

FAST FINITE WIDTH NEURAL TANGENT KERNEL

Anonymous authors

Paper under double-blind review

ABSTRACT

The Neural Tangent Kernel (NTK), defined as the outer product of the neural network (NN) Jacobians, $\Theta_\theta(x_1, x_2) = [\partial f(\theta, x_1)/\partial \theta] [\partial f(\theta, x_2)/\partial \theta]^T$, has emerged as a central object of study in deep learning. In the infinite width limit, the NTK can sometimes be computed analytically and is useful for understanding training and generalization of NN architectures. At finite widths, the NTK is also used to better initialize NNs, compare the conditioning across models, perform architecture search, and do meta-learning. Unfortunately, the finite width NTK is notoriously expensive to compute, which severely limits its practical utility.

We perform the first in-depth analysis of the compute and memory requirements for NTK computation in finite width networks. Leveraging the structure of neural networks, we further propose two novel algorithms that change the *exponent* of the compute and memory requirements of the finite width NTK, dramatically improving efficiency.

We open-source [github.com/iclr2022anon/fast_finite_width_ntk] our two algorithms as general-purpose JAX function transformations that apply to any differentiable computation (convolutions, attention, recurrence, etc.) and introduce no new hyper-parameters.

1 INTRODUCTION

The past few years have seen significant progress towards a theoretical foundation for deep learning. Much of this work has focused on understanding the properties of random functions in high dimensions. One significant line of work (Neal, 1994; Lee et al., 2018; Matthews et al., 2018; Novak et al., 2019; Garriga-Alonso et al., 2019; Hron et al., 2020; Yang, 2019) established that in the limit of infinite width, randomly initialized Neural Networks (NNs) are Gaussian Processes (called the NNGP). Building on this development, Jacot et al. (2018) showed that in function space the dynamics under gradient descent could be computed analytically using the so-called Neural Tangent Kernel (NTK) and Lee et al. (2019) showed that wide neural networks reduce to their linearization in weight space throughout training. A related set of results (Belkin et al., 2019; Spigler et al., 2019) showed that the ubiquitous bias-variance decomposition breaks down as high-dimensional models enter the so-called interpolating regime. Together these results describe learning in the infinite width limit and help explain the impressive generalization capabilities of NNs.

Insights from the wide network limit have had significant practical impact. The conditioning of the NTK has been shown to significantly impact trainability and generalization in NNs (Schoenholz et al., 2017; Xiao et al., 2018; 2020). This notion inspired initialization schemes like Fixup (Zhang et al., 2019), MetaInit (Dauphin & Schoenholz, 2019), and Normalizer Free networks (Brock et al., 2021a;b) and has enabled efficient neural architecture search (Park et al., 2020; Chen et al., 2021b). The NTK has additionally given insight into a wide range of phenomena such as: behavior of Generative Adversarial Networks (Franceschi et al., 2021), neural scaling laws (Bahri et al., 2021), and neural irradiance fields (Tancik et al., 2020). Kernel regression using the NTK has further enabled strong performance on small datasets (Arora et al., 2020), and applications such as dataset distillation (Nguyen et al., 2020; 2021) and uncertainty prediction (He et al., 2020; Adlam et al., 2020).

Despite the significant promise of theory based on the NTK, computing the NTK in practice is challenging. In the infinite width limit, the NTK can sometimes be computed analytically. However, it remains intractable for many architectures, and finite width corrections can be important to describe actual NNs used in practice. The NTK matrix can be computed for finite width networks as the outer

product of Jacobians using forward or reverse mode automatic differentiation (AD),

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{O} \times \mathbf{O}} := \underbrace{\left[\frac{\partial f(\theta, x_1)}{\partial \theta} \right]}_{\mathbf{O} \times \mathbf{P}} \underbrace{\left[\frac{\partial f(\theta, x_2)}{\partial \theta} \right]^T}_{\mathbf{P} \times \mathbf{O}}, \quad (1)$$

where f is the forward pass NN function producing outputs in $\mathbb{R}^{\mathbf{O}}$, $\theta \in \mathbb{R}^{\mathbf{P}}$ are all trainable parameters, and x_1 and x_2 are two inputs to the network. If inputs are batches of sizes \mathbf{N}_1 and \mathbf{N}_2 , the NTK is an $\mathbf{N}_1 \mathbf{O} \times \mathbf{N}_2 \mathbf{O}$ matrix.

Unfortunately, evaluating Eq. (1) is often infeasible due to time and memory requirements.

In this paper, we perform the first in-depth analysis of the compute and memory requirements for the NTK as in Eq. (1). Noting that forward and reverse mode AD are two extremes of a wide range of AD strategies (Naumann, 2004; 2008), we explore other methods for computing the NTK leveraging the structure of NNs used in practice. We propose two novel methods for computing the NTK that exploit different orderings of the computation. We describe the compute and memory requirements of our techniques in fully-connected (FCN) and convolutional (CNN) settings, and show that one is asymptotically more efficient in both settings. We compute the NTK over a wide range of NN architectures and demonstrate that these improvements are robust in practice. We open-source implementations of both methods as JAX function transformations.

2 RELATED WORK

The finite width NTK (denoted as simply NTK throughout this work) has been used extensively in many recent works, but to our knowledge implementation details and compute costs were rarely made public. Below we draw comparison to some of these works, but we stress that it only serves as a sanity check to make sure our contribution is valuable relative to the scale of problems that have been attempted (none of these works had efficient NTK computation as their central goal).

In order to compare performance of models based on the NTK and the infinite width NTK, Arora et al. (2019a, Table 2) compute the NTK of up to 20-layer, 128-channel CNN in a binary CIFAR-2 classification setting. In an equivalent setting with the same hardware (NVIDIA V100), we are able to compute the NTK of a 2048-channel CNN, i.e. a network with at least 256 times more parameters.

To demonstrate the stability of the NTK during training for wide networks, Lee et al. (2019, Figure S6) compute the NTK of up to 3-layer 2^{12} -wide or 1-layer 2^{14} -wide FCNs. In the same setting with the same hardware (NVIDIA V100), we can reach widths of at least 2^{14} and 2^{18} respectively, i.e. handle networks with at least 16 times more parameters.

To investigate convergence of a WideResNet WRN-28- k (Zagoruyko & Komodakis, 2016) to its infinite width limit, Novak et al. (2020, Figure 2) evaluate the NTK of this model with widening factor k up to 32. In matching setting and hardware, we are able to reach the widening factor of at least 64, i.e. work with models at least 4 times larger.

To meta-learn NN parameters for transfer learning in a MAML-like (Finn et al., 2017) setting, Zhou et al. (2021, Table 7) replace the inner training loop with NTK-based inference. They use up to 5-layer, 200-channel CNNs on MiniImageNet (Oreshkin et al., 2018) with scalar outputs and batch size 25. In same setting we achieve at least 512 channels, i.e. support models at least 6 times larger.

Park et al. (2020, §4.1) use the NTK to predict the generalization performance of architectures in the context of Neural Architecture Search (Zoph & Le, 2017, NAS); however, the authors comment on its high computational burden and ultimately use a different proxy. In another NAS setting, Chen et al. (2021a, §3.1.1) use the condition number of NTK to predict a model’s trainability. Remark-ing its prohibitive cost, Chen et al. (2021b, Table 1) also use the NTK to evaluate the trainability of several ImageNet (Deng et al., 2009) models such as ResNet 50/152 (He et al., 2016), Vision Transformer (Dosovitskiy et al., 2021) and MLP-Mixer (Tolstikhin et al., 2021). However, in all of the above cases the authors only evaluate a pseudo-NTK, i.e. an NTK of a scalar-valued function,¹ which impacts the quality of the respective trainability/generalization proxy.

¹Precisely, computing the Jacobian only for a single logit or the sum of all 1000 class logits. The result is not the full NTK, but rather a single diagonal block or the sum of its 1000 diagonal blocks (finite width NTK is a dense matrix, not block-diagonal).

Method	Time	Memory	Use when
Jacobian contraction	$\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \mathbf{W}^2$	$\mathbf{N} \mathbf{O} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{L} \mathbf{W}^2$	Don't
NTK-vector products	$\mathbf{N}^2 \mathbf{O}^2 \mathbf{W} + \mathbf{N}^2 \mathbf{L} \mathbf{O} \mathbf{W}^2$	$\mathbf{N} \mathbf{O} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{L} \mathbf{W}^2$	$\mathbf{O} > \mathbf{W}$ or $\mathbf{N} = 1$
Structured derivatives	$\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \mathbf{W} + \mathbf{N} \mathbf{L} \mathbf{O} \mathbf{W}^2$	$\mathbf{N} \mathbf{O} \mathbf{W} + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{L} \mathbf{W}^2$	$\mathbf{O} < \mathbf{W}$ or $\mathbf{L} = 1$

Table 1: **Asymptotic time and memory cost of computing the NTK for an FCN.** Costs are for a pair of batches of inputs of size \mathbf{N} each, and for \mathbf{L} -deep, \mathbf{W} -wide FCN with \mathbf{O} outputs. Resulting NTK has shape $\mathbf{N} \mathbf{O} \times \mathbf{N} \mathbf{O}$. **NTK-vector products** allow a reduction of the time complexity, while **Structured derivatives** reduce both time and memory complexity. **Note:** presented are asymptotic cost estimates; in practice, all methods incur large constant multipliers (e.g. at least 3x for time; see §3.1). However, this generally does not impact the relative performance of different methods. See §3.6 for discussion, Table 7 for CNN, and Table 2 for more generic cost analysis.

