

# THE RECURRENT NEURAL TANGENT KERNEL

Sina Alemohammad, Zichao Wang, Randall Balestriero, Richard G. Baraniuk

Department of Electrical and Computer Engineering

Rice University

{sa86, zw16, rb42, richb}@rice.edu

## ABSTRACT

The study of deep neural networks (DNNs) in the infinite-width limit, via the so-called *neural tangent kernel* (NTK) approach, has provided new insights into the dynamics of learning, generalization, and the impact of initialization. One key DNN architecture remains to be kernelized, namely, the recurrent neural network (RNN). In this paper we introduce and study the *Recurrent Neural Tangent Kernel* (RNTK), which provides new insights into the behavior of overparametrized RNNs. A key property of the RNTK should greatly benefit practitioners is its ability to compare inputs of different length. To this end, we characterize how the RNTK weights different time steps to form its output under different initialization parameters and nonlinearity choices. A synthetic and 56 real-world data experiments demonstrate that the RNTK offers significant performance gains over other kernels, including standard NTKs, across a wide array of data sets.

## 1 INTRODUCTION

The overparameterization of modern deep neural networks (DNNs) has resulted in not only remarkably good generalization performance on unseen data (Novak et al., 2018; Neyshabur et al., 2019; Belkin et al., 2019) but also guarantees that gradient descent learning can find the global minimum of their highly nonconvex loss functions (Du et al., 2019b; Allen-Zhu et al., 2019b;a; Zou et al., 2018; Arora et al., 2019b). From these successes, a natural question arises: What happens when we take overparameterization to the limit by allowing the width of a DNN’s hidden layers to go to infinity? Surprisingly, the analysis of such an (impractical) DNN becomes analytically tractable. Indeed, recent work has shown that the training dynamics of (infinite-width) DNNs under gradient flow is captured by a constant kernel called the *Neural Tangent Kernel* (NTK) that evolves according to a linear ordinary differential equation (ODE) (Jacot et al., 2018; Lee et al., 2019; Arora et al., 2019a).

Every DNN architecture and parameter initialization produces a distinct NTK. The original NTK was derived from the Multilayer Perceptron (MLP)(Jacot et al., 2018) and was soon followed by kernels derived from Convolutional Neural Networks (CNTK) (Arora et al., 2019a; Yang, 2019a), Residual DNNs (Huang et al., 2020), and Graph Convolutional Neural Networks (GNTK) (Du et al., 2019a). In (Yang, 2020a), a general strategy to obtain the NTK of any architecture is provided.

***In this paper, we extend the NTK concept to the important class of overparametrized Recurrent Neural Networks (RNNs), a fundamental DNN architecture for processing sequential data. We show that RNN in its infinite-width limit converges to a kernel that we dub the *Recurrent Neural Tangent Kernel* (RNTK). The RNTK provides high performance for various machine learning tasks, and an analysis of the properties of the kernel provides useful insights into the behavior of RNNs in the following overparametrized regime. In particular, we derive and study the RNTK to answer the following theoretical questions:***

***Q: Can the RNTK extract long-term dependencies between two data sequences?*** RNNs are known to underperform at learning long-term dependencies due to the gradient vanishing or exploding (Bengio et al., 1994). Attempted ameliorations have included orthogonal weights (Arjovsky et al., 2016; Jing et al., 2017; Henaff et al., 2016) and gating such as in Long Short-Term Memory (LSTM) (Hochreiter & Schmidhuber, 1997) and Gated Recurrent Unit (GRU) (Cho et al., 2014) RNNs. We demonstrate that the RNTK can detect long-term dependencies with proper initialization of the hyperparameters, and moreover, we show how the dependencies are extracted through time via different hyperparameter choices.

**Q: Do the recurrent weights of the RNTK reduce its representation power compared to other NTKs?** An attractive property of an RNN that is shared by the RNTK is that it can deal with sequences of different lengths via weight sharing through time. This enables the reduction of the number of learnable parameters and thus more stable training at the cost of reduced representation power. We prove the surprising fact that employing tied vs. untied weights in an RNN *does not* impact the analytical form of the RNTK.

**Q: Does the RNTK generalize well?** A recent study has revealed that the use of an SVM classifier with the NTK, CNTK, and GNTK kernels outperforms other classical kernel-based classifiers and trained finite DNNs on small data sets (typically fewer than 5000 training samples) (Lee et al., 2020; Arora et al., 2019a; 2020; Du et al., 2019a). We extend these results to RNTKs to demonstrate that the RNTK outperforms a variety of classic kernels, NTKs and finite RNNs for time series data sets in both classification and regression tasks. Carefully designed experiments with data of varying lengths demonstrate that the RNTK’s performance accelerates beyond other techniques as the difference in lengths increases. Those results extend the empirical observations from (Arora et al., 2019a; 2020; Du et al., 2019a; Lee et al., 2020) into finite DNNs, NTK, CNTK, and GNTK comparisons by observing that their performance-wise ranking depends on the employed DNN architecture.

We summarize our contributions as follows:

**[C1]** We derive the analytical form for the RNTK of an overparametrized RNN at initialization using rectified linear unit (ReLU) and error function (erf) nonlinearities for arbitrary data lengths and number of layers (Section 3.1).

**[C2]** We prove that the RNTK remains constant during (overparametrized) RNN training and that the dynamics of training are simplified to a set of ordinary differential equations (ODEs) (Section 3.2).

**[C3]** When the input data sequences are of equal length, we show that the RNTKs of weight-tied and weight-untied RNNs converge to the same RNTK (Section 3.3).

**[C4]** Leveraging our analytical formulation of the RNTK, we empirically demonstrate how correlations between data at different times are weighted by the function learned by an RNN for different sets of hyperparameters. We also offer practical suggestions for choosing the RNN hyperparameters for deep information propagation through time (Section 3.4).

**[C5]** We demonstrate that the RNTK is eminently practical by showing its superiority over classical kernels, NTKs, and finite RNNs in exhaustive experiments on time-series classification and regression with both synthetic and 56 real-world data sets (Section 4).

## 2 BACKGROUND AND RELATED WORK

**Notation.** We denote  $[n] = \{1, \dots, n\}$ , and  $I_d$  the identity matrix of size  $d$ .  $[A]_{ij}$  represents the  $(i, j)$ -th entry of a matrix, and similarly  $[a]_i$  represents the  $i$ -th entry of a vector. We use  $\phi(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  to represent the activation function that acts coordinate wise on a vector and  $\phi'$  to denote its derivative. We will often use the rectified linear unit (ReLU)  $\phi(x) = \max(0, x)$  and error function (erf)  $\phi(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-z^2} dz$  activation functions.  $\mathcal{N}(\mu, \Sigma)$  represents the multidimensional Gaussian distribution with the mean vector  $\mu$  and the covariance matrix  $\Sigma$ .

**Recurrent Neural Networks (RNNs).** Given an input sequence data  $\mathbf{x} = \{x_t\}_{t=1}^T$  of length  $T$  with data at time  $t$ ,  $x_t \in \mathbb{R}^m$ , a *simple RNN* (Elman, 1990) performs the following recursive computation at each layer  $\ell$  and each time step  $t$

$$\mathbf{g}^{(\ell;t)}(\mathbf{x}) = \mathbf{W}^{(\ell)} \mathbf{h}^{(\ell;t-1)}(\mathbf{x}) + \mathbf{U}^{(\ell)} \mathbf{h}^{(\ell-1;t)}(\mathbf{x}) + \mathbf{b}^{(\ell)}, \quad \mathbf{h}^{(\ell;t)}(\mathbf{x}) = \phi(\mathbf{g}^{(\ell;t)}(\mathbf{x})),$$

where  $\mathbf{W}^{(\ell)} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b}^{(\ell)} \in \mathbb{R}^n$  for  $\ell \in [L]$ ,  $\mathbf{U}^{(1)} \in \mathbb{R}^{n \times m}$  and  $\mathbf{U}^{(\ell)} \in \mathbb{R}^{n \times n}$  for  $\ell \geq 2$  are the RNN parameters.  $\mathbf{g}^{(\ell;t)}(\mathbf{x})$  is the pre-activation vector at layer  $\ell$  and time step  $t$ , and  $\mathbf{h}^{(\ell;t)}(\mathbf{x})$  is the after-activation (hidden state). For the input layer  $\ell = 0$ , we define  $\mathbf{h}^{(0;t)}(\mathbf{x}) := \mathbf{x}_t$ ,  $\mathbf{h}^{(\ell;0)}(\mathbf{x})$  as the initial hidden state at layer  $\ell$  that must be initialized to start the RNN recursive computation.

The output of an  $L$ -hidden layer RNN with linear read out layer is achieved via

$$f(\mathbf{x}) = \mathbf{V} \mathbf{h}^{(L;T)}(\mathbf{x}),$$

where  $\mathbf{V} \in \mathbb{R}^{d \times n}$ . Figure 1 visualizes an RNN unrolled through time.

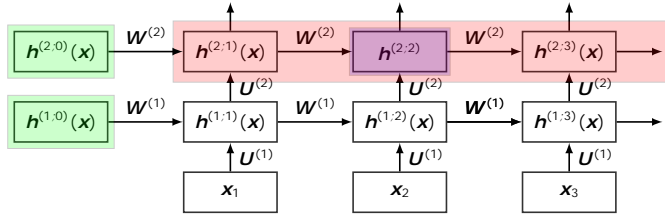


Figure 1: Visualization of a simple RNN that highlights a cell (purple), a layer (red) and the initial hidden state of each layer (green). (Best viewed in color.)

**Neural Tangent Kernel (NTK).** Let  $f(\mathbf{x}) \in \mathbb{R}^d$  be the output of a DNN with parameters  $\theta$ . For two input data sequences  $\mathbf{x}$  and  $\mathbf{x}^\ell$ , the NTK is defined as (Jacot et al., 2018)

$$\widehat{\mathcal{N}}_s(\mathbf{x}, \mathbf{x}^\ell) = \mathbf{h}^T r_s f_s(\mathbf{x}), \mathbf{r}_s f_s(\mathbf{x}^\ell) i,$$

where  $f_s$  and  $\theta_s$  are the network output and parameters during training at time  $s$ .<sup>1</sup> Let  $\mathcal{X}$  and  $\mathcal{Y}$  be the set of training inputs and targets,  $\ell(\hat{y}, y) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$  be the loss function, and  $L = \frac{1}{|\mathcal{X}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}} \ell(f_s(\mathbf{x}), \mathbf{y})$  be the empirical loss. The evolution of the parameters  $\theta_s$  and output of the network  $f_s$  on a test input using gradient descent with infinitesimal step size (a.k.a gradient flow) with learning rate  $\eta$  is given by

$$\frac{\partial \theta_s}{\partial s} = \eta \mathbf{r}_s f_s(\mathcal{X})^T \mathbf{r}_s f_s(\mathcal{X}) L \quad (1)$$

$$\frac{\partial f_s(\mathbf{x})}{\partial s} = \eta \mathbf{r}_s f_s(\mathbf{x}) \mathbf{r}_s f_s(\mathcal{X})^T \mathbf{r}_s f_s(\mathcal{X}) L = \widehat{\mathcal{N}}_s(\mathbf{x}, \mathcal{X}) \mathbf{r}_s f_s(\mathcal{X}) L. \quad (2)$$

Generally,  $\widehat{\mathcal{N}}_s(\mathbf{x}, \mathbf{x}^\ell)$ , hereafter referred to as the empirical NTK, changes over time during training, making the analysis of the training dynamics difficult. When  $f_s$  corresponds to an infinite-width MLP, (Jacot et al., 2018) showed that  $\widehat{\mathcal{N}}_s(\mathbf{x}, \mathbf{x}^\ell)$  converges to a limiting kernel at initialization and stays constant during training, i.e.,

$$\lim_{n \uparrow} \widehat{\mathcal{N}}_s(\mathbf{x}, \mathbf{x}^\ell) = \lim_{n \uparrow} \widehat{\mathcal{N}}_0(\mathbf{x}, \mathbf{x}^\ell) := \widehat{\mathcal{N}}(\mathbf{x}, \mathbf{x}^\ell) \partial s,$$

which is equivalent to replacing the outputs of the DNN by their first-order Taylor expansion in the parameter space (Lee et al., 2019). With a mean-square error (MSE) loss function, the training dynamics in (1) and (2) simplify to a set of linear ODEs, which coincides with the training dynamics of kernel ridge regression with respect to the NTK when the ridge term goes to zero. A nonzero ridge regularization can be conjured up by adding a regularization term  $\frac{\lambda}{2} k \theta_s = \theta_0 k_2^2$  to the empirical loss (Hu et al., 2020).

### 3 THE RECURRENT NEURAL TANGENT KERNEL

We are now ready to derive the RNTK. We first prove the convergence of an RNN at initialization to the RNTK in the infinite-width limit and discuss various insights it provides. We then derive the convergence of an RNN after training to the RNTK. Finally, we analyze the effects of various hyperparameter choices on the RNTK. Proofs of all of our results are provided in the Appendices.

#### 3.1 RNTK FOR AN INFINITE-WIDTH RNN AT INITIALIZATION

First we specify the following parameter initialization scheme that follows previous work on NTKs (Jacot et al., 2018), which is crucial to our convergence results:

$$\mathbf{W}^{(\cdot)} = \frac{\sigma_w}{\sqrt{n}} \mathbf{W}^{(\cdot)}, \quad \mathbf{U}^{(1)} = \frac{\sigma_u}{\sqrt{m}} \mathbf{U}^{(1)}, \quad \mathbf{U}^{(\cdot)} = \frac{\sigma_u}{\sqrt{n}} \mathbf{U}^{(\cdot)} (\ell \geq 2), \quad \mathbf{V} = \frac{\sigma_v}{\sqrt{n}} \mathbf{V}, \quad \mathbf{b}^{(\cdot)} = \sigma_b \mathbf{b}^{(\cdot)}, \quad (3)$$

where

$$[\mathbf{W}^{(\cdot)}]_{ij}, [\mathbf{U}^{(\cdot)}]_{ij}, [\mathbf{V}]_{ij}, [\mathbf{b}^{(\cdot)}]_i \sim N(0, 1). \quad (4)$$

We will refer to (3) and (4) as the *NTK initialization*. The choices of the hyperparameters  $\sigma_w$ ,  $\sigma_u$ ,  $\sigma_v$  and  $\sigma_b$  can significantly impact RNN performance, and we discuss them in detail in Section

<sup>1</sup>We use  $s$  to denote time here, since  $t$  is used to index the time steps of the RNN inputs.

3.4. For the initial (at time  $t = 0$ ) hidden state at each layer  $\ell$ , we set  $\mathbf{h}^{(\cdot;0)}(\mathbf{x})$  to an i.i.d. copy of  $\mathcal{N}(0, \sigma_h)$  (Wang et al., 2018). For convenience, we collect all of the learnable parameters of the RNN into  $\theta = \text{vect}[\mathcal{F}\mathcal{F}\mathbf{W}^{(\cdot)}, \mathbf{U}^{(\cdot)}, \mathbf{b}^{(\cdot)}g_{=1}^L, \mathbf{V}g]$ .

