UNDERSTANDING NEURAL TANGENT KERNEL DYNAM ICS THROUGH ITS TRACE EVOLUTION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027 028 029

030

Paper under double-blind review

Abstract

The Neural Tangent Kernel (NTK) has emerged as a valuable tool for analyzing the training and generalization properties of neural networks. While the behavior of the NTK in the infinite-width limit is well understood, a comprehensive investigation is still required to comprehend its dynamics during training in the finite-width regime. In this paper, we present a detailed exploration of the NTK's behavior through the examination of its trace during training.

By conducting experiments on standard supervised classification tasks, we observe that the NTK trace typically exhibits an increasing trend and stabilizes when the network achieves its highest accuracy on the training data. Additionally, we investigate the phenomenon of "grokking", which has recently garnered attention, as it involves an intriguing scenario where the test accuracy suddenly improves long after the training accuracy plateaus. To shed light on this phenomenon, we employ the NTK trace to monitor the training dynamics during grokking. Furthermore, we utilize the NTK trace to gain insights into the training dynamics of semisupervised learning approaches, including the employment of exponential moving average mechanisms. Through these investigations, we demonstrate that the NTK, particularly its trace, remains a powerful and valuable tool for comprehending the training dynamics of modern finite-width neural networks.

1 INTRODUCTION

Deep neural networks have proven to be highly successful in a wide range of machine learning tasks, particularly in the domain of image classification. Understanding the dynamics of these networks during training and their generalization capabilities is of paramount importance to further advance their performance and interpretability. In recent years, the field of deep learning has seen the rise of a crucial analytical tool known as the Neural Tangent Kernel (NTK) (Jacot et al., 2018), primarily in the infinite-width limit under suitable initialization scales and infinitesimal learning rate. Empirically, researchers adopt neural networks of finite-width, therefore the behavior of the NTK in the finite-width regime necessitates a more comprehensive investigation.

Notably, in the infinite width limit, Jacot et al. (2018) show that when the network depth is fixed 040 the NTK is determined by the initialization and is fixed throughout the whole training process. 041 Theoretically, Hanin & Nica (2019); Littwin et al. (2021) show that the NTK at initialization is also 042 influenced by the network depth even in the infinite width limit. Recently, Fort et al. (2020); Loo et al. 043 (2022) empirically study the NTK of finite width neural network during the whole training process. 044 By defining a "kernel distance" function, they show that after a rather quick chaotic data-dependent kernel learning process, the network learns a kernel that has the linearized training accuracy that matches the performance of the neural network. However, the defined "kernel distance" $S(K_1, K_2)$ 046 (Fort et al., 2020; Loo et al., 2022) is not a distance function in mathematics as $S(K_1, K_2) = 0$ does 047 not imply $K_1 = K_2$. But interestingly, one can show that if in addition $Tr(K_1) = Tr(K_2)$ then 048 $K_1 = K_2$. Moreover, the NTK trace has an efficient approximation method, making it a feasible and orthogonal approach to prior works (Fort et al., 2020; Loo et al., 2022). 050

To understand the behavior of the NTK in finite-width neural networks, we conduct experiments that
 evaluate the NTK trace along the training process on 3 settings: standard supervised classification
 tasks, grokking, and semi-supervised learning. The latter 2 settings grokking and semi-supervised
 learning are closely related to standard supervised classification, which we will briefly discuss as

follows. Grokking is first observed in (Power et al., 2022), where a *supervised* classification task
is trained to perform arithmetic tasks. It refers to a sudden and unexpected great improvement in
test accuracy long after the training accuracy has reached a plateau. Semi-supervised learning uses
a small amount of labeled data and a large amount of unlabeled data in training. It usually adopts
the pseudo-labeling strategy, where during training the network will assign pseudo labels to some
unlabeled data and thus transform it into a *supervised* learning problem.

- Our contributions can be summarized as follows:
 - 1. We present an efficient approximation for NTK trace and link it to the margin of a kernel SVM problem.
 - 2. On standard supervised classification tasks, we find the NTK trace usually rises and stabilizes when top training accuracy is met. We also observe that in the grokking settings, when the NTK trace is stabilized the test accuracy attains its maximum.
 - 3. The NTK trace also helps us understand training dynamics in semi-supervised learning approaches, including the use of the exponential moving average (EMA) mechanism.
- 068 069 070

071

083

095 096

097

062

063

064

065

066

067

2 RELATED WORK

072 **Neural Tangent Kernel.** Jacot et al. (2018) provide neural tangent kernel (NTK) as a tool to 073 understand the dynamics of neural networks training under suitable initialization and small learning 074 rates by considering the infinite-width limit. In this limiting regime, the NTK becomes deterministic 075 and fixed during training when the width is taken to infinity (Zou & Gu, 2019; Ji & Telgarsky, 2019; 076 Chen et al., 2019). As modern neural networks are in the finite-width regime, researchers have 077 gained more and more interest in the behavior of NTK in modern finite-width architectures (Shan & 078 Bordelon, 2021; Atanasov et al., 2021; Seleznova et al., 2024; Hanin & Nica, 2019; Littwin et al., 2021; Fort et al., 2020; Loo et al., 2022). For example, Hanin & Nica (2019); Littwin et al. (2021) 079 have given finite-width (and depth) correction of neural network's NTK at initialization. And Fort et al. (2020); Loo et al. (2022) study the evolution of NTK during modern finite-width neural network 081 training using a "kernel distance". In this paper, we study dynamics of NTK using its trace.

