

# Fast-NTK: Parameter-Efficient Unlearning for Large-Scale Models

Guihong Li, Hsiang Hsu, Chun-Fu Chen, and Radu Marculescu

**Abstract**—The rapid growth of machine learning has spurred legislative initiatives such as “the Right to be Forgotten,” allowing users to request data removal. In response, “machine unlearning” proposes the selective removal of unwanted data without the need for retraining from scratch. While the Neural-Tangent-Kernel-based (NTK-based) unlearning method excels in performance, it suffers from significant computational complexity, especially for large-scale models and datasets. Our work introduces “Fast-NTK,” a novel NTK-based unlearning algorithm that significantly reduces the computational complexity by incorporating parameter-efficient fine-tuning methods, such as fine-tuning batch normalization layers in a CNN or visual prompts in a vision transformer. Our experimental results demonstrate scalability to much larger neural networks and datasets (e.g., 88M parameters; 5k images), surpassing the limitations of previous full-model NTK-based approaches designed for smaller cases (e.g., 8M parameters; 500 images). Notably, our approach maintains a performance comparable to the traditional method of retraining on the retain set alone. Fast-NTK can thus enable for practical and scalable NTK-based unlearning in deep neural networks.

**Index Terms**—Machine Unlearning, Neural Tangent Kernel, Parameter-Efficient Fine-Tuning.

## I. INTRODUCTION

THE surge in machine learning applications has prompted legislative actions, notably “the Right to be Forgotten,” allowing individuals to request the removal of their online information [1]. However, the privacy challenge remains as erasing data from databases may persist in machine learning models, particularly in deep neural networks (DNNs), which are recognized for their efficient training data memorization [2]. To address this issue, “machine unlearning” has emerged to enable selective removal of unwanted “forget samples” without the need of retraining the model from scratch [3].

Among various unlearning algorithms [4–9], neural-tangent-kernel-based (NTK-based) unlearning stands out for its state-of-the-art performance [10, 11]. However, NTK-based unlearning algorithms are challenging due to the need of computing kernel matrices with respect to all samples and model weights. This computational complexity grows polynomially with the number of samples and model weights, thus resulting in intensive computation costs and memory consumption. Consequently, the effectiveness of NTK-based unlearning algorithms is often limited only to small-scale models and datasets (e.g., 8M parameters and 500 images).

Guihong Li and Radu Marculescu are with the University of Texas at Austin, Austin, TX 78712 USA. This work was done during an internship at JPMorgan Chase & Co. (e-mails: lgh@utexas.edu and radum@utexas.edu).

Hsiang Hsu and Chun-Fu Chen are with JPMorgan Chase & Co., USA (e-mails: hsiang.hsu@jpmchase.com and richard.cf.chen@jpmchase.com).

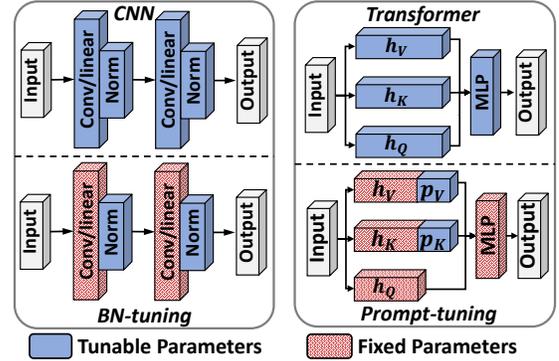


Fig. 1: A schematic representation of parameter-efficient fine-tuning and unlearning. For CNNs (left), instead of updating the entire model, we conduct the fine-tuning and NTK-based unlearning on the batch normalization (BN) layers. For transformers (right), we only modify the prompts ( $p_K$  and  $p_V$ ) appended to the entire model.

In this letter, we draw inspiration from recent strides in parameter-efficient fine-tuning (PEFT) [12–15] and leverage the NTK-based unlearning algorithms—specifically, the computation of kernel matrices—to work with a limited set of important parameters, such as those in batch normalization layers and visual prompts. We term this approach “Fast-NTK,” as shown in Figure 1. Unlike the conventional application of NTK-based unlearning algorithms across all model weights, Fast-NTK significantly reduces the parameter count (cf. Table I) of the standard implementation of the entire model. Remarkably, our experimental results, e.g., vision transformers (ViTs) on the ImageNet-R dataset, demonstrate indistinguishable performance compared to the commonly-used baseline that retrains the model from scratch only on the remaining data. Consequently, we believe our approach can enable a practical and scalable paradigm for the NTK-based unlearning approaches.<sup>1</sup>

