol{z}_l
angle}}
ight]_{i \in [N]}.$$

Where W_K, W_Q, W_V are the key, query, and value projection matrices respectively. We can simplify the presentation by replacing $W_0^{\dagger}W_K$ with a single parameterizing matrix for query-key projections denoted by $W_{QK} \in \mathbb{R}^{D_e \times D_e}$, and absorbing W_V into the weights of the feedforward layer. This provides us with a simplified parameterization of attention, which we denote by $f_{\text{Attn}}(\boldsymbol{p}; \boldsymbol{W}_{\text{QK}})$. This simplification is standard in theoretical works (see e.g. [LIPO23, ACDS23, ZFB24, WWHL24]). 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 $f_{\mathtt{Attn}}(\boldsymbol{p}; \boldsymbol{W}_{\mathrm{OK}}^{(h)})$. The output of the multi-head attention at position i is then given by

$$f_{\mathtt{Attn}}^{(H)}(\boldsymbol{p};\boldsymbol{W}_{\mathrm{QK}}^{(1)},\ldots,\boldsymbol{W}_{\mathrm{QK}}^{(H)})_{i} = (f_{\mathtt{Attn}}(\boldsymbol{p};\boldsymbol{W}_{\mathrm{QK}}^{(1)})_{i},\ldots,f_{\mathtt{Attn}}(\boldsymbol{p};\boldsymbol{W}_{\mathrm{QK}}^{(H)})_{i})^{\top} \in \mathbb{R}^{HD_{e}}.$$

We will denote by $\Theta_{\rm QK} = ({m W}_{\rm OK}^{(1)}, \dots, {m W}_{\rm OK}^{(H)})$ 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 $oldsymbol{h} \in \mathring{\mathbb{R}}^{HD_e}$, the output of the network is given by

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

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

$$\hat{y}_{\mathtt{TR}}(oldsymbol{p};oldsymbol{\Theta}_{\mathtt{TR}})_i = f_{\mathtt{2NN}}(f_{\mathtt{Attn}}^{(H)}(oldsymbol{p};oldsymbol{\Theta}_{\mathtt{QK}})_i;oldsymbol{\Theta}_{\mathtt{2NN}})$$

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

A.2 Proof of Theorem 3

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

Theorem 10. Let $\hat{\mathbf{\Theta}} := \arg\min_{\mathbf{\Theta} \in \Theta_{\mathtt{TR}}} \hat{R}_n^{\mathtt{TR}}(\mathbf{\Theta})$, where

$$\Theta_{\mathrm{TR}} \coloneqq \Big\{ \|\boldsymbol{a}_{\mathrm{2NN}}\|_2 \leq r_a / \sqrt{m}, \|(\boldsymbol{W}_{\mathrm{2NN}}, \boldsymbol{b}_{\mathrm{2NN}})\|_F \leq r_w \sqrt{m}, \left\|\boldsymbol{W}_{\mathrm{QK}}^{(h)}\right\|_{2.1} \leq \alpha \ \forall h \in [H] \Big\}.$$

Suppose H=q, $m=m_g$, and $\alpha=\tilde{\Theta}(1)$ (given in Lemma 11). Then, under Assumptions 1, 2 and 3, with probability at least $1-n^{-c}$ for some absolute constant c>0, we have

$$R^{\text{TR}}(\hat{\mathbf{\Theta}}) \leq \mathcal{O}(\varepsilon_{\text{NN}}^2) + \tilde{\mathcal{O}}\left(C_1\sqrt{\frac{(m_gq(d+q) + r_z^6r_a^2r_w^2q^2 \wedge q(q^2+d^2))}{n}}\right), \tag{A.1}$$

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 qSTR models.

Lemma 11. 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

$$\sup_{\left\{\left\|\boldsymbol{x}_{j}\right\|_{2} \leq r_{x},\,\forall j \in [N]\right\}}\left|g(\boldsymbol{x}_{t_{i1}},\ldots,\boldsymbol{x}_{t_{iq}}) - \hat{y}_{\mathtt{TR}}(\boldsymbol{p};\boldsymbol{\Theta}_{\mathtt{TR}})_{i}\right| \leq 2\sqrt{\varepsilon_{\mathtt{2NN}}},$$

and

$$\left\|\boldsymbol{a}_{2\mathrm{NN}}\right\|_{2} \leq \frac{r_{a}}{\sqrt{m}}, \quad \left\|(\boldsymbol{W}_{2\mathrm{NN}},\boldsymbol{b}_{2\mathrm{NN}})\right\|_{F} \leq \sqrt{m}r_{w}, \quad \left\|\boldsymbol{W}_{\mathrm{QK}}^{(h)}\right\|_{2,1} \leq \frac{2d_{e}q}{d}\log\left(\frac{2r_{a}r_{w}r_{x}N\sqrt{q}}{\varepsilon_{2\mathrm{NN}}}\right),$$

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

$$f_{ ext{Attn}}(oldsymbol{p}; oldsymbol{W}_{ ext{OK}}^{(h)})_i pprox oldsymbol{z}_{t_{ih}}.$$

Note that to do so, for each key token z_j , we only need to compute $\langle \omega_{t_{ih}}, \omega_j \rangle$. Therefore, most entries in $W_{\mathrm{QK}}^{(h)}$ can be zero. We only require a block of $d_e \times d_e$, which corresponds to comparing ω_j and $\omega_{t_{ih}}$ when comparing query z_i and key z_j . Thus, we let

$$\boldsymbol{W}_{\mathrm{QK}}^{(h)} = \begin{pmatrix} \mathbf{0}_{(d+hd_e)\times d} & \mathbf{0}_{(d+hd_e)\times d_e} & \mathbf{0}_{(d+hd_e)\times qd_e} \\ \mathbf{0}_{d_e\times d} & \alpha \mathbf{I}_{d_e} & \mathbf{0}_{d_e\times qd_e} \\ \mathbf{0}_{(q-h)d_e\times d} & \mathbf{0}_{(q-h)d_e\times d_e} & \mathbf{0}_{(q-h)d_e\times qd_e} \end{pmatrix}$$
(A.2)

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

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

for every matrix A. We will specifically choose 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 a_{2NN} , W_{2NN} , b_{2NN} , we let $a_{2NN} = a_g$ and $b_{2NN} = b_g$ from Assumption 2, and define W_{2NN} by extending W_g with zero entries such that

$$egin{aligned} oldsymbol{W}_{ exttt{2NN}}egin{pmatrix} oldsymbol{z}_{t_{i1}} \ dots \ oldsymbol{z}_{t_{iq}} \end{pmatrix} = oldsymbol{W}_{g}egin{pmatrix} oldsymbol{x}_{t_{i1}} \ dots \ oldsymbol{x}_{t_{iq}} \end{pmatrix}. \end{aligned}$$

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

$$\left|g(\boldsymbol{x}_{t_{i1}},\ldots,\boldsymbol{x}_{t_{iq}}) - \hat{y}_{\texttt{TR}}(\boldsymbol{p};\boldsymbol{\Theta}_{\texttt{TR}})_i\right| \leq \sqrt{\varepsilon_{\texttt{2NN}}} + \varepsilon_{\texttt{Attn}},$$

where we recall

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

and

$$\begin{split} \varepsilon_{\texttt{Attn}} &= \left| f_{\texttt{2NN}}((\boldsymbol{z}_{t_{i1}}, \dots, \boldsymbol{z}_{t_{iq}}); \boldsymbol{\Theta}_{\texttt{2NN}}) - f_{\texttt{2NN}}(f_{\texttt{Attn}}^{(q)}(\boldsymbol{p}; \boldsymbol{\Theta}_{\texttt{QK}}); \boldsymbol{\Theta}_{\texttt{2NN}}) \right| \\ &\leq r_{a} r_{w} \sqrt{\sum_{h=1}^{q} \left\| \boldsymbol{A} f_{\texttt{Attn}}(\boldsymbol{p}; \boldsymbol{W}_{\texttt{QK}}^{(h)})_{i} - \boldsymbol{A} \boldsymbol{z}_{t_{ih}} \right\|_{2}^{2}} \\ &\leq 2 r_{a} r_{w} r_{x} N \sqrt{q} e^{-\alpha d/(2q)}, \end{split}$$

where we recall $Az_j = x_j$. Thus, with

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

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

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

Lemma 12. Suppose $\{p^{(i)}\}_{i=1}^n$ are n input prompts (not necessarily independent) drawn from the input distribution, with tokens denoted by $\{(x_j^{(i)})_{j=1}^N\}_{i=1}^n$. Under Assumption 1, for any $r_x > 0$ we have

$$\mathbb{P}\bigg(\max_{i\in[n],j\in[N]}\bigg\|\boldsymbol{x}_{j}^{(i)}\bigg\|_{2}\geq r_{x}\bigg)\leq nNe^{-r_{x}^{2}/(2C_{x}ed)}.$$

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

$$\mathbb{P}\left(\max_{i\in[n],j\in[N]} \left\| \boldsymbol{x}_{j}^{(i)} \right\|_{2} \ge r_{x}\right) \le \sqrt{\frac{1}{nN}}.$$

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

$$\mathbb{P}\left(\max_{i,j}\left\|\boldsymbol{x}_{j}^{(i)}\right\|_{2} \geq r_{x}\right) \leq \frac{\mathbb{E}\left[\max_{i,j}\left\|\boldsymbol{x}_{j}^{(i)}\right\|_{2}^{p}\right]}{r_{x}^{p}} \leq \frac{\mathbb{E}\left[\sum_{i,j}\left\|\boldsymbol{x}_{j}^{(i)}\right\|_{2}^{p}\right]}{r_{x}^{p}} \leq \frac{Nn(C_{x}pd)^{p/2}}{r_{x}^{p}}.$$

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

$$\mathbb{P}\left(\max_{i,j} \left\| \boldsymbol{x}_{j}^{(i)} \right\|_{2} \ge r_{x}\right) \le nNe^{-r_{x}^{2}/(2C_{x}ed)},$$

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

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

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

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

$$\mathcal{F}_{2NN} = \{ \boldsymbol{h} \mapsto f_{2NN}(\boldsymbol{h}; \boldsymbol{\Theta}_{2NN}) : \boldsymbol{\Theta}_{2NN} \in \boldsymbol{\Theta}_{2NN} \},$$

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

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

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

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

We choose the radius $\sqrt{H}r_z$ for defining d_∞^{2NN} since on the event of Lemma 12, 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}_{2NN} .

Lemma 13. Suppose f_{2NN} is L_f Lipschitz for every $f_{2NN} \in \mathcal{F}_{2NN}$. Then, for any ϵ_{2NN} , $\epsilon_{Attn} > 0$, on the event of Lemma 12 we have

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

Proof. The proof simply follows from the triangle inequality, namely

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

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

Lemma 14. Suppose $\|\operatorname{vec}(\boldsymbol{\Theta}_{\mathtt{RNN}})\|_2 \leq R$ and $\|\boldsymbol{z}_j^{(i)}\|_2 \leq R$ for all $i \in [n]$ and $j \in [N]$. Then,

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

This is a special case of Lemma 30, proved in Appendix B.

For the next step, define the distance

$$d_{\infty}^{\mathrm{QK}}(\boldsymbol{\Theta}_{\mathrm{QK}},\boldsymbol{\Theta}_{\mathrm{QK}}')\coloneqq\sup_{\boldsymbol{p},j}\left\|\boldsymbol{\Theta}_{\mathrm{QK}}^{\top}\boldsymbol{z}_{j}-\boldsymbol{\Theta'}_{\mathrm{QK}}^{\top}\boldsymbol{z}_{j}\right\|_{2}$$

on Θ_{QK} , where we recall $\Theta_{QK} = (\boldsymbol{W}_{QK}^{(1)}, \dots, \boldsymbol{W}_{QK}^{(H)}) \in \mathbb{R}^{D_e \times HD_e}$. 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 15. Suppose $\left\| \boldsymbol{z}_{j}^{(i)} \right\|_{2} \leq r_{z}$ for all $i \in [n]$ and $j \in [N]$. Then,

$$\log \mathcal{C}(\mathcal{F}_{\mathtt{Attn}}, d_{\infty}^{\mathtt{Attn}}, \epsilon) \leq \log \mathcal{C}\bigg(\Theta_{\mathrm{QK}}, d_{\infty}^{\mathrm{QK}}, \frac{\epsilon}{2r_{z}^{2}}\bigg).$$

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

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

where we used Lemma 39 for the last inequality. Moreover, by [EGKZ22, Corollary A.7],

$$\begin{aligned} \left\| \operatorname{softmax} \left(\boldsymbol{z}_{j}^{\top} \boldsymbol{W}_{\mathrm{QK}}^{(h)} \boldsymbol{Z}^{\top} \right)^{\top} - \operatorname{softmax} \left(\boldsymbol{z}_{j}^{\top} \hat{\boldsymbol{W}}_{\mathrm{QK}}^{(h)} \boldsymbol{Z}^{\top} \right) \right\|_{1} &\leq 2 \left\| \boldsymbol{Z} \boldsymbol{W}^{(h)}_{\mathrm{QK}}^{\top} \boldsymbol{z}_{j} - \boldsymbol{Z} \hat{\boldsymbol{W}}_{\mathrm{QK}}^{(h)}_{\mathrm{QK}}^{\top} \boldsymbol{z}_{j} \right\|_{\infty} \\ &\leq 2 \left\| \boldsymbol{Z}^{\top} \right\|_{2,\infty} \left\| \boldsymbol{W}^{(h)}_{\mathrm{QK}}^{\top} \boldsymbol{z}_{j} - \hat{\boldsymbol{W}}_{\mathrm{QK}}^{(h)}_{\mathrm{QK}}^{\top} \boldsymbol{z}_{j} \right\|_{2}. \end{aligned}$$

Consequently,

$$\Delta \leq 4r_z^4 \sup_{\boldsymbol{p}, j \in S_n} \sum_{h \in [H]} \left\| \boldsymbol{W}_{\text{QK}}^{(h)} \boldsymbol{z}_j - \hat{\boldsymbol{W}}^{(h)} \boldsymbol{\zeta}_j \boldsymbol{z}_j \right\|_2^2$$
$$= 4r_z^4 \sup_{\boldsymbol{p}, j \in S_n} \left\| \boldsymbol{\Theta}_{\text{QK}}^{\top} \boldsymbol{z}_j - \hat{\boldsymbol{\Theta}}_{\text{QK}}^{\top} \boldsymbol{z}_j \right\|_2^2,$$

which completes the proof.