Neural network training dynamics. Researchers have found that there has usually been some 084 implicit optimization process during the supervised training dynamics of neural networks (Nacson 085 et al., 2019; Xu et al., 2021; Soudry et al., 2018; Banburski et al., 2019; Lyu & Li, 2019; Blanc et al., 2020). For example, Lyu & Li (2019) show the network implicitly maximizes the margin during 087 the training using sgd under loss like cross-entropy loss and Blanc et al. (2020) study the implicit 088 regularization of the network under MSE loss using sgd. Very recently, Power et al. (2022) find cases 089 in supervised learning that the network generalizes long after it overfits the training data, which they termed the phenomenon grokking. There have been attempts to understand the training dynamics 091 of grokking: Liu et al. (2022a) study the dynamics using physics-inspired effective theory, Nanda et al. (2023) study the process using trigonometric series, and Tan & Huang (2023) study the process 092 using the network's robustness. In this paper, we study the dynamics of neural networks and their 093 associated NTK using the trace of NTK. 094

3 PRELIMINARY

We will first introduce some basic knowledge to better understand NTK. Suppose the training dataset as $\mathcal{X} = \{(x_i, y_i)\}_{i=1}^N$ with scalar labels $\{y_i\}_{i=1,...,N} \subset \mathbb{R}$ and a loss function $l : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$. Then the empirical training loss, defined on functions $f : \mathbb{R}^{n_1} \to \mathbb{R}$, is given by

$$l(f) = \sum_{i=1}^{N} l(f(x_i), y_i).$$

105 When the network $f(\cdot; w) : \mathbb{R}^{n_1} \to \mathbb{R}$ is trained by minimizing the loss l(f) via gradient flow, 106 then the network parameters $(w(t))_{t\geq 0}$ will change according to the following ordinary differential 107 equation:

$$\partial_t w(t) = -\nabla_w l(f(\cdot; w(t))). \tag{1}$$

Therefore during training, the output function of the neural network for any $x \in \mathbb{R}^{n_1}$ follows another differential equation given in terms of the Neural Tangent Kernel (NTK) (Lee et al., 2019):

$$\partial_t f(x; w(t)) = -\sum_{i=1}^N K(x, x_i; w) \,\partial_z l(z, y_i) \bigg|_{z=f(x_i; w(t))}$$

where $K(\cdot, \cdot; w)$ is the Neural Tangent Kernel defined as follows:

$$K(x,y;w) = \sum_{k=1}^{d} \partial_{w_k} f(x;w) \partial_{w_k} f(y;w) = \langle \nabla f(x;w), \nabla f(y;w) \rangle.$$

The above equation shows that the NTK determines the training dynamics of $f(\cdot; w(t))$ in the function spaces $\mathbb{R}^{n_1} \to \mathbb{R}$ during gradient flow training.

For a neural network $f(x;w) \in \mathbb{R}^{K}$ with vector outputs (logits) with parameter $w \in \mathbb{R}^{d}$, we will define the NTK as follows. Its Jacobian matrix can be computed as $J(x,w) \in \mathbb{R}^{K \times d}$. Suppose the dataset as $\mathcal{X} = \{(x_i, y_i)\}_{i=1}^{N}$ and denote $J(\mathcal{X}, w) = [J^{\top}(x_1, w), \cdots, J^{\top}(x_N, w)]^{\top} \in \mathbb{R}^{NK \times d}$. Then the empirical neural tangent kernel on \mathcal{X} can be computed as $K(\mathcal{X}, \mathcal{X}; w) = J(\mathcal{X}, w)J^{\top}(\mathcal{X}, w)$. All the proofs of this paper's theorems can be found in Appendix A.

4 EVOLUTION OF NTK (TRACE)

We are interested in evaluating the dynamics of NTK during network training, as the NTK is a matrix, we intend to evaluate its trace (a scalar) $\frac{\text{Tr}(K(\mathcal{X},\mathcal{X};w))}{|\mathcal{X}|K|}$ during training to reflects its dynamic.

Jacot et al. (2018) show that NTK is constant along training in the infinite-width limit. As we are interested in the dynamics of finite-width NTK, we will abbreviate each time-step's $K(\mathcal{X}, \mathcal{X}; w(t))$ as $K_t(\mathcal{X}, \mathcal{X})$. Note that directly calculating the NTK matrix for the whole dataset is computational and memory inefficient, we will provide the following theorem 4.1 that can efficiently approximate the tendency of the trace of NTK.

138 Theorem 4.1.
$$\mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0, I_d)} \| f(\mathcal{X}; w + \Delta) - f(\mathcal{X}; w) \|_F^2 \sim \epsilon^2 \operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; w)).$$

140 We can immediately obtain an approximate calculation of $\operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; w))$ as $\operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; w)) \approx \frac{1}{\epsilon^2} \|f(\mathcal{X}; w + \Delta) - f(\mathcal{X}; w)\|_F^2$, where $\Delta \sim \epsilon \mathcal{N}(0, I_d)$. We plot the evolution of NTK traces on 3 142 widely adopted benchmarks in Figure 1 with the backbone ResNet18 and sgd with momentum. More 143 training details and effect of different learning rate, momentum, and initialization can be found in 144 Appendix B. We can see that the trend of NTK trace is similar among datasets as the trace continues to 145 increase until the training accuracy reaches its maximum. This shows that the NTK matrix stabilizes 146 until training reaches the optimal point.