## II. BACKGROUND AND RELATED WORK

Consider a training dataset  $\mathcal{D}$  that can be divided into two disjoint subsets: a forget set  $\mathcal{D}_f$  which is the target for unlearning, and a retain set  $\mathcal{D}_r$  which collects the remaining samples. The objective of machine unlearning is to eliminate the knowledge from the forget samples in  $\mathcal{D}_f$  of a model trained with  $\mathcal{D}$ , while minimizing the performance degradation of the retain samples in  $\mathcal{D}_r$  [16]. One intuitive strategy is to retrain the entire model from scratch, utilizing only the samples in  $\mathcal{D}_r$ . However, this process can be time-consuming, particularly when dealing with large-scale datasets and models. Consequently, current research endeavors to directly erase the

<sup>1</sup>Codes to reproduce our experiments will be made public.

knowledge associated with the forget samples from the model, without necessitating a complete retraining.

There exist three distinct strategies for accomplishing machine unlearning: data partitioning [4, SISA], mimicking differential privacy [5], and adjusting the model weights [6–9]. This letter specifically delves into the intricacies of updating the model weights, hence targeting machine unlearning through the computation of Neural Tangent Kernels (NTKs) [17, 18].

Consider a neural network  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $\theta \in \mathbb{R}^d$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are the support sets of the input and output, respectively. The NTK matrix of the two datasets  $\mathcal{D}_1$  and  $\mathcal{D}_2$  is defined as  $\Theta(\mathcal{D}_1, \mathcal{D}_2) \triangleq \nabla_\theta f_\theta(\mathcal{D}_1) \nabla_\theta f_\theta(\mathcal{D}_2)^\top$ . Let  $\theta$  and  $\theta_r$  be the weights from training with the entire training set  $\mathcal{D}$  and the retain set  $\mathcal{D}_r$  alone, respectively. Note that directly obtaining  $\theta_r$  from  $\theta$  is the goal of machine unlearning by updating model weights. By linearizing the outputs of  $f_\theta$ , we can approximate  $\theta$  and  $\theta_r$  in closed forms, and directly move the model weights from  $\theta$  to  $\theta_r$  by an optimal one-shot update:

$$\theta_r = \theta + \mathbf{P} \nabla_\theta f_\theta(\mathcal{D}_f)^\top \mathbf{M} \mathbf{V}, \quad (1)$$

where  $\mathbf{P} = \mathbf{I} - \nabla_\theta f_\theta(\mathcal{D}_r)^\top \Theta(\mathcal{D}_r, \mathcal{D}_r)^{-1} \nabla_\theta f_\theta(\mathcal{D}_r)$  is the matrix that projects the gradients of the samples to forget  $\nabla_\theta f_\theta(\mathcal{D}_f)$  to a space that is orthogonal to the space spanned by the gradients of all retain samples;  $\mathbf{M} = [\Theta(\mathcal{D}_f, \mathcal{D}_f) - \Theta(\mathcal{D}_r, \mathcal{D}_f)^\top \Theta(\mathcal{D}_r, \mathcal{D}_r)^{-1} \Theta(\mathcal{D}_r, \mathcal{D}_f)]^{-1}$  and  $\mathbf{V} = (\mathbf{y}_f - f_\theta(\mathcal{D}_f)) + \Theta(\mathcal{D}_r, \mathcal{D}_f)^\top \Theta(\mathcal{D}_f, \mathcal{D}_f)^{-1} (\mathbf{y}_r - f_\theta(\mathcal{D}_r))$  are the re-weighting matrices, while  $\mathbf{y}_f$  and  $\mathbf{y}_r$  are the ground truth labels for the forget set and retain set, respectively.

Despite NTK-based unlearning showcasing state-of-the-art performance in comparison to other methods [19], there are concerns regarding its numerical instability and scalability for models with many parameters [10, 11]. The inherent computational complexity has spurred efforts to enhance the efficiency of NTK-based unlearning algorithms, especially in large-scale setups. One approach to mitigate the computational costs involves the utilization of sketching techniques to approximate tensor products associated with NTK [20]. This method not only scales linearly with data sparsity, but also efficiently truncates the Taylor series of arc-cosine kernels. Additionally, improvements in the spectral approximation of the kernel matrix are achieved through leveraging the score sampling, or introducing a distribution that efficiently generates random features by approximating scores of arc-cosine kernels. Further strides in computational efficiency are made by novel algorithms employing mixed-order or high-order automatic differentiation [21]. It is important to note that these methods are often tailored to specific types of deep neural networks, thus limiting their widespread applicability. Moreover, their efficiency may still fall short for some larger deep networks [21]. Consequently, our objective is to propose a parameter-efficient and practical implementation of NTK-based unlearning methods, as discussed next.