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

Lemma 16. Suppose $\Theta_{QK} = \{\|\Theta_{QK}\|_{2,1} \leq R_{2,1}, \|\Theta_{QK}\|_F \leq R_F\}$ and $\|z_j^{(i)}\|_2 \leq r_z$ for all $i \in [n]$ and $j \in [N]$. Then,

$$\log \mathcal{C}\big(\Theta_{\mathrm{QK}}, d_{\infty}^{\mathrm{QK}}, \epsilon\big) \lesssim \min \Bigg(\frac{r_z^2 R_{2,1}^2 \log(2HD_e^2)}{\epsilon^2}, HD_e^2 \log \Big(1 + \frac{2R_F r_z}{\epsilon}\Big)\Bigg).$$

Proof. The first estimate comes from Maurey's sparsification lemma [BFT17, Lemma 3.2], while the second estimate is based on the inequality

$$\left\| \mathbf{\Theta}_{\mathrm{QK}}^{\top} \mathbf{z}_{j} - \hat{\mathbf{\Theta}}_{\mathrm{QK}}^{\top} \mathbf{z}_{j} \right\|_{2} \leq r_{z} \left\| \mathbf{\Theta}_{\mathrm{QK}} - \hat{\mathbf{\Theta}}_{\mathrm{QK}} \right\|_{\mathrm{F}},$$

and covering $\Theta_{\rm QK}$ with the Frobenius norm, see e.g. Lemma 41.

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

Proposition 17. Suppose $\|\boldsymbol{a}_{2NN}\|_2 \leq r_{m,a}$, $\|(\boldsymbol{W}_{2NN}, \boldsymbol{b}_{2NN})\|_F \leq R_{m,w}$, and $\|\boldsymbol{W}_{QK}^{(h)}\|_{2,1} \leq r_{QK}$ for all $h \in [H]$. Further assume $\|\boldsymbol{z}_j^{(i)}\|_2 \leq r_z$ for all $i \in [n]$ and $j \in [N]$. Let $R := \max(r_{m,a}, R_{m,w}, r_z)$. Then,

 $\log \mathcal{C}(\mathcal{F}_{TR}, d_{\mathcal{F}}, \epsilon) \lesssim m_q H D_e \log(1 + R/\epsilon)$

$$+ \min \left(\frac{r_z^6 r_{m,a}^2 R_{m,w}^2 H^2 r_{QK}^2 \log(HD_e^2)}{\epsilon^2}, HD_e^2 \log \left(1 + \frac{\sqrt{H} r_{QK} r_z^3 r_{m,a} R_{m,w}}{\epsilon} \right) \right).$$

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 13. Moreover, we have

 $R_F \le \sqrt{H}r_{\rm QK}$ and $R_{2,1} \le Hr_{\rm QK}$ in Lemma 16. The rest follows from combining the statements of the previous lemmas.

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

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

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

$$\mathcal{L}_{\tau} := \{ (\boldsymbol{p}, \boldsymbol{y}, j) \mapsto (f_{TR}(\boldsymbol{p})_j - y_j)^2 \wedge \tau : f_{TR} \in \mathcal{F}_{TR} \}, \tag{A.3}$$

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

Lemma 18. Suppose $\max_{i \in [n], j \in [N]} \left\| \mathbf{z}_{j}^{(i)} \right\|_{2} \leq r_{z}$. For the loss class \mathcal{L}_{τ} given by (A.3), we have

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

where $C_1 = m_g H D_e$, $C_2 = r_z^6 r_{m,a}^2 R_{m,w}^2 H^2 r_{OK}^2$, 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 $\ell(\boldsymbol{p}, \boldsymbol{y}, j) = (f(\boldsymbol{p})_j - y_j)^2 \wedge \tau$ and $\ell'(\boldsymbol{p}, \boldsymbol{y}, j) = (f'(\boldsymbol{p})_j - y_j)^2 \wedge \tau$. Then, for any $\alpha \geq 0$, by a standard chaining argument,

$$\begin{split} \hat{\mathfrak{R}}_{n}(\mathcal{L}_{\tau}) &\lesssim \alpha + \int_{\alpha}^{\tau} \sqrt{\frac{\log \mathcal{C}(\mathcal{L}, d_{\infty}^{\mathcal{L}}, \epsilon)}{n}} \mathrm{d}\epsilon. \\ &\lesssim \alpha + \int_{\alpha}^{\tau} \sqrt{\frac{\log \mathcal{C}(\mathcal{F}, d_{\infty}^{\mathsf{TR}}, \epsilon/(2\sqrt{\tau}))}{n}} \\ &\lesssim \alpha + \int_{\alpha}^{\tau} \sqrt{\frac{C_{1} \log(R\sqrt{\tau}/\epsilon)}{n}} \mathrm{d}\epsilon + \left\{ \int_{\alpha}^{\tau} \sqrt{\frac{\tau C_{2} \log(HD_{e}^{2})}{n\epsilon^{2}}} \mathrm{d}\epsilon \right\} \wedge \left\{ \int_{\alpha}^{\tau} \sqrt{\frac{C_{3} \log(1 + C_{4}\sqrt{\tau}/\epsilon)}{n}} \mathrm{d}\epsilon \right\} \\ &\lesssim \alpha + \sqrt{\frac{\tau^{2} C_{1} \log(R\sqrt{\tau}/\alpha)}{n}} + \left\{ \sqrt{\frac{\tau C_{2} \log(HD_{e}^{2})}{n}} \log\left(\frac{\tau}{\alpha}\right) \right\} \wedge \left\{ \sqrt{\frac{\tau^{2} C_{3} \log(1 + C_{4}\sqrt{\tau}/\alpha)}{n}} \right\}, \end{split}$$

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

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 19. Let $\hat{\Theta} = \arg\min_{\Theta \in \Theta_{TR}} \hat{R}_n^{TR}(\Theta)$, where Θ_{TR} is described in Proposition 17. Define $r_z = \sqrt{r_x^2 + d(1+1/q)}$ where r_x is defined in Lemma 12. Let C_1 , C_2 , and C_3 be defined as in Lemma 18. Then, with probability at least $1 - \delta - (nN)^{-1/2}$ over S_n , we have

$$R_{\tau}^{\mathtt{TR}}(\hat{\mathbf{\Theta}}) - \hat{R}_{n}^{\mathtt{TR}}(\hat{\mathbf{\Theta}}) \leq \tilde{\mathcal{O}}\Bigg(\tau\sqrt{\frac{(C_{1} + C_{2} \wedge C_{3})}{n}}\Bigg) + \mathcal{O}\Bigg(\tau\sqrt{\frac{\log(1/\delta)}{n}}\Bigg),$$

where $R_{ au}^{\mathtt{RNN}}(\hat{\mathbf{\Theta}}) \coloneqq \mathbb{E}_{\mathbf{p},j,y}\Big[(\hat{y}_{\mathtt{TR}}(\mathbf{p};\hat{\mathbf{\Theta}})_{j} - y_{j})^{2} \wedge \tau\Big]$

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

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

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

Lemma 20. 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

$$R^{ extsf{TR}}(\hat{oldsymbol{\Theta}}) - R_{ au}^{ extsf{TR}}(\hat{oldsymbol{\Theta}}) \leq \sqrt{rac{1}{n}}.$$

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

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

Moreover,

$$\mathbb{E}\left[\Delta_y^4\right]^{1/2} \le 2 \,\mathbb{E}\left[y_j^4\right]^{1/2} + 2 \,\mathbb{E}\left[\hat{y}(\boldsymbol{p}; \hat{\boldsymbol{\Theta}})_j^4\right]^{1/2}.$$

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

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

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

$$\mathbb{E}\left[\max_{l \in [N]} \|\boldsymbol{x}_{l}\|_{2}^{4}\right] \leq \mathbb{E}\left[\max_{l \in [N]} \|\boldsymbol{x}_{l}\|_{2}^{4r}\right]^{1/r} \leq \mathbb{E}\left[\sum_{l=1}^{N} \|\boldsymbol{x}_{l}\|_{2}^{4r}\right]^{1/r}$$

$$\leq N^{1/r} \, \mathbb{E}\left[\|\boldsymbol{x}_{1}\|_{2}^{4r}\right]^{1/r}$$

$$\lesssim N^{1/r} C_{x}^{2} d^{2} r^{2}$$

$$\lesssim (C_{x} d \log(N))^{2},$$

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

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

We now turn to bounding the probability. We have

$$\mathbb{P}(\Delta_y \ge \sqrt{\tau}) \le \mathbb{P}(|y_j| \ge \frac{\sqrt{\tau}}{2}) + \mathbb{P}(|\hat{y}(\boldsymbol{p}; \hat{\boldsymbol{\Theta}})_j| \ge \frac{\sqrt{\tau}}{2}) \\
\le \exp(-\Omega(\tau^{1/s})) + N \exp(-\Omega(\frac{\tau}{Hr_{m,a}^2 R_{m,w}^2 r_z^2})),$$

where the second inequality follows from sub-Weibull concentration bounds for y and Lemma 12. Choosing $\tau = \Theta(\kappa^2 \log(\kappa^2 N \sqrt{n}) + \log(\kappa^2 \sqrt{n})^s)$ completes the proof.

Proof of Theorem 10. The theorem follows immediately from the approximation guarantee of Lemma 11, the generalization bound of Corollary 19, and the truncation control of Lemma 20. \Box

A.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 $H \ge q$ on the number of heads for all d, at the expense of additional restrictions on the attention matrix.

Proposition 21. Consider the qSTR data model. Suppose d=1 and $y_i=\frac{1}{\sqrt{q}}\sum_{j=1}^q(x_{t_{ij}}^2-\mathbb{E}[x_{t_{ij}}^2])$. Assume $x_i\sim\mathcal{N}(0,\sigma_i^2)$ independently, such that $\sigma_i=1$ for $i<\mathcal{N}/2$ and $\sigma_i=0$ for $i\geq\mathcal{N}/2$. Further, assume the attention weights between the data and positional encoding parts of the tokens are fixed at zero, i.e. $\mathbf{W}_{\mathrm{QK}}^{(h)}=\begin{pmatrix} \mathbf{W}_{x}^{(h)} & \mathbf{0}_{d\times(q+1)d_e}\\ \mathbf{0}_{(q+1)d_e\times d} & \mathbf{W}_{\omega}^{(h)} \end{pmatrix}$ where $\mathbf{W}_{x}^{(h)}\in\mathbb{R}^{d\times d}$ and $\mathbf{W}_{\omega}^{(h)}\in\mathbb{R}^{d\times d}$ 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 $\mathbf{\Theta}_{\mathrm{TR}}$, we have

$$\frac{1}{N} \mathbb{E} \Big[\| \boldsymbol{y} - \hat{\boldsymbol{y}}_{\texttt{TR}}(\boldsymbol{p}; \boldsymbol{\Theta}_{\texttt{TR}}) \|_2^2 \Big] \geq 1 - \frac{H}{q}.$$

Note that in our approximation constructions for learning qSTR, we always fixed the attention weights between data and positional components to be zero, which is why we assume the same in Proposition 21.

Proof of Proposition 21. We will simply choose $t_i = (1, ..., q)$ deterministically for $i \ge \frac{N}{2}$ and draw t_i from an arbitrary distribution for i < N/2. Note that we have

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

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

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

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

$$\alpha_{ij}^{(h)} = \frac{e^{\left\langle \boldsymbol{z}_{i}, \boldsymbol{W}_{\text{QK}}^{(h)} \boldsymbol{z}_{j} \right\rangle}}{\sum_{l=1}^{N} e^{\left\langle \boldsymbol{z}_{i}, \boldsymbol{W}_{\text{QK}}^{(h)} \boldsymbol{z}_{j} \right\rangle}},$$

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

$$R^{TR}(\boldsymbol{\Theta}_{TR}) \ge \frac{1}{N} \sum_{i=N/2}^{N} \mathbb{E}\left[\left(y_{i} - \tilde{\phi}\left(\boldsymbol{A}^{(i)}\boldsymbol{x}_{1:q}, (\boldsymbol{z}_{l})_{l=q+1}^{N}\right)\right)^{2}\right]$$

$$\ge \frac{1}{Nq} \sum_{i=N/2}^{N} \mathbb{E}\left[\operatorname{Var}\left(\|\boldsymbol{x}_{1:q}\|^{2} |\boldsymbol{V}^{(i)}\boldsymbol{x}_{1:q}\right)\right]$$
(A.4)

where $\boldsymbol{V}^{(i)} \in \mathbb{R}^{H \times q}$ is a matrix whose rows form an orthonormal basis of $\operatorname{span}(\boldsymbol{\alpha}_i^{(1)}, \dots, \boldsymbol{\alpha}_i^{(H)})$ where $\boldsymbol{\alpha}_i^{(h)} = (\alpha_{i1}^{(h)}, \dots, \alpha_{iq}^{(h)})^{\top} \in \mathbb{R}^q$ (note that $\boldsymbol{V}^{(i)}$ 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 \boldsymbol{z}_l is independent of $\boldsymbol{x}_{1:q}$ for $l \geq q+1$, and the fact that best predictor of y_i (in L_2 error) given $\boldsymbol{A}^{(i)}\boldsymbol{x}_{1:q}$ is $\mathbb{E}\Big[y_i \,|\, \boldsymbol{V}^{(i)}\boldsymbol{x}_{1:q}\Big]$.