156 157

111 112

113 114

127 128

129

139

147

148

149

150 151 152

153

154

158

Figure 1: Evolution of NTK trace on different datasets using ResNet18.

In this paper, we mainly consider using cross-entropy loss to train a classification problem, it may have different outcomes when facing other settings. When given a dataset $\mathcal{X} = \{(x_i, y_i)\}_{i=1}^N$, the loss is $\frac{1}{N} \sum_{i=1}^N H(\text{onehot}(y_i), \text{prob}(x_i))$, where we use $H(p, q) = -p \log q$ to denote the cross-entropy loss between two probability distributions p and q, and the generated probability of a sample x under 162 model f(x; w) is defined as $\operatorname{prob}(x; w) == \left[\frac{e^{f_1(x;w)}}{\sum_{j=1}^K e^{f_j(x;w)}}, \cdots, \frac{e^{f_K(x;w)}}{\sum_{j=1}^K e^{f_j(x;w)}}\right]^\top$, which we usually abbreviate it as $\operatorname{prob}(x)$.

Another question remains, what property of the neural network does the NTK trace value Tr($K_t(\mathcal{X}, \mathcal{X})$)) reflect during training? We will use the setting in (Lyu & Li, 2019) to gain some insights into the above question. We will summarize the key technical assumptions in (Lyu & Li, 2019) that will be useful in this paper, for more detailed discussions of other assumptions, please refer to the initial paper (Lyu & Li, 2019).

- 1. The problem is a binary classification problem, i.e. $y_i \in \{+1, -1\}$.
- 2. For each sample (x_i, y_i) , the loss is the logistic loss (binary cross-entropy loss), i.e. $l((x_i, y_i)) = \log(1 + e^{-y_i f(x_i; w)}).$
- 173 174

182 183

185 186

170

171 172

Then we have the following theorem which connects the dynamics of the neural network to the solution of an NTK kernel SVM.

Theorem 4.2 ((Lyu & Li, 2019)). Under some technical assumptions and assuming the network is trained by gradient flow, then any limit point \bar{w} of $\{\frac{w(t)}{\|w(t)\|}|t > 0\}$ is along the max-margin direction for a hard-margin SVM with kernel $K(x, y) = \langle \nabla f(x; \bar{w}), \nabla f(y; \bar{w}) \rangle$.

181 This means that for some $\beta > 0$, $\beta \overline{w}$ is the optimal solution to the following optimization problem:

$$\min \quad \frac{1}{2} \|w\|^2$$

s.t. $y_i \langle w, \nabla f(x_i; \bar{w}) \rangle \ge 1 \quad i = 1, 2, \cdots, N$ (2)

Considering the dual of optimization problem (2), we will obtain an optimization problem as follows:

$$\max \sum_{i=1}^{N} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(x_{i}, x_{j}; \bar{w})$$

s.t. $\alpha \ge 0$ (3)

Motivated by (Seleznova et al., 2024), we will make the following simplification assumptions:

197

199

200 201 202

203

Collect the dual variables in α which corresponds to $y_i = +1$ as α^1 and $y_i = -1$ as α^2 . The optimization problem (3) simplifies as follows:

 $\max -\frac{x}{2} (\sum_{i} \alpha_{i}^{1})^{2} + (\sum_{i} \alpha_{i}^{1}) + -\frac{x}{2} (\sum_{i} \alpha_{i}^{2})^{2} + (\sum_{i} \alpha_{i}^{2})$ s.t. $\alpha^{1}, \alpha^{2} \ge 0$ (4)

It is easy to see the optimal objective value for the optimization problem (4) is $\frac{1}{x}$. By noticing that Under the assumption "If $y_i = y_j$, then $K(x_i, x_j; \bar{w}) = x$ ", we can deduce that Nx = $Tr(K(\mathcal{X}, \mathcal{X}; \bar{w}))$. Therefore, we find that $1/\frac{Tr(K(\mathcal{X}, \mathcal{X}; \bar{w}))}{N}$ is a good approximate optimal objective value for (2) by utilizing the duality of (2) and (3). Therefore the results in Figure 1 show that the network is implicitly maximizing the margin of kernel SVM during training.

Notably, (Fort et al., 2020; Loo et al., 2022) also study the dynamics of NTK in the finite-width neural network through the following "kernel distance" between two kernel matrices: $S(K_1, K_2) = 1 - \frac{\operatorname{Tr}(K_1^\top K_2)}{\|K_1\|_F \|K_2\|_F}$.

213 (Fort et al., 2020; Loo et al., 2022) define the kernel velocity as $v(t) = \frac{S(K(\mathcal{X},\mathcal{X};w(t)),K(\mathcal{X},\mathcal{X};w(t+dt)))}{dt}$.

If the total training time is T_1 , (Fort et al., 2020; Loo et al., 2022) find the following facts:

216 1. There exists a time T_0 , where v(t) = 0 ($\forall t > T_0$).

218

225

226 227 228

229

230

231

232

233

235

237 238

239 240

241 242 243

244

245 246

247

2.
$$S(K(\mathcal{X}, \mathcal{X}; w(T_0)), K(\mathcal{X}, \mathcal{X}; w(T_1))) = 0.$$

From the above facts, one may find that $S(K_{t+1}(\mathcal{X}, \mathcal{X}), K_t(\mathcal{X}, \mathcal{X}))) = 0$ $(t \ge T_0)$. However, the defined "kernel distance" is not a distance function in mathematics as $S(K_1, K_2) = 0$ does not imply $K_1 = K_2$. Note we are interested in the dynamics of $K_t(\mathcal{X}, \mathcal{X}))$, the following theorem 4.3 shows that $\text{Tr}(K_t(\mathcal{X}, \mathcal{X})))$ together with the "kernel distance" can characterize whether the NTK stabilize or not. Figure 2 together with Figure 1 further validates that the "kernel distance" and NTK trace give complementary characterizations of the dynamics of NTK.