### III. PROPOSED METHOD

The major barrier in NTK-based unlearning arises from the computation of the Jacobian matrix  $\nabla_\theta f_\theta(\mathcal{D})$ , defined

in Eq. (1), with dimensions  $|\mathcal{Y}| |\mathcal{D}_f| \times d$ . In the context of deep neural networks, the parameter count  $d$  spans a vast range, from millions to trillions [22, 23]. This abundance of parameters poses a formidable challenge due to the prohibitive costs in computation and storage, and has indeed been a primary impediment in applying NTK-based unlearning algorithm on large scale models. To mitigate the computational and storage burdens, the concept of PEFT has been proposed in Hously et al. [24]. PEFT selectively fine-tunes only a small subset of (additional) model parameters. Recent empirical findings indicate that state-of-the-art PEFT techniques achieve performance comparable to that of full fine-tuning (i.e., tuning all parameters) [13].

Drawing inspiration from PEFT, we extend the approach to NTK-based unlearning by selectively focusing on a subset of model parameters—this combined technique is referred to as “Fast-NTK.” As illustrated in Fig. 1, in the case of convolutional neural networks (CNNs), our approach involves fine-tuning the batch normalization (BN) layers, which has proven to be an effective strategy for adapting the model to new data domains [12, 15]. Meanwhile, for vision transformers (ViTs), success is achieved by fine-tuning several prompts appended to the attention blocks [13, 14, 25]. To elaborate, given a pre-trained CNN or ViT, we perform fine-tuning on the downstream dataset  $\mathcal{D}$  by using BN-based adjustments (for CNNs) or prompt-based modifications (for ViTs). Subsequently, when provided with a forget set  $\mathcal{D}_f$ , we execute NTK-based unlearning using Eq. (1) exclusively on the fine-tuned parameters. This streamlined Fast-NTK approach significantly reduces the parameters subjected to fine-tuning, down to a range of 0.05% ~ 4.88% of the full model parameters. Remarkably, Fast-NTK achieves a performance comparable to tuning all parameters, as demonstrated in the next section. For an analysis on the parameter reduction, see Appendix A.

## IV. EMPIRICAL RESULTS

### A. Setup

Our method starts with the CNNs and ViTs pre-trained on the CIFAR-100 and ImageNet-1K datasets respectively. We fine-tune these models on the CIFAR-10 [26] and ImageNet-R [27] datasets and then assess the performance of FAST-NTK. In the case of CIFAR-10, we designate one class as  $\mathcal{D}_f$ , while considering the remaining classes as  $\mathcal{D}_r$ . Similarly, for the ImageNet-R dataset, we randomly choose one class as  $\mathcal{D}_f$  and select either 19 or 49 classes from the 200 classes as  $\mathcal{D}_r$  (i.e., resulting in 20 or 50 total classes in  $\mathcal{D}$ ) to demonstrate the scalability of our approach.

We consider the following three metrics. First, we measure accuracy on both  $\mathcal{D}_r$  and  $\mathcal{D}_f$ —an unlearning algorithm should maintain high accuracy on  $\mathcal{D}_r$  while minimizing accuracy on  $\mathcal{D}_f$ . Second, we calculate accuracy on a hold-out set to ensure consistent performance on unseen data. Note that the hold-out set may contain samples from classes present in both  $\mathcal{D}_f$  and  $\mathcal{D}_r$ . The accuracy on the hold-out set should remain unaffected by the unlearning algorithm. Third, we incorporate relearning time [11], representing the number of epochs required to achieve a training loss below 0.05 on the forget set. Here, the

TABLE I: BN-based FAST-NTK on CNNs with CIFAR-10. All metrics are averaged over 5 runs with different seeds.