Next, thanks to the structural property of $\boldsymbol{W}_{\mathrm{QK}}^{(h)}$ 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, $\boldsymbol{V}^{(i)}$ is independent of $\boldsymbol{x}_{1:q}$. Therefore,

$$m{x}_{1:q} \, | \, m{V}^{(i)} m{x}_{1:q} \sim \mathcal{N}(m{V}^{(i)}^{ op} m{V}^{(i)} m{x}_{1:q}, m{\mathrm{I}}_q - m{V}^{(i)}^{ op} m{V}^{(i)}).$$

By Lemma 40, we have $\operatorname{Var}(\|\boldsymbol{x}_{1:q}\|^2 | \boldsymbol{V}^{(i)}\boldsymbol{x}_{1:q}) = 2(q-H)$, which combined with (A.4) completes the proof.

We now present the similarly structured proof of Proposition 4.

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

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

Recall $z_i = (x_i^\top, \text{enc}(i, t_i)^\top)$. Fix some $i \ge N/2$, and define

$$\tilde{lpha}_{ij}^{(h)} = e^{\left\langle \operatorname{enc}(i, \boldsymbol{t}_i), \boldsymbol{W}_{\mathrm{QK}}^{(h, e, x)} \boldsymbol{x}_j \right\rangle + \left\langle \operatorname{enc}(i, \boldsymbol{t}_i), \boldsymbol{W}_{\mathrm{QK}}^{(h, e, e)} \operatorname{enc}(j, \boldsymbol{t}_j) \right\rangle}$$

where we use the notation

$$oldsymbol{W}_{ ext{QK}}^{(h)} = egin{pmatrix} oldsymbol{W}_{ ext{QK}}^{(h,x,x)} & oldsymbol{W}_{ ext{QK}}^{(h,x,e)} \ oldsymbol{W}_{ ext{QK}}^{(h,e,x)} & oldsymbol{W}_{ ext{QK}}^{(h,e,e)} \end{pmatrix},$$

for the query-key matrix of each head. Recall that $x_i = 0$ 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 21 that we denote the feedforward layer by $\phi: \mathbb{R}^{HD_e} \to \mathbb{R}$. With this notation, we have

$$\begin{split} \hat{y}_{\text{TR}}(\boldsymbol{p};\boldsymbol{\Theta}_{\text{TR}})_{i} &= \phi(\sum_{j=1}^{N} \alpha_{ij}^{(1)} \boldsymbol{z}_{j}, \dots, \sum_{j=1}^{N} \alpha_{ij}^{(H)} \boldsymbol{z}_{j}) \\ &= \tilde{\phi}\Big(\sum_{j=1}^{q} \alpha_{ij}^{(1)} \boldsymbol{x}_{j}, \dots, \sum_{j=1}^{q} \alpha_{ij}^{(H)} \boldsymbol{x}_{j}, (\tilde{\alpha}_{ij}^{(h)})_{h=1, j=1}^{h=H, j=N}, (\boldsymbol{z}_{j})_{j=l+1}^{N}\Big). \end{split}$$

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

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

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

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

which yields $\left< m{lpha}_i^{(h,r)}, m{x}_{1:q} \right> = \sum_{j=1}^q lpha_{ij}^{(h)} x_{jr}$, and $m{w}_{i,j}^{(h)} \in \mathbb{R}^{qd}$ such that

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

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

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

Once again, by Lemma 40, we conclude that $var(\|\boldsymbol{x}_{1:q}\|^2 | \boldsymbol{V}^{(i)}\boldsymbol{x}_{1:q}) \ge 2(qd - H(q+d))$, which completes the proof.

B Details and Proofs of Section 4

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

B.1 Complete Setup of RNNs

When introducing RNNs in Section 4, we used L_h -layer deep feedforward networks to implement the transitions $f_h^{\rightarrow}(\cdot; \Theta_h^{\rightarrow})$ and $f_h^{\leftarrow}(\cdot; \Theta_h^{\leftarrow})$. These transitions are given by

$$f_h^{\rightarrow}(\cdot; \boldsymbol{\Theta}_h^{\rightarrow}) = \boldsymbol{W}_{L_h}^{\rightarrow} \sigma(\boldsymbol{W}_{L_h-1}^{\rightarrow} \dots \sigma(\boldsymbol{W}_2^{\rightarrow} \sigma(\boldsymbol{W}_1^{\rightarrow}(\cdot) + \boldsymbol{b}_1^{\rightarrow}) + \boldsymbol{b}_2^{\rightarrow}) \dots + \boldsymbol{b}_{L_h-1}^{\rightarrow}), \tag{B.1}$$

with $\Theta_h^{\rightarrow} = (\boldsymbol{W}_1^{\rightarrow}, \boldsymbol{b}_1^{\rightarrow}, \dots, \boldsymbol{W}_{L_h-1}^{\rightarrow}, \boldsymbol{b}_{L_h-1}^{\rightarrow}, \boldsymbol{W}_{L_h}^{\rightarrow})$ and a similar equation for $f^{\leftarrow}(\cdot; \Theta_h^{\leftarrow})$. Recall that the output of the RNN is denoted by

$$\hat{m{y}}_{\mathtt{RNN}}(m{p};m{\Theta}_{\mathtt{RNN}}) = (f_y(m{h}_1^{
ightarrow},m{h}_1^{\leftarrow},m{z}_1;m{\Theta}_y),\dots,f_y(m{h}_N^{
ightarrow},m{h}_N^{\leftarrow},m{z}_N;m{\Theta}_y)) \in \mathbb{R}^N.$$

We now define the constraint set of this architecture. Let

$$\Theta_{\text{RNN}} = \left\{ \boldsymbol{\Theta} : \left\| \operatorname{vec}(\boldsymbol{\Theta}) \right\|_{2} \le R, \left\| \boldsymbol{W}_{L_{h}}^{\rightarrow} \right\|_{\operatorname{op}} \dots \left\| \boldsymbol{W}_{1,h}^{\rightarrow} \right\|_{\operatorname{op}} \le \alpha_{N}, \left\| \boldsymbol{W}_{L_{h}}^{\leftarrow} \right\|_{\operatorname{op}} \dots \left\| \boldsymbol{W}_{1,h}^{\leftarrow} \right\|_{\operatorname{op}} \le \alpha_{N} \right\},$$
(B.2)

where $W_{1,h}^{\rightarrow}$ contains the first d_h columns of 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 $W_{1,h}^{\rightarrow}$ by a factor of $\alpha_N/\prod_{l=2}^{L_h} \|W_l^{\rightarrow}\|_{op}$ 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. [CLZ20].

For Theorem 5 we only require $\alpha_N \leq N^{-1}$. In particular, we can choose $\alpha_N = 0$ and fix $\boldsymbol{W}_{1,h}^{\rightarrow} = \boldsymbol{W}_{1,h}^{\leftarrow} = \boldsymbol{0}$, which would simplify the parameterization of the network. Namely, in our construction f^{\rightarrow} and f^{\leftarrow} do not need to depend on $\boldsymbol{h}^{\rightarrow}$ and $\boldsymbol{h}^{\leftarrow}$ respectively.

B.2 Overview of the Proof of Theorem 5

The following is the roadmap we will take for the proof of Section 4.1. The goal here is to implement a bi-directional RNN in such a way that

$$oldsymbol{h}_i^{
ightarrow} pprox ig(oldsymbol{x}_{t_1} \mathbb{1}[t_1 < i], \ldots, oldsymbol{x}_{t_q} \mathbb{1}[t_q < i]ig),$$

and

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

Throughout this section, we will use the notation

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

We can obtain the hidden states above through the following updates

$$\boldsymbol{h}_{i+1}^{\rightarrow} = \boldsymbol{h}_{i}^{\rightarrow} + \Psi(\boldsymbol{x}_{i}, \boldsymbol{\omega_{t}}, \boldsymbol{\omega}_{i}),$$

and

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

where

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

where we recall $\omega_t = (\omega_{t_1}, \dots, \omega_{t_n})$, and σ is ReLU. As a result, our network must approximate

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

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

$$\|f_h^{\rightarrow}(h_i^{\rightarrow}, z_i; \Theta_h^{\rightarrow}) - \Psi(x_i, \omega_t, \omega_i)\|_2 \le \varepsilon,$$
 (B.3)

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

$$\left\|oldsymbol{h}_i^{
ightarrow} - \sum_{j=1}^{i-1} \Psi(oldsymbol{x}_j, oldsymbol{\omega_t}, oldsymbol{\omega_j})
ight\|_2 \lesssim N arepsilon.$$

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

$$\|f_h^{\rightarrow}(\boldsymbol{h}_i^{\rightarrow}, \boldsymbol{z}_i; \boldsymbol{\Theta}_h^{\rightarrow})_l - \Psi(\boldsymbol{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i)_l\|_2 \le \begin{cases} 0 & t_l \neq i \\ \varepsilon & t_l = i. \end{cases}$$
(B.4)

Note that

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

Since for each $l \in [q]$, $t_l = j$ is possible for at most one $j \in [N]$, (B.4) implies

$$\left\|oldsymbol{h}_i^{
ightarrow} - \sum_{j=1}^{i-1} \Psi(oldsymbol{x}_j, oldsymbol{\omega_t}, oldsymbol{\omega_t})
ight\|_2 \leq \sqrt{q} arepsilon,$$

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

We can implement f_h^{\rightarrow} to satisfy (B.3) 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 [Bac17, Proposition 6]), and the third performs coordinate-wise product between 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 (B.4), 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 h_i^{\rightarrow} and h_i^{\leftarrow} , we first calculate

$$\begin{split} \boldsymbol{h}_i &= f_{hh}(\boldsymbol{h}_i^{\rightarrow}, \boldsymbol{h}_i^{\leftarrow}, \boldsymbol{x}_i, \boldsymbol{\omega}_i, \boldsymbol{\omega}_t) \\ &\approx \boldsymbol{h}_i^{\rightarrow} + \boldsymbol{h}_i^{\leftarrow} + \Psi(\boldsymbol{x}_i, \boldsymbol{\omega}_t, \boldsymbol{\omega}_i) \\ &\approx (\boldsymbol{x}_{t_1}, \dots, \boldsymbol{x}_{t_q}). \end{split}$$

Finally, y_i can be generated from h_i by applying the two-layer neural network from Assumption 2 that approximates $y_i = g(x_t)$.

Note that the construction above has a complexity 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

$$oldsymbol{z}_i = (oldsymbol{x}_i, oldsymbol{\omega}_i, oldsymbol{\omega}_{t_{i1}}, \dots, oldsymbol{\omega}_{t_{iq}})^ op,$$

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

B.3 Approximations

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

$$egin{pmatrix} egin{pmatrix} m{h} \ m{x} \ m{\omega}_i \ m{\omega}_{t_1} \ dots \ m{\omega}_{t_q} \end{pmatrix} \mapsto egin{pmatrix} m{x} \ \langle m{\omega}_i, m{\omega}_{t_1}
angle \ dots \ \langle m{\omega}_i, m{\omega}_{t_q}
angle \end{pmatrix} \mapsto egin{pmatrix} 2m{x} \sigma(\langle m{\omega}_i, m{\omega}_{t_1}
angle - 1/2) \ dots \ 2m{x} \sigma(\langle m{\omega}_i, m{\omega}_{t_q}
angle - 1/2) \end{pmatrix}.$$

The first mapping will be provided by

$$\boldsymbol{\chi}_1 = \boldsymbol{A}_1 \sigma(\boldsymbol{W}_1 \boldsymbol{\chi}_0 + \boldsymbol{b}_1),$$

where $\boldsymbol{\chi}_0 = (\boldsymbol{h}^\top, \boldsymbol{x}^\top, \boldsymbol{\omega}_i^\top, \boldsymbol{\omega}_{t_1}^\top, \dots, \boldsymbol{\omega}_{t_q}^\top)^\top \in \mathbb{R}^{d_h + d + (q+1)d_e}, \, \boldsymbol{W}_1 \in \mathbb{R}^{m_1 \times (d_h + d + (q+1)d_e)}, \, \boldsymbol{b}_1 \in \mathbb{R}^{m_1}, \, \text{and} \, \boldsymbol{A}_1 \in \mathbb{R}^{(d+q) \times m_1}, \, \text{with} \, m_1 \, \text{as the width of the first layer. We will use the notation}$

$$\boldsymbol{\chi}_1 = (\boldsymbol{\chi}_1^{\boldsymbol{x}}, \boldsymbol{\chi}_1^{\boldsymbol{\omega}}(1), \dots, \boldsymbol{\chi}_1^{\boldsymbol{\omega}}(q))$$

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

$$\boldsymbol{\chi}_2 = \boldsymbol{A}_2 \sigma(\boldsymbol{W}_2 \boldsymbol{\chi}_1 + \boldsymbol{b}_2),$$

where $W_2 \in \mathbb{R}^{m_2 \times (d+q)}$, $b_2 \in \mathbb{R}^{m_2}$, and $A_2 \in \mathbb{R}^{dq \times m_2}$. We will similarly use the notation $\chi_2 = (\chi_2(1), \dots, \chi_2(q))$, where our goal is to have $\chi_2(l) \approx 2x\sigma(\langle \omega_i, \omega_{t_l} \rangle - 1/2)$. To implement the first mapping, we rely on the following lemma.