The derivation of the RNTK at initialization is based on the correspondence between Gaussian initialized, infinite-width DNNs and Gaussian Processes (GPs), known as the DNN-GP. In this setting every coordinate of the DNN output tends to a GP as the number of units/neurons in the hidden layer (its width) goes to infinity. The corresponding DNN-GP kernel is computed as

$$K(\mathbf{x}, \mathbf{x}^\theta) = \mathbb{E}_N [[f(\mathbf{x})]_i \cdot [f(\mathbf{x}^\theta)]_i], \quad \forall i \geq [d]. \quad (5)$$

First introduced for a single-layer, fully-connected neural network by (Neal, 1995), recent works on NTKs have extended the results for various DNN architectures (Lee et al., 2018; Duvenaud et al., 2014; Novak et al., 2019; Garriga-Alonso et al., 2019; Yang, 2019b), where in addition to the output, all pre-activation layers of the DNN tends to a GPs in the infinite-width limit. In the case of RNNs, each coordinate of the RNN pre-activation  $\mathbf{g}^{(\cdot;t)}(\mathbf{x})$  converges to a centered GP depending on the inputs with kernel

$$\mathbf{g}^{(\cdot;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \mathbb{E}_N [[\mathbf{g}^{(\cdot;t)}(\mathbf{x})]_i \cdot [\mathbf{g}^{(\cdot;t^\theta)}(\mathbf{x}^\theta)]_i] \quad \forall i \geq [n]. \quad (6)$$

As per (Yang, 2019a), the gradients of random infinite-width DNNs computed during backpropagation are also Gaussian distributed. In the case of RNNs, every coordinate of the vector  $\mathbf{r}^{(\cdot;t)}(\mathbf{x}) := \frac{1}{n} \mathbf{r}_{\mathbf{g}^{(\cdot;t)}(\mathbf{x})} f(\mathbf{x})$  converges to a GP with kernel

$$\mathbf{r}^{(\cdot;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \mathbb{E}_N [[\mathbf{r}^{(\cdot;t)}(\mathbf{x})]_i \cdot [\mathbf{r}^{(\cdot;t^\theta)}(\mathbf{x}^\theta)]_i] \quad \forall i \geq [n]. \quad (7)$$

Both convergences occur independently of the coordinate index  $i$  and for inputs of possibly different lengths, i.e.,  $T \neq T^\theta$ . With (6) and (7), we now prove that an infinite-width RNN at initialization converges to the limiting RNTK.

**Theorem 1** *Let  $\mathbf{x}$  and  $\mathbf{x}^\theta$  be two data sequences of potentially different lengths  $T$  and  $T^\theta$ , respectively. Without loss of generality, assume that  $T \geq T^\theta$ , and let  $\tau := T^\theta - T$ . Let  $n$  be the number of units in the hidden layers, the empirical RNTK for an  $L$ -layer RNN with NTK initialization converges to the following limiting kernel as  $n \rightarrow \infty$*

$$\lim_{n \rightarrow \infty} \hat{\mathbf{K}}_0(\mathbf{x}, \mathbf{x}^\theta) = \mathbf{K}(\mathbf{x}, \mathbf{x}^\theta) = \mathbf{K}^{(L;T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \quad \mathbf{I}_d, \quad (8)$$

where

$$\mathbf{K}^{(L;T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \left( \sum_{i=1}^L \sum_{t=1}^T \left( \mathbf{r}^{(\cdot;t+)}(\mathbf{x}, \mathbf{x}^\theta) \cdot \mathbf{r}^{(\cdot;t+)}(\mathbf{x}, \mathbf{x}^\theta) \right) \right) + K(\mathbf{x}, \mathbf{x}^\theta), \quad (9)$$

with  $K(\mathbf{x}, \mathbf{x}^\theta)$ ,  $\mathbf{r}^{(\cdot;t+)}(\mathbf{x}, \mathbf{x}^\theta)$ , and  $\mathbf{K}^{(L;T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta)$  defined in (5)–(7).

**Remarks.** Theorem 1 holds generally for any two data sequences, including different lengths ones. This highlights the RNTK’s ability to produce a similarity measure  $\mathbf{K}(\mathbf{x}, \mathbf{x}^\theta)$  even if the inputs are of different lengths, without resorting to heuristics such as zero padding the inputs to the to the max length of both sequences. Dealing with data of different length is in sharp contrast to common kernels such as the classical radial basis functions, polynomial kernels, and current NTKs. We showcase this capability below in Section 4.

To visualize Theorem 1, we plot in the left plot in Figure 2 the convergence of a single layer, sufficiently wide RNN to its RNTK with the two simple inputs  $\mathbf{x} = [1, -1, 1]g$  of length 3 and  $\mathbf{x}^\theta = [\cos(\alpha), \sin(\alpha)]g$  of length 2, where  $\alpha = [0, 2\pi]$ . For an RNN with a sufficiently large hidden state ( $n = 1000$ ), we see clearly that it converges to the RNTK ( $n \rightarrow \infty$ ).

**RNTK Example for a Single-Layer RNN.** We present a concrete example of Theorem 1 by showing how to recursively compute the RNTK for a single-layer RNN; thus we drop the layer index for notational simplicity. **We compute and display the RNTK for the general case of a multi-layer RNN in Appendix B.3.** To compute the RNTK  $\mathbf{K}^{(T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta)$ , we need to compute the GP

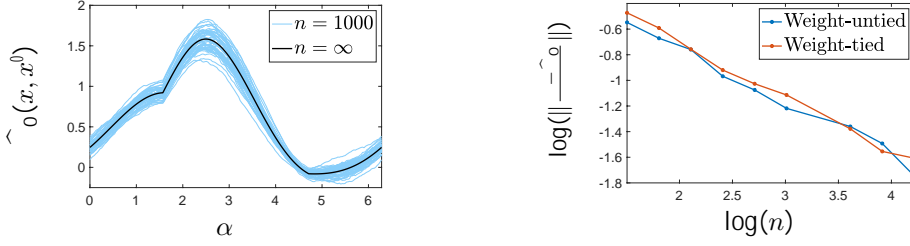


Figure 2: Empirical demonstration of a wide, single-layer RNN converging to its limiting RNTK. **Left:** convergence for a pair of different-length inputs  $\mathbf{x} = \{1; -1; 1\}$  and  $\mathbf{x}' = \{\cos(\alpha); \sin(\alpha)\}$ , with varying  $\alpha = [0; 2\pi]$ . The vertical axis corresponds to the RNTK values for different values of  $\alpha$ . **Right:** convergence of weight-tied and weight-untied single layer RNN to the same limiting RNTK with increasing width (horizontal axis). The vertical axis corresponds to the average of the log-normalized error between the empirical RNTK computed using finite RNNs and the RNTK for 50 Gaussian normal signals of length  $T = 5$ .

kernels  $\mathcal{K}^{(t;t^+)}(\mathbf{x}, \mathbf{x}^\theta)$  and  $\mathcal{K}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta)$ . We first define the operator  $\mathcal{V}[\mathbf{K}]$  that depends on the nonlinearity  $\phi(\cdot)$  and a positive semi-definite matrix  $\mathbf{K} \succeq \mathbb{R}^{2 \times 2}$

$$\mathcal{V}[\mathbf{K}] = \mathbb{E}[\phi(\mathbf{z}_1) \cdot \phi(\mathbf{z}_2)], \quad (\mathbf{z}_1, \mathbf{z}_2) \sim \mathcal{N}(0, \mathbf{K}). \quad (10)$$

Following (Yang, 2019a), we obtain the analytical recursive formula for the GP kernel  $\mathcal{K}^{(t;t^+)}(\mathbf{x}, \mathbf{x}^\theta)$  for a single layer RNN as

$$\mathcal{K}^{(1;1)}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_w^2 \sigma_h^2 \mathbb{1}_{(x=x^\theta)} + \frac{\sigma_u^2}{m} h \mathbf{x}_1, \mathbf{x}_1^\theta + \sigma_b^2 \quad (11)$$

$$\mathcal{K}^{(t;1)}(\mathbf{x}, \mathbf{x}^\theta) = \frac{\sigma_u^2}{m} h \mathbf{x}_t, \mathbf{x}_1^\theta + \sigma_b^2 \quad t > 1 \quad (12)$$

$$\mathcal{K}^{(1;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \frac{\sigma_u^2}{m} h \mathbf{x}_1, \mathbf{x}_{t^\theta}^\theta + \sigma_b^2 \quad t^\theta > 1 \quad (13)$$

$$\mathcal{K}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_w^2 \mathcal{V}[\mathbf{K}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta)] + \frac{\sigma_u^2}{m} h \mathbf{x}_t, \mathbf{x}_{t^\theta}^\theta + \sigma_b^2 \quad t, t^\theta > 1 \quad (14)$$

$$\mathcal{K}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_v^2 \mathcal{V}[\mathbf{K}^{(T+1;T^\theta+1)}(\mathbf{x}, \mathbf{x}^\theta)], \quad (15)$$

where

$$\mathbf{K}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \begin{bmatrix} \mathcal{K}^{(t-1;t-1)}(\mathbf{x}, \mathbf{x}) & \mathcal{K}^{(t-1;t^\theta-1)}(\mathbf{x}, \mathbf{x}^\theta) \\ \mathcal{K}^{(t-1;t^\theta-1)}(\mathbf{x}, \mathbf{x}^\theta) & \mathcal{K}^{(t^\theta-1;t^\theta-1)}(\mathbf{x}^\theta, \mathbf{x}^\theta) \end{bmatrix}. \quad (16)$$

Similarly, we obtain the analytical recursive formula for the GP kernel  $\mathcal{K}^{(t;t^+)}(\mathbf{x}, \mathbf{x}^\theta)$  as

$$\mathcal{K}^{(T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_v^2 \mathcal{V}_o[\mathbf{K}^{(T+1;T^\theta+1)}(\mathbf{x}, \mathbf{x}^\theta)] \quad (17)$$

$$\mathcal{K}^{(t;t^+)}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_w^2 \mathcal{V}_o[\mathbf{K}^{(t+1;t^++1)}(\mathbf{x}, \mathbf{x}^\theta)] \quad \mathcal{K}^{(t+1;t^++1)}(\mathbf{x}, \mathbf{x}^\theta) \quad t \geq [T-1] \quad (18)$$

$$\mathcal{K}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = 0 \quad t^\theta \neq t. \quad (19)$$

For  $\phi = \text{ReLU}$  and  $\phi = \text{erf}$ , we provide analytical expressions for  $\mathcal{V}[\mathbf{K}]$  and  $\mathcal{V}_o[\mathbf{K}]$  in Appendix B.5. These yield an explicit formula for the RNTK that enables fast and point-wise kernel evaluations. For other activation functions, one can apply the Monte Carlo approximation to obtain  $\mathcal{V}[\mathbf{K}]$  and  $\mathcal{V}_o[\mathbf{K}]$  (Novak et al., 2019).

### 3.2 RNTK FOR AN INFINITE-WIDTH RNN DURING TRAINING

We prove that an infinitely-wide RNN, not only at initialization but also *during* gradient descent training, converges to the limiting RNTK at initialization.

**Theorem 2** *Let  $n$  be the number of units of each RNN's layer. Assume that  $(X, X)$  is positive definite on  $X$  such that  $\lambda_{\min}((X, X)) > 0$ . Let  $\eta := 2(\lambda_{\min}((X, X)) + \lambda_{\max}((X, X)))^{-1}$ . For an  $L$ -layer RNN with NTK initialization as in (3), (4) trained under gradient flow (recall (1) and (2)) with  $\eta < \eta$ , we have with high probability*

$$\sup_s \frac{k\theta_s}{\rho} \frac{\theta_0 k_2}{n}, \sup_s \widehat{k}_s(X, X) \quad \widehat{k}_0(X, X) k_2 = \mathcal{O}\left(\frac{1}{n}\right).$$

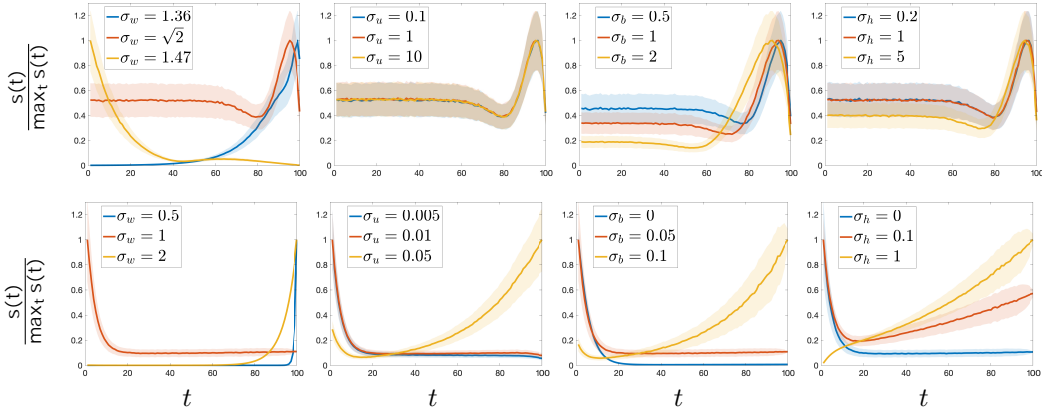


Figure 3: Per time step  $t$  (horizontal axis) sensitivity analysis (vertical axis) of the RNTK for the ReLU (top row) and erf (bottom row) activation functions for various weight noise hyperparameters. We also experiment with different RNTK hyperparameters in each of the subplots, given by the subplot internal legend. Clearly, the ReLU (top-row) provides a more stable kernel across time steps (highlighted by the near constant sensitivity through time). On the other hand, erf (bottom row) sees a more erratic behavior either focusing entirely on early time-steps or on the latter ones.

**Remarks.** Theorem 2 states that the training dynamics of an RNN in the infinite-width limit as in (1), (2) are governed by the RNTK derived from the RNN at its initialization. Intuitively, this is due to the NTK initialization (3), (4) which positions the parameters near a local minima, thus minimizing the amount of update that needs to be applied to the weights to obtain the final parameters.

### 3.3 RNTK FOR AN INFINITE-WIDTH RNN WITHOUT WEIGHT SHARING

We prove that, in the infinite-width limit, an RNN without weight sharing (untied weights), i.e., using independent new weights  $\mathbf{W}^{(\cdot:t)}$ ,  $\mathbf{U}^{(\cdot:t)}$  and  $\mathbf{b}^{(\cdot:t)}$  at each time step  $t$ , converges to the same RNTK as an RNN with weight sharing (tied weights). First, recall that it is a common practice to use weight-tied RNNs, i.e., in layer  $\ell$ , the weights  $\mathbf{W}^{(\cdot)}$ ,  $\mathbf{U}^{(\cdot)}$  and  $\mathbf{b}^{(\cdot)}$  are the same across all time steps  $t$ . This practice conserves memory and reduces the number of learnable parameters. We demonstrate that, when using untied-weights, the RNTK formula remains unchanged.

**Theorem 3** *For inputs of the same length, an RNN with untied weights converges to the same RNTK as an RNN with tied weights in the infinite-width ( $n \rightarrow \infty$ ) regime.*

**Remarks.** Theorem 3 implies that weight-tied and weight-untied RNNs have similar behaviors in the infinite-width limit. It also suggests that existing results on the simpler, weight-untied RNN setting may be applicable for the more general, weight-tied RNN. The plot on the right side of Figure 2 empirically demonstrates the convergence of both the weight-tied and weight-untied RNNs to the RNTK with increasing hidden layer size  $n$ ; moreover, the convergence rates are similar.

### 3.4 INSIGHTS INTO THE ROLES OF THE RNTK’S HYPERPARAMETERS

Our analytical form for the RNTK is fully determined by a small number of hyperparameters, which contains the various weight variances collected into  $S = \{\sigma_w, \sigma_u, \sigma_b, \sigma_h\}$  and the activation function.<sup>2</sup> In standard supervised-learning settings, one often performs cross-validation to select the hyperparameters. However, since kernel methods become computationally intractable for large datasets, we seek a more computationally friendly alternative to cross-validation. Here we conduct a novel exploratory analysis that provides new insights into the impact of the RNTK hyperparameters on the RNTK output and suggests a simple method to select them a priori in a deliberate manner.