Dataset	Architectures	MobileNetV2			ResNet-110		
		#Images per class	100	200	500	100	200
CIFAR-10	#Params ratio (%)	4.88	4.88	4.88	0.51	0.51	0.51
Accuracy on $\mathcal{D}_r$	FULL	74.42±2.17	78.54±0.62	84.12±0.24	66.87±1.03	72.28±1.39	77.22±1.27
	MAX LOSS	71.13±1.91	68.24±1.40	14.12±1.59	56.64±2.17	49.17±2.19	13.49±1.89
	RANDOM LABEL	69.58±2.21	69.02±1.72	66.94±2.84	58.76±1.58	66.58±1.71	72.58±2.42
	<b>Fast-NTK</b>	70.80±2.04	73.70±0.68	80.76±0.40	65.60±4.36	71.04±1.65	76.84±0.21
	RETRAIN	75.56±2.36	79.50±0.55	85.27±0.26	69.02±1.65	74.13±1.54	78.98±0.43
Accuracy on $\mathcal{D}_f$	FULL	68.40±5.28	75.00±4.17	84.80±2.07	67.20±3.06	73.70±1.81	75.20±0.98
	MAX LOSS	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	RANDOM LABEL	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	<b>Fast-NTK</b>	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	RETRAIN	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
Accuracy on Hold-Out set	FULL	65.00±1.11	71.63±1.25	77.91±0.21	54.12±0.72	62.29±1.12	71.02±0.71
	MAX LOSS	58.14±0.95	58.28±1.22	12.51±1.27	43.18±0.34	41.61±2.06	12.06±2.81
	RANDOM LABEL	57.04±1.15	63.36±0.48	69.27±0.15	43.38±0.91	52.53±0.64	61.18±0.60
	<b>Fast-NTK</b>	58.54±0.88	63.96±1.67	69.96±3.64	50.80±5.57	59.88±1.59	60.58±0.63
	RETRAIN	60.50±1.15	66.41±0.46	71.80±0.19	50.36±1.17	57.57±0.71	65.58±1.51
#Relearning Epochs	MAX LOSS	28.80±0.40	22.20±0.40	77.20±6.01	24.00±0.89	25.20±2.48	22.00±0.93
	RANDOM LABEL	19.80±0.40	10.00±0.00	4.00±0.00	10.80±0.40	6.00±0.00	3.00±0.49
	<b>Fast-NTK</b>	21.00±0.63	10.80±0.40	4.00±0.00	12.40±0.80	6.00±0.00	2.80±0.40
	RETRAIN	21.20±0.40	11.00±0.00	4.80±0.40	12.60±0.49	6.20±0.40	3.00±0.00

value 0.05 is manually chosen, and could be chosen to other values. Relearning time serves as a measure of the difficulty in recovering knowledge from the forget set. If the model fails to achieve a loss below 0.05 within 100 epochs, we denote it as ‘>100’.

We compare FAST-NTK against the following baselines:

- FULL: The original model fine-tuned on  $\mathcal{D} = \mathcal{D}_f \cup \mathcal{D}_r$  without unlearning, serving as the reference model.
- MAX LOSS [28]: This baseline maximizes the training loss with respect to the ground truth labels of the samples in the forget set  $\mathcal{D}_f$ .
- RANDOM LABEL [29, 30]: This baseline minimizes the training loss by assigning uniformly random labels to the samples in the forget set  $\mathcal{D}_f$ .
- RETRAIN: The model trained only on the retain set  $\mathcal{D}_r$ .

Among these baselines, RETRAIN is commonly referred to as the **golden baseline**. This designation stems from its lack of prior knowledge about the samples in the forget set  $\mathcal{D}_f$ , making it an ideal reference point for comparing any unlearning algorithms. By evaluating FAST-NTK against RETRAIN, we aim to ensure that the unlearned model closely approximates the ideal scenario. This comparison helps ascertain that the unlearning process effectively eliminates unwanted data without causing significant performance degradation on  $\mathcal{D}_r$ . Essentially, an ideal unlearned model should exhibit indistinguishability in terms of the specified evaluation metrics to the golden baseline RETRAIN (see [3, Section 3.2]).

## B. Evaluation of Fast-NTK

We perform BN-based fine-tuning on MobileNet-v2 and ResNet-110 using a subset of the CIFAR-10 dataset, followed by unlearning algorithms that involves forgetting the class labeled “0.” To showcase the scalability of our approach, we vary the number of images per class (#IPC). The results in Table I reveal that our method requires less than **4.88%** of the parameters involved in tuning the entire model, making the unlearning process practical and achievable for these large models. Notably, FAST-NTK exhibits negligible or no accuracy degradation on the retain set compared to the golden baseline RETRAIN. In contrast, the accuracy on the forget set is indistinguishable from RETRAIN (drops to “0”) across various setups, with a similar number of relearning epochs needed as RETRAIN. Compared to the other baselines, MAX LOSS and RANDOM LABEL, FAST-NTK effectively preserves accuracy on the retrain set  $\mathcal{D}_r$  and the general test set, highlighting the robustness and efficiency of our proposed technique for CNNs.