Method	Time	Memory	Use when
Jacobian contraction	$\mathbf{N} \mathbf{O} [\mathbf{F}\mathbf{P}] + \mathbf{N}^2 \mathbf{O}^2 \mathbf{P}$	$\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{O} [\mathbf{Y}^k + \mathbf{P}^l] + \mathbf{N} \mathbf{Y} + \mathbf{P}$	$\mathbf{P} \ll \mathbf{Y}$, small \mathbf{O} , exotic primitives
NTK-vector products	$\mathbf{N}^2 \mathbf{O} [\mathbf{F}\mathbf{P}]$	$\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{O} [\mathbf{Y}^k + \mathbf{P}^l] + \mathbf{N} \mathbf{Y} + \mathbf{P}$	$\mathbf{F}\mathbf{P} < \mathbf{O}\mathbf{P}$, large \mathbf{O} , small \mathbf{N}
Structured derivatives	$\mathbf{N} \mathbf{O} [\mathbf{F}\mathbf{P}] + \mathbf{N} \mathbf{O} \mathbf{G} + \mathbf{N} [\mathbf{J} - \mathbf{O}\mathbf{P}]$	$\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{O} \mathbf{Y}^k + \mathbf{N} \mathbf{J}^k + \mathbf{N} \mathbf{Y} + \mathbf{P}$	$\mathbf{F}\mathbf{P} > \mathbf{O}\mathbf{P}$, large \mathbf{O} , large \mathbf{N}

Table 2: **Asymptotic time and memory cost estimates of computing the NTK for a generic function.** \mathbf{P} stands for the number of all parameters in the network, \mathbf{Y} stands for size of all pre-activations in the network, $\mathbf{F}\mathbf{P}$ stands for forward pass, and \mathbf{G} and \mathbf{J} depend on the structure of $\mathbf{F}\mathbf{P}$ (§B). For example, FCNs usually have a cheap $\mathbf{F}\mathbf{P} \leq \mathbf{O}\mathbf{P}$, as it consists of a single matrix multiply with the parameter matrix, and therefore **NTK-vector products** are recommended. CNNs, notably when the number of output pixels \mathbf{D} is large, have a costly $\mathbf{F}\mathbf{P} \geq \mathbf{O}\mathbf{P}$, since it amounts to \mathbf{D} matrix multiplies with the parameters, and therefore **Structured derivatives** are preferred. For precise analysis, see Table 1 for FCN and Table 7 for CNN.

In contrast, in this work we can compute the full 1000×1000 NTK on the same models (1000 classes), i.e. perform a task 1000 times more costly.

Finally, we remark that in all of the above settings, scaling up by increasing width or by working with the true NTK (vs the pseudo-NTK) should lead to improved downstream task performance due to better infinite width/linearization approximation or higher-quality trainability/generalization proxy respectively, which makes our work especially relevant to modern research.

3 EFFICIENT FINITE WIDTH NTKS IN A SIMPLIFIED SETTING

To gain intuition for the problem, we start by analyzing and improving the cost of computing the NTK for a simple FCN. See §F for an equivalent analysis of CNNs. We summarize the resulting complexities for FCN in Table 1, CNN in Table 7, and a general takeaway in Table 2.

Setting. Consider an \mathbf{L} -layer FCN $f(\theta, x) = \theta^{\mathbf{L}} \phi(\theta^{\mathbf{L}-1} \dots \theta^1 \phi(\theta^0 x) \dots) \in \mathbb{R}^{\mathbf{O}}$, where \mathbf{O} is the number of logits. We denote individual weight matrices as θ^l with shapes $\mathbf{W} \times \mathbf{W}$ (except for top-layer $\theta^{\mathbf{L}}$ of shape $\mathbf{O} \times \mathbf{W}$), where \mathbf{W} is the width of the network, and write the set of all parameters as $\theta = \text{vec}[\theta^0, \dots, \theta^{\mathbf{L}}] \in \mathbb{R}^{\mathbf{L}\mathbf{W}^2 + \mathbf{O}\mathbf{W}}$. We further define $x^l := \phi(y^{l-1})$ as post-activations (with $x^0 := x$), and $y^l := \theta^l x^l$ as pre-activations with $y^{\mathbf{L}} = f(\theta, x)$. See Fig. 5 for a visual schematic of these quantities. For simplicity, we assume that inputs x also have width \mathbf{W} , and $\mathbf{O} = \mathcal{O}(\mathbf{L}\mathbf{W})$, i.e. the number of logits is dominated by the product of width and depth. In §L we repeat the same derivations without the latter assumption, and arrive at qualitatively identical conclusions.

The NTK of f evaluated at two inputs x_1 and x_2 is an $\mathbf{O} \times \mathbf{O}$ matrix defined as

$$\Theta_{\theta} := \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T = \sum_{l=0}^{\mathbf{L}} \frac{\partial f(\theta, x_1)}{\partial \theta^l} \frac{\partial f(\theta, x_2)}{\partial \theta^l}^T =: \sum_{l=0}^{\mathbf{L}} \Theta_{\theta}^l \in \mathbb{R}^{\mathbf{O} \times \mathbf{O}}, \quad (2)$$

where we have defined Θ_{θ}^l to be the summands. We omit dependence on x_1, x_2 , and f for brevity.

In §3.1 and §3.2 we describe the cost of several fundamental AD operations that we will use as building blocks throughout the text. We borrow the nomenclature introduced by Autograd (Maclau-

rin et al.) and describe Jacobian-vector products (JVP), vector-Jacobian products (VJP), as well as the cost of computing the Jacobian $\partial f(\theta, x)/\partial \theta$.

In §3.3, we describe the baseline complexity of evaluating the NTK, by computing two Jacobians and contracting them. This approach is used in most (likely all) prior works, and scales poorly with the NN width \mathbf{W} and output size \mathbf{O} .

In §3.4 we present our first contribution, that consists in observing that many intermediate operations on weights performed by NNs possess a certain structure, that can allow linear algebra simplifications of the NTK expression, leading to a cheaper contraction and smaller memory footprint.

In §3.5 we present our second contribution, where we rephrase the NTK computation as instantiating itself row-by-row by applying the NTK-vector product function to columns of an identity matrix. As we will show, this trades off Jacobian contraction for more forward passes, which proves beneficial in many (but not all) settings.

3.1 JACOBIAN-VECTOR PRODUCTS AND VECTOR-JACOBIAN PRODUCTS

We begin by defining Jacobian-vector products and vector-Jacobian products:

$$\text{JVP}_{(f, \theta, x)} : \theta_t \in \mathbb{R}^{\mathbf{LW}^2 + \mathbf{OW}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta} \theta_t \in \mathbb{R}^{\mathbf{O}}, \quad (3)$$

$$\text{VJP}_{(f, \theta, x)} : f_c \in \mathbb{R}^{\mathbf{O}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta}^T f_c \in \mathbb{R}^{\mathbf{LW}^2 + \mathbf{OW}}. \quad (4)$$

The JVP can be understood as pushing forward a tangent vector in weight space to a tangent vector in the space of outputs; by contrast the VJP pulls back a cotangent vector in the space of outputs to a cotangent vector in weight space. These elementary operations correspond to forward and reverse mode AD respectively and serve as a basis for typical AD computations such as gradients, Jacobians, Hessians, etc.

Time and memory costs of JVP and VJP are asymptotically equivalent to the cost of the forward pass (**FP**), except for VJP additionally requires storing all intermediate activations. (see §N and Fig. 6).

For the case of FCNs, the time cost² of both operations is therefore

$$[\mathbf{FP}] = [\text{cost of all intermediate layers}] + [\text{cost of the top layer}] = [\mathbf{LW}^2] + [\mathbf{OW}] \sim \mathbf{LW}^2.$$

For a single input, the memory cost of computing both the JVP and the VJP are respectively,

$$[\text{size of all weights}] + [\text{size of activations at a single layer}] = [\mathbf{LW}^2 + \mathbf{OW}] + [\mathbf{W} + \mathbf{O}] \sim \mathbf{LW}^2,$$

$$[\text{size of all weights}] + [\text{size of activations in all layers}] = [\mathbf{LW}^2 + \mathbf{OW}] + [\mathbf{LW} + \mathbf{O}] \sim \mathbf{LW}^2.$$

Despite the fact that the VJP requires more memory to store intermediate activations, we see that for FCNs both computations are dominated by the cost of storing the weights.

Batched inputs. If x is a batch of inputs of size \mathbf{N} , the time cost of JVP and VJP increases linearly to \mathbf{NLW}^2 . The memory cost is slightly more nuanced. Since weights can be shared across inputs, the memory cost of the JVP and VJP are respectively,

$$\begin{aligned} & [\text{size of all weights}] + \mathbf{N} [\text{size of activations at a single layer}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N} [\mathbf{W} + \mathbf{O}] \sim \mathbf{LW}^2 + \mathbf{NW} + \mathbf{NO}, \\ & [\text{size of all weights}] + \mathbf{N} [\text{size of activations in all layers}] + \mathbf{N} [\text{size of all weight matrices}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N} [\mathbf{LW} + \mathbf{O}] + \mathbf{N} [\mathbf{LW}^2 + \mathbf{OW}] \sim \mathbf{NLW}^2. \end{aligned}$$

The cost of the VJP is dominated by the cost of storing the cotangents in weight space. However, for the purposes of computing the NTK, we will be contracting Jacobians layerwise and so we will only need to store one cotangent weight matrix, $\partial f/\partial \theta^l$, at a time. Thus, for the purposes of this work we end up with the following costs:

- JVP costs \mathbf{NLW}^2 time and $\mathbf{LW}^2 + \mathbf{NW} + \mathbf{NO}$ memory.
- VJP costs \mathbf{NLW}^2 time and $\mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NW}^2 + \mathbf{NOW}$ memory.

²To declutter notation, we omit the \mathcal{O} symbol to indicate asymptotic complexity in this work.

3.2 JACOBIAN COMPUTATION

For neural networks, the Jacobian is most often computed by evaluating the VJP on rows of the identity matrix $I_{\mathbf{O}}$, i.e.

$$[\partial f(\theta, x) / \partial \theta]^T = [\partial f(\theta, x) / \partial \theta]^T I_{\mathbf{O}} \in \mathbb{R}^{(\mathbf{LW}^2 + \mathbf{OW}) \times \mathbf{O}}. \quad (5)$$

It follows that computing the Jacobian takes \mathbf{O} evaluations of the VJP. However, as mentioned in §3.1, we only need to store one $\partial f / \partial \theta^l$ at a time, while the weights and intermediate activations are reused across evaluations. Thus, time and memory costs to compute the Jacobian are respectively,

$$\begin{aligned} & \mathbf{ON}([\text{cost of all intermediate layers}] + [\text{cost of the top layer}]) \\ &= \mathbf{ON}([\mathbf{LW}^2] + [\mathbf{OW}]) \sim \mathbf{NLOW}^2 + \mathbf{NO}^2\mathbf{W}, \\ & [\text{size of all weights}] + \mathbf{N}[\text{size of activations in all layers}] + \mathbf{ON}[\text{size of a single weight matrix}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LW} + \mathbf{O}] + \mathbf{ON}[\mathbf{W}^2 + \mathbf{OW}] \sim \mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W}. \end{aligned}$$

Therefore, asymptotically,

Jacobian costs $\mathbf{NLOW}^2 + \mathbf{NO}^2\mathbf{W}$ time and $\mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W}$ memory.