Lemma 22. Let σ be the ReLU activation. For any $\varepsilon > 0$ and positive integer d_e , there exists $m = \mathcal{O}(d_e^3(\log(d_e/\varepsilon)/\varepsilon)^2)$, $\boldsymbol{a} \in \mathbb{R}^m$, $\boldsymbol{W} \in \mathbb{R}^{m \times 2d_e}$, and $\boldsymbol{b} \in \mathbb{R}^m$, such that

$$\sup_{\boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \in \mathbb{S}^{d_e-1}} \left| \langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle - \boldsymbol{a}^\top \sigma \bigg(\boldsymbol{W} \begin{pmatrix} \boldsymbol{\omega}_1 \\ \boldsymbol{\omega}_2 \end{pmatrix} + \boldsymbol{b} \bigg) \right| \leq \varepsilon,$$

and

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

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 42, for every $\varepsilon_j > 0$, there exists $m_j \leq \mathcal{O}((1/\varepsilon_j \log(1/\varepsilon_j))^2)$, $\boldsymbol{a}_j \in \mathbb{R}^{m_j}$, $\boldsymbol{W}_j \in \mathbb{R}^{m_j \times 2d_e}$, and $\boldsymbol{b}_j \in \mathbb{R}^{m_j}$, such that

$$\sup_{|e_{1j}| \leq 1, |e_{2j}| \leq 1} \left| e_{1j} e_{2j} - \sum_{l=1}^m a_{jl} \sigma \left(\left\langle \boldsymbol{w}_{jl}, (\boldsymbol{\omega}_1^\top, \boldsymbol{\omega}_2^\top)^\top \right\rangle + b_{jl} \right) \right| \leq \varepsilon_j,$$

 $\|m{a}_j\|_2 \leq \mathcal{O}ig((\log(1/arepsilon_j)/arepsilon_j)^{3/2}/\sqrt{m_j}ig), \ \|m{b}_j\|_\infty \leq 1, \ ext{and} \ \|m{w}_{jl}\|_1 \leq 1.$ Specifically, the only nonzero coordinates of $m{w}_{jl}$ are the jth 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 $\boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^m$ and $\boldsymbol{W} \in \mathbb{R}^{m \times 2d_e}$ by concatenating $(\boldsymbol{a}_j), (\boldsymbol{b}_j)$, and (\boldsymbol{W}_j) respectively. The resulting network satisfies

$$\sup_{\boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \in \mathbb{S}^{d_e-1}} \left| \langle \boldsymbol{\omega}_1, \boldsymbol{\omega}_2 \rangle - \boldsymbol{a}^\top \sigma \bigg(\boldsymbol{W} \begin{pmatrix} \boldsymbol{\omega}_1 \\ \boldsymbol{\omega}_2 \end{pmatrix} + \boldsymbol{b} \bigg) \right| \leq \varepsilon,$$

while $\|\boldsymbol{a}\|_{2} \leq \mathcal{O}\left(d_{e}^{5/2}(\log(d_{e}/\varepsilon)/\varepsilon)^{3/2}/\sqrt{m}\right)$, $\|\boldsymbol{b}\|_{\infty} \leq 1$, and $\|\boldsymbol{W}^{\top}\|_{1,\infty} \leq 1$, completing the proof.

We can now specify A_1 , W_1 , and b_1 in our construction.

Lemma 23. 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 $A_1 \in \mathbb{R}^{(d+q)\times m_1}$, $W_1 \in \mathbb{R}^{m_1\times (d_h+d+(q+1)d_e)}$, and $b_1 \in \mathbb{R}^{m_1}$, given by Equations (B.5) to (B.9), such that

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

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

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

Proof. We define the decompositions

$$W_1 = \begin{pmatrix} W_{11} \\ W_{12} \end{pmatrix}, \quad b_1 = \begin{pmatrix} b_{11} \\ b_{12} \end{pmatrix}, \quad A_1 = \begin{pmatrix} A_{11} \\ A_{12} \end{pmatrix},$$
 (B.5)

where $\boldsymbol{W}_{11} \in \mathbb{R}^{2d \times (d_h + d + d_e)}$, $\boldsymbol{W}_{12} \in \mathbb{R}^{q\bar{m}_1 \times (d_h + d + d_e)}$, $\boldsymbol{b}_{11} \in \mathbb{R}^{2d}$, $\boldsymbol{b}_{12} \in \mathbb{R}^{q\bar{m}_1}$, $\boldsymbol{A}_{11} \in \mathbb{R}^{d \times m_1}$, and $\boldsymbol{A}_{12} \in \mathbb{R}^{q \times m_1}$. Let $\boldsymbol{v}_1, \dots, \boldsymbol{v}_d$ 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

$$\boldsymbol{W}_{11} = \begin{pmatrix} \mathbf{0}_{d_h} & \boldsymbol{v}_1^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \mathbf{0}_{d_h} & -\boldsymbol{v}_1^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \vdots & \vdots & & \\ \mathbf{0}_{d_h} & \boldsymbol{v}_d^\top & \mathbf{0}_{(q+1)d_e}^\top \\ \mathbf{0}_{d_h} & -\boldsymbol{v}_d^\top & \mathbf{0}_{(q+1)d_e}^\top \end{pmatrix}, \tag{B.6}$$

as well as

$$\boldsymbol{b}_{1} = \mathbf{0}_{2d}, \quad \text{and} \quad \boldsymbol{A}_{11} = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & \mathbf{0}_{q\bar{m}_{1}}^{\top} \\ 0 & 0 & 1 & -1 & \dots & 0 & \mathbf{0}_{q\bar{m}_{1}}^{\top} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & 0 & 1 & -1 & \mathbf{0}_{q\bar{m}_{1}}^{\top} \end{pmatrix}$$
(B.7)

Notice that $\|\boldsymbol{W}_{11}^{\top}\|_{1,\infty}=1$ and $\|\boldsymbol{A}_{11}^{\top}\|_{1,\infty}=2$. To implement the inner product part of the mapping, we take the construction of weights, biases, and second layer weights from Lemma 22, and rename them as $\tilde{\boldsymbol{W}}_1 \in \mathbb{R}^{\tilde{m}_1 \times 2d_e}$, $\tilde{\boldsymbol{b}}_1 \in \mathbb{R}^{\tilde{m}_1}$, and $\tilde{\boldsymbol{a}}_1 \in \mathbb{R}^{\tilde{m}_1}$. Let us introduce the decomposition $\tilde{\boldsymbol{W}}_1 = (\tilde{\boldsymbol{W}}_{11} \quad \tilde{\boldsymbol{W}}_{12})$, where $\tilde{\boldsymbol{W}}_{11}, \tilde{\boldsymbol{W}}_{12} \in \mathbb{R}^{\tilde{m}_1 \times d_e}$. With this decomposition, we can separate the projections applied to the first and second vectors in Lemma 22. We can then define

$$\boldsymbol{W}_{12} = \begin{pmatrix} \mathbf{0}_{\bar{m}_{1} \times (d_{h} + d)} & \tilde{W}_{11} & \tilde{W}_{12} & \mathbf{0}_{\bar{m}_{1} \times d_{e}} & \dots & \mathbf{0}_{\bar{m}_{1} \times d_{e}} \\ \mathbf{0}_{\bar{m}_{1} \times (d_{h} + d)} & \tilde{W}_{11} & \mathbf{0}_{\bar{m}_{1} \times d_{e}} & \tilde{W}_{12} & \dots & \mathbf{0}_{\bar{m}_{1} \times d_{e}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{\bar{m}_{1} \times (d_{h} + d)} & \tilde{W}_{11} & \mathbf{0}_{\bar{m}_{1} \times d_{e}} & \mathbf{0}_{\bar{m}_{1} \times d_{e}} & \dots & \tilde{W}_{12} \end{pmatrix},$$
(B.8)

as well as

$$\boldsymbol{b}_{12} = \begin{pmatrix} \tilde{\boldsymbol{b}}_1 \\ \vdots \\ \tilde{\boldsymbol{b}}_1 \end{pmatrix}, \quad \text{and} \quad \boldsymbol{A}_{12} = \begin{pmatrix} \mathbf{0}_{2d}^{\top} & \tilde{\boldsymbol{a}}_{1}^{\top} & \mathbf{0}_{\bar{m}_{1}}^{\top} & \dots & \mathbf{0}_{\bar{m}_{1}}^{\top} \\ \mathbf{0}_{2d}^{\top} & \mathbf{0}_{\bar{m}_{1}}^{\top} & \tilde{\boldsymbol{a}}_{1}^{\top} & \dots & \mathbf{0}_{\bar{m}_{1}}^{\top} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{2d}^{\top} & \mathbf{0}_{\bar{m}_{1}}^{\top} & \dots & \mathbf{0}_{\bar{m}_{1}}^{\top} & \tilde{\boldsymbol{a}}_{1}^{\top} \end{pmatrix}. \tag{B.9}$$

From Lemma 22, we have $\left\| \boldsymbol{W}_{12}^{\top} \right\|_{1,\infty} \leq 1, \left\| \boldsymbol{b}_{12} \right\|_{\infty} \leq 1$, and

$$\left\|\boldsymbol{A}_{12}^{\top}\right\|_{1,\infty} = \left\|\tilde{\boldsymbol{a}}_{1}\right\|_{1} \leq \mathcal{O}(d_{e}^{5/2}(\log(d_{e}/\varepsilon)/\varepsilon)^{3/2}),$$

which completes the proof.

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 22.

Lemma 24. Let σ be the ReLU activation. Suppose $|h| \leq r_{\infty}^h$, $|x| \leq r_{\infty}^x$ and $|z| \leq 1$. Let $R \coloneqq \sqrt{1 + r_{\infty}^x{}^2 + r_{\infty}^h{}^2}$. For any $\varepsilon > 0$, there exists $m = \mathcal{O}(R^6(\log(R/\varepsilon)/\varepsilon)^3)$, $\boldsymbol{a} \in \mathbb{R}^m$, $\boldsymbol{W} \in \mathbb{R}^{m \times 2}$, and $\boldsymbol{b} \in \mathbb{R}^m$, such that

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

and

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

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

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

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 42 with dimension 3 when $r_\infty^h\neq 0$ and dimension 2 otherwise.

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

Lemma 25. Suppose $\|\boldsymbol{\chi}_1^{\boldsymbol{x}}\|_{\infty} \leq r_x$ and $\max_l |\boldsymbol{\chi}^{\boldsymbol{\omega}}(l)| \leq 1$. Let $R \coloneqq \sqrt{1+r_x^2}$. Then, for every $\varepsilon > 0$ and absolute constant $\delta \in (0,1)$, there exists $\bar{m}_2 \leq \mathcal{O}(R^4(\log(R/\varepsilon)/\varepsilon)^{3/2})$, $m_2 \coloneqq qd\bar{m}_2$, and $\boldsymbol{A}_2 \in \mathbb{R}^{d_h \times m_2}$, $\boldsymbol{W}_2 \in \mathbb{R}^{m_2 \times (d+q)}$, and $\boldsymbol{b}_2 \in \mathbb{R}^{m_2}$ given by Equations (B.10) and (B.11) such that

$$\|\boldsymbol{\chi}_2(l) - 2\boldsymbol{\chi}_1^{\boldsymbol{x}}\sigma(\boldsymbol{\chi}_1^{\boldsymbol{\omega}}(l) - 1/2)\|_{\infty} \le \varepsilon,$$

for all such χ_1 and $l \in [q]$, where we recall $\chi_2 = A_2 \sigma(W_2 \chi_1 + b_2)$. Moreover, we have

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

Proof. Let $\tilde{W}=(\tilde{w}_{21} \quad \tilde{w}_{22})$, \tilde{b} , and \tilde{a} be the weights obtained from Lemma 24, where $\tilde{w}_{21}, \tilde{w}_{22}, \tilde{b}, \tilde{a} \in \mathbb{R}^{\bar{m}_2}$. To construct W_2 and b_2 , we let

$$\boldsymbol{W}_{2} = \begin{pmatrix} \boldsymbol{W}_{2}(1,1) \\ \vdots \\ \boldsymbol{W}_{2}(1,d) \\ \vdots \\ \boldsymbol{W}_{2}(q,1) \\ \vdots \\ \boldsymbol{W}_{2}(q,d) \end{pmatrix}, \quad \boldsymbol{b}_{22} = \begin{pmatrix} \boldsymbol{b}_{2}(1,1) \\ \vdots \\ \boldsymbol{b}_{2}(1,d) \\ \vdots \\ \boldsymbol{b}_{2}(q,1) \\ \vdots \\ \boldsymbol{b}_{2}(q,d) \end{pmatrix}. \tag{B.10}$$

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

$$W_2(l,j) = (\mathbf{0}_{\bar{m}_2 \times (j-1)} \quad \tilde{w}_{21} \quad \mathbf{0}_{\bar{m}_2 \times (d-j)} \quad \mathbf{0}_{\bar{m}_2 \times (l-1)} \quad \tilde{w}_{22} \quad \mathbf{0}_{\bar{m}_2 \times (q-l)}),$$

and $b_2(l,j) = \tilde{b}_2$. Consequently, $\|\boldsymbol{W}_2^\top\|_{1,\infty} \leq 1$ and $\|\boldsymbol{b}_2\|_{\infty} \leq 1$. Finally, we have

$$\mathbf{A}_{2} = \begin{pmatrix} \tilde{\mathbf{a}}_{2}^{\top} & \mathbf{0}_{\bar{m}_{2}}^{\top} & \dots & \mathbf{0}_{\bar{m}_{2}}^{\top} \\ \mathbf{0}_{\bar{m}_{2}}^{\top} & \tilde{\mathbf{a}}_{2}^{\top} & \dots & \mathbf{0}_{\bar{m}_{2}}^{\top} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}_{\bar{m}_{2}}^{\top} & \dots & \mathbf{0}_{\bar{m}_{2}}^{\top} & \tilde{\mathbf{a}}_{2}^{\top} \end{pmatrix}.$$
(B.11)