Theorem 4.3. Suppose $S(K_{t+1}(\mathcal{X}, \mathcal{X}), K_t(\mathcal{X}, \mathcal{X}))) = 0$ $(t \ge T_0)$. Suppose $T > T_0$, then $K_i(\mathcal{X}, \mathcal{X}) = K_j(\mathcal{X}, \mathcal{X})$ $(i, j \ge T)$ iff $\operatorname{Tr}(K_i(\mathcal{X}, \mathcal{X})) = \operatorname{Tr}(K_j(\mathcal{X}, \mathcal{X}))$ $(i, j \ge T)$.



Figure 2: Accuracy and NTK "kernel distance".

5 UNDERSTANDING DYNAMICS OF OTHER TRAINING SCENARIOS USING NTK TRACE

5.1 LAYERWISE DYNAMICS DURING TRAINING

After investigating the behavior of the NTK trace in the previous section, we now focus on understanding the behavior of inner layers in the neural network. To do this, we depict the architecture of the network in Figure 3. We specifically examine the outputs of three layers: the front layer (Layer 1 in Figure 3), which is close to the input picture, the middle layer (Layer 3 in Figure 3), and the final layer, whose output is the logit. For simplicity, we denote these layers as f_1 , f_2 , and f_3 , respectively.

In addition to calculating the NTK trace for the logits, as described in theorem 4.1, we also compute $\frac{\mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} |f(\mathcal{X};w+\Delta) - f(\mathcal{X};w)|_F^2}{\epsilon^2 |\mathcal{X}|K_i} \quad (i = 1, 2), \text{ where } K_i \text{ is the dimension of the layer output. We present}$ the results in Figure 4.

From Figure 4, we observe that each layer exhibits a distinct evaluation pattern. The last layer stabilizes the latest, while the middle layer stabilizes the fastest. This discrepancy suggests that different layers play different roles in the neural network. As the last layer stabilizes slowly, it is crucial to track the NTK trace of the logits to gain a comprehensive understanding of the entire training dynamics.

262 263 5.2 THE DYNAMICS OF GROKKING

In this section, we consider a special setting in supervised classification, called grokking (Liu et al., 2022b; Nanda et al., 2023). It refers to a strange phenomenon that the test accuracy suddenly increases
long after training accuracy reaches 100%, Figures 5 and 6 depict two scenarios for this phenomenon.
To conduct our experiments, we consider two canonical grokking settings. The first setting focuses
on the MNIST image classification task (Liu et al., 2022b). For this task, we adopt a ReLU MLP
architecture with a width of 200 and a depth of 3. The loss function used is mean squared error (MSE) loss. To optimize the network, we employ the AdamW optimizer with a learning rate of



Figure 4: The evolution of NTK trace among different layers, on both training and testing datasets.

298 0.001, and we use a batch size of 200 during the training process. In the second setting, we explore 299 the modulo addition dataset (Nanda et al., 2023). Here, we train transformers to perform addition 300 modulo P = 113. The input format follows the pattern of a + b =, where a and b are encoded as 301 *P*-dimensional one-hot vectors. The output c is extracted from the special token =. Our model for 302 this task consists of a one-layer ReLU transformer. The token embeddings have hidden dimension 303 d = 128, and we incorporate learned positional embeddings. The model comprises 4 attention heads with a dimension of $\frac{d}{4} = 32$, and an MLP with 512 hidden units. We use a training dataset containing 304 30% of all possible input pairs, and the loss function employed is cross-entropy loss. 305

In our analysis, we observe that the estimated NTK trace, given by $\frac{1}{\epsilon^2} |f(\mathcal{X}; w + \Delta) - f(\mathcal{X}; w)|_F^2$, bears a resemblance to the robustness metric (Tan & Huang, 2023). Specifically, we can define the logit perturb distance on the training data as $\frac{1}{\epsilon^2} |f(\mathcal{X} + \hat{\Delta}; w) - f(\mathcal{X}; w)|_F^2$, where $\hat{\Delta}$ represents a Gaussian perturbation applied to the input data. We consider both the NTK trace and the logit perturb distance and present them in Figures 5 and 6.

311 Our analysis reveals that the NTK traces on both the train and test datasets stabilize only when the 312 test accuracy approaches its maximum value. This finding aligns with the observations made in 313 standard supervised classification, as discussed in section 4, indicating that the underlying dynamics 314 of the network still changes when the training accuracy reaches 100% and stabilizes only after the 315 test accuracy saturates. Note the experiments in Figure 5 utilizes cross-entropy loss and the NTK 316 trace continuous to increase before the test accuracy first reaches 100%, making it similar to the 317 observation in section 4. Interestingly, we also observe that the trend of the NTK trace and the logit perturb distance exhibits several similarities, which could be attributed to the similarities in their 318 respective defining formulas. 319

320 321

322

296 297

5.3 SEMI-SUPERVISED LEARNING DYNAMICS

We will first present a simple review of a canonical semi-supervised learning method FixMatch. Given a batch of *B* labelled samples $\mathcal{X}_l = \{(x_i^1, y_i^1)\}_{i=1}^B$, there will also be μB unlabelled samples



Figure 5: The evolution of NTK trace and logit perturbation distance on Modular Addition dataset.