3.3 JACOBIAN CONTRACTION

We now analyze the cost of computing the NTK, starting with the direct computation as the product of two Jacobians. Consider a single summand from Eq. (2):

$$\underbrace{\Theta_{\theta}^l}_{\mathbf{O} \times \mathbf{O}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta^l}}_{\mathbf{O} \times (\mathbf{W} \times \mathbf{W})} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta^l}}_{(\mathbf{W} \times \mathbf{W}) \times \mathbf{O}}. \quad (6)$$

The time cost of this contraction is $\mathbf{O}^2\mathbf{W}^2$, and the memory necessary to instantiate each factor and the result is $\mathbf{OW}^2 + \mathbf{O}^2$. Repeating the above operation for each θ^l , we arrive at $\mathbf{LO}^2\mathbf{W}^2$ time cost and unchanged memory, due to being able to process summands sequentially.

Batched inputs. If we consider x_1 and x_2 to be input batches of size \mathbf{N} , then the resulting NTK is a matrix of shape $\mathbf{NO} \times \mathbf{NO}$, and the time cost becomes $\mathbf{N}^2\mathbf{LO}^2\mathbf{W}^2$, while memory grows to $[\text{NTK matrix size}] + [\text{factors size}] = \mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}^2$.

What remains is to account for the cost of computing and storing individual cotangents $\partial f / \partial \theta^l$, which is exactly the cost of computing the Jacobian (§3.2). Adding the costs up we obtain

Jacobian contraction costs $\mathbf{N}^2\mathbf{LO}^2\mathbf{W}^2$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{LW}^2 + \mathbf{NLW}$ memory.

3.4 LEVERAGING STRUCTURED DERIVATIVES FOR COMPUTING THE NTK

We can rewrite Θ_{θ}^l in Eq. (6) using the chain rule and our pre- and post-activation notation as:

$$\Theta_{\theta}^l = \left[\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \frac{\partial y_{x_1}^l}{\partial \theta^l} \right] \left[\frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l} \frac{\partial y_{x_2}^l}{\partial \theta^l} \right]^T = \underbrace{\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l}}_{\mathbf{O} \times \mathbf{W}} \underbrace{\frac{\partial y_{x_1}^l}{\partial \theta^l}}_{\mathbf{W} \times (\mathbf{W} \times \mathbf{W})} \underbrace{\frac{\partial y_{x_2}^l}{\partial \theta^l}}_{(\mathbf{W} \times \mathbf{W}) \times \mathbf{W}} \underbrace{\frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}}_{\mathbf{W} \times \mathbf{O}}. \quad (7)$$

At face value, rewriting Eq. (6) this way is unhelpful as it appears to have introduced additional costly contractions. However, recall that $y^l = \theta^l x^l$, and therefore

$$\frac{\partial y_{x_1}^l}{\partial \theta^l} = I_{\mathbf{W}} \otimes x_1^{lT}, \quad \frac{\partial y_{x_2}^l}{\partial \theta^l} = I_{\mathbf{W}} \otimes x_2^{lT}, \quad (8)$$

where \otimes is the Kronecker product. Plugging Eq. (8) into Eq. (7) we obtain (see §G)

$$\Theta_{\theta}^l = \begin{pmatrix} \underbrace{x_1^l}_{1 \times W}^T & \underbrace{x_2^l}_{W \times 1} \end{pmatrix} \begin{bmatrix} \underbrace{\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l}}_{O \times W} & \underbrace{\frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}}_{W \times O}^T \end{bmatrix}, \quad (9)$$

and observe that it takes only O^2W time and $OW + O^2$ memory. Accounting for depth, time cost becomes LO^2W , while memory does not change since the summands can be processed sequentially.

Batched inputs. The time cost grows quadratically with the batch size N up to N^2LO^2W , while the memory cost increases to $N^2O^2 + NOW$ to store the resulting NTK and $\partial f(\theta, x) / \partial y_x^l$ factors.

Finally, we need to account for the cost of computing the derivatives, $\partial f / \partial y^l$, and post-activations, x^l . Notice that both x^l and $\partial f / \partial y^l$ arises naturally when computing the **Jacobian** as the primals and cotangents in layer l respectively. However, since we do not need to compute the weight space cotangents explicitly (i.e. we cut the backpropagation algorithm short) the memory cost will be,

$$\begin{aligned} & [\text{size of all weights}] + N [\text{size of activations in all layers}] \\ &= [LW^2 + OW] + N [LW + O] \sim LW^2 + NLW. \end{aligned}$$

The extra time cost is asymptotically the cost of O forward passes, $NLOW^2$ which is the same as the **Jacobian**. However, as we will see in experiments, in practice we'll often compute the NTK faster than the Jacobian due to not computing the weight space cotangents $\partial f / \partial \theta^l$. Altogether,

By leveraging **Structured derivatives** in NN computations, we have reduced the cost of NTK to $N^2LO^2W + NLOW^2$ time and $N^2O^2 + NOW + LW^2 + NLW$ memory.

The key insight was to leverage the constant block-diagonal structure of the pre-activation derivatives $\partial y^l / \partial \theta^l$. This idea is quite general; as we discuss in §4 and detail in the appendix, similar structure exists for many common operations such as convolutions, pooling, and arithmetic.

We highlight that these computational improvements do not emerge automatically in AD. While JAX and other libraries leverage structures analogous to Eq. (8) to efficiently compute single evaluations of the VJP and JVP, this knowledge is lost once the (structureless) Jacobian $\partial f(\theta, x_1) / \partial \theta^l$ is instantiated, and cannot be taken advantage of in the following contraction with $\partial f(\theta, x_2) / \partial \theta^l$. We discuss how we impose this structure to compute the NTK for general neural networks in §4.

3.5 NTK VIA NTK-VECTOR PRODUCTS

Computing the **Jacobian contraction** using **Jacobian** first instantiates the Jacobian using VJPs and then performs a contraction. **Structured derivatives** use a similar strategy, but speed-up the contraction and avoid explicitly instantiating the weight space cotangents. Here we avoid performing a contraction altogether at the cost of extra VJP/JVP calls; this ends up being beneficial for FCNs.

We introduce the linear function performing the NTK-vector product: $\Theta VP : v \in \mathbb{R}^O \mapsto \Theta_{\theta} v \in \mathbb{R}^O$. Applying this function to O columns of the identity matrix I_O allows us to compute the NTK, i.e. $\Theta_{\theta} I_O = \Theta_{\theta}$. The cost of evaluating the NTK in this fashion is equal to O times the cost of a single NTK-vector product evaluation $\Theta VP(v)$. We now expand $\Theta VP(v) = \Theta_{\theta} v$ as

$$\frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)], \quad (10)$$

where we have observed that, if contracted from right to left, the NTK-vector product can be expressed as a composition of a JVP and VJP of the underlying function f . The cost of this operation is asymptotically equivalent to the cost of the **Jacobian**, since it consists of O VJPs followed by O (cheaper) JVPs. Therefore it costs $LOW^2 + O^2W$ time and $LW^2 + OW^2 + O^2W$ memory.

Batched inputs. In the batched setting Eq. (10) is repeated for each pair of inputs, and therefore time increases by a factor of N^2 to become $N^2LOW^2 + N^2O^2W$. However, the memory cost grows only linearly in N (except for the cost of storing the NTK of size N^2O^2), since intermediate activations and derivatives necessary to compute the JVP and VJP can be computed for each batch x_1 and x_2

separately; these quantities are then reused for every pairwise combination resulting in a memory equivalent to the **Jacobian**, i.e. $\mathbf{N}^2\mathbf{O}^2 + (\mathbf{LW}^2 + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLW})$, resulting in

NTK computation as a sequence of **NTK-vector products** costs $\mathbf{N}^2\mathbf{LOW}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{W}$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}^2 + \mathbf{LW}^2 + \mathbf{NLW}$ memory.

3.6 SUMMARY

Structured derivatives and **NTK-vector products** allow a reduction in the time cost of NTK computation in different ways, and **Structured derivatives** also reduce memory requirements. **Structured derivatives** are beneficial for wide networks, with large \mathbf{W} , and **NTK-vector products** are beneficial for networks with large outputs \mathbf{O} . We confirm our predictions with FLOPs measurements in Fig. 1.

We further confirm our methods can provide orders of magnitude speed-ups and memory savings on all major hardware platforms in Fig. 1 (right) and Fig. 3. However, we notice that our wall-clock time measurements often deviate from predictions due to unaccounted constant overheads of various methods, hardware specifics, padding, and the (largely black-box) behavior of the **XLA** compiler. Notably, in practice, we find **Structured derivatives** almost always outperform **NTK-vector products**.

Finally, we evaluate our methods in the wild, and confirm computational benefits on full ImageNet models in Fig. 2 (ResNets, He et al. (2016)) and Fig. 4 (WideResNets, Zagoruyko & Komodakis (2016); Vision Transformers and Transformer-ResNet hybrids Dosovitskiy et al. (2021); Steiner et al. (2021); and MLP-Mixers Tolstikhin et al. (2021)). Computing the full $\mathbf{O} \times \mathbf{O} = 1000 \times 1000$ NTK for many of these models on modern accelerators is only possible with **Structured derivatives**.

4 STRUCTURED DERIVATIVES FOR GENERIC FUNCTIONS

Here we generalize the idea of leveraging structure in subexpressions presented in §3.4. This section (and our implementation) is not specific to NNs and applies to any differentiable function.