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

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

Proposition 26. Let $z=(x,\omega_i,\omega_{t_1},\ldots,\omega_{t_q})$. Then, for every $\varepsilon>0$, there exists a feedforward network with $L_h=4$ layers given by

$$f^{
ightarrow}(oldsymbol{h},oldsymbol{z};oldsymbol{\Theta}_h^{
ightarrow}) = oldsymbol{W}_{L_h}\sigma\Bigl(\dots\sigma\bigl(oldsymbol{W}_2\sigmaig(oldsymbol{W}_1(oldsymbol{h}^{ op},oldsymbol{z}^{ op})^{ op}+oldsymbol{b}_1ig)+oldsymbol{b}_2ig)\dots\Bigr)$$

where $W_i \in \mathbb{R}^{m_i \times m_{i-1}}$, $b_i \in \mathbb{R}^m_i$ for $i \in \{2, \dots, L_h - 1\}$, $W_1 \in \mathbb{R}^{m_1 \times d_h + d + (q+1)d_e}$, $b_1 \in \mathbb{R}^{m_1}$, and $W_{L_h} \in \mathbb{R}^{d_h \times m_{L_h} - 1}$ that satisfies the following:

1. If
$$t_l = i$$
, then

$$\left\|f^{\rightarrow}(\boldsymbol{h}, \boldsymbol{z}; \hat{\boldsymbol{\Theta}}_{h}^{\rightarrow})_{l} - \boldsymbol{x}\right\|_{2} \leq \varepsilon$$

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

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

Proof. Let $\tilde{A}_1 \in \mathbb{R}^{(d+q) \times m_1}$, $\tilde{W}_1 \in \mathbb{R}^{m_1 \times (d_h + d + (q+1)d_e)}$, $\tilde{b}_1 \in \mathbb{R}^{m_1}$ be given by Lemma 23 with error parameter ε_1 and $\tilde{A}_2 \in \mathbb{R}^{d_h \times m_2}$, $\tilde{W}_2 \in \mathbb{R}^{m_2 \times (d+q)}$, $\tilde{b}_2 \in \mathbb{R}^{m_2}$ be given by Lemma 25 with error parameter ε_2 . Recall that

$$oldsymbol{\chi}_1 = ilde{oldsymbol{A}}_1 \sigma ig(ilde{oldsymbol{W}}_1 oldsymbol{\chi}_0 + ilde{oldsymbol{b}}_1ig), \quad oldsymbol{\chi}_2 = ilde{oldsymbol{A}}_2 \sigma ig(ilde{oldsymbol{W}}_2 oldsymbol{\chi}_1 + ilde{oldsymbol{b}}_2ig).$$

By the triangle inequality.

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

where $\bar{\chi}_1 = (x^\top, \langle \omega_i, \omega_{t_1} \rangle, \dots, \langle \omega_i, \omega_{t_n} \rangle)^\top$. By letting $\varepsilon_2 = \varepsilon/4$, we obtain

$$m_2, \|\tilde{\boldsymbol{A}}_2\|_{\mathrm{F}}, \|\tilde{\boldsymbol{W}}_2\|_{\mathrm{F}}, \|\tilde{\boldsymbol{b}}_2\|_2 \leq \operatorname{poly}(r_x, D_e, \varepsilon^{-1}).$$

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

$$m_1, \|\tilde{\boldsymbol{A}}_2\|_{F}, \|\tilde{\boldsymbol{W}}_2\|_{F}, \|\tilde{\boldsymbol{b}}_2\|_{2} \leq \operatorname{poly}(r_x, D_e, \varepsilon^{-1}).$$

Let

$$oldsymbol{W}_2 = ilde{oldsymbol{W}}_2 ilde{oldsymbol{A}}_1, \hspace{1cm} oldsymbol{W}_1 = ilde{oldsymbol{W}}_1, \hspace{1cm} oldsymbol{b}_1 = ilde{oldsymbol{b}}_1, \hspace{1cm} oldsymbol{b}_2 = ilde{oldsymbol{b}}_2.$$

Then,

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

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

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 $\Psi(\boldsymbol{x},\boldsymbol{t},i)_l = \mathbf{0}_d$. 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 $\tilde{\boldsymbol{W}}_3 \in \mathbb{R}^{2d_h \times d_h}$, $\boldsymbol{b}_3 \in \mathbb{R}^{2d_h}$, and $\boldsymbol{W}_4 \in \mathbb{R}^{d_h \times 2d_h}$ via

$$\tilde{\boldsymbol{W}}_{3} = \begin{pmatrix} \boldsymbol{v}_{1}^{\top} \\ -\boldsymbol{v}_{1}^{\top} \\ \vdots \\ \boldsymbol{v}_{d}^{\top} \\ -\boldsymbol{v}_{d}^{\top} \end{pmatrix}, \quad \boldsymbol{b}_{3} = -\frac{\varepsilon}{2} \mathbf{1}_{2d_{h}}, \quad \boldsymbol{W}_{4} = \begin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & -1 & \dots & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 1 & -1 \end{pmatrix}.$$

As a result, $\chi_3 = W_4 \sigma(\tilde{W}_3 \chi_2 + b_3)$ satisfies

$$|(\chi_3)_j - (\chi_2)_j| \le \begin{cases} 0 & |(\chi_2)_j| \le \varepsilon/2\\ \varepsilon/2 & |(\chi_2)_j| > \varepsilon/2 \end{cases}, \quad \forall j \in [d_h]. \tag{B.12}$$

We thus make two observations. First, $\|\boldsymbol{\chi}_3 - \boldsymbol{\chi}_2\|_{\infty} \leq \varepsilon/2$, and consequently $\|\boldsymbol{\chi}_3(l) - \Psi(\boldsymbol{x}, \boldsymbol{t}, i)_l\|_{\infty} \leq \varepsilon$ for all $l \in [q]$. Second, when $t_l \neq i$, we have $\Psi(\boldsymbol{x}, \boldsymbol{t}, i)_l = \mathbf{0}_d$ and $|\chi_2(l)_j| \leq \varepsilon/2$ for all $j \in [d]$ since $\|\boldsymbol{\chi}_2(l) - \Psi(\boldsymbol{x}, \boldsymbol{t}, i)_l\|_{\infty} \leq \varepsilon/2$. Consequently, by the first case in (B.12), we have $\chi_3(l)_j = 0$ for all $j \in [d]$. We can summarize these two observations as follows

$$\|\boldsymbol{\chi}_3(l) - \Psi(\boldsymbol{x}, \boldsymbol{t}, i)_l\|_{\infty} \le \begin{cases} 0 & t_l \neq i \\ \varepsilon & t_l = i \end{cases}$$

which completes the proof.

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

Corollary 27. Let f_h^{\rightarrow} be given by the construction in Proposition 26, and suppose $r_h \geq \sqrt{q}(r_x + \sqrt{d\varepsilon})$. Then, h_i^{\rightarrow} satisfies the following guarantees for all $i \in [N]$ and $l \in [q]$:

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

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

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

- If $t_l \geq i+1$, then $\boldsymbol{h}_i^{\rightarrow}(l) = \boldsymbol{0}_d$ and $f_h^{\rightarrow}(\boldsymbol{h}_i^{\rightarrow}, \boldsymbol{z}_i; \boldsymbol{\Theta}_h^{\rightarrow}) = \boldsymbol{0}_d$ by Proposition 26.
- If $t_l < i < i+1$, then $\|\boldsymbol{h}_i^{\rightarrow}(l) \boldsymbol{x}_{t_l}\|_{\infty} \leq \varepsilon$ by induction hypothesis, and $f_h^{\rightarrow}(\boldsymbol{h}_j^{\rightarrow}, \boldsymbol{z}_j; \boldsymbol{\Theta}_h^{\rightarrow}) = \boldsymbol{0}_d$.
- Finally, if $t_l = i < i+1$, then $\boldsymbol{h}_i^{\rightarrow}(l) = 0$ and $\|f_h^{\rightarrow}(\boldsymbol{h}_i^{\rightarrow}, \boldsymbol{z}_i; \boldsymbol{\Theta}_h^{\rightarrow}) \boldsymbol{x}_{t_l}\|_{\infty} \leq \varepsilon$.

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

By symmetry, the same construction for f_h^{\leftarrow} would yield a similar guarantee on h_i^{\leftarrow} .

The last step is to design $f_u(\mathbf{h}^{\rightarrow}, \mathbf{h}^{\leftarrow}, \mathbf{z}; \mathbf{\Theta}_u)$ such that

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

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

Proposition 28. Suppose g satisfies Assumption 2. Then there exist RNN weights Θ_{RNN} with $vec(\Theta_{RNN}) \in \mathbb{R}^p$ (i.e. with p parameters) and $r_h \geq \sqrt{q}r_x + \sqrt{\varepsilon_{2NN}}/(r_ar_w)$, such that

$$\sup_{i \in [N]} \left| g(\boldsymbol{x}_{t_1}, \dots, \boldsymbol{x}_{t_q}) - \hat{y}(\boldsymbol{p}; \boldsymbol{\Theta}_{\text{RNN}})_i \right|^2 \le 4\varepsilon_{\text{2NN}} \tag{B.13}$$

for all $t \in [N]^q$ and $||x_j||_2 \le r_x$ for all $j \in [N]$. Additionally, we have

 $\|\operatorname{vec}(\boldsymbol{\Theta}_{\mathtt{RNN}})\|_{2} \leq \operatorname{poly}(r_{x}, D_{e}, r_{w}, r_{a}, \varepsilon_{\mathtt{2NN}}^{-1}), \quad p \leq \operatorname{poly}(r_{x}, D_{e}, m_{g}, r_{w}, r_{a}, \varepsilon_{\mathtt{2NN}}^{-1}),$ (B.14) and f_{h}^{\rightarrow} , f_{h}^{\leftarrow} do not depend on h^{\rightarrow} and h^{\leftarrow} , namely the first d_{h} columns of $\boldsymbol{W}_{1}^{\rightarrow}$ and $\boldsymbol{W}_{1}^{\leftarrow}$ that are multiplied by h^{\rightarrow} and 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 $(v_j)_{j=1}^{d_h}$ 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

$$m{W}_{ ext{id}} = egin{pmatrix} m{v}_1^{ op} \ -m{v}_1^{ op} \ dots \ m{v}_{ ext{id}}^{ op} \ -m{v}_{ ext{id}}^{ op} \end{pmatrix}, \quad m{b}_{ ext{id}} = m{0}_{2d_h}, m{A}_{ ext{id}} = egin{pmatrix} 1 & -1 & 0 & 0 & \dots & 0 & 0 \ 0 & 0 & 1 & -1 & \dots & 0 & 0 \ 0 & 0 & 0 & 0 & \dots & 1 & -1 \end{pmatrix},$$

where $\boldsymbol{W}_{\mathrm{id}} \in \mathbb{R}^{2d_h \times d_h}$, $\boldsymbol{b}_{\mathrm{id}} \in \mathbb{R}^{2d_h}$, and $\boldsymbol{A}_{\mathrm{id}} \in \mathbb{R}^{d_h \times 2d_h}$. Let $\boldsymbol{W}_1, \boldsymbol{b}_1, \tilde{\boldsymbol{A}}_1, \tilde{\boldsymbol{W}}_2, \boldsymbol{b}_2, \tilde{\boldsymbol{A}}_2$ be given as in the proof of Proposition 26, for achieving an L_{∞} error of $\tilde{\varepsilon}$, to be fixed later. Recall $\boldsymbol{z}_i = (\boldsymbol{x}_i^{\top}, \boldsymbol{\omega}_i^{\top}, \boldsymbol{\omega}_{t_1}^{\top}, \dots, \boldsymbol{\omega}_{t_q}^{\top})^{\top}$. In the following, we remove the zero columns of \boldsymbol{W}_1 corresponding to the \boldsymbol{h} part of the input (see Lemma 23), which does not change the resulting function. Our construction can then be denoted by

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

$$m{W}_{ ext{add}} = egin{pmatrix} m{v}_1^ op & m{v}_1^ op & m{v}_1^ op & m{v}_1^ op \ -m{v}_1^ op & -m{v}_1^ op & -m{v}_1^ op \ dots & dots & dots \ m{v}_{d_h}^ op & m{v}_{d_h}^ op & m{v}_{d_h}^ op \ -m{v}_{d_h}^ op & -m{v}_{d_h}^ op & -m{v}_{d_h}^ op \end{pmatrix}, \quad m{b}_{ ext{add}} = m{0}_{2d_h}, \quad m{A}_{ ext{add}} = m{A}_{ ext{id}},$$

where $W_{\text{add}} \in \mathbb{R}^{2d_h \times 3d_h}$, $b_{\text{add}} \in \mathbb{R}^{2d_h}$, $A_{\text{add}} \in \mathbb{R}^{d_h \times 2d_h}$.