To visualize the role of the RNTK hyperparameters, we introduce the *sensitivity*  $s(t)$  of the RNTK of two input sequences  $\mathbf{x}$  and  $\mathbf{x}^\theta$  with respect to the input  $\mathbf{x}_t$  at time  $t$

$$s(t) = k \nabla_{\mathbf{x}_t} (\mathbf{x}, \mathbf{x}^\theta) k_2. \tag{20}$$

<sup>2</sup>From (11) to (18) we emphasize that  $\sigma$  merely scales the RNTK and does not change its overall behavior.

Table 1: Summary of time series classification results on 56 real-world data sets. The RNTK outperforms classical kernels, the NTK, and trained RNNs across all metrics. See Appendix A for detailed description of the metrics.

	RNTK	NTK	RBF	Polynomial	Gaussian RNN	Identity RNN	GRU
Acc. mean "	<b>80.15%</b> <b>15.99%</b>	77.74% 16.61%	78.15% 16.59%	77.69% 16.40%	55.98% 26.42%	63.08% 19.02%	69.50% 22.67%
P90 "	<b>92.86%</b>	85.71%	87.60%	82.14%	28.57%	42.86%	60.71%
P95 "	<b>80.36%</b>	66.07%	75.00%	67.86%	17.86%	21.43%	46.43%
PMA "	<b>97.23%</b>	94.30%	94.86%	94.23%	67.06%	78.22%	84.31%
Friedman Rank #	<b>2.38</b>	3.00	2.89	3.46	5.86	5.21	4.21

Here,  $s(t)$  indicates how sensitive the RNTK is to the data at time  $t$ , i.e.,  $\mathbf{x}_t$ , in presence of another data sequence  $\mathbf{x}^j$ . Intuitively, large/small  $s(t)$  indicates that the RNTK is relatively sensitive/insensitive to the input  $\mathbf{x}_t$  at time  $t$ .

The sensitivity is crucial to understanding to which extent the RNTK prediction is impacted by the input at each time step. In the case where some time indices have a small sensitivity, then any input variation in those corresponding times will not alter the RNTK output and thus will produce a metric that is invariant to those changes. This situation can be beneficial or detrimental based on the task at hand. Ideally, and in the absence of prior knowledge on the data, one should aim to have a roughly constant sensitivity across time in order to treat all time steps equally in the RNTK input comparison.

Figure 3 plots the normalized sensitivity  $s(t)/\max_t(s(t))$  for two data sequences of the same length  $T = 100$ , with  $s(t)$  computed numerically for  $\mathbf{x}_t, \mathbf{x}_t^j \sim \mathcal{N}(0, 1)$ . We repeated the experiments 10000 times; the mean of the sensitivity is shown in Figure 3. Each of the plots shows the changes of parameters  $S_{\text{ReLU}} = \{2, 1, 0, 0\}$  for  $\phi = \text{ReLU}$  and  $S_{\text{erf}} = \{1, 0.01, 0.05, 0\}$  for  $\phi = \text{erf}$ .

From Figure 3 we first observe that both ReLU and erf show similar per time step sensitivity measure  $s(t)$  behavior around the hyperparameters  $S_{\text{ReLU}}$  and  $S_{\text{erf}}$ . If one varies any of the weight variance parameters, the sensitivity exhibits a wide range of behavior, and in particular with erf. We observe that  $\sigma_w$  has a major influence on  $s(t)$ . For ReLU, a small decrease/increase in  $\sigma_w$  can lead to over-sensitivity of the RNTK to data at the last/first times steps, whereas for erf, any changes in  $\sigma_w$  leads to over-sensitivity to the last time steps.

Another notable observation is the importance of  $\sigma_h$ , which is usually set to zero for RNNs. (Wang et al., 2018) showed that a non-zero  $\sigma_h$  acts as a regularization that improves the performance of RNNs with the ReLU nonlinearity. From the sensitivity perspective, a non-zero  $\sigma_h$  results in reducing the importance of the first time steps of the input. We also see the same behavior in erf, but with stronger changes as  $\sigma_h$  increases. Hence whenever one aims at reinforcing the input pairwise comparisons, such parameters should be favored.

This sensitivity analysis provides a practical tool for RNTK hyperparameter tuning. In the absence of knowledge about the data, hyperparameters should be chosen to produce the least time varying sensitivity. If given a priori knowledge, hyperparameters can be selected that direct the RNTK to the desired time-steps.

## 4 EXPERIMENTS

We now empirically validate the performance of the RNTK compared to classic kernels, NTKs, and trained RNNs on both classification and regression tasks using a large number of time series data sets. Of particular interest is the capability of the RNTK to offer high performance even on inputs of different lengths.

**Time Series Classification.** The first set of experiments considers time series inputs of the same lengths from 56 datasets in the UCR time-series classification data repository (Dau et al., 2019). We restrict ourselves to selected data sets with fewer than 1000 training samples and fewer than 1000 time steps ( $T$ ) as kernel methods become rapidly intractable for larger datasets. We compare the RNTK with a variety of other kernels, including the Radial Basis Kernel (RBF), polynomial kernel, and NTK (Jacot et al., 2018), as well as finite RNNs with Gaussian, identity (Le et al., 2015) initialization, and GRU (Cho et al., 2014). We use  $\phi = \text{ReLU}$  for both the RNTKs and NTKs. For each kernel, we train a C-SVM (Chang & Lin, 2011) classifier, and for each finite RNN we use gradient descent training. For model hyperparameter tuning, we use 10-fold cross-validation. Details on the data sets and experimental setup are available in Appendix A.1.

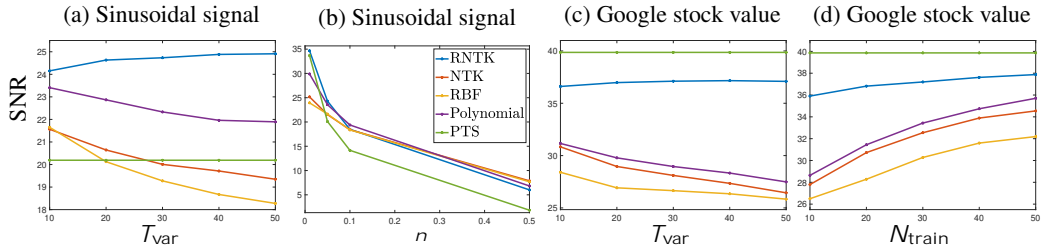


Figure 4: Performance of the RNTK on the synthetic sinusoid and real-world Google stock price data sets compared to three other kernels. We vary the input lengths (a,c), the input noise level (b), and training set size (d). We compute the average SNR by repeating each experiment 1000 times. The RNTK clearly outperforms all of the other kernels under consideration. Figure 4b suggests that the RNTK performs better when input noise level is low demonstrating one case where time recurrence from RNTK might be sub-optimal as it collects and accumulates the high noise from each time step as opposed to other kernels treating each independently.

We summarize the classification results over all 56 datasets in Table 1; detailed results on each data set is available in Appendix A.2. We see that the RNTK outperforms not only the classical kernels but also the NTK and trained RNNs in all metrics. The results demonstrate the ability of RNTK to provide increased performances compare to various other methods (kernels and RNNs). The superior performance of RNTK compared to other kernels, including NTK, can be explained by the internal recurrent mechanism present in RNTK, allowing time-series adapted sample comparison. In addition, RNTK also outperforms RNN and GRU. As the datasets we consider are relative small in size, finite RNNs and GRUs that typically require large amount of data to succeed do not perform well in our setting. An interesting future direction would be to compare RNTK to RNN/GRU on larger datasets.

**Time Series Regression.** We now validate the performance of the RNTK on time series inputs of *different* lengths on both synthetic data and real data. For both scenarios, the target is to predict the next time-step observation of the randomly extracted windows of different length using kernel ridge regression.

We compare the RNTK to other kernels, the RBF and polynomial kernels and the NTK. We also compare our results with a data independent predictor that requires no training, that is simply to predict the next time step with previous time step (PTS).

For the synthetic data experiment, we simulate 1000 samples of one period of a sinusoid and add white Gaussian noise with default  $\sigma_n = 0.05$ . From this fixed data, we extract training set size  $N_{train} = 20$  segments of uniform random lengths in the range of  $[T_{xed}, T_{xed} + T_{var}]$  with  $T_{xed} = 10$ . We use standard kernel ridge regression for this task. The test set is comprised of  $N_{test} = 5000$  obtained from other randomly extracted segments, again of varying lengths. For the real data, we use 975 days of the Google stock value in the years 2014–2018. As in the simulated signal setup above, we extract  $N_{train}$  segments of different lengths from the first 700 days and test on the  $N_{test}$  segments from days 701 to 975. Details of the experiment are available in Appendix A.2.

We report the predicted signal-to-noise ratio (SNR) for both datasets in Figures 4a and 4c for various values of  $T_{var}$ . We vary the noise level and training set size for fixed  $T_{var} = 10$  in Figures 4b and 4d. As we see from Figures 4a and 4c, the RNTK offers substantial performance gains compared to the other kernels, due to its ability to naturally deal with variable length inputs. Moreover, the performance gap increases with the amount of length variation of the inputs  $T_{var}$ . Figure 4d demonstrates that, unlike the other methods, the RNTK maintains its performance even when the training set is small. Finally, Figure 4c demonstrates that the impact of noise in the data on the regression performance is roughly the same for all models but becomes more important for RNTK with a large  $\sigma_n$ ; this might be attributed to the recurrent structure of the model allowing for a time propagation and amplification of the noise for very low SNR. These experiments demonstrate the distinctive advantages of the RNTK over classical kernels, and NTKs for input data sequences of varying lengths.

In the case of PTS, we expect the predictor to outperform kernel methods when learning from the training samples is hard, due to noise in the data or small training size which can lead to over fitting. In Figure 4a RNTK and Polynomial kernels outperforms PTS for all values of  $T_{var}$ , but for larger  $T_{var}$ , NTK and RBF under perform PTS due to the increasing detrimental effect of zero padding.



For the Google stock value, we see a superior performance of PTS with respect to all other kernel methods due to the nature of those data heavily relying on close past data. However, RNTK is able to reduce the effect of over-fitting, and provide the closest results to PTS among all kernel methods we employed, with increasing performance as the number of training samples increase.

## 5 CONCLUSIONS

In this paper, we have derived the RNTK based on the architecture of a simple RNN. We have proved that, at initialization, after training, and without weight sharing, any simple RNN converges to the same RNTK. This convergence provides new insights into the behavior of infinite-width RNNs, including how they process different-length inputs, their training dynamics, and the sensitivity of their output at every time step to different nonlinearities and initializations. We have highlighted the RNTK’s practical utility by demonstrating its superior performance on time series classification and regression compared to a range of classical kernels, the NTK, and trained RNNs. There are many avenues for future research, including developing RNTKs for gated RNNs such as the LSTM (Hochreiter & Schmidhuber, 1997) and investigating which of our theoretical insights extend to finite RNNs.

## ACKNOWLEDGMENTS

This work was supported by NSF grants CCF-1911094, IIS-1838177, and IIS-1730574; ONR grants N00014-18-12571, N00014-20-1-2787, and N00014-20-1-2534; AFOSR grant FA9550-18-1-0478; and a Vannevar Bush Faculty Fellowship, ONR grant N00014-18-1-2047.

## REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *In Advances in Neural Information Processing Systems*, pp. 6155–6166, 2019a.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. In *Advances in Neural Information Processing Systems*, pp. 6676–6688, 2019b.
- Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In *International Conference on Machine Learning*, pp. 1120–1128, 2016.
- Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pp. 8141–8150, 2019a.
- Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*, 2019b.
- Sanjeev Arora, Simon S. Du, Zhiyuan Li, Ruslan Salakhutdinov, Ruosong Wang, and Dingli Yu. Harnessing the power of infinitely wide deep nets on small-data tasks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkl8sJBYvH>.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. URL <https://www.pnas.org/content/116/32/15849>.
- Yoshua. Bengio, Patrice. Simard, and Paolo. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Networks*, 5(2):157–166, 1994.
- Erwin Bolthausen. An iterative construction of solutions of the tap equations for the sherrington–kirkpatrick model. *Communications in Mathematical Physics*, 325(1):333–366, 2014.
- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.

- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Youngmin Cho and Lawrence K Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, pp. 342–350, 2009.
- Hoang Anh Dau, Eamonn Keogh, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, Yanping Chen, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, Gustavo Batista, and ML Hexagon. The UCR time series classification archive, 2019. URL [https://www.cs.ucr.edu/~eamonn/time\\_series\\_data\\_2018/](https://www.cs.ucr.edu/~eamonn/time_series_data_2018/).
- Simon S. Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, pp. 5724–5734, 2019a.
- Simon S. Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pp. 1675–1685, 2019b.
- David Duvenaud, Oren Rippel, Ryan Adams, and Zoubin Ghahramani. Avoiding pathologies in very deep networks. In *Artificial Intelligence and Statistics*, pp. 202–210, 2014.
- Jeffrey L. Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, 1990. URL <https://www.sciencedirect.com/science/article/pii/036402139090002E>.
- M. Fernández-Delgado, M.S. Sirsat, E. Cernadas, S. Alawadi, S. Barro, and M. Febrero-Bande. An extensive experimental survey of regression methods. *Neural Networks*, 111:11 – 34, 2019.
- Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional networks as shallow gaussian processes. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bklfsi0cKm>.
- Mikael Henaff, Arthur Szlam, and Yann LeCun. Recurrent orthogonal networks and long-memory tasks. In *International Conference on Machine Learning*, pp. 2034–2042. PMLR, 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8): 1735–1780, 1997.
- Wei Hu, Zhiyuan Li, and Dingli Yu. Simple and effective regularization methods for training on noisily labeled data with generalization guarantee. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=Hke3gyHYwH>.
- Kaixuan Huang, Yuqing Wang, Molei Tao, and Tuo Zhao. Why do deep residual networks generalize better than deep feedforward networks?—a neural tangent kernel perspective. *arXiv preprint arXiv:2002.06262*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Li Jing, Yichen Shen, Tena Dubcek, John Peurifoy, Scott Skirlo, Yann LeCun, Max Tegmark, and Marin Soljačić. Tunable efficient unitary neural networks (eunn) and their application to rnns. In *International Conference on Machine Learning*, pp. 1733–1741, 2017.
- Quoc V Le, Navdeep Jaitly, and Geoffrey E Hinton. A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941*, 2015.
- Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=B1EA-M-0Z>.

- Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in Neural Information Processing Systems*, pp. 8572–8583, 2019.
- Jaehoon Lee, Samuel S Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *arXiv preprint arXiv:2007.15801*, 2020.
- Radford M Neal. *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto, 1995.
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=BygfgHAcYX>.
- Roman Novak, Yasaman Bahri, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=HJC2SszZCW>.
- Roman Novak, Lechao Xiao, Yasaman Bahri, Jaehoon Lee, Greg Yang, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-dickstein. Bayesian deep convolutional networks with many channels are gaussian processes. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Blg30j0qF7>.
- Zichao Wang, Randall Balestriero, and Richard Baraniuk. A max-affine spline perspective of recurrent neural networks. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=BJej72AqF7>.
- Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019a.
- Greg Yang. Tensor programs I: Wide feedforward or recurrent neural networks of any architecture are gaussian processes. *arXiv preprint arXiv:1910.12478*, 2019b.
- Greg Yang. Tensor programs II: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020a.
- Greg Yang. Tensor programs III: Neural matrix laws. *arXiv preprint arXiv:2009.10685*, 2020b.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. *arXiv preprint arXiv:1811.08888*, 2018.

## A EXPERIMENT DETAILS

### A.1 TIME SERIES CLASSIFICATION

**Kernel methods settings.** We used RNTK, RBF, polynomial and NTK (Jacot et al., 2018). For data pre-processing, we normalized the norm of each  $\mathbf{x}$  to 1. For training we used C-SVM in LIBSVM library (Chang & Lin, 2011) and for hyperparameter selection we performed 10-fold validation for splitting the training data into 90% training set and 10% validation test. We then choose the best performing set of hyperparameters on all the validation sets, retrain the models with the best set of hyperparameters on the entire training data and finally report the performance on the unseen test data. The performance of all kernels on each data set is shown in table 2.