Consider two differentiable functions defined on a common input domain:

$$f_i : (\theta^0, \dots, \theta^{\mathbf{L}}) \in \mathbb{R}^{\mathbf{P}_0 \times \dots \times \mathbf{P}_{\mathbf{L}}} \mapsto f_i(\theta^0, \dots, \theta^{\mathbf{L}}) \in \mathbb{R}^{\mathbf{O}_i} \quad (i \in \{1, 2\}).$$

For NNs, typically $(\theta^0, \dots, \theta^{\mathbf{L}})$ correspond to trainable parameters in layers $0, \dots, \mathbf{L}$, and $f_i(\theta^0, \dots, \theta^{\mathbf{L}}) = f(\theta^0, \dots, \theta^{\mathbf{L}}, x_i)$, x_i being network inputs, $\mathbf{O}_i = \mathbf{O}$ being the number of outputs (logits, classes). The NTK is defined as

$$\Theta_{\theta}(f_1, f_2) := \sum_{l=0}^{\mathbf{L}} \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T \in \mathbb{R}^{\mathbf{O}_1 \times \mathbf{O}_2}. \quad (11)$$

Assume the following decompositions of f_i into computational graphs made of primitives y_i :

$$f_i(\theta^0, \dots, \theta^{\mathbf{L}}) = \tilde{f}_i(y_i^1(\theta^0, \dots, \theta^{\mathbf{L}}), \dots, y_i^{\mathbf{K}_i}(\theta^0, \dots, \theta^{\mathbf{L}})) \quad (i \in \{1, 2\}). \quad (12)$$

with $y_i^k(\theta^0, \dots, \theta^{\mathbf{L}}) \in \mathbb{R}^{\mathbf{Y}_i^k}$. In common NNs, y_i^k would correspond to pre-activations evaluated on inputs x_i in layer k_i , and, without weight sharing, typically $\mathbf{K}_1 = \mathbf{K}_2 = \mathbf{L}$. However, we do not impose any relationship between the number of parameter variables \mathbf{L} and number of primitives \mathbf{K}_1 and \mathbf{K}_2 , allowing arbitrary weight sharing. We can then use the chain rule in Eq. (2) to obtain:

$$\Theta_{\theta}(f_1, f_2) = \sum_{l, k_1, k_2}^{\mathbf{L}, \mathbf{K}_1, \mathbf{K}_2} \left(\frac{\partial \tilde{f}_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial \tilde{f}_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T = \sum_{l, k_1, k_2}^{\mathbf{L}, \mathbf{K}_1, \mathbf{K}_2} \frac{\partial \tilde{f}_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l}^T \frac{\partial \tilde{f}_2}{\partial y_2^{k_2}}^T. \quad (13)$$

All methods from §3 perform the sum of contractions in Eq. (13) one way or another. **Jacobian contraction** uses VJPs to implicitly contract each summand “outside-in”, i.e. it first computes $\partial f_i / \partial \theta^l$ terms with VJPs followed by their contraction. As discussed in §3.3, this costs $\mathbf{NO}[\mathbf{FP}] + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$.

NTK-vector products use both JVPs and VJPs to contract “Right-to-left”, i.e. first compute $\partial f_2 / \partial \theta^l$ as an implicit contraction of $\partial f_2 / \partial y_2$ with $\partial y_2 / \partial \theta^l$ via VJP, followed by an implicit contraction of the result with $\partial y_1 / \partial \theta^l$ via a JVP, followed by another implicit contraction with $\partial f_1 / \partial y_1$ with another JVP. Per §3.5 this costs $\mathbf{N}^2\mathbf{O}[\mathbf{FP}]$.

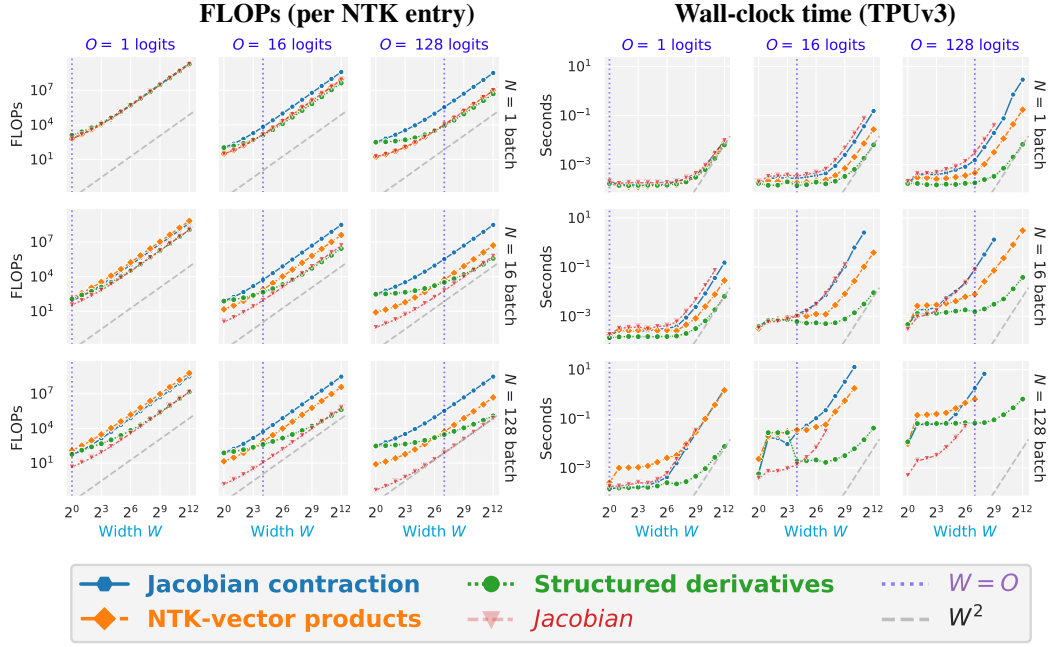


Figure 1: **FLOPs (left) and wall-clock time (right) of computing the NTK for a 10-layer ReLU FCN.** As predicted by Table 1, our methods almost always outperform **Jacobian contraction**, allowing orders of magnitude speed-ups and memory improvements for realistic problem sizes. **FLOPs per NTK entry:** We confirm several specific predictions: (1) **NTK-vector products** are the best performing method for $N = 1$, and have cost equivalent to **Jacobian** for any width W or output size O (top row); (2) **NTK-vector products** offer an O -fold improvement over **Jacobian contraction** (left to right columns); (3) **NTK-vector products** are equivalent to **Jacobian contraction** for $O = 1$ (leftmost column); (4) **Structured derivatives** outperform **NTK-vector products** iff $O < W$ ($O = W$ are plotted as pale vertical lines, which is where **Structured derivatives** and **NTK-vector products** intersect); (5) **Structured derivatives** approach the cost of **Jacobian** in the limit of large width W (left to right). (6) All methods, as expected, scale quadratically with width W . **Wall-clock runtime:** In real applications, given hardware-specific constraints, padding, and delicate interplay with the XLA compiler, we observe that: (1) **NTK-vector products** improve upon **Jacobian contraction** for $O > 1$, but the effect is not perfectly robust (see bottom row for small W and Fig. 3, notably GPU platforms); (2) **Structured derivatives** robustly outperform all other methods, including simply computing the **Jacobian**, as discussed in §3.4; (3) **Structured derivatives** have lower memory footprint, and reach up to 8x larger widths (bottom right; missing points indicate out-of-memory), i.e. can process models up to 64x larger than other methods, as discussed in §3.4; (4) All methods have a smaller memory footprint than **Jacobian** (see §3.1). **More:** Fig. 3 for other hardware platforms, §H for details.

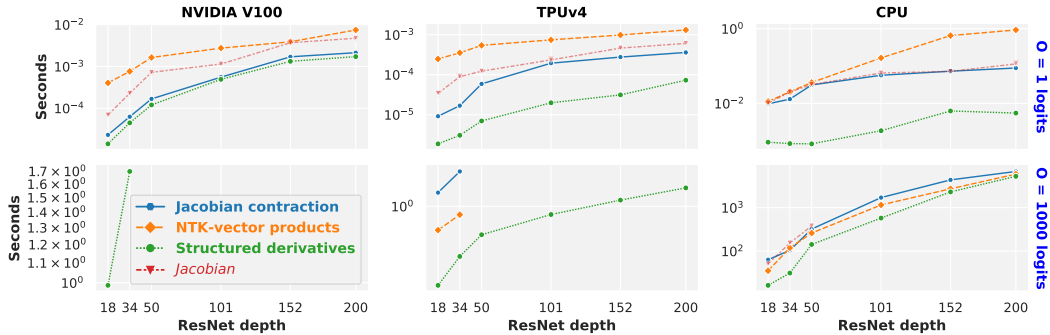


Figure 2: **Wall-clock time cost of computing an NTK for several ResNet sizes on a pair of ImageNet inputs.** **Structured derivatives** allow the NTK to be computed faster and for larger models (see bottom row – missing points indicate out-of-memory). **NTK-vector products**, as predicted in §3.6 and Table 2, are advantageous for large O (bottom row), but are suboptimal when the cost of the forward pass is large relative to the number of parameters, e.g. when there is a lot of weight sharing (see Table 7 and Table 2), which is the case for convolutions. See Fig. 4 for more ImageNet models, §F for analysis of CNN NTK computational complexity, and §H for experimental details.

However, recall from §3.4, while JVPs and VJPVs themselves are computationally optimal, higher-order computations like their contraction (Jacobian contraction) or composition (NTK-vector products) are generally not. The idea of Structured derivatives is to design rules for efficient computation of such contractions, similarly to how JAX has rules for efficient JVPs and VJPVs. From Eq. (13), in the general case this requires hand-made rules for all pairwise combinations of primitives y_1 and y_2 . Due to quadratic scaling in the number of primitives, we restrict the current implementation to rules that operate on individual primitives y . This still provides substantial computational benefit.

Specifically, our rules identify a few simple types of structure (e.g. block diagonal, constant-block diagonal, tiling) in $\partial y / \partial \theta^l$, that allow us to simplify the contraction in Eq. (13). In practice this amounts to replacing the inner terms $\partial y_1^{k_1} / \partial \theta^l$ and $\partial y_2^{k_2} / \partial \theta^l$ with their (much) smaller subarrays, and modifying the instructions passed to `np.einsum` that contracts all 4 terms. In §C we provide specific descriptions of our rules and their impact on the computational complexity of Eq. (13).

In Table 1 and Table 7 we show that our rules are asymptotically better than Jacobian contraction for matrix multiplications and convolutions, and verify that they are practically beneficial in a much wider set of operations used by contemporary ImageNet models in Fig. 2 and Fig. 4.

For both Structured derivatives and NTK-vector products a fully general and rigorous comparison of complexities is not feasible since it will rely upon specifics of the model architecture and the pairs of primitives, y_1 and y_2 , present in the network. Nonetheless, we can offer heuristics that suggest when each method will be beneficial. The time complexity of Structured derivatives has the form of $\text{NO}[\text{FP}] + \text{NOG} + \text{N}[\text{J} - \text{OP}]$, where G is related to the cost of contraction, and J to the cost of computing $\partial y / \partial \theta^l$ (exact values depend on the structure present in y_1 and y_2). This is guaranteed to be no worse than Jacobian contraction for FCNs and CNNs. From Table 2, the performance of NTK-vector products relative to Jacobian contraction ultimately depends on the cost of the forward pass through the network, $[\text{FP}]$, relative to OP . In practice this amounts to best performance on models without weight sharing like FCNs.