Figure 6: The evolution of NTK trace and logit perturbation distance on MNIST dataset.

 $\mathcal{X}_u = \{(x_i^2)\}_{i=1}^{\mu B}$ involved in the training. FixMatch loss utilizes two types of augmentations: weak and strong, defined by $\omega(\cdot)$ and $\Omega(\cdot)$. The loss consists of two parts: the supervised loss l_1 and the unsupervised loss l_2 .

$$l_1 = \frac{1}{B} \sum_{i=1}^{B} \mathrm{H}(\mathrm{onehot}(y_i^1), \mathrm{prob}(\omega(x_i^1))).$$
(5)

 $l_2 = \frac{1}{\mu B} \sum_{i=1}^{\mu B} \mathbb{I}(\max(\operatorname{prob}(\omega(x_i^2))) \ge \phi) \operatorname{H}(\operatorname{onehot}(\hat{y}_i^2), \operatorname{prob}(\Omega(x_i^2))),$ (6)

where $\hat{y}_i^2 = \arg \max \operatorname{prob}(\omega(x_i^2))$.

5.3.1 EVALUATING THE DYNAMICS ON LABELED AND UNLABELLED DATA

Semi-supervised learning involves both labeled and unlabeled samples, the key difference is that compared to the supervised case where the model is only trained on labeled samples, the training



Figure 7: The evolution of NTK trace estimations of (weakly-augmented) labeled, (weakly and strongly augmented) unlabeled, and pseudo-labeled samples during training. We consider three configurations with different numbers of labeled samples.



Figure 8: The evolution of test accuracy and NTK trace under different EMA coefficient τ .

process of semi-supervised learning is much more complex. Specifically, it involves both the (weakly 391 augmented-)labeled and strongly augmented (unlabeled-)samples that are pseudo-labeled in the 392 training. Motivated by the discussions in the supervised learning settings, we are interested in the dynamics of NTK trace of the following samples: 1) All the weakly augmented-labeled samples. 2) 394 All The weakly augmented-unlabeled samples. 3) All the strongly augmented-unlabeled samples. 395 4) All the strongly augmented-unlabeled samples that are pseudo-labeled. Denote $\hat{\mathcal{X}}_{\mu}$ as the set of 396 samples that are pseudo-labeled, i.e. those x_i^2 which satisfies $\max(\operatorname{prob}(\omega(x_i^2))) \geq \tau$. Then we are 397 interested in the following 4 quantities: 1) Supervised (labeled) samples (Sup) $\frac{\text{Tr}(K_t(\omega(\mathcal{X}_l),\omega(\mathcal{X}_l)))}{\mathcal{B} \mathcal{U}}$ 398 2) Unsupervised (unlabeled) weakly augmented-samples (Unsup Weak) $\frac{\text{Tr}(K_t(\omega(\mathcal{X}_u),\omega(\mathcal{X}_u)))}{\mu B K}$. 3) 399 Unsupervised (unlabeled) strongly augmented-samples (Unsup Strong) $\frac{\operatorname{Tr}(K_t(\Omega(\mathcal{X}_u),\Omega(\mathcal{X}_u)))}{\mu BK}$. 400 4) 401 Pseudo-labeled samples (Pseudo) $\frac{\text{Tr}(K_t(\Omega(\hat{X}_u), \Omega(\hat{X}_u)))}{\frac{1}{2}K}$ 402

403 In our experimental setup, we have closely followed the configuration used in FixMatch (Sohn et al., 404 2020). Our approach involves the utilization of a sgd optimizer with a momentum of 0.9. Additionally, 405 we have employed a cosine learning rate scheduler, initialized at 0.03, to dynamically adjust the 406 learning rate during training. The total number of training steps is set to 2^{20} , with evaluations conducted at every 5000 steps and makes the total number of evaluation steps to be 210. For the 407 labeled samples, we have chosen a batch size of 64, while maintaining a ratio μ of 7 between the 408 unlabeled and labeled samples. Furthermore, we have set the threshold value ϕ to 0.95 to guide 409 the decision-making process. The weak and strong augmentation functions are used following the 410 RandAugment approach (Cubuk et al., 2020). Lastly, we have employed the WideResNet-28-2 411 architecture as the underlying backbone model for our experimental investigations. 412

Figure 7 demonstrates the NTK trace for varying label quantities, with a range of 40, 250, and 413 4000 labeled samples in total. Notably, the changing patterns of each quantity exhibit remarkable 414 similarities across the different label quantities. To enhance clarity, a logarithmic scale is employed for 415 the evaluation steps. The figure reveals that the accuracy undergoes slow changes after approximately 416 20 steps, which corresponds to around 10% of the total training steps. Simultaneously, all NTK 417 traces experience a significant decrease, with distinct values for each of the four traces. Subsequently, 418 throughout the remaining training process, all NTK traces exhibit an increasing trend and converge 419 towards similar values. This phenomenon can be attributed to the scarcity of labeled data in the 420 early stages of training, which leads to a less stable training process. Furthermore, it indicates 421 that the weakly augmented labeled and unlabeled data share similar NTK traces, implying that the 422 predicted pseudo-labels may have similar distributions to the labeled ones. Additionally, the NTK 423 traces of weakly and strongly augmented samples are similar, suggesting that the network generalizes well on strongly augmented samples. Finally, the NTK traces of strongly augmented samples and 424 pseudo-labeled samples exhibit resemblances, indicating a high efficiency of pseudo-labeling. 425

426

388 389 390

427 428

5.3.2 INVESTIGATING THE DYNAMICS OF EXPONENTIAL MOVING AVERAGE MECHANISM

To enhance the prediction accuracy, in semi-supervised learning an exponential moving average (EMA) technique is usually employed. Specifically, at each step *n*, denote the EMA model as z_n^2 and the initial model as z_n^1 , the update rule for the EMA model with momentum τ is $z_{n+1}^2 = \tau z_{n+1}^2 + (1-\tau) z_n^1$.