Additionally, we extend the same setting to ViTs on the ImageNet-R dataset. As demonstrated in Table II, our approach requires less than **0.4%** of the parameters compared to tuning the entire model, making practical unlearning feasible for these large models. Comparisons with RETRAIN, MAX LOSS, and RANDOM LABEL show that FAST-NTK effectively preserves accuracy on the retain set  $\mathcal{D}_r$  and the general test set, achieving close accuracy to RETRAIN on the retain set. These results confirm the effectiveness and practicality of our unlearning approach for ViTs. Importantly, our method scales up to ViTs, representing a significant advancement compared to previous

TABLE II: Prompt-based FAST-NTK on ViTs with ImageNet-R. All metrics are averaged over 5 runs with different seeds.

Dataset	Architectures #Classes/#IPC	ViT-Tiny		ViT-Small		ViT-Base	
		20/50	50/20	20/50	50/20	20/50	50/20
ImageNet-R	#Params ratio (%)	0.24	0.35	0.12	0.18	0.06	0.09
Accuracy on $\mathcal{D}_r$	FULL	66.40±0.91	65.48±0.85	87.60±0.89	85.56±1.58	36.82±2.55	15.36±1.17
	MAX LOSS	57.71±1.33	51.80±0.70	77.35±1.23	71.17±0.66	24.78±2.74	8.52±0.77
	RANDOM LABEL	58.29±1.70	51.50±0.97	78.51±1.52	71.43±0.20	23.56±2.99	7.60±0.77
	<b>Fast-NTK</b>	66.53±0.63	65.24±0.48	87.03±1.49	85.31±2.14	40.84±1.84	17.40±1.89
	RETRAIN	68.21±1.50	65.92±0.93	87.77±0.42	86.58±0.15	37.45±2.03	16.02±1.02
Accuracy on $\mathcal{D}_f$	FULL	77.20±5.60	56.67±20.95	91.60±3.44	87.50±2.50	56.80±11.91	20.00±5.00
	MAX LOSS	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	RANDOM LABEL	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	<b>Fast-NTK</b>	0.00±0.00	0.00±0.00	0.00±0.00	3.00±2.50	0.00±0.00	0.00±0.00
	RETRAIN	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
Accuracy on Hold-Out set	FULL	47.73±0.45	31.43±1.70	68.06±1.12	52.75±0.75	32.53±2.40	12.05±0.05
	MAX LOSS	41.56±1.05	26.13±0.87	59.41±1.03	45.15±0.05	19.67±3.20	6.55±0.55
	RANDOM LABEL	45.02±1.92	26.03±1.31	64.03±1.23	45.40±0.10	19.96±3.60	6.55±0.25
	<b>Fast-NTK</b>	45.44±0.93	31.17±1.28	64.03±1.03	52.40±0.50	23.26±2.40	9.15±0.35
	RETRAIN	46.54±1.26	30.67±2.25	64.69±0.90	51.60±0.90	30.29±2.24	11.60±0.10
#Relearning Epochs	MAX LOSS	17.00±0.00	13.67±0.47	18.00±0.00	15.00±0.00	>100	>100
	RANDOM LABEL	4.20±0.40	3.67±0.47	6.40±0.49	6.00±0.00	>100	>100
	<b>Fast-NTK</b>	4.40±0.49	4.00±0.00	5.80±0.40	6.00±0.00	>100	>100
	RETRAIN	5.00±0.00	4.67±0.47	6.40±0.49	6.50±0.50	>100	>100

TABLE III: Linear probing on the ImageNet-R dataset. All metrics are averaged over 5 runs with different seeds.

	Network #Classes/#IPC	ViT-Small			ViT-Base		
		20/20	20/50	50/20	20/20	20/50	50/20
Acc on $\mathcal{D}_r$	PRE-TRAINED	60.39±1.27	58.24±1.49	53.70±2.36	99.93±0.11	99.32±0.24	99.87±0.08
	RANDOM-INIT	35.66±1.69	26.40±1.06	22.32±0.60	32.31±0.74	19.30±0.66	17.33±0.69
	<b>Fast-NTK</b>	60.25±3.76	53.71±2.45	47.24±0.00	86.58±2.13	87.66±1.01	87.24±0.00
Acc on $\mathcal{D}_f$	PRE-TRAINED	72.50±9.01	79.50±6.22	66.25±5.45	100.00±0.00	99.00±1.00	98.75±2.17
	RANDOM-INIT	54.53±2.46	33.33±17.00	43.33±11.12	49.40±4.42	15.00±8.16	17.33±7.72
	<b>Fast-NTK</b>	0.00±0.00	0.00±0.00	0.00±0.00	2.50±2.50	0.00±0.00	0.00±0.00

approaches like [11], which are confined only to toy networks and small datasets (e.g., less than 200 samples). For additional results on CIFAR-10, see Appendix B.