Let Θ_h^{\rightarrow} (and similarly Θ_h^{\leftarrow}) be given by Proposition 26 with corresponding error ε_h . Using the shorthand notation $\boldsymbol{x_t} = (\boldsymbol{x_{t_1}}, \dots, \boldsymbol{x_{t_q}}) \in \mathbb{R}^{dq}$ and $\hat{\boldsymbol{x_t}} = \boldsymbol{h_i^{\rightarrow}} + \boldsymbol{h_i^{\leftarrow}} + \boldsymbol{\chi_2}$, we have

$$egin{aligned} \|oldsymbol{h}_i^{
ightarrow} + oldsymbol{h}_i^{
ightharpoonup} - \hat{oldsymbol{x}}_t\|_2 &\leq \left\|oldsymbol{h}_i^{
ightarrow} - \sum_{j=1}^{i-1} \Psi(oldsymbol{x}_j, oldsymbol{t}, j)
ight\|_2 + \left\|oldsymbol{h}_i^{\leftarrow} - \sum_{j=N}^{i+1} \Psi(oldsymbol{x}_j, oldsymbol{t}, j)
ight\|_2 + \|oldsymbol{\chi}_2 - \Psi(oldsymbol{x}_i, oldsymbol{t}, i)\|_2 \\ &\leq \sqrt{qd} (2arepsilon_h + ilde{arepsilon}), \end{aligned}$$

which holds for all input prompts p with $||x_j||_2 \le r_x$ for all $j \in [N]$. Finally, we have

$$\begin{split} \sup_{\|\boldsymbol{x}_{j}\|_{2} \leq r_{x}, \, \forall j \in [N]} & \left| g(\boldsymbol{x}_{t}) - \boldsymbol{a}_{g}^{\top} \sigma(\boldsymbol{W}_{g} \hat{\boldsymbol{x}}_{t} + \boldsymbol{b}_{g}) \right| \leq \sup_{\|\boldsymbol{x}_{j}\|_{2} \leq r_{x}, \, \forall j \in [N]} & \left| g(\boldsymbol{x}_{t}) - \boldsymbol{a}_{g}^{\top} \sigma(\boldsymbol{W}_{g} \boldsymbol{x}_{t} + \boldsymbol{b}_{g}) \right| \\ & + \sup_{\|\boldsymbol{x}_{j}\|_{2} \leq r_{x}, \, \forall j \in [N]} & \left| \boldsymbol{a}_{g}^{\top} \sigma(\boldsymbol{W}_{g} \boldsymbol{x}_{t} + \boldsymbol{b}_{g}) - \boldsymbol{a}_{g}^{\top} \sigma(\boldsymbol{W}_{g} \hat{\boldsymbol{x}}_{t} + \boldsymbol{b}_{g}) \right| \\ & \leq \sqrt{\varepsilon_{2\mathrm{NN}}} + r_{a} r_{w} \sqrt{q d} (2\varepsilon_{h} + \tilde{\varepsilon}). \end{split}$$

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

B.4 Generalization Upper Bounds for RNNs

Recall the state transitions

$$\begin{split} & \boldsymbol{h}_{j+1}^{\rightarrow} = \Pi_{r_h} \big(\boldsymbol{h}_j^{\rightarrow} + f_h^{\rightarrow} (\boldsymbol{h}_j^{\rightarrow}, \boldsymbol{z}_j; \boldsymbol{\Theta}_h^{\rightarrow}) \big) \\ & \boldsymbol{h}_{j-1}^{\leftarrow} = \Pi_{r_h} \big(\boldsymbol{h}_j^{\leftarrow} + f^{\leftarrow} (\boldsymbol{h}^{\leftarrow}, \boldsymbol{z}_j; \boldsymbol{\Theta}^{\leftarrow}) \big). \end{split}$$

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

$$F(\boldsymbol{p};\boldsymbol{\Theta}_{h}^{\rightarrow},\boldsymbol{\Theta}_{h}^{\leftarrow},\boldsymbol{\Theta}_{y})_{j}=f_{y}(\boldsymbol{h}_{j}^{\rightarrow}(\boldsymbol{p};\boldsymbol{\Theta}_{h}^{\rightarrow}),\boldsymbol{h}_{j}^{\leftarrow}(\boldsymbol{p};\boldsymbol{\Theta}_{h}^{\leftarrow}),\boldsymbol{z}_{j};\boldsymbol{\Theta}_{y}).$$

We can now define the function class

$$\mathcal{F}_{\mathtt{RNN}} = \{ \boldsymbol{p}, j \mapsto F(\boldsymbol{p}; \boldsymbol{\Theta}_h^{\rightarrow}, \boldsymbol{\Theta}_h^{\leftarrow}, \boldsymbol{\Theta}_y)_j \, : \, \boldsymbol{\Theta}_h^{\rightarrow}, \boldsymbol{\Theta}_h^{\leftarrow}, \boldsymbol{\Theta}_y \in \boldsymbol{\varTheta}_{\mathtt{RNN}} \}.$$

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

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

We will further use the notation

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

and

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

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

Lemma 29. Suppose for every Θ_h^{\rightarrow} , Θ_h^{\leftarrow} , $\Theta_y \in \Theta_{RNN}$ we have

$$\left\| \boldsymbol{W}_{L_{y}}^{y} \dots \boldsymbol{W}_{1}^{y} \right\|_{op} \leq C_{W}^{y}, \quad \left\| \boldsymbol{W}_{L_{h}}^{\rightarrow} \right\|_{op} \dots \left\| \boldsymbol{W}_{1,h}^{\rightarrow} \right\|_{op} \leq \alpha_{N}, \quad \left\| \boldsymbol{W}_{L_{h}}^{\leftarrow} \right\|_{op} \dots \left\| \boldsymbol{W}_{1,h}^{\leftarrow} \right\|_{op} \leq \alpha_{N},$$
where $\alpha_{N} \leq N^{-1}$. Then,

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

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

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

where

$$\mathcal{E}_1 \coloneqq \sup_{oldsymbol{p},j \in S_n} \left| f_y(oldsymbol{h}_j^{
ightarrow}, oldsymbol{h}_j^{\leftarrow}, oldsymbol{z}_j; oldsymbol{\Theta}_y) - f_y(oldsymbol{h}_j^{
ightarrow}, oldsymbol{h}_j^{\leftarrow}, oldsymbol{z}_j; oldsymbol{\hat{\Theta}}_y)
ight|$$
 $\mathcal{E}_2 \coloneqq \sup_{oldsymbol{p},j \in S_n} \left| f_y(oldsymbol{h}_j^{
ightarrow}, oldsymbol{h}_j^{\leftarrow}, oldsymbol{h}_j^{\leftarrow}, oldsymbol{z}_j; oldsymbol{\hat{\Theta}}_y) - f_y(\hat{oldsymbol{h}}_j^{
ightarrow}, \hat{oldsymbol{h}}_j^{\leftarrow}, oldsymbol{z}_j; oldsymbol{\hat{\Theta}}_y)
ight|.$

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

Using the Lipschitzness of f_y , we obtain

$$\mathcal{E}_{2} \leq \left\| \boldsymbol{W}_{L_{y}}^{y} \dots \boldsymbol{W}_{1}^{y} \right\|_{\text{op}} \left(\sup_{\boldsymbol{p}, j} \left\| \boldsymbol{h}_{j}^{\rightarrow} - \hat{\boldsymbol{h}}_{j}^{\rightarrow} \right\|_{2} + \sup_{\boldsymbol{p}, j} \left\| \boldsymbol{h}_{j}^{\leftarrow} - \hat{\boldsymbol{h}}_{j}^{\leftarrow} \right\|_{2} \right)$$

$$\leq C_{W}^{y} \left(\sup_{\boldsymbol{p}, j} \left\| \boldsymbol{h}_{j}^{\rightarrow} - \hat{\boldsymbol{h}}_{j}^{\rightarrow} \right\|_{2} + \sup_{\boldsymbol{p}, j} \left\| \boldsymbol{h}_{j}^{\leftarrow} - \hat{\boldsymbol{h}}_{j}^{\leftarrow} \right\|_{2} \right).$$

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

$$\begin{split} \sup_{\boldsymbol{p},j} \left\| \boldsymbol{h}_{j}^{\rightarrow} - \hat{\boldsymbol{h}}_{j}^{\rightarrow} \right\|_{2} \leq \sup_{\boldsymbol{p},j} \left\| \boldsymbol{h}_{j-1}^{\rightarrow} - \hat{\boldsymbol{h}}_{j-1}^{\rightarrow} \right\|_{2} + \underbrace{\sup_{\boldsymbol{p},j} \left\| f_{h}^{\rightarrow}(\boldsymbol{h}_{j-1}^{\rightarrow}, \boldsymbol{z}_{j-1}; \hat{\boldsymbol{\Theta}}_{h}^{\rightarrow}) - f_{h}^{\rightarrow}(\hat{\boldsymbol{h}}_{j-1}^{\rightarrow}, \boldsymbol{z}_{j-1}; \hat{\boldsymbol{\Theta}}_{h}^{\rightarrow}) \right\|_{2}}_{=:\mathcal{E}_{h}^{h}} \\ + \sup_{\boldsymbol{p},j} \left\| f_{h}^{\rightarrow}(\boldsymbol{h}_{j-1}^{\rightarrow}, \boldsymbol{z}_{j-1}; \hat{\boldsymbol{\Theta}}_{h}^{\rightarrow}) - f_{h}^{\rightarrow}(\boldsymbol{h}_{j-1}^{\rightarrow}, \boldsymbol{z}_{j-1}; \hat{\boldsymbol{\Theta}}_{h}^{\rightarrow}) \right\|_{2}}_{=:\mathcal{E}_{h}^{h}}. \end{split}$$

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

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

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

$$\begin{split} \sup_{\boldsymbol{p},j} \left\| \boldsymbol{h}_{j}^{\rightarrow} - \hat{\boldsymbol{h}}_{j}^{\rightarrow} \right\|_{2} &\leq (1 + \alpha_{N}) \sup_{\boldsymbol{p},j} \left\| \boldsymbol{h}_{j-1}^{\rightarrow} - \hat{\boldsymbol{h}}_{j-1}^{\rightarrow} \right\|_{2} + d_{\infty}(f_{h}^{\rightarrow}(\cdot;\boldsymbol{\Theta}_{h}^{\rightarrow}), f_{h}^{\rightarrow}(\cdot;\hat{\boldsymbol{\Theta}}_{h}^{\rightarrow})) \\ &\leq \sum_{l=0}^{j-2} (1 + \alpha_{N})^{l} d_{\infty}(f_{h}^{\rightarrow}(\cdot;\boldsymbol{\Theta}_{h}^{\rightarrow}), f^{\rightarrow}(\cdot;\hat{\boldsymbol{\Theta}}_{h}^{\rightarrow})) \\ &\leq \frac{(1 + \alpha_{N})^{j-1} - 1}{\alpha_{N}} d_{\infty}(f_{h}^{\rightarrow}(\cdot;\boldsymbol{\Theta}_{h}^{\rightarrow}), f_{h}^{\rightarrow}(\cdot;\hat{\boldsymbol{\Theta}}_{h}^{\rightarrow})) \\ &\leq eNd_{\infty}(f_{h}^{\rightarrow}(\cdot;\boldsymbol{\Theta}_{h}^{\rightarrow}), f_{h}^{\rightarrow}(\cdot;\hat{\boldsymbol{\Theta}}_{h}^{\rightarrow})). \end{split}$$

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

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

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

$$\mathcal{C}(\mathcal{F}_{\mathrm{NN},L_{h}}^{\rightarrow},\epsilon/(4eC_{w}^{y}N)), \quad \text{and}, \quad \mathcal{C}(\mathcal{F}_{\mathrm{NN},L_{h}}^{\leftarrow},\epsilon/(4eC_{w}^{y}N))$$

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 30. Let

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

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

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

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

Proof. Let $x_0 = x$, $x_l = \sigma(W_l x_{l-1} + b_l)$ for $l \in [L-1]$, and $x_L = W_L x_{L-1}$. Also let (\hat{x}_l) be the corresponding definitions under weights and biases (\hat{W}_l) and (\hat{b}_l) . First, we remark that for $l \in [L-1]$,

$$\|\boldsymbol{x}_{l}\|_{2} \leq \|\boldsymbol{W}_{l}\|_{\text{op}} \|\boldsymbol{x}_{l-1}\|_{2} + \|\boldsymbol{b}_{l}\|_{2}$$

$$\leq \prod_{i=1}^{l} \|\boldsymbol{W}_{i}\|_{\text{op}} \|\boldsymbol{x}_{0}\|_{2} + \sum_{i=0}^{l-1} \|\boldsymbol{b}_{l-i-1}\|_{2} \prod_{j=0}^{i} \|\boldsymbol{W}_{l-j}\|_{\text{op}} + \|\boldsymbol{b}_{l}\|_{2}$$

$$\leq \text{poly}(R),$$
(B.15)

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

$$\begin{aligned} \|\boldsymbol{x}_{l} - \hat{\boldsymbol{x}}_{l}\|_{2} &\leq \|\boldsymbol{W}_{l}\boldsymbol{x}_{l-1} - \hat{\boldsymbol{W}}_{l}\hat{\boldsymbol{x}}_{l-1}\|_{2} + \|\boldsymbol{b}_{l} - \hat{\boldsymbol{b}}_{l}\|_{2} \\ &\leq \|\boldsymbol{W}_{l}\|_{\text{op}} \|\boldsymbol{x}_{l-1} - \hat{\boldsymbol{x}}_{l-1}\|_{2} + \|\hat{\boldsymbol{x}}_{l-1}\|_{2} \|\boldsymbol{W}_{l} - \hat{\boldsymbol{W}}_{l}\|_{\text{op}} + \|\boldsymbol{b}_{l} - \hat{\boldsymbol{b}}_{l}\|_{2} \\ &\leq \text{poly}(R) \Big\{ \|\boldsymbol{x}_{l-1} - \hat{\boldsymbol{x}}_{l-1}\|_{2} + \|\boldsymbol{W}_{l} - \hat{\boldsymbol{W}}_{l}\|_{F} + \|\boldsymbol{b}_{l} - \hat{\boldsymbol{b}}_{l}\|_{2} \Big\}. \end{aligned}$$

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

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

Finally, we have the bound

$$\begin{aligned} \|\boldsymbol{x}_{L} - \hat{\boldsymbol{x}}_{L}\|_{2} &\leq \|\boldsymbol{W}_{L}\|_{\text{op}} \|\boldsymbol{x}_{L-1} - \hat{\boldsymbol{x}}_{L-1}\|_{2} + \|\hat{\boldsymbol{x}}_{L-1}\|_{2} \|\boldsymbol{W}_{L} - \hat{\boldsymbol{W}}_{L}\|_{\text{op}} \\ &\leq \text{poly}(R) \|\text{vec}(\boldsymbol{\Theta}_{\text{NN}}) - \text{vec}(\hat{\boldsymbol{\Theta}}_{\text{NN}})\|_{2}. \end{aligned}$$

Consequently, we have

$$\log \mathcal{C}(\mathcal{F}_{\mathrm{NN},L}, d_{\infty}, \epsilon) \leq \log \mathcal{C}(\{\Theta \in \mathbb{R}^p : \|\Theta\|_2 \leq \mathrm{poly}(d, q)\}, \|\cdot\|_2, \epsilon/\operatorname{poly}(R))$$

$$\leq p \log(1 + \operatorname{poly}(R)/\epsilon),$$

where the last inequality follows from Lemma 41.