For C-SVM we chose the cost function value

$$C \in \{0.01, 0.1, 1, 10, 100\}g$$

and for each kernel we used the following hyperparameter sets

- RNTK: We only used single layer RNTK, we  $\phi = \text{ReLU}$  and the following hyperparameter sets for the variances:

$$\begin{aligned} \sigma_w &\in \{1.34, 1.35, 1.36, 1.37, 1.38, 1.39, 1.40, 1.41, 1.42, \sqrt{2}, 1.43, 1.44, 1.45, 1.46, 1.47\}g \\ \sigma_u &= 1 \\ \sigma_b &\in \{0, 0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.7, 0.9, 1, 2\}g \\ \sigma_h &\in \{0, 0.01, 0.1, 0.5, 1\}g \end{aligned}$$

- NTK: The formula for NTK of  $L$ -layer MLP (Jacot et al., 2018) for  $\mathbf{x}, \mathbf{x}^\ell \in \mathbb{R}^m$  is:

$$\begin{aligned} \sigma_w^{(l)} &= \frac{\sigma_w^2}{m} \|\mathbf{h}\mathbf{x}, \mathbf{x}^\ell\| + \sigma_b^2 \\ K^{(l)}(\mathbf{x}, \mathbf{x}^\ell) &= \sigma_w^{(l)2} \mathbb{V} [K^{(l)}(\mathbf{x}, \mathbf{x}^\ell)] + \sigma_b^2 & \ell \in [L] \\ K^{(l-1)}(\mathbf{x}, \mathbf{x}^\ell) &= \sigma_w^{(l-1)2} \mathbb{V} [K^{(l-1)}(\mathbf{x}, \mathbf{x}^\ell)] & \ell \in [L] \\ \mathbf{K}^{(l)}(\mathbf{x}, \mathbf{x}^\ell) &= \begin{bmatrix} K^{(l-1)}(\mathbf{x}, \mathbf{x}) & K^{(l-1)}(\mathbf{x}, \mathbf{x}^\ell) \\ K^{(l-1)}(\mathbf{x}, \mathbf{x}^\ell) & K^{(l-1)}(\mathbf{x}^\ell, \mathbf{x}^\ell) \end{bmatrix} \\ K(\mathbf{x}, \mathbf{x}^\ell) &= \sigma_w^2 \mathbb{V} [K^{(L+1)}(\mathbf{x}, \mathbf{x}^\ell)] \\ k_{\text{NTK}} &= \sum_{l=1}^L \left( K^{(l)}(\mathbf{x}, \mathbf{x}^\ell) \prod_{l'=0}^{l-1} K^{(l')}(\mathbf{x}, \mathbf{x}^\ell) \right) + K(\mathbf{x}, \mathbf{x}^\ell) \end{aligned}$$

and we used the following hyperparameters

$$\begin{aligned} L &\in [10] \\ \sigma_w &\in \{0.5, 1, \sqrt{2}, 2, 2.5, 3\}g \\ \sigma_b &\in \{0, 0.01, 0.1, 0.2, 0.5, 0.8, 1, 2, 5\}g \end{aligned}$$

- RBF:

$$\begin{aligned} k_{\text{RBF}}(\mathbf{x}, \mathbf{x}^\ell) &= e^{-\alpha \|\mathbf{x} - \mathbf{x}^\ell\|_2^2} \\ \alpha &\in \{0.01, 0.05, 0.1, 0.2, 0.5, 0.6, 0.7, 0.8, 1, 2, 3, 4, 5, 10, 20, 30, 40, 100\}g \end{aligned}$$

- Polynomial:

$$\begin{aligned} k_{\text{Polynomial}}(\mathbf{x}, \mathbf{x}^\ell) &= (r + \|\mathbf{x} - \mathbf{x}^\ell\|_2)^d \\ d &\in [5] \\ r &\in \{0, 0.1, 0.2, 0.5, 1, 2\}g \end{aligned}$$

Table 2: Test accuracy of each model on 56 time series data set from UCR time-series classification data repository (Dau et al., 2019).

Dataset	RNTK	NTK	RBF	POLY	Gaussian RNN	Identity RNN	GRU
Strawberry	<b>98.38</b>	97.57	97.03	96.76	94.32	75.4	91.62
ProximalPhalanxOutlineCorrect	<b>89</b>	87.97	87.29	86.94	82.81	74.57	86.94
PowerCons	97.22	97.22	96.67	91.67	96.11	95	<b>99.44</b>
Ham	70.48	<b>71.63</b>	66.67	71.43	53.33	60	60.95
SmallKitchenAppliances	67.47	38.4	40.27	37.87	60.22	<b>76</b>	71.46
ScreenType	41.6	43.2	<b>43.47</b>	38.4	40	41.06	36.26
MiddlePhalanxOutlineCorrect	57.14	57.14	48.7	64.29	<b>76.28</b>	57.04	74.57
RefrigerationDevices	46.93	37.07	36.53	41.07	36	<b>50.93</b>	46.66
Yoga	<b>84.93</b>	84.63	84.63	84.87	46.43	76.66	61.83
Computers	<b>59.2</b>	55.2	58.8	56.4	53.2	55.2	58.8
ECG5000	93.76	<b>94.04</b>	93.69	93.96	88.4	93.15	93.26
Fish	<b>90.29</b>	84	85.71	88	28	38.28	24
UWaveGestureLibraryX	<b>79.59</b>	78.7	78.48	65.83	55.97	75.34	73.64
UWaveGestureLibraryY	<b>71.56</b>	70.63	70.35	70.32	44.5	65.18	65.38
UWaveGestureLibraryZ	<b>73.95</b>	73.87	72.89	71.94	43.29	67.81	70.32
StarLightCurves	95.94	<b>96.19</b>	94.62	94.44	82.13	86.81	96.15
CricketX	60.51	59.49	62.05	62.56	8.46	<b>63.58</b>	26.41
CricketY	<b>63.85</b>	58.97	60.51	59.74	15.89	59.23	36.15
CricketZ	<b>60.26</b>	59.23	62.05	59.23	8.46	57.94	41.28
DistalPhalanxOutlineCorrect	<b>77.54</b>	<b>77.54</b>	75.36	73.91	69.92	69.56	75
Worms	<b>57.14</b>	50.65	55.84	50.65	35.06	49.35	41.55
SyntheticControl	98.67	96.67	98	97.67	92.66	97.66	<b>99</b>
Herring	56.65	<b>59.38</b>	<b>59.38</b>	<b>59.38</b>	23.28	<b>59.37</b>	<b>59.38</b>
MedicalImages	74.47	73.29	<b>75.26</b>	74.61	48.15	64.86	69.07
SwedishLeaf	90.56	91.04	<b>91.36</b>	90.72	59.2	45.92	91.04
ChlorineConcentration	90.76	77.27	86.35	<b>91.54</b>	65.99	55.75	61.14
SmoothSubspace	<b>96</b>	87.33	92	86.67	94	95.33	92.66
TwoPatterns	94.25	90.45	91.25	93.88	99.7	99.9	<b>100</b>
Faceall	74.14	<b>83.33</b>	83.25	82.43	53.66	70.53	70.65
DistalPhalanxTW	66.19	<b>69.78</b>	66.91	67.37	67.62	64.74	69.06
MiddlePhalanxTW	57.79	<b>61.04</b>	59.74	60.39	58.44	58.44	59.09
FacesUCR	81.66	80.2	80.34	<b>82.98</b>	53.21	75.26	79.46
OliveOil	<b>90</b>	86.67	86.67	83.33	66.66	40	40
UMD	91.67	92.36	97.22	90.97	44.44	71.52	<b>100</b>
nsectEPGRegular	99.6	99.2	99.6	96.79	<b>100</b>	<b>100</b>	98.39
Meat	<b>93.33</b>	<b>93.33</b>	<b>93.33</b>	<b>93.33</b>	0.55	55	33.33
Lightning2	<b>78.69</b>	73.77	70.49	68.85	45.9	70.49	67.21
Lightning7	61.64	60.27	63.01	60.27	23.28	69.86	<b>76.71</b>
Car	<b>83.33</b>	78.83	80	80	23.33	58.33	26.66
GunPoint	<b>98</b>	95.33	95.33	94	82	74.66	80.66
Arrowhead	80.57	<b>83.43</b>	80.57	74.86	48	56	37.71
Coffee	<b>100</b>	<b>100</b>	92.86	92.86	<b>100</b>	42.85	57.14
Trace	96	81	76	76	70	71	<b>100</b>
ECG200	<b>93</b>	89	89	86	86	72	76
plane	<b>98.1</b>	96.19	97.14	97.14	96.19	84.76	96.19
GunPointOldVersusYoung	<b>98.73</b>	97.46	<b>98.73</b>	94.6	53.96	52.38	98.41
GunPointMaleVersusFemale	99.05	<b>99.68</b>	99.37	<b>99.68</b>	68.67	52.53	99.68
GunPointAgeSpan	<b>96.52</b>	94.62	95.89	93.99	47.78	47.78	95.56
FreezerRegularTrain	<b>97.44</b>	94.35	96.46	96.84	76.07	7.5	86.59
SemgHandSubjectCh2	84.22	85.33	86.14	86.67	20	36.66	<b>89.11</b>
WormsTwoClass	<b>62.34</b>	<b>62.34</b>	61.04	59.74	51.94	46.75	57.14
Earthquakes	74.82	74.82	74.82	74.82	65.46	<b>76.97</b>	<b>76.97</b>
FiftyWords	68.57	68.57	<b>69.67</b>	68.79	34.28	60.21	65.27
Beef	90	73.33	83.33	<b>93.33</b>	26.67	46.67	36.67
Adiac	76.63	71.87	73.40	<b>77.75</b>	51.4	16.88	60.61
WordSynonyms	57.99	58.46	61.13	<b>62.07</b>	17.71	45.77	53.76

**Finite-width RNN settings.** We used 3 different RNNs. The first is a ReLU RNN with Gaussian initialization with the same NTK initialization scheme, where parameter variances are  $\sigma_w = \sigma_v = \frac{1}{2}$ .

$\sigma_u = 1$  and  $\sigma_b = 0$ . The second is a ReLU RNN with identity initialization following (Le et al., 2015). The third is a GRU (Cho et al., 2014) with uniform initialization. All models are trained with RMSProp algorithm for 200 epochs. Early stopping is implemented when the validation set accuracy does not improve for 5 consecutive epochs.

We perform standard 5-fold cross validation. For each RNN architecture we used hyperparameters of number of layer, number of hidden units and learning rate as

$$\begin{aligned} L &\in \{1, 2, g\} \\ n &\in \{50, 100, 200, 500\} \\ \eta &\in \{0.01, 0.001, 0.0001, 0.00001\} \end{aligned}$$

**Metrics descriptions** First, only in this paragraph, let  $i \in \{1, 2, \dots, Ng\}$  index a total of  $N$  datasets and  $j \in \{1, 2, \dots, Mg\}$  index a total of  $M$  classifiers. Let  $y_{ij}$  be the accuracy of the  $j$ -th classifier on the  $i$ -th dataset. We reported results on 4 metrics: average accuracy (Acc. mean), P90, P95, PMA and Friedman Rank. P90 and P95 is the fraction of datasets that the classifier achieves at least 90% and 95% of the maximum achievable accuracy for each dataset, i.e.,

$$P90_j = \frac{1}{N} \sum_i \mathbf{1}(y_{ij} \geq 0.9(\max_j y_{ij})). \quad (21)$$

PMA is the accuracy of the classifier on a dataset divided by the maximum achievable accuracy on that dataset, averaged over all datasets:

$$\text{PMA}_j = \frac{1}{N} \sum_i \frac{y_{ij}}{\max_j y_{ij}}. \quad (22)$$

Friedman Rank (Fernández-Delgado et al., 2019) first ranks the accuracy of each classifier on each dataset and then takes the average of the ranks for each classifier over all datasets, i.e.,

$$\text{FR}_j = \frac{1}{N} \sum_i r_{ij}, \quad (23)$$

where  $r_{ij}$  is the ranking of the  $j$ -th classifier on the  $i$ -th dataset.

Note that a better classifier achieves a lower Friedman Rank, Higher P90 and PMA.

**Remark.** In order to provide insight into the performance of RNTK in long time steps setting, we picked two datasets with more than 1000 time steps: SemgHandSubjectCh2 ( $T = 1024$ ) and StarLightCurves ( $T = 1024$ ).

## A.2 TIME SERIES REGRESSION

For time series regression, we used the 5-fold validation of training set and same hyperparameter sets for all kernels. For training we kernel ridge regression with ridge term chosen from

$$\lambda \in \{0, 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.8, 1, 2, 3, 4, 5, 6, 7, 8, 10, 100\}$$

## B PROOFS FOR THEOREMS 1 AND 3: RNTK CONVERGENCE AT INITIALIZATION

### B.1 PRELIMINARY: NETSOP PROGRAMS

Calculation of NTK in any architecture relies on finding the GP kernels that correspond to each pre-activation and gradient layers at initialization. For feedforward neural networks with  $n_1, \dots, n_L$  number of neurons (channels in CNNs) at each layer the form of this GP kernels can be calculated via taking the limit of  $n_1, \dots, n_L$  sequentially one by one. The proof is given by induction, where by conditioning on the previous layers, each entry of the current layer is sum of infinite i.i.d Gaussian random variables, and based on Central Limit Theorem (CLT), it becomes a Gaussian process with

kernel calculated based on the previous layers. Since the first layer is an affine transformation of input with Gaussian weights, it is a Gaussian process and the proof is completed. See (Lee et al., 2018; Duvenaud et al., 2014; Novak et al., 2019; Garriga-Alonso et al., 2019) for a formal treatment. However, due to weight-sharing, sequential limit is not possible and condoning on previous layers does not result in i.i.d. weights. Hence the aforementioned arguments break. To deal with it, in (Yang, 2019a) a proof using Gaussian conditioning trick (Bolthausen, 2014) is presented which allows use of recurrent weights in a network. More precisely, it has been demonstrated that neural networks (without batch normalization) can be expressed and a series of matrix multiplication and (piece wise) nonlinearity application, generally referred as *Netsor programs*. It has been shown that any architecture that can be expressed as *Netsor programs* that converge to GPs as width goes to infinity in the same rate, which a general rule to obtain the GP kernels. For completeness of this paper, we briefly restate the results from (Yang, 2019a) which we will use later for calculation derivation of RNTK.

There are 3 types of variables in *Netsor programs*; *A*-vars, *G*-vars and *H*-vars. *A*-vars are matrices and vectors with i.i.d Gaussian entries, *G*-vars are vectors introduced by multiplication of a vector by an *A*-var and *H*-vars are vectors after coordinate wise nonlinearities is applied to *G*-vars. Generally, *G*-vars can be thought of as pre-activation layers which are asymptotically treated as a Gaussian distributed vectors, *H*-vars as after-activation layers and *A*-vars are the weights. Since in neural networks inputs are immediately multiplied by a weight matrix, it can be thought of as an *G*-var, namely  $\mathbf{g}_{in}$ . Generally *Netsor programs* supports *G*-vars with different dimension, however the asymptotic behavior of a neural networks described by *Netsor programs* does not change under this degree of freedom, as long as they go to infinity at the same rate. For simplicity, let the *G*-vars and *H*-vars have the same dimension  $n$  since the network of interest is RNN and all pre-activation layers have the same dimension. We introduce the *Netsor programs* under this simplification. To produce the output of a neural network, *Netsor programs* receive a set of *G*-vars and *A*-vars as input, and new variables are produced sequentially using the three following operators:

- **Matmul** : multiplication of an *A*-var:  $\mathbf{A}$  with an *H*-var:  $\mathbf{h}$ , which produce a new *G*-var,  $\mathbf{g}$ .

$$\mathbf{g} = \mathbf{A}\mathbf{h} \quad (24)$$

- **Lincomp**: Linear combination of *G*-vars,  $\mathbf{g}^i, 1 \leq i \leq k$ , with coefficients  $a^i \in \mathbb{R}$  which produce of new *G*-var:

$$\mathbf{g} = \sum_{i=1}^k a^i \mathbf{g}^i \quad (25)$$

- **Nonlin**: creating a new *H*-var,  $\mathbf{h}$ , by using a nonlinear function  $\phi : \mathbb{R}^k \rightarrow \mathbb{R}$  that act coordinate wise on a set of *G*-vars,  $\mathbf{g}^i, 1 \leq i \leq k$ :

$$\mathbf{h} = \varphi(\mathbf{g}^1, \dots, \mathbf{g}^k) \quad (26)$$

Any output of the neural network  $y \in \mathbb{R}$  should be expressed as inner product of a new *A*-var which has not been used anywhere else in previous computations and an *H*-var:

$$y = \mathbf{v}^T \mathbf{h}$$

Any other output can be produced by another  $\mathbf{v}^0$  and  $\mathbf{h}^0$  (possibility the same  $\mathbf{h}$  or  $\mathbf{v}$ ).