## V. DISCUSSION

**Risk of using pre-trained models.** It is crucial to emphasize that FAST-NTK starts with a pre-trained model rather than one initialized randomly. Despite the increasing popularity of leveraging pre-trained foundation models [31], these pre-trained models may possess some knowledge of classes from  $\mathcal{D}_f$ . This prior knowledge introduces an inherent risk during the unlearning process, as erasing all information and concepts associated with the classes in  $\mathcal{D}_f$  solely through the use of forget samples becomes a challenging task. To reassess this risk, for the pre-trained models used in our evaluation (PRE-TRAINED), we conduct fine-tuning of the classification head (i.e., linear probing) on  $\mathcal{D}_r \cup \mathcal{D}_f$ , while keeping the

parameters in the remaining layers frozen. We also conduct the linear probing on the randomly initialized model (RANDOM-INIT) and the unlearned model obtained by FAST-NTK (cf. Section IV). As illustrated in Table III, the accuracy of PRE-TRAINED on  $\mathcal{D}_r$  and  $\mathcal{D}_f$  is much higher than RANDOM-INIT (very close to 100%), indicating that the pre-trained model already possesses some level of knowledge about  $\mathcal{D}_r$  and  $\mathcal{D}_f$ . As expected, FAST-NTK effectively removes the knowledge on  $\mathcal{D}_f$  as the accuracy on  $\mathcal{D}_f$  is zero. This finding underscores the need for further investigation into the interplay between unlearning and PEFT on pre-trained models.

**Future work.** Our current implementation to obtain the NTK matrix relies on exact computations. To further improve the efficiency of FAST-NTK, one future direction is to explore approximate computation of the NTK matrix by, e.g., low-rank approximation or factorization.

## DISCLAIMER

This paper was prepared for informational purposes by the Global Technology Applied Research center of JPMorgan Chase & Co. This paper is not a product of the Research Department of JPMorgan Chase & Co. or its affiliates. Neither JPMorgan Chase & Co. nor any of its affiliates makes any explicit or implied representation or warranty and none of them accept any liability in connection with this paper, including, without limitation, with respect to the completeness, accuracy, or reliability of the information contained herein and the potential legal, compliance, tax, or accounting effects thereof. This document is not intended as investment research or investment advice, or as a recommendation, offer, or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction. Guihong Li's and Radu Marculescu's contributions were made as part of Guihong Li's internship at the Global Technology Applied Research center of JPMorgan Chase & Co.