Owing to the nuanced trade-offs between different computational methods in the general case, we release all our implementations as a single function that allows the user to manually select the desired implementation. For convenience, we include an automated setting which will perform FLOPs analysis for each method at compilation time and automatically select the most efficient one.

5 IMPLEMENTATION

Both algorithms are implemented in JAX (Bradbury et al., 2018) as the following function transformation `ntk_fn` : $[f : (\theta, x) \mapsto f(\theta, x)] \mapsto [\Theta : (x_1, x_2, \theta) \mapsto \Theta_\theta(x_1, x_2)]$, i.e. our function accepts any function f with the above signature and returns the efficient NTK kernel function operating on inputs x_1 and x_2 and parameterized by θ . Inputs x , parameters θ , and outputs $f(\theta, x)$ can be arbitrary PyTrees. We rely on many utilities from JAX and Neural Tangents (Novak et al., 2020).

NTK-vector products algorithm is implemented by using JAX core operations such as `vjp`, `jvp`, and `vmap` to map the NTK-vp function to the I_0 matrix and to parallelize the computation over pairwise combinations of N inputs in each batch x_1 and x_2 .

Structured derivatives algorithm is implemented as a Jaxpr interpreter, built on top of the default JAX reverse mode AD interpreter. On a high level, the algorithm performs the sum in Eq. (13). Each summand is a contraction of 4 factors: $\partial \tilde{f}_1 / \partial y_1, \partial y_1 / \partial \theta, \partial y_2 / \partial \theta, \partial \tilde{f}_2 / \partial y_2$.

First, we linearize f to obtain a computational graph constructed out of a limited set (54,³ see Table 5) of linear primitives y^1, \dots, y^K . Then, we can obtain two factors $\partial \tilde{f}_1 / \partial y_1, \partial \tilde{f}_2 / \partial y_2$ as part of a backward pass almost identical to calling `jax.jacobian(f)(θ, x)`. To contract these terms with $\partial y_1 / \partial \theta$ and $\partial y_2 / \partial \theta$, as described above, we query a dictionary of rules which map primitives to a structural description (§C.8); for a given pair of primitives, these rules allow us to analytically simplify the contraction and avoid explicitly instantiating the derivatives.

³JAX leverages a similar approach to implement only 54 transpose rules for linear primitives for reverse mode differentiation instead of 131 VJP rules (Frostig et al., 2021).

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, 2016.
- Ben Adlam, Jaehoon Lee, Lechao Xiao, Jeffrey Pennington, and Jasper Snoek. Exploring the uncertainty properties of neural networks’ implicit priors in the infinite-width limit. In *International Conference on Learning Representations*, 2020.
- 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. Curran Associates, Inc., 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=rkl8sJBVvH>.
- Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *arXiv preprint arXiv:2102.06701*, 2021.
- 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.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. *arXiv preprint arXiv:2101.08692*, 2021a.
- Andrew Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. *arXiv preprint arXiv:2102.06171*, 2021b.
- Wuyang Chen, Xinyu Gong, and Zhangyang Wang. Neural architecture search on imagenet in four gpu hours: A theoretically inspired perspective. In *International Conference on Learning Representations*, 2021a.
- Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets without pretraining or strong data augmentations, 2021b.
- Yann Dauphin and Samuel S Schoenholz. Metainit: Initializing learning by learning to initialize. 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Póczos, Ruosong Wang, and Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*. 2019.

- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1126–1135. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/finn17a.html>.
- Jean-Yves Franceschi, Emmanuel de Bézenac, Ibrahim Ayed, Mickaël Chen, Sylvain Lamprier, and Patrick Gallinari. A neural tangent kernel perspective of gans. *arXiv preprint arXiv:2106.05566*, 2021.
- Roy Frostig, Matthew J Johnson, Dougal Maclaurin, Adam Paszke, and Alexey Radul. Decomposing reverse-mode automatic differentiation. *arXiv preprint arXiv:2105.09469*, 2021.
- Adrià Garriga-Alonso, Laurence Aitchison, and Carl Edward Rasmussen. Deep convolutional networks as shallow gaussian processes. In *International Conference on Learning Representations*, 2019.
- Andreas Griewank and Andrea Walther. *Evaluating Derivatives*. Society for Industrial and Applied Mathematics, second edition, 2008. doi: 10.1137/1.9780898717761. URL <https://epubs.siam.org/doi/abs/10.1137/1.9780898717761>.
- Roger Grosse. Neural net training dynamics, January 2021. URL https://www.cs.toronto.edu/~rgrosse/courses/csc2541.2021/readings/L02_Taylor_approximations.pdf.
- Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SJgndT4KwB>.
- Bobby He, Balaji Lakshminarayanan, and Yee Whye Teh. Bayesian deep ensembles via the neural tangent kernel. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/0b1ec366924b26fc98fa7b71a9c249cf-Abstract.html>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. URL <http://github.com/google/flax>.
- Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin. Haiku: Sonnet for JAX, 2020. URL <http://github.com/deepmind/dm-haiku>.
- Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: NNGP and NTK for deep attention networks. In *International Conference on Machine Learning*, 2020.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Sam Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- Jaehoon Lee, Lechao Xiao, Samuel S. 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*, 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. 2020.
- Dougal Maclaurin, David Duvenaud, and Ryan P Adams. Autograd: Effortless gradients in numpy.

- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- Herman Müntz. Solution directe de l’équation séculaire et de quelques problèmes analogues transcendents. *C. R. Acad. Sci. Paris*, 156:43–46, 1913.
- Uwe Naumann. Optimal accumulation of jacobian matrices by elimination methods on the dual computational graph. *Mathematical Programming*, 99(3):399–421, 2004.
- Uwe Naumann. Optimal jacobian accumulation is np-complete. *Mathematical Programming*, 112(2):427–441, 2008.
- Radford M. Neal. Priors for infinite networks (tech. rep. no. crg-tr-94-1). *University of Toronto*, 1994.
- Timothy Nguyen, Zhourong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-regression. *arXiv preprint arXiv:2011.00050*, 2020.
- Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. *arXiv preprint arXiv:2107.13034*, 2021.
- Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, 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.
- Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A. Alemi, Jascha Sohl-Dickstein, and Samuel S. Schoenholz. Neural tangents: Fast and easy infinite neural networks in python. In *International Conference on Learning Representations*, 2020. URL <https://github.com/google/neural-tangents>.
- Boris N. Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *NeurIPS*, 2018.
- Daniel S Park, Jaehoon Lee, Daiyi Peng, Yuan Cao, and Jascha Sohl-Dickstein. Towards nngp-guided neural architecture search. *arXiv preprint arXiv:2011.06006*, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Jeffrey Pennington and Yasaman Bahri. Geometry of neural network loss surfaces via random matrix theory. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2798–2806. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/pennington17a.html>.
- Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. *International Conference on Learning Representations*, 2017.
- Stefano Spigler, Mario Geiger, Stéphane d’Ascoli, Levent Sagun, Giulio Biroli, and Matthieu Wyart. A jamming transition from under-to over-parametrization affects generalization in deep learning. *Journal of Physics A: Mathematical and Theoretical*, 52(47):474001, 2019.

- Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv preprint arXiv:2106.10270*, 2021.
- Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020.
- Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. Mlp-mixer: An all-mlp architecture for vision, 2021.
- Lechao Xiao, Yasaman Bahri, Jascha Sohl-Dickstein, Samuel Schoenholz, and Jeffrey Pennington. Dynamical isometry and a mean field theory of CNNs: How to train 10,000-layer vanilla convolutional neural networks. In *International Conference on Machine Learning*, 2018.
- Lechao Xiao, Jeffrey Pennington, and Samuel S Schoenholz. Disentangling trainability and generalization in deep learning. In *International Conference on Machine Learning*, 2020.
- Sho Yaida. Non-Gaussian processes and neural networks at finite widths. In *Mathematical and Scientific Machine Learning Conference*, 2020.
- 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*, 2019.
- Greg Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020.
- Greg Yang, Jeffrey Pennington, Vinay Rao, Jascha Sohl-Dickstein, and Samuel S. Schoenholz. A mean field theory of batch normalization. In *International Conference on Learning Representations*, 2019.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.
- Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv preprint arXiv:1901.09321*, 2019.
- Yufan Zhou, Zhenyi Wang, Jiayi Xian, Changyou Chen, and Jinhui Xu. Meta-learning with neural tangent kernels. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ti87Pv50c8>.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. 2017. URL <https://arxiv.org/abs/1611.01578>.