It is assumed that each entry of any *A*-var:  $\mathbf{A} \in \mathbb{R}^{n \times n}$  in the *netsor programs* computations is drawn from  $\mathcal{N}(0, \frac{2}{n})$  and the input *G*-vars are Gaussian distributed. The collection of a specific entry of all *G*-vars of in the *netsor program* converges in probability to a Gaussian vector  $[\mathbf{g}^1]_i, \dots, [\mathbf{g}^k]_i$   $\mathcal{N}(\mu, \Sigma)$  for all  $i \in [n]$  as  $n$  goes to infinity.

Let  $\mu(\mathbf{g}) := \mathbb{E}[[\mathbf{g}]_i]$  be the mean of a *G*-var and  $\text{cov}(\mathbf{g}, \mathbf{g}^0) := \mathbb{E}[[\mathbf{g}]_i \cdot [\mathbf{g}^0]_i]$  be the covariance between any two *G*-vars. The general rule for  $\mu(\mathbf{g})$  is given by the following equations:

$$\mu(\mathbf{g}) = \begin{cases} \mu^{in}(\mathbf{g}) & \text{if } \mathbf{g} \text{ is input} \\ \sum_{i=1}^k a^i \mu(\mathbf{g}^i) & \text{if } \mathbf{g} = \sum_{i=1}^k a^i \mathbf{g}^i \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

For  $\mathbf{g}$  and  $\mathbf{g}^\theta$ , let  $G = \{g^1, \dots, g^r\}$  be the set of  $G$ -vars that has been introduced *before*  $\mathbf{g}$  and  $\mathbf{g}^\theta$  with distribution  $N(\mathbf{0}, \Sigma_G)$ , where  $\Sigma_G \in \mathbb{R}^{|G| \times |G|}$  containing the pairwise covariances between the  $G$ -vars.  $\kappa(\mathbf{g}, \mathbf{g}^\theta)$  is calculated via the following rules:

$$\kappa(\mathbf{g}, \mathbf{g}^\theta) = \begin{cases} \kappa^{\text{in}}(\mathbf{g}, \mathbf{g}^\theta) & \text{if } \mathbf{g} \text{ and } \mathbf{g}^\theta \text{ are inputs} \\ \sum_{i=1}^k a^i \kappa(\mathbf{g}^i, \mathbf{g}^\theta) & \text{if } \mathbf{g} = \sum_{i=1}^k a^i \mathbf{g}^i \\ \sum_{i=1}^k a^i \kappa(\mathbf{g}, \mathbf{g}^i) & \text{if } \mathbf{g}^\theta = \sum_{i=1}^k a^i \mathbf{g}^i \\ \sigma_A^2 \mathbb{E}_{\mathbf{z} \sim N(\mathbf{0}, \Sigma_G)} [\varphi(\mathbf{z})\varphi(\mathbf{z})] & \text{if } \mathbf{g} = \mathbf{A}\mathbf{h} \text{ and } \mathbf{g}^\theta = \mathbf{A}\mathbf{h}^\theta \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

Where  $\mathbf{h} = \varphi(\mathbf{g}^1, \dots, \mathbf{g}^r)$  and  $\mathbf{h}^\theta = \varphi(\mathbf{g}^1, \dots, \mathbf{g}^r)$  are functions of  $G$ -vars in  $G$  from possibly different nonlinearities. This set of rules presents a recursive method for calculating the GP kernels in a network where the recursive formula starts from data dependent quantities  $\kappa^{\text{in}}$  and  $\mu^{\text{in}}$  which are given.

All the above results holds when the nonlinearities are bounded uniformly by  $e^{(\alpha x^2)}$  for some  $\alpha > 0$  and when their derivatives exist.

**Standard vs. NTK initialization.** The common practice (which *netsor programs* uses) is to initialize DNNs weights  $[\mathbf{A}]_{i,j}$  with  $N(0, \frac{\sigma_w^2}{n})$  (known as *standard* initialization) where generally  $n$  is the number of units in the previous layer. In this paper we have used a different parameterization scheme as used in (Jacot et al., 2018) and we factor the standard deviation as shown in 3 and initialize weights with standard standard Gaussian. This approach does not change the the forward computation of DNN, but normalizes the backward computation (when computing the gradients) by factor  $\frac{1}{n}$ , otherwise RNTK will be scales by  $n$ . However this problem can be solved by scaling the step size by  $\frac{1}{n}$  and there is no difference between *NTK* and *standard* initialization (Lee et al., 2019).

## B.2 PROOF FOR THEOREM 1: SINGLE LAYER CASE

We first derive the RNTK in a simpler setting, i.e., a single layer and single output RNN. We then generalize the results to multi-layer and multi-output RNNs. We drop the layer index  $\ell$  to simplify notation. From 3 and 4, the forward pass for computing the output under NTK initialization for each input  $\mathbf{x} = \{x_t, g_{t=1}^T\}$  is given by:

$$\mathbf{g}^{(t)}(\mathbf{x}) = \frac{\sigma_w}{m} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}) + \frac{\sigma_u}{n} \mathbf{U} \mathbf{x}_t + \sigma_b \mathbf{b} \quad (29)$$

$$\mathbf{h}^{(t)}(\mathbf{x}) = \phi(\mathbf{g}^{(t)}(\mathbf{x})) \quad (30)$$

$$f(\mathbf{x}) = \frac{\sigma_v}{n} \mathbf{v}^\top \mathbf{h}^{(T)}(\mathbf{x}) \quad (31)$$

Note that (29), (30) and (31) use all the introduced operators introduced in 24, 25 and 26 given input variables  $\mathbf{W}$ ,  $\mathbf{U}$ ,  $\mathbf{x}_t$ ,  $\mathbf{b}$ ,  $\mathbf{v}$  and  $\mathbf{h}^{(0)}(\mathbf{x})$ .

First, we compute the kernels of forward pass  $\kappa^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta)$  and backward pass  $\mu^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta)$  introduced in (6) and (7) for two input  $\mathbf{x}$  and  $\mathbf{x}^\theta$ . Note that based on (27) the mean of all variables is zero since the inputs are all zero mean. In the forward pass for the intermediate layers we have:

$$\begin{aligned} \kappa^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) &= \kappa(\mathbf{g}^{(t)}(\mathbf{x}), \mathbf{g}^{(t)}(\mathbf{x}^\theta)) \\ &= \left( \frac{\sigma_w}{n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}) + \frac{\sigma_u}{m} \mathbf{U} \mathbf{x}_t + \sigma_b \mathbf{b}, \frac{\sigma_w}{n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}^\theta) + \frac{\sigma_u}{m} \mathbf{U} \mathbf{x}_t^\theta + \sigma_b \mathbf{b} \right) \\ &= \left( \frac{\sigma_w}{n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}), \frac{\sigma_w}{n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}^\theta) \right) + \kappa^{\text{in}} \left( \frac{\sigma_u}{m} \mathbf{U} \mathbf{x}_t, \frac{\sigma_u}{m} \mathbf{U} \mathbf{x}_t^\theta \right) + \kappa^{\text{in}}(\sigma_b \mathbf{b}, \sigma_b \mathbf{b}). \end{aligned}$$



We have used the second and third rule in (28) to expand the formula, We have also used the first and fifth rule to set the cross term to zero, i.e.,

$$\begin{aligned} \left( \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}), \frac{\sigma_u}{\rho_n} \mathbf{U} \mathbf{x}_{t^\theta}^\theta \right) &= 0 \\ \left( \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}), \sigma_b \mathbf{b} \right) &= 0 \\ \left( \frac{\sigma_u}{\rho_m} \mathbf{U} \mathbf{x}_t, \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t^\theta-1)}(\mathbf{x}^\theta) \right) &= 0 \\ \left( \sigma_b \mathbf{b}, \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t^\theta-1)}(\mathbf{x}^\theta) \right) &= 0 \\ \text{in } \left( \frac{\sigma_u}{\rho_m} \mathbf{U} \mathbf{x}_t, \sigma_b \mathbf{b} \right) &= 0 \\ \text{in } \left( \sigma_b \mathbf{b}, \frac{\sigma_u}{\rho_m} \mathbf{U} \mathbf{x}_{t^\theta}^\theta \right) &= 0. \end{aligned}$$

For the non-zero terms we have

$$\begin{aligned} \text{in } (\sigma_b \mathbf{b}, \sigma_b \mathbf{b}) &= \sigma_b^2 \\ \text{in } \left( \frac{\sigma_u}{\rho_m} \mathbf{U} \mathbf{x}_t, \frac{\sigma_u}{\rho_m} \mathbf{U} \mathbf{x}_{t^\theta}^\theta \right) &= \frac{\sigma_u^2}{m} h_{\mathbf{x}_t, \mathbf{x}_{t^\theta}^\theta}^i, \end{aligned}$$

which can be achieved by straight forward computation. If  $t \neq 1$  and  $t^\theta \neq 1$ , by using the forth rule in (28) we have

$$\left( \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t-1)}(\mathbf{x}), \frac{\sigma_w}{\rho_n} \mathbf{W} \mathbf{h}^{(t^\theta-1)}(\mathbf{x}^\theta) \right) = \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim N(0; \mathbf{K}^{(t, t^\theta)}(\mathbf{x}, \mathbf{x}^\theta))} [\phi(\mathbf{z}_1) \phi(\mathbf{z}_2)] = \mathbb{V} [\mathbf{K}^{(t, t^\theta)}(\mathbf{x}, \mathbf{x}^\theta)].$$

With  $\mathbf{K}^{(t, t^\theta)}(\mathbf{x}, \mathbf{x}^\theta)$  defined in (16). Otherwise, it will be zero by the fifth rule (if  $t$  or  $t^\theta = 1$ ).

Here the set of previously introduced  $G$ -vars is  $G = \{f\mathbf{g}^{(t-1)}(\mathbf{x}), \mathbf{U} \mathbf{x}_{g^t=1}, f\mathbf{g}^{(t^\theta-1)}(\mathbf{x}^\theta), \mathbf{U} \mathbf{x}_{g^{t^\theta}=1}^\theta, \mathbf{h}^{(0)}(\mathbf{x}), \mathbf{h}^{(0)}(\mathbf{x}^\theta)\}$ , but the dependency is only on the last layer  $G$ -vars,  $\varphi(f\mathbf{g} : \mathbf{g} \geq G\mathbf{g}) = \phi(\mathbf{g}^{(t-1)}(\mathbf{x}))$ ,  $\varphi(f\mathbf{g} : \mathbf{g} \geq G\mathbf{g}) = \phi(\mathbf{g}^{(t^\theta-1)}(\mathbf{x}^\theta))$ , leading the calculation to the operator defined in (10). As a result

$$\mathbf{K}^{(t, t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) = \sigma_w^2 \mathbb{V} [\mathbf{K}^{(t, t^\theta)}(\mathbf{x}, \mathbf{x}^\theta)] + \frac{\sigma_u^2}{m} h_{\mathbf{x}_t, \mathbf{x}_{t^\theta}^\theta}^i + \sigma_b^2.$$

To complete the recursive formula, using the same procedure for the first layers we have

$$\begin{aligned} \mathbf{K}^{(1;1)}(\mathbf{x}, \mathbf{x}^\theta) &= \sigma_w^2 \sigma_h^2 1_{(x=x^\theta)} + \frac{\sigma_u^2}{m} h_{\mathbf{x}_1, \mathbf{x}_1}^i + \sigma_b^2, \\ \mathbf{K}^{(1; t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) &= \frac{\sigma_u^2}{m} h_{\mathbf{x}_1, \mathbf{x}_{t^\theta}^\theta}^i + \sigma_b^2, \\ \mathbf{K}^{(t;1)}(\mathbf{x}, \mathbf{x}^\theta) &= \frac{\sigma_u^2}{m} h_{\mathbf{x}_t, \mathbf{x}_1}^i + \sigma_b^2. \end{aligned}$$

The output GP kernel is calculated via

$$K(\mathbf{x}, \mathbf{x}^\theta) = \sigma_v^2 \mathbb{V} [\mathbf{K}^{(T+1; T^\theta+1)}(\mathbf{x}, \mathbf{x}^\theta)]$$

The calculation of the gradient vectors  $\mathbf{g}^{(t)}(\mathbf{x}) = \rho_n^{-1} (r_{g^{(t)}(\mathbf{x})} f(\mathbf{x}))$  in the backward pass is given by

$$\begin{aligned} \mathbf{g}^{(T)}(\mathbf{x}) &= \sigma_v \mathbf{V} \phi^\theta(\mathbf{g}^{(T)}(\mathbf{x})) \\ \mathbf{g}^{(t)}(\mathbf{x}) &= \frac{\sigma_w}{\rho_n} \mathbf{W}^\top \left( \phi^\theta(\mathbf{g}^{(t)}(\mathbf{x})) \quad \mathbf{g}^{(t+1)}(\mathbf{x}) \right) \quad t \geq [T-1] \end{aligned}$$

To calculate the backward pass kernels, we rely on the following Corollary from (Yang, 2020b)

**Corollary 1** *In infinitely wide neural networks weights used in calculation of back propagation gradients ( $\mathbf{W}^>$ ) is an i.i.d copy of weights used in forward propagation ( $\mathbf{W}$ ) as long as the last layer weight ( $\mathbf{v}$ ) is sampled independently from other parameters and has mean 0.*

The immediate result of Corollary 1 is that  $\mathbf{g}^{(t)}(\mathbf{x})$  and  $\mathbf{x}^{(t)}$  are two independent Gaussian vector as their covariance is zero based on the fifth rule in (28). Using this result, we have:

$$\begin{aligned} \text{cov}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^0) &= \text{cov}(\mathbf{x}^{(t)}, \mathbf{x}^{(t^0)}) \\ &= \mathbb{E} \left[ [\mathbf{x}^{(t)}]_i \cdot [\mathbf{x}^{(t^0)}]_i \right] \\ &= \sigma_w^2 \mathbb{E} \left[ [\phi^\theta(\mathbf{g}^{(t)}(\mathbf{x}))]_i \cdot [\mathbf{x}^{(t+1)}]_i \cdot [\phi^\theta(\mathbf{g}^{(t^0)}(\mathbf{x}^0))]_i \cdot [\mathbf{x}^{(t^0+1)}(\mathbf{x}^0)]_i \right] \\ &= \sigma_w^2 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(0; \mathbf{K}^{(t+1:t^0+1)}(\mathbf{x}, \mathbf{x}^0))} [\phi^\theta(\mathbf{z}_1) \cdot \phi^\theta(\mathbf{z}_2)] \cdot \mathbb{E} \left[ [\mathbf{x}^{(t+1)}]_i \cdot [\mathbf{x}^{(t^0+1)}(\mathbf{x}^0)]_i \right] \\ &= \sigma_w^2 \mathbb{V}_o[\mathbf{K}^{(t+1:t^0+1)}(\mathbf{x}, \mathbf{x}^0)] \text{cov}^{(t+1:t^0+1)}(\mathbf{x}, \mathbf{x}^0). \end{aligned}$$