## REFERENCES

- [1] G. D. P. Regulation, "General data protection regulation (gdpr)," *Intersoft Consulting, Accessed in October*, vol. 24, no. 1, 2018.
- [2] Y. Wu, Y. Burda, R. Salakhutdinov, and R. B. Grosse, "On the quantitative analysis of decoder-based generative models," in *Proceedings of ICLR*, 2017.
- [3] T. T. Nguyen, T. T. Huynh, P. L. Nguyen, A. W. Liew, H. Yin, and Q. V. H. Nguyen, "A survey of machine unlearning," *CoRR*, vol. abs/2209.02299, 2022.
- [4] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, "Machine unlearning," in *42nd IEEE Symposium on Security and Privacy*. IEEE, 2021, pp. 141–159.
- [5] V. Gupta, C. Jung, S. Neel, A. Roth, S. Sharifi-Malvajerdi, and C. Waites, "Adaptive machine unlearning," in *Advances in NeurIPS*, 2021.
- [6] S. Neel, A. Roth, and S. Sharifi-Malvajerdi, "Descent-to-delete: Gradient-based methods for machine unlearning," in *Proceedings of ALT*. PMLR, 2021.
- [7] A. K. Tarun, V. S. Chundawat, M. Mandal, and M. S. Kankanhalli, "Deep regression unlearning," in *Proceedings of ICML*. PMLR, 2023.
- [8] R. Chourasia and N. Shah, "Forget unlearning: Towards true data-deletion in machine learning," in *Proceedings of ICML*. PMLR, 2023.
- [9] Y. Chen, J. Xiong, W. Xu, and J. Zuo, "A novel online incremental and decremental learning algorithm based on variable support vector machine," *Cluster Computing*, vol. 22, pp. 7435–7445, 2019.
- [10] A. Golatkar, A. Achille, and S. Soatto, "Forgetting outside the box: Scrubbing deep networks of information accessible from input-output observations," in *Proceedings of ECCV*. Springer, 2020.
- [11] —, "Eternal sunshine of the spotless net: Selective forgetting in deep networks," in *Proceedings of CVPR*. IEEE, 2020.
- [12] J. Lin, L. Zhu, W.-M. Chen, W.-C. Wang, C. Gan, and S. Han, "On-device training under 256kb memory," in *Advances in NeurIPS*, 2022.
- [13] Z. Zheng, X. Yue, K. Wang, and Y. You, "Prompt vision transformer for domain generalization," *arXiv preprint arXiv:2208.08914*, 2022.
- [14] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim, "Visual prompt tuning," in *Proceedings of ECCV*. Springer, 2022.
- [15] H.-Y. Chiang, N. Frumkin, F. Liang, and D. Marculescu, "MobileTL: on-device transfer learning with inverted residual blocks," in *Proceedings of the AAAI*, 2023.
- [16] H. Xu, T. Zhu, L. Zhang, W. Zhou, and P. S. Yu, "Machine unlearning: A survey," *ACM Comput. Surv.*, vol. 56, no. 1, aug 2023.
- [17] J. Lee *et al.*, "Wide neural networks of any depth evolve as linear models under gradient descent," in *Advances in NeurIPS*, 2019.
- [18] A. Jacot, F. Gabriel, and C. Hongler, "Neural tangent kernel: Convergence and generalization in neural networks," in *Advances in NeurIPS*, 2018.
- [19] J. Jia, J. Liu, P. Ram, Y. Yao, G. Liu, Y. Liu, P. Sharma, and S. Liu, "Model sparsification can simplify machine unlearning," *CoRR*, vol. abs/2304.04934, 2023.
- [20] A. Zandieh, I. Han, H. Avron, N. Shoham, C. Kim, and J. Shin, "Scaling neural tangent kernels via sketching and random features," in *Advances in NeurIPS*, 2021.
- [21] R. Novak, J. Sohl-Dickstein, and S. S. Schoenholz, "Fast finite width neural tangent kernel," in *Proceedings of ICML*. PMLR, 2022.
- [22] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of ICML*. PMLR, 2021.
- [23] A. Dosovitskiy *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proceedings of ICLR*, 2021.
- [24] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *Proceedings of ICML*. PMLR, 2019.
- [25] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *arXiv preprint arXiv:2304.08485*, 2023.
- [26] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [27] D. Hendrycks *et al.*, "The many faces of robustness: A critical analysis of out-of-distribution generalization," *arXiv preprint arXiv:2006.16241*, 2020.
- [28] A. Halimi, S. Kadhe, A. Rawat, and N. Baracaldo, "Federated unlearning: How to efficiently erase a client in fl?" *CoRR*, vol. abs/2207.05521, 2022.
- [29] L. Graves, V. Nagisetty, and V. Ganesh, "Amnesiac machine learning," in *Proceedings of the AAAI*, 2021, pp. 11 516–11 524.
- [30] Z. Kong and K. Chaudhuri, "Data redaction from conditional generative models," *CoRR*, vol. abs/2305.11351, 2023.
- [31] R. Bommasani *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021.

## APPENDIX

## A. Parameter Reduction of Fast-NTK

**Fine-tune/unlearn CNNs with BN layers.** As shown in Fig. 1, a convolutional layer is typically followed by a batch normalization layer in a CNN. For a typical convolutional layer with  $C_o$  output channels,  $C_i$  input channels, kernel size  $K \times K$ , and  $g$  separable groups, the total number of parameters (weights) in this layer is  $\text{Parameters}_{\text{conv}} = \frac{C_o \times C_i \times K^2}{g}$ . In contrast, for a batch normalization (BN) layer, the only learnable parameters are the scaling ( $\gamma$ ) and shifting ( $\beta$ ) terms for each channel. Hence, the total number of learnable parameters in a BN layer is then  $\text{Parameters}_{\text{BN}} = 2 \times C_o$ . Usually,  $C_i \times K^2 \gg 2$  and  $g = 1$ ; therefore

$$\frac{\text{Parameters}_{\text{conv}}}{\text{Parameters}_{\text{BN}}} = \frac{C_i \times K^2}{2g} \gg 1. \quad (\text{A.2})$$