APPENDIX

A ADDITIONAL FIGURES

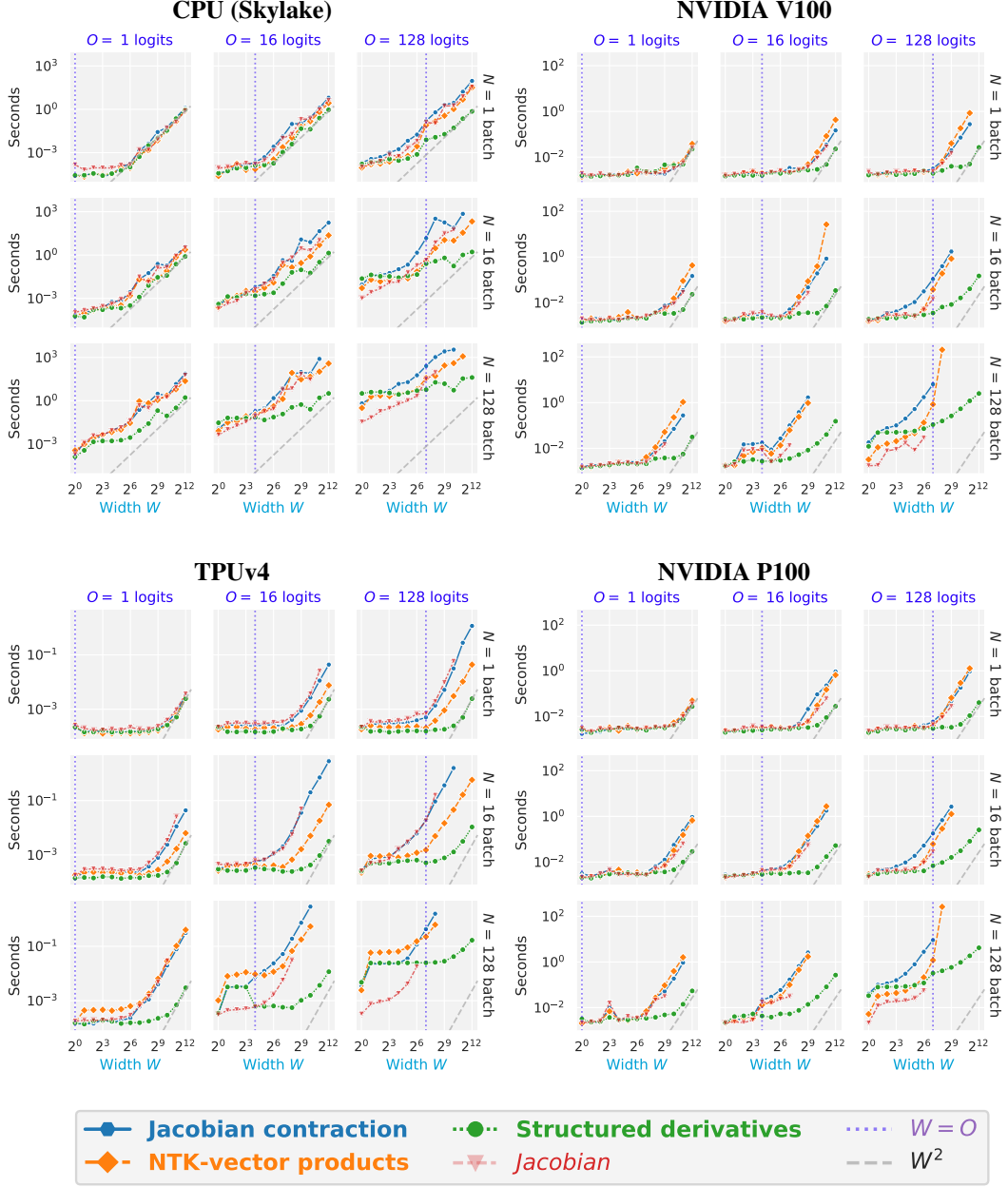


Figure 3: **Wall-clock time of computing the NTK of a 10-layer ReLU FCN on different platforms.** In all settings, **Structured derivatives** allow orders of magnitude improvement in wall-clock time and memory (missing points indicate out-of-memory error). However, we remark that on GPU platforms (right), **NTK-vector products** deliver a robust improvement only for large O (rightmost column), while for $O = 16$ the cost is comparable or even larger than **Jacobian contraction**. See Fig. 1 for FLOPs, TPUv3 platform, and more discussion. See §H for details.

B GLOSSARY

- **N** - batch size of inputs x to the NN $f(\theta, x)$. In a more general setting (§4), the number of functions $f(\theta)$.
 n - batch indices ranging from 1 to **N**.
- **O** - output size (e.g. number of logits) of the NN $f(\theta, x)$ for a single (**N** = 1) input x .
- The NTK matrix has shape **NO** \times **NO**.
- **W** - width of an FCN, or number of channels of a CNN. Individual inputs x are usually assumed to have the same size / number of channels.
- **L** - depth of the network, number of layers. In a more general setting (§4), number of trainable parameter matrices, that are used in a possibly different number of subexpressions in the network.
 l - depth index ranging from 0 to **L**.
- **K** - number of subexpressions (primitives, nodes in the computation graph) of the network $f(\theta, x)$. For NNs without weight sharing, **K** = **L**.
 k - subexpression/primitive index ranging from 1 to **K**.
- **D** - total number of pixels (e.g. 1024 for a 32×32 image; 1 for an FCN) in an input and every intermediate layer of a CNN (SAME or CIRCULAR padding is assumed, to consider the spatial size unchanged from layer to layer).
- **F** - total filter size (e.g. 9 for a 3×3 filter; 1 for an FCN) in a convolutional filter of a CNN (no striding and dilation is assumed for simplicity).
- **Y** - total size of a pre-activation / primitive / subexpression y (e.g. **Y** = **DW** for a layer with **D** pixels and **W** channels; **Y** = **W** for FCN). Depending on the context, can represent size of a single or particular pre-activation in the network, or the size of all pre-activations together.
- **C** - in §4, the size of the axis along which a subexpression derivative $\partial y / \partial \theta$ admits certain structure (**C** can often be equal to **Y** or a significant fraction of it, e.g. **W**).
 c - index along the structured axis, ranging from 1 to **C**.
- **P** - total size of trainable parameters. Depending on the context, can represent the size of a particular weight matrix θ^l in some layer l (e.g. **W**² for width-**W** FCN), or the size of all parameters in the network.
- **FP** - forward pass, cost (time or memory, depending on the context) of evaluating $f(\theta, x)$ on a single (**N** = 1) input x .
- If a variable is present in time or memory complexity analysis with an index such as k or l , it is considered to be the maximum over that index, e.g. $\mathbf{Y}^k = \max_k \mathbf{Y}^k$. This is used in Table 2 for brevity.
- \mathbf{J}_i^k - the (usually negligible) cost of evaluating a single primitive Jacobian $\partial y^k / \partial \theta^l$ subarray (§D), given the structure present in y^k according to §C.
- **G** - in Table 2 and §3.4, a variable related to the cost of contraction of Eq. (12), the precise value of which also depends on the structure present in the computation (see §C and Table 3).

C TYPES OF STRUCTURED DERIVATIVES

Here we continue §4 and list the types of structures in primitive derivatives $\partial y / \partial \theta$ that allow linear algebra simplifications of the NTK expression. Analysis from the following subsections is summarized in Table 3.

Structure of $\partial y / \partial \theta \downarrow$	Outside-in	Left-to-right	Inside-out
None w/ VJPs & JVPs:	$\mathbf{NO}[\mathbf{FP}] + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{O}[\mathbf{FP}]$	Not possible
None w/ explicit matrices	$\mathbf{NOYP} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{OYP} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P} + \mathbf{N}^2\mathbf{OY}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Block-diagonal	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C}^2 + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Constant block-diagonal	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C}^3 + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Input block-tiled	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C} + \mathbf{N}^2\mathbf{OY}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Output block-tiled	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P} + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C}^2 + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$
Block-tiled	$\mathbf{NOYP/C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{P/C} + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{OYP/C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C}^2 + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C}^3 + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$

Table 3: **Asymptotic time complexities of computing the contractions for NTK summands** $\Theta(f_1^{n_1}, f_2^{n_2}) (\theta^0, \dots, \theta^L)_l^{k_1, k_2} \in \mathbb{R}^{\mathbf{O} \times \mathbf{O}}$ in Eq. (14), for all n_1 and n_2 from 1 to \mathbf{N} (resulting in an $\mathbf{NO} \times \mathbf{NO}$ NTK matrix). Time complexity of **Structured derivatives** is the minimum (due to using `np.einsum` with optimal contraction order) of the row corresponding to the structure present in a pair of primitives $y_1^{k_1}$ and $y_2^{k_2}$. How it compares to **Jacobian contraction** and **NTK-vector products** (top row) depends on many variables, including the cost of evaluating the primitive **FP**. See Table 1 and Table 7 for exact comparison for matrix multiplication and convolution. See §B for legend.

C.1 NO STRUCTURE

We first consider the default cost of evaluating a single summand in Eq. (13), denoting individual matrix shapes underneath:

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) := \frac{\partial \tilde{f}_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2 T}}{\partial \theta^l} \frac{\partial \tilde{f}_2^T}{\partial y_2^{k_2}} =: \overbrace{\frac{\partial \tilde{f}_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2^T}{\partial \theta} \frac{\partial \tilde{f}_2^T}{\partial y_2}}^{\mathbf{O}_1 \times \mathbf{O}_2} \quad (14)$$

$\mathbf{O}_1 \times \mathbf{Y}_1 \quad \mathbf{Y}_1 \times \mathbf{P} \quad \mathbf{P} \times \mathbf{Y}_2 \quad \mathbf{Y}_2 \times \mathbf{O}_2$

We have dropped indices l, k_1 and k_2 on the right-hand side of Eq. (14) to avoid clutter, and consider $\theta := \theta^l, y_1 := y_1^{k_1}, y_2 := y_2^{k_2}$ until the end of this section. To simplify exposition, we also assume $\mathbf{O}_1 = \mathbf{O}_2 = \mathbf{O}$ and $\mathbf{Y}_1 = \mathbf{Y}_2 = \mathbf{Y}$.

Under this assumption, there are 3 ways of contracting Eq. (14) that cost

- (a) **Outside-in:** $\mathbf{OYP} + \mathbf{O}^2\mathbf{P}$
- (b) **Left-to-right and right-to-left:** $\mathbf{OYP} + \mathbf{O}^2\mathbf{Y}$.
- (c) **Inside-out-left and inside-out-right:** $\mathbf{Y}^2\mathbf{P} + \mathbf{OY}^2 + \mathbf{O}^2\mathbf{Y}$.

In the next sections, we look at how these costs are reduced given certain structure in $\partial y / \partial \theta$.