If  $T^0 = t^0 = T = t$ , then the the formula will lead to

$$\begin{aligned} \text{cov}^{(T;T^0)}(\mathbf{x}, \mathbf{x}^0) &= \mathbb{E} \left[ [\mathbf{x}^{(T)}]_i \cdot [\mathbf{x}^{(T^0)}]_i \right] \\ &= \sigma_v^2 \mathbb{E} \left[ [\mathbf{v}]_i \cdot [\phi^\theta(\mathbf{g}^{(T)}(\mathbf{x}))]_i \cdot [\mathbf{v}]_i \cdot [\phi^\theta(\mathbf{g}^{(T^0)}(\mathbf{x}^0))]_i \right] \\ &= \mathbb{E} \left[ [\phi^\theta(\mathbf{g}^{(T)}(\mathbf{x}))]_i \cdot [\phi^\theta(\mathbf{g}^{(T^0)}(\mathbf{x}^0))]_i \right] \cdot \mathbb{E} [[\mathbf{v}]_i [\mathbf{v}]_i] \\ &= \sigma_v^2 \mathbb{V}_o[\mathbf{K}^{(T+1:T+1)}(\mathbf{x}, \mathbf{x}^0)]. \end{aligned}$$

Otherwise it will end to either of two cases for some  $t^0 < T$  or  $T^0 = T = \tau$  and by the fifth rule in (28) we have:

$$\begin{aligned} \text{cov}^{(t^0;t^0)}(\mathbf{x}, \mathbf{x}^0) &= \left( \frac{\sigma_w}{n} \mathbf{W}^> \left( \phi^\theta(\mathbf{g}^{(t^0)}(\mathbf{x})) \quad \mathbf{x}^{(t^0+1)}(\mathbf{x}^0) \right), \mathbf{v} \cdot \phi^\theta(\mathbf{g}^{(T^0)}(\mathbf{x})) \right) = 0 \\ \text{cov}^{(T;t^0)}(\mathbf{x}, \mathbf{x}^0) &= \left( \mathbf{v} \cdot \phi^\theta(\mathbf{g}^{(T)}(\mathbf{x})), \frac{\sigma_w}{n} \mathbf{W}^> \left( \phi^\theta(\mathbf{g}^{(t^0)}(\mathbf{x}^0)) \quad \mathbf{x}^{(t^0+1)}(\mathbf{x}^0) \right) \right) = 0. \end{aligned}$$

Without loss of generality, from now on assume  $T^0 < T$  and  $T^0 = T = \tau$ , the final formula for computing the backward gradients becomes:

$$\begin{aligned} \text{cov}^{(T;T^0)}(\mathbf{x}, \mathbf{x}^0) &= \sigma_v^2 \mathbb{V}_o[\mathbf{K}^{(T+1:T+1)}(\mathbf{x}, \mathbf{x}^0)] \\ \text{cov}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^0) &= \sigma_w^2 \mathbb{V}_o[\mathbf{K}^{(t+1:t^0+1)}(\mathbf{x}, \mathbf{x}^0)] \text{cov}^{(t+1:t^0+1)}(\mathbf{x}, \mathbf{x}^0) \quad t \geq [T-1] \\ \text{cov}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^0) &= 0 \quad t^0 = t \notin \tau \end{aligned}$$

Now we have derived the single layer RNTK. Recall that  $\theta = \text{Vect}[f\mathbf{W}, \mathbf{U}, \mathbf{b}, \mathbf{v}g]$  contains all of the network's learnable parameters. As a result, we have:

$$r f(\mathbf{x}) = \text{Vect} \left[ f \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}}, \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}}, \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}}, \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}} g \right].$$

As a result

$$\begin{aligned} \text{cov} r f(\mathbf{x}), r f(\mathbf{x}^0) &= \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}}, \frac{\partial f(\mathbf{x}^0)}{\partial \mathbf{W}} \right\rangle + \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}}, \frac{\partial f(\mathbf{x}^0)}{\partial \mathbf{U}} \right\rangle + \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}}, \frac{\partial f(\mathbf{x}^0)}{\partial \mathbf{b}} \right\rangle \\ &\quad + \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^0)}{\partial \mathbf{v}} \right\rangle \end{aligned}$$

Where the gradients of output with respect to weights can be formulated as the following compact form:

$$\begin{aligned}\frac{\partial f(\mathbf{x})}{\partial \mathbf{W}} &= \sum_{t=1}^T \left( \frac{1}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}) \rangle \cdot \left( \frac{\sigma_w}{n} \mathbf{h}^{(t-1)}(\mathbf{x}) \right)^\top \\ \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}} &= \sum_{t=1}^T \left( \frac{1}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}) \rangle \cdot \left( \frac{\sigma_u}{m} \mathbf{x}_t \right)^\top \\ \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}} &= \sum_{t=1}^T \left( \frac{\sigma_b}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}) \rangle \right) \\ \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}} &= \frac{\sigma_v}{n} \mathbf{h}^{(T)}(\mathbf{x}).\end{aligned}$$

As a result we have:

$$\begin{aligned}\left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{W}} \right\rangle &= \sum_{t^0=1}^{T^0} \sum_{t=1}^T \left( \frac{1}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \rangle \right) \cdot \left( \frac{\sigma_w^2}{n} \langle \mathbf{h}^{(t-1)}(\mathbf{x}), \mathbf{h}^{(t-1)}(\mathbf{x}^\theta) \rangle \right) \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{U}} \right\rangle &= \sum_{t^0=1}^{T^0} \sum_{t=1}^T \left( \frac{1}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \rangle \right) \cdot \left( \frac{\sigma_u^2}{m} \langle \mathbf{x}_t, \mathbf{x}_{t^0}^\theta \rangle \right) \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{b}} \right\rangle &= \sum_{t^0=1}^{T^0} \sum_{t=1}^T \left( \frac{1}{n} \langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \rangle \right) \cdot \sigma_b^2 \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{v}} \right\rangle &= \left( \frac{\sigma_v^2}{n} \langle \mathbf{h}^{(T)}(\mathbf{x}), \mathbf{h}^{(T)}(\mathbf{x}^\theta) \rangle \right).\end{aligned}$$

Remember that for any two  $G$ -var  $E[\langle \mathbf{g} \rangle_i \langle \mathbf{g}^\theta \rangle_i]$  is independent of index  $i$ . Therefore,

$$\begin{aligned}\frac{1}{n} \langle \mathbf{h}^{(t-1)}(\mathbf{x}), \mathbf{h}^{(t-1)}(\mathbf{x}^\theta) \rangle &= \forall [\mathbf{K}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^\theta)] \quad t > 1 \\ \frac{1}{n} \langle \mathbf{h}^{(0)}(\mathbf{x}), \mathbf{h}^{(0)}(\mathbf{x}^\theta) \rangle &= \sigma_h^2.\end{aligned}$$

Hence, by summing the above terms in the infinite-width limit we get

$$\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \rangle = \left( \sum_{t^0=1}^{T^0} \sum_{t=1}^T \mathbf{K}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^\theta) \right) + K(\mathbf{x}, \mathbf{x}^\theta). \quad (32)$$

Since  $\mathbf{K}^{(t;t^0)}(\mathbf{x}, \mathbf{x}^\theta) = 0$  for  $t^0 \neq t$  it is simplified to

$$\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \rangle = \left( \sum_{t=1}^T \mathbf{K}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \right) + K(\mathbf{x}, \mathbf{x}^\theta).$$

**Multi-dimensional output.** For  $f(\mathbf{x}) \in \mathbb{R}^d$ , the  $i$ -th output for  $i \in [d]$  is obtained via

$$[f(\mathbf{x})]_i = \frac{\sigma_v}{n} \mathbf{v}_i^\top \mathbf{h}^{(T)}(\mathbf{x}),$$

where  $\mathbf{v}_i$  is independent of  $\mathbf{v}_j$  for  $i \neq j$ . As a result, for The RNTK  $\langle [f(\mathbf{x})]_i, [f(\mathbf{x}^\theta)]_i \rangle \in \mathbb{R}^d$  for multi-dimensional output we have

$$\langle [f(\mathbf{x})]_i, [f(\mathbf{x}^\theta)]_i \rangle = \left\langle \left( \sum_{t=1}^T \mathbf{K}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \right) + K(\mathbf{x}, \mathbf{x}^\theta), \mathbf{v}_i \right\rangle$$

For  $i = j$ , the kernel is the same as computed in (32) and we denote it as

$$\langle [f(\mathbf{x})]_i, [f(\mathbf{x}^\theta)]_i \rangle = \mathbf{K}^{(T;T)}(\mathbf{x}, \mathbf{x}^\theta).$$

For  $i \notin j$ , since  $v_i$  is independent of  $v_j$ ,  $\frac{\partial}{\partial \mathbf{x}^0} (T; T^0)(\mathbf{x}, \mathbf{x}^0)$  and all the backward pass gradients become zero, so

$$\left\langle r [f(\mathbf{x})]_i, r [f(\mathbf{x}^0)]_j \right\rangle = 0 \quad i \notin j$$

which gives us the following formula

$$\frac{\partial}{\partial \mathbf{x}^0} (T; T^0)(\mathbf{x}, \mathbf{x}^0) = \mathbf{I}_d.$$

This concludes the proof for Theorem 1 for single-layer case.

### B.3 PROOF FOR THEOREM 1: MULTI-LAYER CASE

Now we derive the RNTK for multi-layer RNTK. We will only study single output case and the generalization to multi-dimensional case is identical as the single layer case. The set of equations for calculation of the output of a  $L$ -layer RNN for  $\mathbf{x} = \{x_t\}_{t=1}^T$  are

$$\begin{aligned} \mathbf{g}^{(\cdot; t)}(\mathbf{x}) &= \frac{\sigma_w}{n} \mathbf{W}^{(\cdot)} \mathbf{h}^{(\cdot; t-1)}(\mathbf{x}) + \frac{\sigma_u}{n} \mathbf{U}^{(\cdot)} \mathbf{x}_t + \sigma_b \mathbf{b}^{(\cdot)} & \ell = 1 \\ \mathbf{g}^{(\cdot; t)}(\mathbf{x}) &= \frac{\sigma_w}{n} \mathbf{W}^{(\cdot)} \mathbf{h}^{(\cdot; t-1)}(\mathbf{x}) + \frac{\sigma_u}{n} \mathbf{U}^{(\cdot)} \mathbf{h}^{(\cdot-1; t)}(\mathbf{x}) + \sigma_b \mathbf{b}^{(\cdot)} & \ell > 1 \\ \mathbf{h}^{(\cdot; t)}(\mathbf{x}) &= \phi(\mathbf{g}^{(\cdot; t)}(\mathbf{x})) \\ f(\mathbf{x}) &= \frac{\sigma_v}{n} \mathbf{v}^{\triangleright} \mathbf{h}^{(L; T)}(\mathbf{x}) \end{aligned}$$

The forward pass kernels for the first layer is the same as calculated in B.2. For  $\ell \geq 2$  we have:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}^0} (\cdot; t^0)(\mathbf{x}, \mathbf{x}^0) &= \left( \mathbf{g}^{(\cdot; t)}(\mathbf{x}), \mathbf{g}^{(\cdot; t^0)}(\mathbf{x}^0) \right) \\ &= \left( \frac{\sigma_w}{n} \mathbf{W}^{(\cdot)} \mathbf{h}^{(\cdot; t-1)}(\mathbf{x}), \frac{\sigma_w}{n} \mathbf{W}^{(\cdot)} \mathbf{h}^{(\cdot; t^0-1)}(\mathbf{x}^0) \right) \\ &+ \left( \frac{\sigma_u}{n} \mathbf{U}^{(\cdot)} \mathbf{h}^{(\cdot-1; t)}(\mathbf{x}), \frac{\sigma_u}{n} \mathbf{U}^{(\cdot)} \mathbf{h}^{(\cdot-1; t^0)}(\mathbf{x}^0) \right) + \text{in}(\sigma_b \mathbf{b}^{(\cdot)}, \sigma_b \mathbf{b}^{(\cdot)}) \\ &= (\sigma_w)^2 \mathbf{V} [\mathbf{K}^{(\cdot; t; t^0)}(\mathbf{x}, \mathbf{x}^0)] + (\sigma_u)^2 \mathbf{V} [\mathbf{K}^{(\cdot-1; t+1; t^0+1)}(\mathbf{x}, \mathbf{x}^0)] + (\sigma_b)^2, \end{aligned}$$

where

$$\mathbf{K}^{(\cdot; t; t^0)}(\mathbf{x}, \mathbf{x}^0) = \begin{bmatrix} \Sigma^{(\cdot; t-1; t-1)}(\mathbf{x}, \mathbf{x}) & \Sigma^{(\cdot; t-1; t^0-1)}(\mathbf{x}, \mathbf{x}^0) \\ \Sigma^{(\cdot; t-1; t^0-1)}(\mathbf{x}, \mathbf{x}^0) & \Sigma^{(\cdot; t^0-1; t^0-1)}(\mathbf{x}^0, \mathbf{x}^0) \end{bmatrix},$$

and  $\text{in}$  is defined in (B.2). For the first first time step we have:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}^0} (\cdot; 1; 1)(\mathbf{x}, \mathbf{x}^0) &= (\sigma_w)^2 \sigma_h^2 1_{(x=x^0)} + (\sigma_u)^2 \mathbf{V} [\mathbf{K}^{(\cdot; 2; 2)}(\mathbf{x}, \mathbf{x}^0)] + (\sigma_b)^2, \\ \frac{\partial}{\partial \mathbf{x}^0} (\cdot; t; 1)(\mathbf{x}, \mathbf{x}^0) &= (\sigma_u)^2 \mathbf{V} [\mathbf{K}^{(\cdot; t+1; 2)}(\mathbf{x}, \mathbf{x}^0)] + (\sigma_b)^2, \\ \frac{\partial}{\partial \mathbf{x}^0} (\cdot; 1; t^0)(\mathbf{x}, \mathbf{x}^0) &= (\sigma_u)^2 \mathbf{V} [\mathbf{K}^{(\cdot; 2; t^0+1)}(\mathbf{x}, \mathbf{x}^0)] + (\sigma_b)^2. \end{aligned}$$

And the output layer

$$\frac{\partial}{\partial \mathbf{x}^0} K(\mathbf{x}, \mathbf{x}^0) = \sigma_v^2 \mathbf{V} [\mathbf{K}^{(L; T+1; T^0+1)}(\mathbf{x}, \mathbf{x}^0)].$$

Note that because of using new weights at each layer we get

$$(\mathbf{g}^{(\cdot; t)}(\mathbf{x}), \mathbf{g}^{(\cdot; t^0)}(\mathbf{x}^0)) = 0 \quad \ell \notin \ell^0$$

Now we calculate the backward pass kernels in multi-layer RNTK. The gradients at the last layer is calculated via

$$\frac{\partial}{\partial \mathbf{x}^0} (L; T)(\mathbf{x}) = \sigma_v \mathbf{v}^{\triangleright} \phi'(\mathbf{g}^{(L; T)}(\mathbf{x})).$$

In the last hidden layer for different time steps we have

$$\frac{\partial}{\partial \mathbf{x}^0} (\cdot; t)(\mathbf{x}) = \frac{\sigma_w^L}{n} \left( \mathbf{W}^{(L)} \right)^{\triangleright} \left( \phi'(\mathbf{g}^{(L; t)}(\mathbf{x})) \frac{\partial}{\partial \mathbf{x}^0} (\cdot; t+1)(\mathbf{x}) \right) \quad t \geq [T-1]$$