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

Corollary 31. Suppose $\Theta_{RNN} \subseteq \{ \Theta \in \mathbb{R}^p : \| \operatorname{vec}(\Theta) \|_2 \leq R \}$ and $\| \boldsymbol{z}_j^{(i)} \|_2 \leq R$ for all $i \in [n]$ and $j \in [N]$. Then,

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

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

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

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

Lemma 32. *In the same setting as Corollary 31 and with* $\tau \geq 1$ *, we have*

$$\hat{\mathfrak{R}}_n(\mathcal{L}_{ au}^{\mathtt{RNN}}) \leq \mathcal{O}\!\left(au\sqrt{rac{p\log(RNn au)}{n}}
ight).$$

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

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

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

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

Corollary 33. 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 $\sqrt{3C_x ed \log(nN)} + q + 1 \leq R$. Then, for every $\delta > 0$, with probability at least $1 - \delta - (nN)^{-1/2}$ over the training set, we have

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

Proof. We highlight that for the specified R, Lemma 12 guarantees $\|z_j^{(i)}\|_2 \le R$ 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 32 complete the proof.

Note that $\hat{R}^{\text{RNN}}_{\tau}(\hat{\mathbf{\Theta}}) \leq \hat{R}^{\text{RNN}}_{n}(\hat{\mathbf{\Theta}})$ which is further controlled in the approximation section by Proposition 28. Therefore, the last step is to demonstrate that choosing $\tau = \text{poly}(d,q,\log n)$ suffices to achieve a desirable bound on $R^{\text{RNN}}(\hat{\mathbf{\Theta}})$ through $R^{\text{RNN}}_{\tau}(\hat{\mathbf{\Theta}})$.

Lemma 34. Consider the setting of Corollary 33, and additionally assume $R \ge r_h$. Then, for some $\tau = \text{poly}(R, \log n)$, we have

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

Proof. The proof of this lemma proceeds similarly to the proof of Lemma 20. By defining

$$\Delta_y \coloneqq \left| \hat{y}_{\mathtt{RNN}}(oldsymbol{p}; \hat{oldsymbol{\Theta}})_j - y_j
ight|$$

33

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

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

where

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

and

$$\mathbb{P}\big(\Delta_y > \sqrt{\tau}\big) \leq \mathbb{P}\bigg(|y_j| \geq \frac{\sqrt{\tau}}{2}\bigg) + \mathbb{P}\bigg(\Big|\hat{y}_{\mathtt{RNN}}(\boldsymbol{p}; \hat{\boldsymbol{\Theta}})_j\Big| \geq \frac{\sqrt{\tau}}{2}\bigg)$$

From Assumption 1, we have $\mathbb{E}\big[y_j^4\big]^{1/2}\lesssim 1$ 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 (B.16) for the derivation)

$$\left|\hat{y}_{\mathtt{RNN}}(\boldsymbol{p};\hat{\boldsymbol{\Theta}})_{j}\right| \leq \prod_{l=1}^{L_{y}} \lVert \boldsymbol{W}_{l}^{y} \rVert_{\mathrm{op}} \lVert (\boldsymbol{h}_{j}^{\rightarrow},\boldsymbol{h}_{j}^{\leftarrow},\boldsymbol{z}_{j}) \rVert_{2} + \sum_{i=0}^{L_{y}-1} \left\lVert \boldsymbol{b}_{L_{y}-i-1}^{y} \right\rVert_{2} \prod_{l=0}^{i} \left\lVert \boldsymbol{W}_{L_{y}-l}^{y} \right\rVert_{\mathrm{op}}.$$

As a result,

$$\left|\hat{y}_{\text{RNN}}(\boldsymbol{p}; \hat{\boldsymbol{\Theta}})_{j}\right| \leq \text{poly}(R)(1 + r_{h} + \|\boldsymbol{z}_{j}\|).$$

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

$$\mathbb{E}\Big[\hat{y}_{\mathtt{RNN}}(\boldsymbol{p}; \hat{\boldsymbol{\Theta}})_{j}^{4}\Big]^{1/2} \leq \mathrm{poly}(R).$$

On the other hand, from Lemma 12 (with n = N = 1), we obtain

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

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

We can summarize the above facts into the proof of Theorem 5.

Proof of Theorem 5. From the approximation bound of Proposition 28, we know that for some $R = \text{poly}(d, q, r_a, r_w, \varepsilon_{2NN}^{-1}, \log(nN))$ and the constraint set

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

with any $\alpha_N \leq N^{-1}$, we have $\hat{R}^{\text{RNN}}(\hat{\mathbf{\Theta}}) \lesssim \varepsilon_{\text{2NN}}$. The proof is then completed by letting $r_h = \sqrt{q}r_x + \sqrt{\varepsilon_{\text{2NN}}}/(r_a r_w)$, invoking the generalization bound of Corollary 33, and the bound on truncation error given in Lemma 34, with $R = \text{poly}(d,q,r_a,r_w,\varepsilon_{\text{2NN}}^{-1},\log(nN))$.

B.5 Proof of Proposition 6

The crux of the proof of Proposition 6 is to show the following position, which provides a lower bound on the prediction error at any fixed position in the prompt.

Proposition 35. Consider the same setting as in Proposition 6. There exists an absolute constant c > 0, such that for any fixed $j \in [N]$, if

$$\mathbb{E}\big[(\hat{y}_{RNN}(\boldsymbol{p})_j - y_j)^2\big] \le c,$$

then

$$d_h \geq \Omega\Big(rac{N}{\log(1+\mathfrak{L}^2\|oldsymbol{U}\|_{op}^2)}\Big), \quad ext{and} \quad \|oldsymbol{U}\|_{op}^2 \geq \Omega\Big(rac{N}{\mathfrak{L}^2\log(1+d_h)}\Big).$$

We shortly remark that the statement of Proposition 6 directly follows from that of Proposition 35.

Proof of Proposition 6. Let c be the constant given by Proposition 35. Suppose that

$$rac{1}{N} \mathbb{E} \Big[\| \hat{oldsymbol{y}}_{ exttt{RNN}}(oldsymbol{p}) - oldsymbol{y} \|_2^2 \Big] \leq c.$$

Then,

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

As a result, there exists some $j \in [N]$ such that $\mathbb{E}[(\hat{y}_{RNN}(\boldsymbol{p})_j - y_j)^2] \leq c$. We can then invoke Proposition 35 to obtain lower bounds on d_h and $\|\boldsymbol{U}\|_{op}$, completing the proof of Proposition 6. \square

We now present the proof of Proposition 35.

Proof of Proposition 35. Let $h_j = (U^{\rightarrow} h_i^{\rightarrow}, U^{\leftarrow} h_i^{\leftarrow}) \in \mathbb{R}^{2d_h}$, and define

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

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

Let $p^{(1)}, \ldots, p^{(P)}$ be an i.i.d. sequence of prompts, then modify them to share the jth input token, i.e. $x_j^{(i)} = x_j^{(1)}$ 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 $p^{(i)}$. Similarly, define

$$m{g}^{(i)} \coloneqq (m{g}(m{x}_1^{(i)}), \dots, m{g}(m{x}_{i-1}^{(i)}), m{g}(m{x}_{i+1}^{(i)}), \dots, m{g}(m{x}_N^{(i)}))^{ op} \in \mathbb{R}^{N-1}$$

for each prompt. We also let $\boldsymbol{h}^{(i)}_{j}^{\rightarrow}$, $\boldsymbol{h}^{(i)}_{j}^{\leftarrow}$ be the corresponding hidden states obtained from passing these prompts through the RNN, and define $\boldsymbol{h}_{j}^{(i)}$ using them. Note that $\boldsymbol{g}^{(1)},\ldots,\boldsymbol{g}^{(P)}$ 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

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

and

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

where $\delta \in (0,1)$ will be chosen later. In other words, E_1 is the event in which $\boldsymbol{g}^{(i)}$ 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 ${m g}^{(i)}-{m g}^{(k)}\stackrel{(d)}{=}\sqrt{2}{m g}$ where ${m g}\sim\mathcal{N}(0,{f I}_{N-1}).$ By a union bound we have

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

for all $\varepsilon_g \leq c\sqrt{2}$, where c>0 is an absolute constant such that $c\sqrt{N-1} \leq \mathbb{E}[\|\boldsymbol{g}\|]$, 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 $\mathbb{P}\big(E_1^C\big) \leq P^2 e^{-c^2(N-1)/8}$.

To lower bound $\mathbb{P}(E_2)$, consider a random prompt-label pair p, y and the corresponding g. Note that in the prompt p, the index t_i is drawn independently of the rest of p, and has a uniform distribution

in [N]. Let $p[t_j \mapsto k]$ denote a modification of p where we set t_j equal to k, and let $y[t_j \mapsto k]$ be the labels corresponding to this modified prompt. We then have

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

As a result, via a Markov inequality, we obtain

$$\begin{split} \mathbb{P}\bigg(\frac{1}{N}\|\Phi(\boldsymbol{h}_{j})-\boldsymbol{g}\|_{2}^{2} &\geq \frac{\varepsilon^{2}}{\delta^{2}}\bigg) = \mathbb{P}\bigg(\mathbb{E}_{t_{j}}\big[(\hat{y}_{\mathrm{RNN}}(\boldsymbol{p})_{j}-y_{j})^{2}\big] \geq \frac{\varepsilon^{2}}{\delta^{2}}\bigg) \\ &\leq \frac{\delta^{2}\,\mathbb{E}\big[(\hat{y}_{\mathrm{RNN}}(\boldsymbol{p})_{j}-y_{j})^{2}\big]}{\varepsilon^{2}} \\ &\leq \delta^{2}. \end{split}$$

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

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

Note that $(z^{(i)})$ are i.i.d. since $h_i^{(i)}$ and $g^{(i)}$ do not depend on x_j . Then, by Hoeffding's inequality,

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

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

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

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

$$\begin{split} \left\| \Phi(\boldsymbol{h}_{j}^{(i)}) - \Phi(\boldsymbol{h}_{j}^{(k)}) \right\|_{2} &\geq \left\| \boldsymbol{g}^{(i)} - \boldsymbol{g}^{(k)} \right\|_{2} - \left\| \Phi(\boldsymbol{h}_{j}^{(i)}) - \boldsymbol{g}^{(i)} \right\|_{2} - \left\| \Phi(\boldsymbol{h}_{j}^{(k)}) - \boldsymbol{g}^{(k)} \right\|_{2} \\ &\geq \varepsilon_{g} \sqrt{N - 1} - \frac{2\varepsilon\sqrt{N}}{\delta} =: \mathfrak{L}\sqrt{N}\varepsilon_{h}. \end{split}$$

Note that from the Lipschitzness of f_y , we have $\left\|\Phi(\boldsymbol{h}_j^{(i)}) - \Phi(\boldsymbol{h}_j^{(k)})\right\|_2 \leq \frac{\mathfrak{L}\sqrt{N}}{r_h} \left\|\boldsymbol{h}_j^{(i)} - \boldsymbol{h}_j^{(k)}\right\|_2$. As a result, the set $\left\{\boldsymbol{h}_j^{(i)}: i \in \mathcal{I}\right\}$ is an $r_h \varepsilon_h$ -packing for $\{\boldsymbol{h}: \|\boldsymbol{h}\|_2 \leq \sqrt{2}\|\boldsymbol{U}\|_{\text{op}} r_h\}$. Using Lemma 41, the log packing number can be bounded by

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

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,

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

and

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

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

B.6 Proof of Theorem 7

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

Lemma 36. Suppose $\{x^{(i)}\}_{i=1}^n \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ and let $y^{(i)} = \langle u, x_{t_i} \rangle$ for arbitrary $t_i \in [N]$ and $u \in \mathbb{S}^{d-1}$. 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 $a, b \in \mathbb{R}^m$ and $W \in \mathbb{R}^{m \times d}$, such that

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

and

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

Proof. The proof of Lemma 36 is an immediate consequence of two lemmas.

- 1. Lemma 37 shows that the inputs $x^{(1)}, \ldots, x^{(n)}$ can be projected to sufficiently separated scalar values with a unit vector v.
- 2. Lemma 38 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 $\epsilon = \Omega(1/n^2)$ as given by Lemma 37.

The only missing piece is to upper bound $\|y\|_2$ appearing in the final bound of Lemma 38. To that end, we apply the following Markov inequality,

$$\mathbb{P}\Big(\|\boldsymbol{y}\|_{2}^{2} \geq 6n\Big) \leq \frac{\mathbb{E}\Big[\|\boldsymbol{y}\|_{2}^{2}\Big]}{6n} \leq \frac{1}{6}.$$

As the statement of Lemma 37 holds with probability at least $\frac{1}{3}$, this suggests that the statement of Lemma 36 holds with probability at least $\frac{1}{6}$, concluding the proof.