C.2 BLOCK-DIAGONAL

Assume $\partial y / \partial \theta = \oplus_{c=1}^{\mathbf{C}} \partial y^c / \partial \theta_c$, where \oplus stands for **direct sum of matrices**, i.e. $\partial y / \partial \theta$ is a block-diagonal matrix made of blocks $\{\partial y^c / \partial \theta_c\}_{c=1}^{\mathbf{C}}$, where $\partial y^c / \partial \theta_c$ have shapes $(\mathbf{Y/C}) \times (\mathbf{P/C})$. Here $\{y^c\}_{c=1}^{\mathbf{C}}$ and $\{\theta_c\}_{c=1}^{\mathbf{C}}$ are **partitions** of y and θ respectively. In NNs this structure is present in binary bilinear operations (on θ and another argument) such as multiplication, division, batched matrix multiplication, or depthwise convolution. Then Eq. (14) can be re-written as

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \frac{\partial \tilde{f}_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2^T}{\partial \theta} \frac{\partial \tilde{f}_2^T}{\partial y_2} \quad (15)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\oplus_{c=1}^{\mathbf{C}} \frac{\partial y_1^c}{\partial \theta_c} \right) \left(\oplus_{c=1}^{\mathbf{C}} \frac{\partial y_2^c}{\partial \theta_c} \right)^T \frac{\partial \tilde{f}_2^T}{\partial y_2} \quad (16)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\oplus_{c=1}^{\mathbf{C}} \left[\frac{\partial y_1^c}{\partial \theta_c} \frac{\partial y_2^{c T}}{\partial \theta_c} \right] \right) \frac{\partial \tilde{f}_2^T}{\partial y_2} \quad (17)$$

$$= \sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_1}{\partial y_1^c} \left[\frac{\partial y_1^c}{\partial \theta_c} \frac{\partial y_2^{c T}}{\partial \theta_c} \right] \frac{\partial \tilde{f}_2^T}{\partial y_2^c}, \quad (18)$$

where we have applied the block matrix identity

$$[A^1, \dots, A^{\mathbf{C}}]^T (\oplus_{c=1}^{\mathbf{C}} B^c) [D^1, \dots, D^{\mathbf{C}}] = \sum_{c=1}^{\mathbf{C}} A^c B^c D^c. \quad (19)$$

We now perform a complexity analysis similar to Eq. (14):

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \sum_{c=1}^{\mathbf{C}} \overbrace{\begin{matrix} \frac{\partial \tilde{f}_1}{\partial y_1^c} & \frac{\partial y_1^c}{\partial \theta_c} & \frac{\partial y_2^c}{\partial \theta_c}^T & \frac{\partial \tilde{f}_2}{\partial y_2^c}^T \\ \underbrace{\hspace{1.5cm}}_{\mathbf{O} \times (\mathbf{Y}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{P}/\mathbf{C}) \times (\mathbf{Y}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{Y}/\mathbf{C}) \times \mathbf{O}} \end{matrix}}^{\mathbf{O} \times \mathbf{O}}$$

In this case complexities of the three methods become

1. **Outside-in:** $\mathbf{OYP}/\mathbf{C} + \mathbf{O}^2\mathbf{P}$.
2. **Left-to-right and right-to-left:** $\mathbf{OYP}/\mathbf{C} + \mathbf{O}^2\mathbf{Y}$.
3. **Inside-out-left and inside-out-right:** $\mathbf{Y}^2\mathbf{P}/\mathbf{C}^2 + \mathbf{OY}^2/\mathbf{C} + \mathbf{O}^2\mathbf{Y}$.

C.3 CONSTANT BLOCK-DIAGONAL

Assume $\frac{\partial y}{\partial \theta} = I_{\mathbf{C}} \otimes \frac{\partial y^1}{\partial \theta_1}$, and $\frac{\partial y^1}{\partial \theta_1}$ has shape $(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})$. In NNs, this is present in fully-connected, convolutional, locally-connected, attention, and many other layers that contain a matrix multiplication along some axis. This is also present in all unary elementwise linear operations on θ like transposition, negation, reshaping and many others.

This is a special case of §C.2 with $\frac{\partial y^c}{\partial \theta_c} = \frac{\partial y^1}{\partial \theta_1}$ for any c . Here a similar analysis applies, yielding

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \sum_{c=1}^{\mathbf{C}} \overbrace{\begin{matrix} \frac{\partial \tilde{f}_1}{\partial y_1^c} & \frac{\partial y_1^1}{\partial \theta_1} & \frac{\partial y_2^1}{\partial \theta_1}^T & \frac{\partial \tilde{f}_2}{\partial y_2^c}^T \\ \underbrace{\hspace{1.5cm}}_{\mathbf{O} \times (\mathbf{Y}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{P}/\mathbf{C}) \times (\mathbf{Y}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{Y}/\mathbf{C}) \times \mathbf{O}} \end{matrix}}^{\mathbf{O} \times \mathbf{O}}$$

and the same contraction complexities as in §C.2, except for the **Inside-out** order, where the inner contraction term costs only $\mathbf{Y}^2\mathbf{P}/\mathbf{C}^3$, since it is only contracted once instead of \mathbf{C} times as in **Block-diagonal**.

C.4 INPUT BLOCK-TILED

Assume $\frac{\partial y}{\partial \theta} = \mathbb{1}_{(1, \mathbf{C})} \otimes \frac{\partial y}{\partial \theta_1}$, where $\mathbb{1}_{(1, \mathbf{C})}$ is an **all-ones matrix** of shape $1 \times \mathbf{C}$, and $\frac{\partial y}{\partial \theta_1}$ has shape $\mathbf{Y} \times (\mathbf{P}/\mathbf{C})$. In this case

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \frac{\partial \tilde{f}_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (20)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbb{1}_{(1, \mathbf{C})} \otimes \frac{\partial y_1}{\partial \theta_1} \right) \left(\mathbb{1}_{(1, \mathbf{C})} \otimes \frac{\partial y_2}{\partial \theta_1} \right)^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (21)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbf{C} \mathbb{1}_{(1, 1)} \otimes \left[\frac{\partial y_1}{\partial \theta_1} \frac{\partial y_2}{\partial \theta_1}^T \right] \right) \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (22)$$

$$= \mathbf{C} \frac{\partial \tilde{f}_1}{\partial y_1} \left[\frac{\partial y_1}{\partial \theta_1} \frac{\partial y_2}{\partial \theta_1}^T \right] \frac{\partial \tilde{f}_2}{\partial y_2}^T. \quad (23)$$

The matrix shapes are

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \mathbf{C} \overbrace{\begin{matrix} \frac{\partial \tilde{f}_1}{\partial y_1} & \frac{\partial y_1}{\partial \theta_1} & \frac{\partial y_2}{\partial \theta_1}^T & \frac{\partial \tilde{f}_2}{\partial y_2}^T \\ \underbrace{\hspace{1.5cm}}_{\mathbf{O} \times \mathbf{Y}} & \underbrace{\hspace{1.5cm}}_{\mathbf{Y} \times (\mathbf{P}/\mathbf{C})} & \underbrace{\hspace{1.5cm}}_{(\mathbf{P}/\mathbf{C}) \times \mathbf{Y}} & \underbrace{\hspace{1.5cm}}_{\mathbf{Y} \times \mathbf{O}} \end{matrix}}^{\mathbf{O} \times \mathbf{O}}$$

Which leads to the following resulting complexities:

1. **Outside-in:** $\text{OYP/C} + \text{O}^2\text{P}$.
2. **Left-to-right and right-to-left:** $\text{OYP/C} + \text{O}^2\text{Y}$.
3. **Inside-out and inside-out-right:** $\text{Y}^2\text{P/C} + \text{OY}^2 + \text{O}^2\text{Y}$.

C.5 OUTPUT BLOCK-TILED

Assume $\frac{\partial y}{\partial \theta} = \mathbb{1}_{(\mathbf{C},1)} \otimes \frac{\partial y^1}{\partial \theta}$, where $\frac{\partial y^1}{\partial \theta}$ has shape $(\mathbf{Y/C}) \times \mathbf{P}$. This occurs during broadcasting or broadcasted arithmetic operations. In this case

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \frac{\partial \tilde{f}_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (24)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C},1)} \otimes \frac{\partial y_1^1}{\partial \theta} \right) \left(\mathbb{1}_{(\mathbf{C},1)} \otimes \frac{\partial y_2^1}{\partial \theta} \right)^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (25)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \left[\frac{\partial y_1^1}{\partial \theta} \frac{\partial y_2^1}{\partial \theta}^T \right] \right) \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (26)$$

$$= \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_1}{\partial y_1^c} \right) \left[\frac{\partial y_1^1}{\partial \theta_1} \frac{\partial y_2^1}{\partial \theta_1}^T \right] \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_2}{\partial y_2^c} \right)^T, \quad (27)$$

where we have used a block matrix identity

$$[A^1, \dots, A^{\mathbf{C}}]^T (\mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes B) [D^1, \dots, D^{\mathbf{C}}] = \left(\sum_{c=1}^{\mathbf{C}} A^c \right) B \left(\sum_{c=1}^{\mathbf{C}} D^c \right).$$

Finally, denoting the shapes,

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \overbrace{\left(\underbrace{\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_1}{\partial y_1^c}}_{\mathbf{O} \times (\mathbf{Y/C})} \underbrace{\frac{\partial y_1^1}{\partial \theta}}_{(\mathbf{Y/C}) \times \mathbf{P}} \underbrace{\frac{\partial y_2^1}{\partial \theta}^T}_{\mathbf{P} \times (\mathbf{Y/C})} \underbrace{\left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_2}{\partial y_2^c} \right)^T}_{(\mathbf{Y/C}) \times \mathbf{O}} \right)}^{\mathbf{O} \times \mathbf{O}},$$

complexities of the three methods become (notice we add an OY term to perform the sums)

1. **Outside-in:** $\text{OYP/C} + \text{O}^2\text{P} + \text{OY}$.
2. **Left-to-right:** $\text{OYP/C} + \text{O}^2\text{Y/C} + \text{OY}$.
3. **Inside-out:** $\text{Y}^2\text{P/C}^2 + \text{OY}^2/\text{C}^2 + \text{O}^2\text{Y/C} + \text{OY}$.

C.6 BLOCK-TILED

Assume $\frac{\partial y}{\partial \theta} = \mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \frac{\partial y^1}{\partial \theta_1}$, where $\frac{\partial y^1}{\partial \theta_1}$ has shape $(\mathbf{Y/C}) \times (\mathbf{P/C})$. This occurs for instance when y is a constant. In this case

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \frac{\partial \tilde{f}_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (28)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \frac{\partial y_1^1}{\partial \theta_1} \right) \left(\mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \frac{\partial y_2^1}{\partial \theta_1} \right)^T \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (29)$$

$$= \frac{\partial \tilde{f}_1}{\partial y_1} \left(\mathbf{C} \mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \left[\frac{\partial y_1^1}{\partial \theta_1} \frac{\partial y_2^1}{\partial \theta_1}^T \right] \right) \frac{\partial \tilde{f}_2}{\partial y_2}^T \quad (30)$$

$$= \mathbf{C} \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_1}{\partial y_1^c} \right) \left[\frac{\partial y_1^1}{\partial \theta_1} \frac{\partial y_2^1}{\partial \theta_1}^T \right] \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_2}{\partial y_2^c} \right)^T, \quad (31)$$

This results in the following contraction:

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \mathbf{C} \overbrace{\left(\underbrace{\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_1}{\partial y_1^c}}_{\mathbf{O} \times (\mathbf{Y}/\mathbf{C})} \underbrace{\frac{\partial y_1^1}{\partial \theta_1}}_{(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})} \underbrace{\frac{\partial y_2^1}{\partial \theta_1}}_{(\mathbf{P}/\mathbf{C}) \times (\mathbf{Y}/\mathbf{C})} \underbrace{\left(\sum_{c=1}^{\mathbf{C}} \frac{\partial \tilde{f}_2}{\partial y_2^c} \right)^T}_{(\mathbf{Y}/\mathbf{C}) \times \mathbf{O}} \right)}^{\mathbf{O} \times \mathbf{O}},$$

with final complexities of

1. **Outside-in:** $\mathbf{OYP}/\mathbf{C}^2 + \mathbf{O}^2\mathbf{P} + \mathbf{OY}$.
2. **Left-to-right:** $\mathbf{OYP}/\mathbf{C}^2 + \mathbf{O}^2\mathbf{Y}/\mathbf{C}^2 + \mathbf{OY}$.
3. **Inside-out:** $\mathbf{Y}^2\mathbf{P}/\mathbf{C}^3 + \mathbf{OY}^2/\mathbf{C}^2 + \mathbf{O}^2\mathbf{Y}/\mathbf{C} + \mathbf{OY}$.