We plot the accuracy and NTK trace of EMA model z_n^2 on the test dataset under different τ in Figure 8 447 using FixMatch on CIFAR-10 with 40 labels. The best accuracy for different τ is listed as: $\tau = 0.99$, 448 accuracy = 94.73. $\tau = 0.999$, accuracy = 94.97. $\tau = 0.9999$, accuracy = 95.13. The evolution of 449 test NTK is similar to the observation in section 5.3.1 with a first chaotic phase and quickly turns into 450 an increasing phase. 451

For the behavior of accuracy under different momentum parameters τ in EMA, we will present an 452 intuitive theoretical understanding of the phenomenon as follows using the gradient flow. 453

Theorem 5.1. Denote the parameter for the initial model along training as $z^{1}(t)$ and the parameter for the EMA updated model with momentum τ as $z_{\tau}^2(t)$. Assume $z_{\tau}^2(t)$ also satisfies the gradient 455 flow equation (1) and the training loss is a cross-entropy loss for input pair (x_i, y_i) , where y_i is a 456 ground-truth label or pseudo-label. Then we have the following bound on the average error:

$$\frac{1-\tau}{\tau} \|z^{1}(t) - z_{\tau}^{2}(t)\| \leq \sqrt{\frac{\operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; z_{\tau}^{2}(t)))}{N}} \sqrt{\frac{\sum_{i=1}^{N} \|onehot(y_{i}) - prob(x_{i}; z_{\tau}^{2}(t))\|^{2}}{N}}.$$
 (7)

We therefore plot the l_2 distance between initial model and EMA model parameters $||z^1(t) - z_{\tau}^2(t)||$ 462 in Figure 9 (a) and the ratio $\frac{1-\tau}{\tau} ||z^1(t) - z_{\tau}^2(t)|| / \sqrt{\frac{\operatorname{Tr}(K(\mathcal{X},\mathcal{X};z_{\tau}^2(t)))}{N}}$ in Figure 9 (b). Recall that the 463 464 test accuracy corresponds to $\tau = 0.99, 0.999, 0.9999$ is 94.73, 94.97, 95.13. From Figure 9 (b), we 465 can see the ratio decreases after the initial unstable training phase. The ratio has a decreasing order 466 with the size of τ and is getting close at the end of training, this matches the fact that test accuracy is 467 close but with decreasing order with the size of τ . These observations show that the inequality (7) 468 offers valuable insight into the dynamic of the EMA model.

469 470 471

446

454

6 CONCLUSION

472 Through the analysis of the NTK trace, we gained valuable insights into the network's training 473 progress and its correlation with accuracy. Our findings in standard supervised image classification 474 settings demonstrated that the NTK trace typically exhibited an increasing trend, eventually stabilizing 475 when the network achieved its highest accuracy on the training data. This observation underscores the 476 intrinsic connection between the NTK dynamics and the network's training performance. Furthermore, 477 our exploration of the phenomenon referred to as "grokking" yielded further insights. Through the utilization of the NTK trace, we closely monitored the training dynamics in grokking scenarios and 478 made a noteworthy observation: the test accuracy reaches its peak when the NTK trace stabilizes. 479 In the final phase of our study, we investigated the training dynamics of semi-supervised learning, 480 with a particular emphasis on examining the effectiveness of exponential moving average (EMA) 481 mechanisms. By leveraging the NTK trace, we gained a deeper and more accurate understanding 482 of the intricate behavior displayed in semi-supervised learning scenarios, as well as the influence 483 of EMA on the training process. This analysis provided valuable insights into the dynamics of 484 semi-supervised learning and shed light on the role of EMA in enhancing training performance. One 485 limitation of our work is not analyzing the whole property of NTK, just its trace.

486 REFERENCES

510

516

526

527

528

538

Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021. 2

- Andrzej Banburski, Qianli Liao, Brando Miranda, Lorenzo Rosasco, Fernanda De La Torre, Jack
 Hidary, and Tomaso Poggio. Theory iii: Dynamics and generalization in deep networks. *arXiv* preprint arXiv:1903.04991, 2019. 2
- Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural
 networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*, pp.
 483–513. PMLR, 2020. 2
- Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep relu networks? *arXiv preprint arXiv:1911.12360*, 2019. 2
- Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 702–703, 2020.
- Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy, and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. *Advances in Neural Information Processing Systems*, 33:5850–5861, 2020. 1, 2, 4
- Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. *arXiv preprint arXiv:1909.05989*, 2019. 1, 2
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
 1, 2, 3
- Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. *arXiv preprint arXiv:1909.12292*, 2019. 2
- Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha SohlDickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
 under gradient descent. Advances in neural information processing systems, 32, 2019. 3
- Etai Littwin, Tomer Galanti, and Lior Wolf. On random kernels of residual architectures. In Uncertainty in Artificial Intelligence, pp. 897–907. PMLR, 2021. 1, 2
- Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams.
 Towards understanding grokking: An effective theory of representation learning. *Advances in Neural Information Processing Systems*, 35:34651–34663, 2022a. 2
 - Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data. *arXiv preprint arXiv:2210.01117*, 2022b. 5
- Noel Loo, Ramin Hasani, Alexander Amini, and Daniela Rus. Evolution of neural tangent kernels
 under benign and adversarial training. *Advances in Neural Information Processing Systems*, 35:
 11642–11657, 2022. 1, 2, 4
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks.
 arXiv preprint arXiv:1906.05890, 2019. 2, 4
- Mor Shpigel Nacson, Suriya Gunasekar, Jason Lee, Nathan Srebro, and Daniel Soudry. Lexicographic
 and depth-sensitive margins in homogeneous and non-homogeneous deep models. In *International Conference on Machine Learning*, pp. 4683–4692. PMLR, 2019. 2
- 539 Neel Nanda, Lawrence Chan, Tom Liberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023. 2, 5, 6