**Fine-tune/unlearn ViTs with Prompts.** In a ViT, the embedding layer transforms the input image into a sequence-like feature representation with the embedding dimension of  $E$ . Then the representation is processed by several transformer block consisting of a multi-head self-attention (MSA) block and two multi-layer perceptron (MLP) layers to obtain the outputs. Within each block, each MLP layer has  $E \times rE$ , where  $r$  is usually 4; so two MLP layers have  $8E^2$  parameters. Besides, each MSA has three weight matrices of size  $\frac{E}{m} \times E$ . Hence, MSA has  $3E^2$  parameters, and in total,  $\text{Parameters}_{\text{Block}} = 11E^2$ . As shown in Fig. 1, the prompt-based fine-tuning inserts the prompt parameters  $p_K$  and  $p_V$  to the Key and Value  $h_K$  and  $h_V$  of an MSA. As a contrast to tuning the entire MSA, the prompt-based method fine-tunes only  $L_p \times E$  parameters. Typically, the embedding dimensions  $E$  is much higher than the prompt length  $L_p$  (in our experimental setup,  $L_p = 10$ ); therefore:

$$\frac{\text{Parameters}_{\text{MSA}}}{\text{Parameters}_{\text{prompt}}} = \frac{11E}{L_p} \gg 1. \quad (\text{A.3})$$

## B. Additional Results on CIFAR-10.

We provide the results of Fast-NTK with ViTs on the CIFAR-10 dataset. Again, our approach requires less than 0.2% of the parameters and outperforms other baselines.

TABLE A.1: Prompt-based FAST-NTK on ViTs with CIFAR-10. All metrics are averaged over 5 runs with different seeds.

Dataset	Architectures	ViT-Small			ViT-Base		
		100	200	500	100	200	500
CIFAR-10	#Images per class						
	#Params ratio (%)	0.11	0.11	0.11	0.05	0.05	0.05
Accuracy on $\mathcal{D}_r$	FULL	95.78±0.52	94.93±1.06	94.78±0.48	84.18±1.09	85.67±0.62	87.07±0.24
	MAX LOSS	87.18±1.19	86.53±0.79	83.47±0.40	78.04±0.60	84.26±0.83	87.39±0.08
	RANDOM LABEL	93.87±0.86	93.72±0.55	93.32±0.35	76.56±0.83	83.83±0.82	87.28±0.17
	<b>Fast-NTK</b>	93.91±0.77	94.84±1.25	94.59±0.03	87.60±1.16	89.13±0.51	89.30±0.12
	RETRAIN	96.02±0.43	95.71±0.53	94.29±0.20	84.36±1.16	86.47±0.58	88.19±0.32
Accuracy on $\mathcal{D}_f$	FULL	97.00±1.55	96.20±1.36	95.73±1.67	84.80±6.31	90.40±1.11	92.00±0.00
	MAX LOSS	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	RANDOM LABEL	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
	<b>Fast-NTK</b>	0.20±0.40	0.20±0.24	0.00±0.00	0.00±0.00	0.00±0.00	0.20±0.20
	RETRAIN	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00	0.00±0.00
Accuracy on Hold-Out set	FULL	86.62±1.42	87.51±0.93	89.29±0.18	82.06±0.77	84.95±0.95	86.78±0.28
	MAX LOSS	71.90±0.45	73.38±0.72	73.09±0.24	68.14±0.71	74.60±0.86	77.73±0.21
	RANDOM LABEL	78.12±0.91	79.01±0.59	80.30±0.07	66.80±1.23	74.34±0.83	77.73±0.35
	<b>Fast-NTK</b>	77.78±1.14	78.87±0.97	80.63±0.23	70.62±2.10	75.37±0.45	78.47±0.15
	RETRAIN	78.94±1.11	79.61±0.44	80.64±0.10	73.76±0.43	76.56±0.87	78.59±0.27
#Relearning Epochs	MAX LOSS	9.20±0.40	8.00±0.00	6.00±0.00	>100	>100	47.50±0.50
	RANDOM LABEL	2.20±0.40	1.20±0.40	1.00±0.00	>100	>100	45.00±1.00
	<b>Fast-NTK</b>	2.60±0.49	1.00±0.00	1.00±0.00	>100	>100	53.50±1.50
	RETRAIN	2.60±0.49	1.40±0.49	1.00±0.00	>100	>100	46.50±0.50