In the last time step for different hidden layers we have

$$\mathbf{h}^{(\ell; T)}(\mathbf{x}) = \frac{\sigma_u^{\ell+1}}{n} \left( \mathbf{U}^{(\ell+1)} \right)^{\top} \left( \phi^{\ell}(\mathbf{g}^{(\ell; T)}(\mathbf{x})) \quad \mathbf{h}^{(\ell+1; T)}(\mathbf{x}) \right) \quad \ell \geq [L-1]$$

At the end for the other layers we have

$$\begin{aligned} \mathbf{h}^{(\ell; t)}(\mathbf{x}) &= \frac{\sigma_w^{\ell}}{n} \left( \mathbf{W}^{(\ell)} \right)^{\top} \left( \phi^{\ell}(\mathbf{g}^{(\ell; t)}(\mathbf{x})) \quad \mathbf{h}^{(\ell+1; t)}(\mathbf{x}) \right) \\ &\quad + \frac{\sigma_u^{\ell+1}}{n} \left( \mathbf{U}^{(\ell+1)} \right)^{\top} \left( \phi^{\ell}(\mathbf{g}^{(\ell; t)}(\mathbf{x})) \quad \mathbf{h}^{(\ell+1; t)}(\mathbf{x}) \right) \quad \ell \geq [L-1], t \geq [T-1] \end{aligned}$$

The recursive formula for the  $\mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta})$  is the same as the single layer, and it is non-zero for  $t^{\theta} = T^{\theta} - T = \tau$ . As a result we have

$$\begin{aligned} \mathbf{h}^{(\ell; T; T+)}(\mathbf{x}, \mathbf{x}^{\theta}) &= \sigma_v^2 \mathbf{V}^{\circ} \left[ \mathbf{K}^{(\ell; T+1; T+)} \right](\mathbf{x}, \mathbf{x}^{\theta}) \\ \mathbf{h}^{(\ell; t; t+)}(\mathbf{x}, \mathbf{x}^{\theta}) &= (\sigma_w^{\ell})^2 \mathbf{V}^{\circ} \left[ \mathbf{K}^{(\ell; t+1; t+)} \right](\mathbf{x}, \mathbf{x}^{\theta}) \cdot \mathbf{h}^{(\ell+1; t+1+)}(\mathbf{x}, \mathbf{x}^{\theta}) \quad t \geq [T-1] \\ \mathbf{h}^{(\ell; t; t^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta}) &= 0 \quad t^{\theta} = t \notin \tau \end{aligned} \quad (33)$$

Similarly by using the same course of arguments used in the single layer setting, for the last time step we have

$$\mathbf{h}^{(\ell+1; T+)}(\mathbf{x}, \mathbf{x}^{\theta}) = (\sigma_u^{\ell+1})^2 \mathbf{V}^{\circ} \left[ \mathbf{K}^{(\ell+1; T+1; T+)} \right](\mathbf{x}, \mathbf{x}^{\theta}) \cdot \mathbf{h}^{(\ell+2; T+1; T+)}(\mathbf{x}, \mathbf{x}^{\theta}) \quad \ell \geq [L-1]$$

For the other layers we have

$$\begin{aligned} \mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta}) &= (\sigma_w^{\ell})^2 \mathbf{V}^{\circ} \left[ \mathbf{K}^{(\ell; t+1; t+)} \right](\mathbf{x}, \mathbf{x}^{\theta}) \cdot \mathbf{h}^{(\ell+1; t+1+)}(\mathbf{x}, \mathbf{x}^{\theta}) \\ &\quad + (\sigma_u^{\ell+1})^2 \mathbf{V}^{\circ} \left[ \mathbf{K}^{(\ell; t+1; t+)} \right](\mathbf{x}, \mathbf{x}^{\theta}) \cdot \mathbf{h}^{(\ell+2; t+1; t^{\theta}+1)}(\mathbf{x}, \mathbf{x}^{\theta}). \end{aligned}$$

For  $t^{\theta} = t \notin \tau$  the recursion continues until it reaches  $\mathbf{h}^{(\ell; T; t^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta}), t^{\theta} < T^{\theta}$  or  $\mathbf{h}^{(\ell; t^{\theta}; T^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta}), t^{\theta} < T$  and as a result based on (33) we get

$$\mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta}) = 0 \quad t^{\theta} = t \notin \tau \quad (34)$$

For  $t^{\theta} = t = \tau$  it leads to  $\mathbf{h}^{(\ell; T; T^{\theta})}(\mathbf{x}, \mathbf{x}^{\theta})$  and has a non-zero value.

Now we derive RNTK for multi-layer:

$$\begin{aligned} \text{Tr} \left[ \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}^{(\ell)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{W}^{(\ell)}} \right] &= \sum_{i=1}^L \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}^{(i)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{W}^{(i)}} \right\rangle + \sum_{i=1}^L \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}^{(i)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{U}^{(i)}} \right\rangle \\ &\quad + \sum_{i=1}^L \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}^{(i)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{b}^{(i)}} \right\rangle + \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{v}} \right\rangle, \end{aligned}$$

where

$$\begin{aligned} \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}^{(\ell)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{W}^{(\ell)}} \right\rangle &= \sum_{t^{\theta}=1}^{T^{\theta}} \sum_{t=1}^T \left( \frac{1}{n} \left\langle \mathbf{h}^{(\ell; t)}(\mathbf{x}), \mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \cdot \left( \frac{(\sigma_w^{\ell})^2}{n} \left\langle \mathbf{h}^{(\ell+1; t)}(\mathbf{x}), \mathbf{h}^{(\ell+1; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}^{(\ell)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{U}^{(\ell)}} \right\rangle &= \sum_{t^{\theta}=1}^{T^{\theta}} \sum_{t=1}^T \left( \frac{1}{n} \left\langle \mathbf{h}^{(\ell; t)}(\mathbf{x}), \mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \cdot \left( \frac{(\sigma_u^{\ell})^2}{m} \langle \mathbf{h}_{\mathbf{x}_t, \mathbf{x}_{t^{\theta}}}^{\ell} \rangle \right) \quad \ell = 1 \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}^{(\ell)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{U}^{(\ell)}} \right\rangle &= \sum_{t^{\theta}=1}^{T^{\theta}} \sum_{t=1}^T \left[ \left( \frac{1}{n} \left\langle \mathbf{h}^{(\ell; t)}(\mathbf{x}), \mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \right. \\ &\quad \left. \cdot \left( \frac{(\sigma_u^{\ell})^2}{n} \left\langle \mathbf{h}^{(\ell+1; t)}(\mathbf{x}), \mathbf{h}^{(\ell+1; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \right] \quad \ell > 1 \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}^{(\ell)}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{b}^{(\ell)}} \right\rangle &= \sum_{t^{\theta}=1}^{T^{\theta}} \sum_{t=1}^T \left( \frac{1}{n} \left\langle \mathbf{h}^{(\ell; t)}(\mathbf{x}), \mathbf{h}^{(\ell; t^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \cdot (\sigma_b)^2 \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^{\theta})}{\partial \mathbf{v}} \right\rangle &= \left( \frac{\sigma_v^2}{n} \left\langle \mathbf{h}^{(T)}(\mathbf{x}), \mathbf{h}^{(T^{\theta})}(\mathbf{x}^{\theta}) \right\rangle \right) \end{aligned}$$

Summing up all the terms and replacing the inner product of vectors with their expectations we get

$$\langle \mathbb{E} f(\mathbf{x}), \mathbb{E} f(\mathbf{x}^\theta) \rangle = \sum_{i=1}^L \sum_{t=1}^T \sum_{t^\theta=1}^{T^\theta} \langle \mathbb{E}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta), \mathbb{E}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \rangle + K(\mathbf{x}, \mathbf{x}^\theta).$$

By (34), we can simplify to

$$\langle \mathbb{E}^{(L;T;T^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \rangle = \left( \sum_{i=1}^L \sum_{t=1}^T \langle \mathbb{E}^{(i;t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta), \mathbb{E}^{(i;t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \rangle \right) + K(\mathbf{x}, \mathbf{x}^\theta).$$

For multi-dimensional output it becomes

$$\langle \mathbb{E} \mathbf{f}(\mathbf{x}), \mathbb{E} \mathbf{f}(\mathbf{x}^\theta) \rangle = \langle \mathbb{E}^{(L;T;T^\theta)}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}^\theta)) \rangle \mathbf{I}_d.$$

This concludes the proof for Theorem 1 for the multi-layer case.

#### B.4 PROOF FOR THEOREM 3: WEIGHT-UNTIED RNTK

The architecture of a weight-untied single layer RNN is

$$\begin{aligned} \mathbf{g}^{(t)}(\mathbf{x}) &= \frac{\sigma_w}{m} \mathbf{W}^{(t)} \mathbf{h}^{(t-1)}(\mathbf{x}) + \frac{\sigma_u}{n} \mathbf{U}^{(t)} \mathbf{x}_t + \sigma_b \mathbf{b}^{(t)} \\ \mathbf{h}^{(t)}(\mathbf{x}) &= \phi(\mathbf{g}^{(t)}(\mathbf{x})) \\ \mathbf{f}(\mathbf{x}) &= \frac{\sigma_v}{n} \mathbf{V} \mathbf{h}^{(T)}(\mathbf{x}) \end{aligned}$$

Where we use new weights at each time step and we index it by time. Like previous sections, we first derive the forward pass kernels for two same length data  $\mathbf{x} = \mathbb{E} \mathbf{x}_t \mathbb{E} \mathbf{x}_{t^\theta}^T, \mathbf{x}^\theta = \mathbb{E} \mathbf{x}_{t^\theta} \mathbb{E} \mathbf{x}_t^T$

$$\begin{aligned} \langle \mathbb{E}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \rangle &= \sigma_w^2 \mathbb{V}[\mathbf{K}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta)] + \frac{\sigma_u^2}{m} \langle \mathbf{h}_{\mathbf{x}_t, \mathbf{x}_{t^\theta}^\theta} \rangle + \sigma_b^2. \\ \langle \mathbb{E}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \rangle &= 0 \quad t \neq t^\theta \end{aligned}$$

Since we are using same weight at the same time step,  $\langle \mathbb{E}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \rangle$  can be written as a function of the previous kernel, which is exactly as the weight-tied RNN. However for different length, it becomes zero as a consequence of using different weights, unlike weight-tied which has non-zero value. The kernel of the first time step and output is also the same as weight-tied RNN. For the gradients we have:

$$\begin{aligned} \langle \mathbb{E}^{(T)}(\mathbf{x}) \rangle &= \sigma_v \mathbf{V} \phi^\theta(\mathbf{g}^{(T)}(\mathbf{x})) \\ \langle \mathbb{E}^{(t)}(\mathbf{x}) \rangle &= \frac{\sigma_w}{n} \langle \mathbf{W}^{(t+1)} \rangle \langle \phi^\theta(\mathbf{g}^{(t)}(\mathbf{x})) \rangle \langle \mathbb{E}^{(t+1)}(\mathbf{x}) \rangle \quad t \geq [T-1] \end{aligned}$$

For  $t^\theta = t$  we have:

$$\begin{aligned} \langle \mathbb{E}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \rangle &= \sigma_w^2 \mathbb{V}[\langle \mathbf{K}^{(t+1;t+1)}(\mathbf{x}, \mathbf{x}^\theta) \rangle \langle \mathbb{E}^{(t+1;t+1)}(\mathbf{x}, \mathbf{x}^\theta) \rangle] \\ \langle \mathbb{E}^{(t;t)}(\mathbf{x}, \mathbf{x}^\theta) \rangle &= \sigma_v^2 \mathbb{V}[\langle \mathbf{K}^{(T+1;T+1)}(\mathbf{x}, \mathbf{x}^\theta) \rangle]. \end{aligned}$$

Due to using different weights for  $t \neq t^\theta$ , we can immediately conclude that  $\langle \mathbb{E}^{(t;t^\theta)}(\mathbf{x}, \mathbf{x}^\theta) \rangle = 0$ . This set of calculation is exactly the same as the weight-tied case when  $\tau = T - T = 0$ .

Finally, with  $\theta = \text{Vect}[\mathbb{E} \mathbf{W}^{(t)}, \mathbf{U}^{(t)}, \mathbf{b}^{(t)} \mathbb{E} \mathbf{g}_{t=1}^T, \mathbf{v} \mathbb{g}]$  we have

$$\begin{aligned} \langle \mathbb{E} f(\mathbf{x}), \mathbb{E} f(\mathbf{x}^\theta) \rangle &= \sum_{t=1}^T \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{W}^{(t)}} \right\rangle + \sum_{t=1}^T \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{U}^{(t)}} \right\rangle \\ &\quad + \sum_{t=1}^T \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{b}^{(t)}} \right\rangle + \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{v}} \right\rangle \end{aligned}$$

with

$$\begin{aligned} \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{W}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{W}^{(t)}} \right\rangle &= \left( \frac{1}{n} \left\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \right\rangle \right) \cdot \left( \frac{\sigma_w^2}{n} \left\langle \mathbf{h}^{(t-1)}(\mathbf{x}), \mathbf{h}^{(t-1)}(\mathbf{x}^\theta) \right\rangle \right) \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{U}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{U}^{(t)}} \right\rangle &= \left( \frac{1}{n} \left\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \right\rangle \right) \cdot \left( \frac{\sigma_u^2}{m} \sum_i \mathbf{h}_{\mathbf{x}_t, \mathbf{x}_t^i} \right) \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{b}^{(t)}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{b}^{(t)}} \right\rangle &= \left( \frac{1}{n} \left\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \right\rangle \right) \cdot \sigma_b^2 \\ \left\langle \frac{\partial f(\mathbf{x})}{\partial \mathbf{v}}, \frac{\partial f(\mathbf{x}^\theta)}{\partial \mathbf{v}} \right\rangle &= \left( \frac{\sigma_v^2}{n} \left\langle \mathbf{h}^{(\tau)}(\mathbf{x}), \mathbf{h}^{(\tau)}(\mathbf{x}^\theta) \right\rangle \right). \end{aligned}$$

As a result we obtain

$$\langle \nabla f(\mathbf{x}), \nabla f(\mathbf{x}^\theta) \rangle = \left( \sum_{t=1}^{\tau} \left\langle \mathbf{h}^{(t)}(\mathbf{x}), \mathbf{h}^{(t)}(\mathbf{x}^\theta) \right\rangle \right) + \mathcal{K}(\mathbf{x}, \mathbf{x}^\theta),$$

same as the weight-tied RNN when  $\tau = 0$ . This concludes the proof for Theorem 3.

### B.5 ANALYTICAL FORMULA FOR $\mathbb{V}[\mathbf{K}]$

For any positive definite matrix  $\mathbf{K} = \begin{bmatrix} K_1 & K_3 \\ K_3 & K_2 \end{bmatrix}$  we have:

- $\phi = \text{ReLU}$  (Cho & Saul, 2009)

$$\mathbb{V}[\mathbf{K}] = \frac{1}{2\pi} \left( c(\pi - \arccos(c)) + \sqrt{1 - c^2} \right) \sqrt{K_1 K_2},$$

$$\mathbb{V}_o[\mathbf{K}] = \frac{1}{2\pi} (\pi - \arccos(c)).$$

where  $c = K_3 / \sqrt{K_1 K_2}$

- $\phi = \text{erf}$  (Neal, 1995)

$$\mathbb{V}[\mathbf{K}] = \frac{2}{\pi} \arcsin \left( \frac{2K_3}{\sqrt{(1+2K_1)(1+2K_3)}} \right),$$

$$\mathbb{V}_o[\mathbf{K}] = \frac{4}{\pi \sqrt{(1+2K_1)(1+2K_2)} - 4K_3^2}.$$

## C PROOF FOR THEOREM 2: RNTK CONVERGENCE AFTER TRAINING

To prove theorem 2, we use the strategy used in (Lee et al., 2019) which relies on the the local lipschitzness of the network Jacobian  $\mathbf{J}(\theta, X) = \nabla f(\mathbf{x}) \in \mathbb{R}^{j \times d}$  at initialization.