540 541 542	Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: General- ization beyond overfitting on small algorithmic datasets. <i>arXiv preprint arXiv:2201.02177</i> , 2022. 2
543 544 545	Mariia Seleznova, Dana Weitzner, Raja Giryes, Gitta Kutyniok, and Hung-Hsu Chou. Neural (tangent kernel) collapse. <i>Advances in Neural Information Processing Systems</i> , 36, 2024. 2, 4
546 547	Haozhe Shan and Blake Bordelon. A theory of neural tangent kernel alignment and its influence on training. <i>arXiv preprint arXiv:2105.14301</i> , 2021. 2
549 550 551 552	Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Do- gus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. <i>Advances in neural information processing systems</i> , 33:596–608, 2020. 8
553 554 555	Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. <i>Journal of Machine Learning Research</i> , 19(70):1–57, 2018. 2
556 557 558	Zhiquan Tan and Weiran Huang. Understanding grokking through a robustness viewpoint. <i>arXiv</i> preprint arXiv:2311.06597, 2023. 2, 6
559 560	Tengyu Xu, Yi Zhou, Kaiyi Ji, and Yingbin Liang. When will gradient methods converge to max- margin classifier under relu models? <i>Stat</i> , 10(1):e354, 2021. 2
562 563 564	Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. <i>Advances in neural information processing systems</i> , 32, 2019. 2
566 567	
569 570	
571 572 573	
574 575 576	
577 578 579	
580 581	
582 583 584	
585 586 587	
588 589 590	
591 592 593	

⁵⁹⁴ A PROOFS OF THEOREMS

Theorem A.1. $\mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0, I_d)} \| f(\mathcal{X}; w + \Delta) - f(\mathcal{X}; w) \|_F^2 \sim \epsilon^2 \operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; w)).$

Proof. From Taylor expansion,

$$\begin{split} & \mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} \| f(\mathcal{X}; w + \Delta) - f(\mathcal{X}; w) \|_F^2 = \sum_{i=1}^N \mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} \| f(x_i; w + \Delta) - f(x_i; w) \|_F^2 \sim \\ & \sum_{i=1}^N \mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} \| J(x_i, w) \Delta \|_F^2 = \sum_{i=1}^N \mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} \operatorname{Tr}(J(x_i, w) \Delta \Delta^\top J^\top(x_i, w)) = \\ & \sum_{i=1}^N \operatorname{Tr}(\mathbb{E}_{\Delta \sim \epsilon \mathcal{N}(0,I_d)} J(x_i, w) \Delta \Delta^\top J^\top(x_i, w)) = \epsilon^2 \sum_{i=1}^N \operatorname{Tr}(J(x_i, w) J^\top(x_i, w)) = \\ & \epsilon^2 \operatorname{Tr}(J(\mathcal{X}, w) J^\top(\mathcal{X}, w)) = \epsilon^2 \operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; w)). \end{split}$$

Theorem A.2. Suppose $S(K_{t+1}(\mathcal{X}, \mathcal{X}), K_t(\mathcal{X}, \mathcal{X}))) = 0$ $(t \ge T_0)$. Suppose $T > T_0$, then $K_i(\mathcal{X}, \mathcal{X}) = K_j(\mathcal{X}, \mathcal{X})$ $(i, j \ge T)$ iff $\operatorname{Tr}(K_i(\mathcal{X}, \mathcal{X})) = \operatorname{Tr}(K_j(\mathcal{X}, \mathcal{X}))$ $(i, j \ge T)$.

Theorem A.3. Denote the parameter for the initial model along training as $z^1(t)$ and the parameter for the EMA updated model with momentum τ as $z_{\tau}^2(t)$. Assume $z_{\tau}^2(t)$ also satisfies the gradient flow equation (1) and the training loss is a cross-entropy loss for input pair (x_i, y_i) , where y_i is a ground-truth label or pseudo-label. Then we have the following bound

$$\frac{1-\tau}{\tau} \|z^{1}(t) - z_{\tau}^{2}(t)\| \leq \sqrt{\frac{\operatorname{Tr}(K(\mathcal{X}, \mathcal{X}; z_{\tau}^{2}(t)))}{N}} \sqrt{\frac{\sum_{i=1}^{N} \|onehot(y_{i}) - prob(x_{i}; z_{\tau}^{2}(t))\|^{2}}{N}}.$$
 (8)