Lemma 37. Suppose $\{x^{(i)}\}_{i=1}^n \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$. Then, with probability at least 1/3, there exists some $v \in \mathbb{S}^{d-1}$ (dependent on $\{x^{(i)}\}$) such that for all $i \neq j$,

$$\left| \boldsymbol{v}^{\top} \boldsymbol{x}^{(i)} - \boldsymbol{v}^{\top} \boldsymbol{x}^{(j)} \right| = \Omega\left(\frac{1}{n^2}\right).$$
 (B.17)

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

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

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

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

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

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

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

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

Lemma 38. Consider some $z = (z^{(1)}, \dots, z^{(n)})^{\top} \in \mathbb{R}^n$ and $y = (y^{(1)}, \dots, y^{(n)})^{\top} \in \mathbb{R}^n$, such that $|z^{(i)} - z^{(j)}| \ge \epsilon$ for all $i \ne j$. For simplicity, assume $\epsilon \le 1$. 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

$$\|\boldsymbol{a}\|_{2}^{2} + \|\boldsymbol{w}\|_{2}^{2} + \|\boldsymbol{b}\|_{2}^{2} = \mathcal{O}\left(\frac{\|\boldsymbol{y}\|_{2}\sqrt{n + \|\boldsymbol{z}\|_{2}^{2}}}{\epsilon}\right).$$
 (B.18)

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

$$g(t) = \sum_{i=1}^{n} a'_{i} \sigma(w'_{i} t - b'_{i}) = y^{(1)} \sigma(t - z^{(1)} + 1) + \left(\frac{y^{(2)} - y^{(1)}}{z^{(2)} - z^{(1)}} - y^{(1)}\right) \sigma(t - z^{(1)})$$

$$+ \sum_{i=3}^{n} \left(\frac{y^{(i)} - y^{(i-1)}}{z^{(i)} - z^{(i-1)}} - \frac{y^{(i-1)} - y^{(i-2)}}{z^{(i-1)} - z^{(i-2)}}\right) \sigma(t - z^{(i-1)}).$$

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 $w_i'=1$, $\|\boldsymbol{b}'\|_2^2\lesssim \|\boldsymbol{z}\|_2^2+1$, and $\|\boldsymbol{a}'\|_2^2\lesssim \|\boldsymbol{y}\|_2^2/\epsilon^2$. For $\alpha=\left((\|\boldsymbol{z}\|_2^2+n)\epsilon^2/\|\boldsymbol{y}\|_2^2\right)^{1/4}$, let $\boldsymbol{u}=\alpha\boldsymbol{u}', \boldsymbol{w}=\boldsymbol{w}'/\alpha$, and $\boldsymbol{b}=\boldsymbol{b}'/\alpha$. By homogeneity, the neural network with weights $(\boldsymbol{u},\boldsymbol{w},\boldsymbol{b})$ has identical outputs to that of $(\boldsymbol{u}',\boldsymbol{w}',\boldsymbol{b}')$ and satisfies (B.18), completing the proof.

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 $Uh = (U^{\rightarrow}h^{\rightarrow}, U^{\leftarrow}h^{\leftarrow})$, f_y is \mathcal{L} -Lipschitz with

$$\mathfrak{L} = \left\|oldsymbol{W}_{L_y}
ight\|_{ ext{op}} \left\|oldsymbol{W}_{L_y-1}
ight\|_{ ext{op}} \ldots \left\|oldsymbol{W}_2
ight\|_{ ext{op}}.$$

Using the AM-GM inequality,

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

As a result, we have $\mathfrak{L}\|U\|_{\text{op}} \leq e^{N^c/2}$. By invoking Proposition 26, to obtain population risk less than some absolute constant $c_3 > 0$, we need

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

This implies $n \ge d_h \ge \Omega(N^{1-c})$. 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 $d_h \ge n$.

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 $\Theta_h^{\rightarrow} = 0$, $\Theta_h^{\leftarrow} = 0$, U = 0, and describe the construction of $W_{L_y}, \ldots, W_2, W_y$, and (b_l) in the following. Using the construction of Lemma 36, we can let

$$oldsymbol{W}_y = egin{pmatrix} oldsymbol{W} & oldsymbol{0}_{n imes d_E} \ oldsymbol{0}_{(m-n) imes d} & oldsymbol{0}_{(m-n) imes d_E} \end{pmatrix}, \quad oldsymbol{b}_1 = egin{pmatrix} oldsymbol{b} \ oldsymbol{0}_{m-n} \end{pmatrix}, \quad oldsymbol{W}_2 = egin{pmatrix} oldsymbol{a}^ op & oldsymbol{0}^ op_{m-n} \ oldsymbol{0}_{(m-n) imes n} & oldsymbol{0}_{m-n} \ oldsymbol{0}_{(m-n) imes n} & oldsymbol{0}_{(m-n) imes n} \end{pmatrix},$$

where $W \in \mathbb{R}^{n \times d}$, and $a, b \in \mathbb{R}^n$ are given by Lemma 36. Then,

$$\boldsymbol{W}_{2}^{\top} \sigma(\boldsymbol{W}_{y} \boldsymbol{x}_{j^{(i)}}^{(i)} + \boldsymbol{b}_{y}) = (y_{j^{(i)}}^{(i)}, -y_{j^{(i)}}^{(i)}, 0, \dots, 0)^{\top}.$$

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

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

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

$$f_y(\boldsymbol{U}^{
ightarrow}\boldsymbol{h}_j^{
ightarrow}, \boldsymbol{U}^{\leftarrow}\boldsymbol{h}_j^{\leftarrow}, \boldsymbol{z}_{j^{(i)}}^{(i)}; \boldsymbol{\Theta}_y) = y_{j^{(i)}}^{(i)}$$

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 $Uh = (U^{\rightarrow}h^{\rightarrow}, U^{\leftarrow}h^{\leftarrow})$, f_u is \mathfrak{L} -Lipschitz with

$$\mathfrak{L} = \left\| \boldsymbol{W}_{L_y} \right\|_{\text{op}} \left\| \boldsymbol{W}_{L_y - 1} \right\|_{\text{op}} \dots \left\| \boldsymbol{W}_2 \right\|_{\text{op}}.$$

From Proposition 6, using the fact that $\|\cdot\|_{op} \leq \|\cdot\|_{F}$ and the AM-GM inequality, we obtain

$$\frac{1}{L_y} \| \operatorname{vec}(\boldsymbol{\Theta}) \|_2^2 \ge \left(\mathfrak{L}^2 \| \boldsymbol{U} \|_{\operatorname{op}}^2 \right)^{1/L_y} \ge \Omega \left(\left(\frac{N}{\log d_h} \right)^{1/L_y} \right)$$

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 $c_1>0$, $\left\|\operatorname{vec}(\hat{\Theta}_{\varepsilon})\right\|_2$ will always be less than the lower bound above, with some absolute constant probability $c_2>0$ over the training set.

C Auxiliary Lemmas

Lemma 39. Suppose $A \in \mathbb{R}^{d_1 \times d_2}$ and $B \in \mathbb{R}^{d_2 \times d_3}$. Then, for all $r, s \geq 1$ and $p, q \geq 1$ such that 1/p + 1/q = 1, we have

$$\|AB\|_{r,s} \le \|A\|_{r,p} \|B\|_{q,s}.$$

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

$$\|m{A}m{b}\|_r = \left\|\sum_{j=1}^{d_2} b_j m{A}_{:,j}
ight\| \ \le \sum_{j=1}^{d_2} |b_j| \|m{A}_{:,j}\|_r \le \|m{A}\|_{r,p} \|m{b}\|_q,$$

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

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

The next lemma follows from standard Gaussian integration.

Lemma 40. Suppose
$$\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
. Then $\operatorname{Var}(\|\boldsymbol{x}\|^2) = 2\operatorname{tr}(\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma}) + 4\boldsymbol{\mu}^{\top}\boldsymbol{\Sigma}\boldsymbol{\mu}$.

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

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

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

Lemma 42 ([Bac17, Propositions 1 and 6]). Suppose $f: \mathbb{R}^d \to \mathbb{R}$ satisfies $|f(x)| \leq LR$ and $|f(x) - f(x')| \leq L||x - x'||_2$ for all $x, x' \in \mathbb{R}^d$ with $||x|||_2 \leq R$ and $||x'||_2 \leq R$ and some

constants L, R > 0. Then, for every $\varepsilon > 0$, 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

$$\sup_{\|\boldsymbol{x}\|_2 \le R} |f(\boldsymbol{x}) - \boldsymbol{a}^{\top} \sigma(\boldsymbol{W} \boldsymbol{x} + \boldsymbol{b})| \le \varepsilon.$$

Additionally, we have

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

D Proof of Theorem 9

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

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

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

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

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

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

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

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

$$\mathbb{E}_{\boldsymbol{u},y,j,\boldsymbol{p}}[(y_j - f(t_1, \boldsymbol{W}\boldsymbol{x})_j)^2] \ge \mathbb{E}_{\boldsymbol{u},t_1,\boldsymbol{V}\boldsymbol{x}}[\operatorname{Var}(y \mid \boldsymbol{u},t_1,\boldsymbol{V}\boldsymbol{x})] = \mathbb{E}_{\boldsymbol{u},t_1,\boldsymbol{V}\boldsymbol{x}}[\operatorname{Var}(\langle \boldsymbol{P}_{t_1}\boldsymbol{u},\boldsymbol{x}\rangle \mid \boldsymbol{u},t_1,\boldsymbol{V}\boldsymbol{x})],$$

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

fact that via training, f and W can depend on u, but the prediction depends on x only through Vx. Consequently, we replace the prediction of the FFN by the best predictor having access to u, t_1 , and Vx. Note that t_1 , u, and Vx are jointly independent, and the joint distribution $(\langle P_{t_1}u, x \rangle, Vx)$ is

given by
$$\mathcal{N}\left(0,\begin{pmatrix} 1 & VP_{t_1}u \\ u^{\top}P_{t_1}^{\top}V^{\top} & \mathbf{I}_n \end{pmatrix}\right)$$
, thus we have

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

In particular,

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

and

$$\mathbb{E}_{\boldsymbol{u},t_1}[\operatorname{Var}(\langle \boldsymbol{P}_{t_1}\boldsymbol{u}, \boldsymbol{x}\rangle \mid \boldsymbol{u}, t_1, \boldsymbol{V}\boldsymbol{x})] = 1 - \frac{1}{Nd} \sum_{t_1=1}^{N} \sum_{i=1}^{n} \left\| \boldsymbol{P}_{t_1}^{\top} \boldsymbol{v}^{(i)} \right\|^2$$
$$= 1 - \frac{1}{Nd} \sum_{i=1}^{n} \left\| \boldsymbol{v}^{(i)} \right\|^2 = 1 - \frac{n}{Nd}.$$

E Experimental Details and Additional Results

In this section, we provide the details of our experimental setup, as well as additional results on the effect of q in Figure 3.

Architectures. We use a Transformer composed of a multihead attention layer with q heads, where each heads observes the entire $d+(q+1)d_e$ -dimensional input token, followed by a fully connected ReLU layer with width 100. For the RNN, we use a simple bidirectional RNN with a hidden state size $500 \times q$, and a linear readout layer. For the FFN, we use a depth-3 fully connected ReLU network, where the first layer has width Ndq and the second layer has width 1000. The output layer of the FFN has width N to match the input sequence.

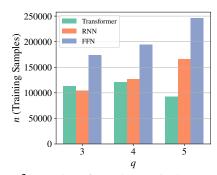


Figure 3: Number of samples required to get to test MSE loss 0.88 while training with online AdamW for the quadratic qSTR model explained in Appendix E with N=7. The gap increases with larger q. A closer theoretical analysis capturing the effect of large q can be an interesting direction for future work.

Optimization. For Figures 1 and 3 we use online AdamW with weight decay 0.1, where in Figure 1 we use a learning rate of 10^{-3} and in Figure 3 we use a learning rate of 10^{-4} . Each optimization step uses an independent batch size of 64 samples, and we track the test MSE loss using an independent set of 10,000 samples. For Figure 2 we use AdamW with weight decay 0.2 and learning rate 10^{-3} on a fixed training set of 50,000 samples.

Data Generating Model. In all experiments, we sample $x \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{Nd})$. For Figures 1 and 2 we have q=1 and define $g(\boldsymbol{x}_1)=\langle \boldsymbol{u},\boldsymbol{x}_1\rangle$ for a unit-norm \boldsymbol{u} uniformly sampled from the unit sphere. For Figure 3 we let $g(\boldsymbol{x}_1,\ldots,\boldsymbol{x}_q)=\frac{1}{\sqrt{q}}\sum_{i=1}^q \operatorname{He}_2(\langle \boldsymbol{u}_i,\boldsymbol{x}_i\rangle)$ where $\operatorname{He}_2(z)=(z^2-1)/\sqrt{2}$ is the normalized second Hermite polynomial. We use a non-linear g as this is a more challenging setting where e.g. Transformers require q heads by Theorem 4.

The code to reproduce all our experiments is provided at: https://github.com/mousavih/transformers-separation.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We prove all the claims made in the abstract and introduction in this paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations throughout the paper and mostly in Section 6.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Our Definition 2 of the qSTR model and our Assumptions 1, 2 and 3 are clearly stated.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA].

Justification: Our contributions are theoretical and our limited numerical simulations are for illustration purposes only. We will include a link to the GitHub repository of our code in the de-anonymized version.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: Our contributions are theoretical and our limited numerical simulations are for illustration purposes only. We will include a link to the GitHub repository of our code in the de-anonymized version.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: Our contributions are theoretical and our limited numerical simulations are for illustration purposes only. We will include a link to the GitHub repository of our code in the de-anonymized version.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA].

Justification: Our contributions are theoretical and our limited numerical simulations are for illustration purposes only.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: Our contributions are theoretical and our limited numerical simulations are for illustration purposes only.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have followed the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA].

Justification: Our contributions are theoretical and do not have immediate societal impacts. Guidelines:

• The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our contributions are theoretical.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: Our contributions are theoretical and we do not use any such assets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: Our contributions are theoretical and we do not introduce new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Our contributions are theoretical and we do not perform this type of experiment.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our contributions are theoretical and we do not perform this type of experiment. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We did not use LLMs.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.