**Definition 1** The Jacobian of a neural network is local lipschitz at NTK initialization  $(\theta_0 \sim \mathcal{N}(0, 1))$  if there is constant  $K > 0$  for every  $C$  such that

$$\begin{cases} \|\mathbf{J}(\theta, X) - \mathbf{J}(\theta_0, X)\|_F < K \\ \|\mathbf{J}(\theta, X)\|_F < K \|\theta - \theta_0\|, \end{cases} \quad \forall \theta, \theta_0 \in B(\theta_0, R)$$

where

$$B(\theta, R) := \{\theta : \|\theta - \theta_0\| < R\}.$$

**Theorem 4** Assume that the network Jacobian is local lipschitz with high probability and the empirical NTK of the network converges in probability at initialization and it is positive definite over the input set. For  $\epsilon > 0$ , there exists  $N$  such that for  $n > N$  when applying gradient flow with  $\eta < 2(\lambda_{\min}(\hat{\mathbf{K}}(X, X)) + \lambda_{\max}(\hat{\mathbf{K}}(X, X)))^{-1}$  with probability at least  $(1 - \epsilon)$  we have:

$$\sup_s \frac{k\theta_s}{n}, \sup_s \frac{\theta_0 k_2}{n}, \sup_s \hat{k}_s(X, X) - \hat{k}_0(X, X) = \mathcal{O}\left(\frac{1}{n}\right).$$

Proof: See (Lee et al., 2019)

Theorem 4 holds for any network architecture and any cost function and it was used in (Lee et al., 2019) to show the stability of NTK for MLP during training.

Here we extend the results for RNTK by proving that the Jacobian of a multi-layer RNN under NTK initialization is local lipschitz with high probability.

To prove it, first, we prove that for any two points  $\theta, \theta' \in B(\theta_0, R)$  there exists constant  $K_1$  such that

$$\|k\mathbf{g}^{(\cdot;t)}(\mathbf{x})\|_2, \|k\delta^{(\cdot;t)}(\mathbf{x})\|_2 \leq K_1 \rho_{-n} \quad (35)$$

$$\|k\mathbf{g}^{(\cdot;t)}(\mathbf{x}) - \mathbf{g}^{(\cdot;t)}(\mathbf{x})\|_2, \|k\delta^{(\cdot;t)}(\mathbf{x}) - \delta^{(\cdot;t)}(\mathbf{x})\|_2 \leq k\theta - \theta' \leq K_1 \rho_{-n} k\theta - \theta' k. \quad (36)$$

To prove (35) and (36) we use the following lemmas.<sup>3</sup>

**Lemma 1** Let  $A \in \mathbb{R}^{n \times m}$  be a random matrix whose entries are independent standard normal random variables. Then for every  $t > 0$ , with probability at least  $1 - e^{-ct^2}$  for some constant  $c$  we have:

$$\|A\|_2 \leq \sqrt{\frac{\rho_{-m}}{m} + \frac{\rho_{-n}}{n}} + t.$$

**Lemma 2** Let  $a \in \mathbb{R}^n$  be a random vector whose entries are independent standard normal random variables. Then for every  $t > 0$ , with probability at least  $1 - e^{-ct^2}$  for some constant  $c$  we have:

$$\|a\|_2 \leq \sqrt{\frac{\rho_{-n}}{n} + \frac{\rho_{-t}}{t}}.$$

Setting  $t = \rho_{-n}$  for any  $\theta \in B(\theta_0, R)$ . With high probability, we get:

$$\|k\mathbf{W}^{(\cdot)}\|_2, \|k\mathbf{U}^{(\cdot)}\|_2 \leq 3\sqrt{\frac{\rho_{-n}}{n}}, \|k\mathbf{b}^{(\cdot)}\|_2 \leq 2\sqrt{\frac{\rho_{-n}}{n}}, \|k\mathbf{h}^{(\cdot;0)}(\mathbf{x})\|_2 \leq 2\sigma_h \sqrt{\frac{\rho_{-n}}{n}}.$$

We also assume that there exists some finite constant  $C$  such that

$$|\phi(x)| \leq C|x|, \quad |\phi(x) - \phi(x^0)| \leq C|x - x^0|, \quad |\phi^0(x)| \leq C, \quad |\phi^0(x) - \phi^0(x^0)| \leq C|x - x^0|.$$

The proof is obtained by induction. From now on assume that all inequalities in (35) and (36) holds with some  $k$  for the previous layers. We have

$$\begin{aligned} \|k\mathbf{g}^{(\cdot;t)}(\mathbf{x})\|_2 &= k\sqrt{\frac{\sigma_w}{n}} \|\mathbf{W}^{(\cdot)} \mathbf{h}^{(\cdot;t-1)}(\mathbf{x}) + \sqrt{\frac{\sigma_u}{n}} \|\mathbf{U}^{(\cdot)} \mathbf{h}^{(\cdot;t-1)}(\mathbf{x}) + \sigma_b \mathbf{b}^{(\cdot)}\|_2 \\ &\leq \sqrt{\frac{\sigma_w}{n}} \|k\mathbf{W}^{(\cdot)}\|_2 k\phi(\|\mathbf{g}^{(\cdot;t-1)}(\mathbf{x})\|_2) + \sqrt{\frac{\sigma_u}{n}} \|k\mathbf{U}^{(\cdot)}\|_2 k\phi(\|\mathbf{g}^{(\cdot;t-1)}(\mathbf{x})\|_2) + \sigma_b \|k\mathbf{b}^{(\cdot)}\|_2 \\ &\leq (3\sigma_w C k + 3\sigma_u C k + 2\sigma_b) \rho_{-n}. \end{aligned}$$

And the proof for (35) and (36) is completed by showing that the first layer is bounded

$$\begin{aligned} \|k\mathbf{g}^{(1;1)}(\mathbf{x})\|_2 &= k\sqrt{\frac{\sigma_w^1}{n}} \|\mathbf{W}^{(1)} \mathbf{h}^{(1;0)}(\mathbf{x}) + \sqrt{\frac{\sigma_u^1}{m}} \|\mathbf{U}^{(1)} \mathbf{x}_1 + \sigma_b^1 \mathbf{b}^{(1)}\|_2 \\ &\leq (3\sigma_w^1 \sigma_h + \sqrt{\frac{3\sigma_u^1}{m}} \|k\mathbf{x}_1\|_2 + 2\sigma_b^1) \rho_{-n}. \end{aligned}$$

For the gradient of first layer we have

$$\begin{aligned} \|k^{(L;T)}(\mathbf{x})\|_2 &= k\sigma_v \|\mathbf{v} \cdot \phi^0(\mathbf{g}^{(L;T)}(\mathbf{x}))\|_2 \\ &\leq \sigma_v \|k\mathbf{v}\|_2 k\phi^0(\|\mathbf{g}^{(L;T)}(\mathbf{x})\|_2) k_1 \\ &= 2\sigma_v C \rho_{-n}. \end{aligned}$$

And similarly we have

$$\|k^{(\cdot;t)}(\mathbf{x})\|_2 \leq (3\sigma_w C k^0 + 3\sigma_u C k^0) \rho_{-n}.$$

<sup>3</sup>See [math.uci.edu/~rvershyn/papers/HDP-book/HDP-book.pdf](http://math.uci.edu/~rvershyn/papers/HDP-book/HDP-book.pdf) for proofs



For  $\theta, \theta \geq B(\theta_0, R)$  we have

$$\begin{aligned} \|\mathbf{g}^{(1:1)}(\mathbf{x}) - \mathbf{g}^{(1:1)}(\mathbf{x})\|_2 &= \left\| \frac{\sigma_w^1}{n} (\mathbf{W}^{(1)} - \mathbf{W}^{(1)}) \mathbf{h}^{(1:0)}(\mathbf{x}) + \frac{\sigma_u^1}{m} (\mathbf{U}^{(1)} - \mathbf{U}^{(1)}) \mathbf{h}^{(1:0)}(\mathbf{x}) \right\|_2 \\ &\leq \left( 3\sigma_w^1 \sigma_h + \frac{3\sigma_u^1}{m} k_{\mathbf{x}_1} k_2 \right) k \theta \leq \theta k_2 \frac{\rho_-}{n}. \end{aligned}$$

$$\begin{aligned} \|\mathbf{g}^{(\cdot:t)}(\mathbf{x}) - \mathbf{g}^{(\cdot:t)}(\mathbf{x})\|_2 &\leq \|\phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x})) - \phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x}))\|_2 \frac{\sigma_w}{n} \|\mathbf{W}^{(\cdot)} - \mathbf{W}^{(\cdot)}\|_2 \\ &\quad + \|\frac{\sigma_w}{n} \mathbf{W}^{(\cdot)}\|_2 \|\phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x})) - \phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x}))\|_2 \\ &\quad + \|\phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x}))\|_2 \|\frac{\sigma_u}{n} (\mathbf{U}^{(\cdot)} - \mathbf{U}^{(\cdot)})\|_2 \\ &\quad + \|\frac{\sigma_u}{n} \mathbf{U}^{(\cdot)}\|_2 \|\phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x})) - \phi(\mathbf{g}^{(\cdot:t-1)}(\mathbf{x}))\|_2 + \sigma_b k \|\mathbf{b}^{(\cdot)} - \mathbf{b}^{(\cdot)}\|_2 \\ &\leq (k\sigma_w + 3\sigma_w C k + k\sigma_u + 3\sigma_u C k + \sigma_b) k \theta \leq \theta k_2 \frac{\rho_-}{n}. \end{aligned}$$

For gradients we have

$$\begin{aligned} \|\mathbf{g}^{(L:T)}(\mathbf{x}) - \mathbf{g}^{(L:T)}(\mathbf{x})\|_2 &\leq \sigma_v k \phi^0(\mathbf{g}^{(L:T)}(\mathbf{x})) \|\mathbf{v} - \mathbf{v}\|_2 + \sigma_v k \|\mathbf{v}\|_2 \|\phi^0(\mathbf{g}^{(L:T)}(\mathbf{x})) - \phi^0(\mathbf{g}^{(L:T)}(\mathbf{x}))\|_2 \\ &\leq (\sigma_v C + 2\sigma_v C k) k \theta \leq \theta k_2 \frac{\rho_-}{n}. \end{aligned}$$

And similarly using same techniques we have

$$\|\mathbf{g}^{(\cdot:t)}(\mathbf{x}) - \mathbf{g}^{(\cdot:t)}(\mathbf{x})\|_2 \leq (\sigma_w C + 3\sigma_w C k + \sigma_u C + 3\sigma_u C k) k \theta \leq \theta k_2 \frac{\rho_-}{n}.$$

As a result, there exists  $K_1$  that is a function of  $\sigma_w, \sigma_u, \sigma_b, L, T$  and the norm of the inputs.

Now we prove the local Lipchitzness of the Jacobian

$$\begin{aligned} \|\mathbf{J}(\theta, \mathbf{x})\|_F &\leq \sum_{l=2}^L \sum_{t=1}^T \left( \frac{1}{n} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \frac{\sigma_w}{n} \mathbf{h}^{(l:t-1)}(\mathbf{x}) \right) \right\|_F \right. \\ &\quad \left. + \frac{1}{n} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \frac{\sigma_u}{n} \mathbf{h}^{(l:t-1)}(\mathbf{x}) \right) \right\|_F + \frac{1}{n} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \sigma_b \mathbf{b}^{(l:t-1)}(\mathbf{x}) \right) \right\|_F \right) \\ &\quad + \sum_{t=1}^T \left( \frac{1}{n} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \frac{\sigma_w}{n} \mathbf{h}^{(1:t-1)}(\mathbf{x}) \right) \right\|_F \right. \\ &\quad \left. + \frac{1}{nm} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \frac{\sigma_u}{n} \mathbf{x}_t \right) \right\|_F + \frac{1}{n} \left\| \frac{\partial}{\partial \mathbf{x}} \left( \sigma_b \mathbf{b}^{(1:t-1)}(\mathbf{x}) \right) \right\|_F \right) + \frac{\sigma_v}{n} k \|\mathbf{h}^{(L:T)}(\mathbf{x})\|_F \\ &\quad \left( \sum_{l=2}^L \sum_{t=1}^T (K_1^2 C \sigma_w + K_1^2 C \sigma_u + \sigma_b K_1) \right. \\ &\quad \left. + \sum_{t=1}^T (K_1^2 C \sigma_w + \frac{K_1 \sigma_u}{m} k_{\mathbf{x}_t} k_2 + \sigma_b K_1) + \sigma_v C K_1 \right). \end{aligned}$$

And for  $\theta, \theta \in B(\theta_0, R)$  we have

$$\begin{aligned}
k\mathcal{J}(\theta, \mathbf{x}) - \mathcal{J}(\theta, \mathbf{x})k_F & \leq \sum_{l=2}^L \sum_{t=1}^T \left( \frac{1}{n} \left\| \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \left( \sigma_w^l \mathbf{h}^{(l;t-1)}(\mathbf{x}) \right)^{\otimes} - \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \left( \sigma_w^l \mathbf{h}^{(l;t-1)}(\mathbf{x}) \right)^{\otimes} \right\|_F \right. \\
& + \frac{1}{n} \left\| \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \left( \sigma_u^l \mathbf{h}^{(l;t-1)}(\mathbf{x}) \right)^{\otimes} - \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \left( \sigma_u^l \mathbf{h}^{(l;t-1)}(\mathbf{x}) \right)^{\otimes} \right\|_F \\
& + \frac{1}{n} \left\| \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \cdot \sigma_b^l - \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \cdot \sigma_b^l \right\|_F \\
& + \sum_{t=1}^T \left( \frac{1}{n} \left\| \langle \cdot \rangle^{(1;t)}(\mathbf{x}) \left( \sigma_w^1 \mathbf{h}^{(1;t-1)}(\mathbf{x}) \right)^{\otimes} - \langle \cdot \rangle^{(1;t)}(\mathbf{x}) \left( \sigma_w^1 \mathbf{h}^{(1;t-1)}(\mathbf{x}) \right)^{\otimes} \right\|_F \right. \\
& + \frac{1}{nm} \left\| \langle \cdot \rangle^{(1;t-1)}(\mathbf{x}) \left( \sigma_u^1 \mathbf{x}_t \right)^{\otimes} - \langle \cdot \rangle^{(1;t-1)}(\mathbf{x}) \left( \sigma_u^1 \mathbf{x}_t \right)^{\otimes} \right\|_F \\
& + \frac{1}{n} \left\| \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \cdot \sigma_b^l - \langle \cdot \rangle^{(l;t)}(\mathbf{x}) \cdot \sigma_b^l \right\|_F \Big) + \frac{\sigma_v}{n} k\mathbf{h}^{(L;T)}(\mathbf{x}) - \mathbf{h}^{(L;T)}(\mathbf{x})k_F \\
& \left( \sum_{l=2}^L \sum_{t=1}^T (4K_1^2 C \sigma_w^l + 4K_1^2 C \sigma_u^l + \sigma_b^l K_1) \right. \\
& \left. + \sum_{t=1}^T (4K_1^2 C \sigma_w^1 + \frac{K_1 \sigma_u^1}{m} k\mathbf{x}_t k_2 + \sigma_b^1 K_1) + \sigma_v C K_1 \right) k\theta - \theta k_2.
\end{aligned}$$

The above proof can be generalized to the entire dataset by a straightforward application of the union bound. This concludes the proof for Theorem 2.