Proof. The update formula for EMA model is as follows:

$$z_{\tau}^{2}(t+dt) = \tau z_{\tau}^{2}(t) + (1-\tau)z^{1}(t).$$

Then we can obtain the following ordinary differential equation

$$\partial_t z_\tau^2(t) = \frac{1 - \tau}{\tau} (z^1(t) - z_\tau^2(t)).$$
(9)

Assume $z_{\tau}^2(t)$ satisfies equation (1), then equation (9) simplifies to the following

$$-\nabla_w l(f(\mathcal{X}; z_\tau^2(t))) = \frac{1-\tau}{\tau} (z^1(t) - z_\tau^2(t)).$$
(10)

Using the expression of l, one may derive the following

$$\frac{1}{N}\sum_{i=1}^{N}J^{\top}(x_i, z_{\tau}^2(t))(\text{onehot}(y_i) - \text{prob}(x_i; z_{\tau}^2(t)) = \frac{1-\tau}{\tau}(z^1(t) - z_{\tau}^2(t)), \quad (11)$$

where $\operatorname{prob}(x_i; z_{\tau}^2(t)) = \left[\frac{e^{f_1(x_i; z_{\tau}^2(t))}}{\sum_{j=1}^K e^{f_j(x_i; z_{\tau}^2(t))}}, \cdots, \frac{e^{f_K(x_i; z_{\tau}^2(t))}}{\sum_{j=1}^K e^{f_j(x_i; z_{\tau}^2(t))}}\right]^{\top}.$

Then by the triangular inequality for norm and Cauchy–Schwarz inequality, we can obtain the following bound

$$\frac{1-\tau}{\tau} \|z^{1}(t) - z_{\tau}^{2}(t)\| \leq \frac{1}{N} \sum_{i=1}^{N} \|J^{\top}(x_{i}, z_{\tau}^{2}(t))\|_{F} \|\text{onehot}(y_{i}) - \text{prob}(x_{i}; z_{\tau}^{2}(t))\| \\
\leq \sqrt{\frac{\text{Tr}(K(\mathcal{X}, \mathcal{X}; z_{\tau}^{2}(t)))}{N}} \sqrt{\frac{\sum_{i=1}^{N} \|\text{onehot}(y_{i}) - \text{prob}(x_{i}; z_{\tau}^{2}(t))\|^{2}}{N}}. \quad (12)$$

Remark: Similar arguments also hold for MSE loss.

648 В MORE EMPIRICAL RESULTS ON STANDARD SUPERVISED CLASSIFICATION 649

650 The default settings for the experiment are as follows: the dataset used is CIFAR10, and the model 651 employed is ResNet18. The parameter ϵ for NTK trace calculation is set to 0.01. The learning rate is 652 initialized to 0.1, and momentum is set to 0.9. Weight decay is set to $5 * 10^{-4}$. The batch size for 653 training is set to 128, and the total number of training epochs is 200. Additionally, a cosine learning 654 rate decay strategy is employed during training. For CIFAR100 dataset, the settings are the same as CIFAR10. For ImageNet dataset, the parameter ϵ for NTK trace calculation is set to 0.001, the 655 learning rate is initialized to 0.01 and weight decay is set to 10^{-4} with 300 epochs training. All 656 experiments can be performed under 8 GeForce RTX 4090 in less than 1 day. 657

658 659

667

669

671 672 673

674

675 676

677

682 683

684 685

686

687 688 689

690 691

692

693

694

695 696

697

(a) mt = 0.1 (94.5)

B.1 DIFFERENT LEARNING RATES AND MOMENTUMS

660 In this section, we conduct two sets of experiments. We first vary the learning rate (lr) from 1 to 661 10^{-4} while keeping the momentum (mt) fixed at 0.9, and we also conduct experiments by fixing 662 the learning rate at 0.1 and vary the momentum from 0.1 to 0.9. We plot the NTK (Neural Tangent 663 Kernel) trace as a function of training epochs and we record the best test accuracy in the captions 664 after lr or mt. 665

Recall that sgd with momentum has update rule: w(t+1) = w(t) - lr * v(t+1), where v(t+1) = w(t) - lr * v(t+1). 666 $\nabla l(w(t)) + mt * v(t)$. When the learning rate or momentum is very small, the network will be hard to jump out of the local minima. This explains the phenomenon in Figures 10 and 11 where the NTK 668 trace will first undergo a certain phase of decreasing and test accuracy will be smaller in the low lr or mt regime due to this optimization insufficiency. 670





(c) mt = 0.5 (95.0)

(d) mt = 0.7 (95.5)

(e) mt = 0.9 (95.3)

B.2 DIFFERENT ARCHITECTURES IN A MODEL FAMILY

(b) mt = 0.3 (94.8)

In the following experiment, we test the dynamic behavior of the NTK trace on a family of different models, namely ResNet18, ResNet34, ResNet50 and ResNet101. We fix the learning rate at 0.1 and momentum at 0.9, and plot the NTK trace as a function of training epochs for each model. From Figure 12, we can see the trend of NTK is similar among models.

B.3 DIFFERENT INITIALIZATION STRATEGIES

698 In the following experiment, we test the influence of different initialization methods (Kaiming, 699 Orthogonal, Xavier, Standard) on the NTK trace dynamics. We fix the learning rate at 0.1 and momentum at 0.9, and plot the NTK trace as a function of training epochs for models initialized 700 with different methods. From Figure 13, we can see the trend of NTK is similar among different 701 initializations.







Figure 16: The evolution of NTK trace estimations under different optimizers and architectures, on both training and testing datasets.