C.7 BATCHED NTK COST ANALYSIS

For simplicity, we have considered evaluating the NTK $\Theta(f_1, f_2)$ on a single pair of functions f_1 and f_2 . In practice one is almost always interested in computing the NTK for all pairs of functions $f_1^{n_1}$ and $f_2^{n_2}$ from two batches $\{f_1^{n_1}\}_{n_1=1}^{\mathbf{N}_1}$ and $\{f_2^{n_2}\}_{n_2=1}^{\mathbf{N}_2}$, resulting in a $\mathbf{N}_1\mathbf{O}_1 \times \mathbf{N}_2\mathbf{O}_2$ NTK matrix. In common NNs, this corresponds to having batches of \mathbf{N}_1 and \mathbf{N}_2 inputs x_1 and x_2 respectively, and having $f_1^{n_i}(\theta^0, \dots, \theta^{\mathbf{L}}) := f(\theta^0, \dots, \theta^{\mathbf{L}}, x_i^{n_i})$. In this case the same argument as in previous section follows (given identical assumptions for all n_1 and n_2), but the cost of contractions involving terms from different batches grow by a multiplicative factor of $\mathbf{N}_1\mathbf{N}_2$, while all other costs grow by a factor of \mathbf{N}_1 or \mathbf{N}_2 . To declutter notation we consider $\mathbf{N}_1 = \mathbf{N}_2 = \mathbf{N}$, and summarize resulting batched costs in Table 3.

C.8 COMPLEX STRUCTURE COST ANALYSIS

In previous sections and in §4, we have considered $\partial y_1/\partial \theta$ and $\partial y_2/\partial \theta$ admitting the same, and at most one kind of structure. While this is a common case, in general these derivatives may admit multiple types of structures along multiple axes (for instance, addition is **Constant block-diagonal** along non-broadcasted axes, and **Output block-tiled** along the broadcasted axes), and $\partial y_1/\partial \theta$ and $\partial y_2/\partial \theta$ may have different types of structures and respective axes, if the same weight θ is used in multiple different subexpressions of different kind. In such cases, equivalent optimizations are possible (and are implemented in the code) along the largest common subsets of axes for each type of structure that $\partial y_1/\partial \theta$ and $\partial y_2/\partial \theta$ have.

For example, let θ be a matrix in $\mathbb{R}^{\mathbf{W} \times \mathbf{W}}$, y_1 be multiplication by a scalar $y_1(\theta) = 2\theta$, and y_2 be matrix-vector multiplication $y_2(\theta) = \theta x$, $x \in \mathbb{R}^{\mathbf{W}}$. In this case $\partial y_1/\partial \theta = 2I_{\mathbf{W}} \otimes I_{\mathbf{W}}$, i.e. it is **Constant block-diagonal** along axes both 1 and 2. $\partial y_2/\partial \theta = I_{\mathbf{W}} \otimes x^T$, i.e. it is also **Constant block-diagonal**, but only along axis 1. Hence, the NTK term containing $\partial y_1/\partial \theta$ and $\partial y_2/\partial \theta$ will be computed with **Constant block-diagonal** simplification along axis 1. There are probably more computationally optimal ways of processing different structure combinations, as well as more types of structures that could be leveraged for NTK computation, and we intend to investigate it in future work.

C.9 EXAMPLE

In §4 and previous sections we have demonstrated how structure in primitive derivatives $\partial y/\partial \theta$ can be leveraged to reduce the cost of computing NTK. In this section we will consider a simple example of applying the framework of structured derivatives to FCNs to reproduce Table 1. See §F for equivalent application for CNNs.

As in §3, we consider a deep FCN with width \mathbf{W} and \mathbf{O} outputs. We assume the network is deep and/or wide enough to ignore the size of inputs x , and we ignore biases. In this case the number of parameters is quadratic in width $\mathbf{P} \sim \mathbf{W}^2$, and intermediate primitive outputs have the same size as the width, $\mathbf{Y} = \mathbf{W}$. As in §3.4, we recognize that individual primitives $y^{k,n}(\theta_k) = \theta_k x^{k,n}$, as matrix multiplications ($\theta_k \in \mathbb{R}^{\mathbf{W} \times \mathbf{W}}$, $x^{k,n} \in \mathbb{R}^{\mathbf{W}}$) admit the **Constant block-diagonal** structure ($\partial y^{k,n}/\partial \theta_k = I_{\mathbf{W}} \otimes x^{k,nT}$) with $\mathbf{C} = \mathbf{Y} = \mathbf{W}$. Finally, **FP** costs \mathbf{W}^2 , and $\mathbf{J} = \mathbf{YP}/\mathbf{C}^2 = \mathbf{W}$, i.e.

Structure of $\partial y / \partial \theta \downarrow$	Outside-in	Left-to-right	Inside-out
None w/ JVPs and VJPs	$N^2 O^2 W^2$	$N^2 O W^2$	Not possible
None	$N W^3 O + N^2 O^2 W^2$	$N^2 O W^3 + N^2 O^2 W$	$N^2 W^4 + N O W^2 + N^2 O^2 W$
Constant block-diagonal	$N^2 O^2 W^2$	$N^2 O W^2 + N^2 O^2 W$	$N^2 O^2 W$

Table 4: **Asymptotic time complexities of computing a single fully-connected layer NTK contribution.** See §C.9 for discussion, Table 3 for a more general setting, Table 1 for the case of deep networks, and §B for detailed legend.

is negligible. Substituting all these equalities into Table 3 we get a simplified Table 4, that confirms the benefits of **NTK-vector products** and **Structured derivatives** for FCNs.

D JACOBIAN RULES FOR STRUCTURED DERIVATIVES

Here we discuss computing primitive $\partial y / \partial \theta$ derivatives as part of our implementation in §5. We provide 4 options to compute them through arguments `j_rules` and `fwd`:

1. **Forward mode**, `fwd = True`, is equivalent to `jax.jacfwd`, forward mode Jacobian computation, performed by applying the JVP to **P** columns of the I_P identity matrix. Best for $P < Y$.
2. **Reverse mode**, `fwd = False`, is equivalent to `jax.jacrev`, reverse mode Jacobian computation, performed by applying the VJP to **Y** columns of the I_Y identity matrix. Best for $P > Y$.
3. **Automatic mode**, `fwd = None`, selects forward or reverse mode for each primitive based on parameters and output shapes.
4. **Rule mode**, `j_rules = True`, queries a dictionary of Jacobian rules (similar to the dictionary of structure rules) with our custom implementations of primitive Jacobians, instead of computing them through VJPs or JVPs. The reason for introducing custom rules follows our discussion in §3.4: while JAX has computationally optimal VJP and JVP rules, respective Jacobian computations are not guaranteed to be most efficient. In practice, we find our rules to be most often faster, however this effect is not perfectly consistent (can occasionally be slower) and often negligible, requiring further investigation.

The default setting is `j_rules = True`, `fwd = None`, i.e. a custom Jacobian implementation is preferred, and, if absent, Jacobian is computed in forward or reverse mode based on parameters and output sizes. Note that in all settings, structure of $\partial y / \partial \theta$ is used to compute only the smallest Jacobian subarray necessary, and therefore most often inputs to VJP/JVP will be smaller identity matrices $I_{P/C}$ or $I_{Y/C}$ respectively, and all methods will return a smaller Jacobian matrix of size $(Y/C) \times (P/C)$. We denote the (usually negligible) memory costs of these sub-arrays as **J**. If for any reason (for example debugging) you want the whole $\partial y / \partial \theta$ Jacobians computed, you can set the `a_rules=False`, i.e. disable structure rules.

E KNOWN ISSUES

We will continue improving our function transformations in various ways after release, and welcome bug reports and feature requests. Below are the missing features / issues at the time of submission:

1. No support for complex differentiation.
2. Not tested on functions with advanced JAX primitives like parallel collectives (`jax.lax.psum`, `jax.lax.pmean`, etc.), gradient checkpointing (`jax.remat`), compiled loops (`jax.lax.scan`; Python loops are supported).
3. Our current implementation of **NTK-vector products** relies on XLA’s common subexpression elimination (CSE) in order to reuse computation across different pairs of inputs x_1 and x_2 , and, as shown in Fig. 1 and Fig. 3, can have somewhat unpredictable wall-clock

Transposable primitive in <code>jax.ad.primitive_transposes</code>	Constant block-diagonal	Block-diagonal	Output block-tiled
add	✓		✓
add_any	✓		✓
all_gather			
all_to_all			
broadcast_in_dim	✓		✓
call			
complex	✓		
concatenate			
conj	✓		
conv_general_dilated			
convert_element_type	✓		
cumsum			
custom_lin			
custom_linear_solve			
device_put	✓		
div	✓	✓	
dot_general	✓	✓	
dynamic_slice			
dynamic_update_slice			
fft			
gather			
imag	✓		
linear_call			
mul	✓	✓	
named_call			
neg	✓		
pad	✓		
pdot			
ppermute			
psum			
real	✓		
reduce_sum	✓		
reduce_window_sum	✓		
remat_call			
reshape	✓		
rev	✓		
scatter			
scatter-add			
scatter-mul			
select			
select_and_gather_add			
select_and_scatter_add			
sharding_constraint			
sharding_constraint			
slice			
squeeze	✓		
sub	✓		✓
transpose	✓		
triangular_solve			
while			
xla_call			
xla_pmap			
xmap			
zeros_like	✓		

Table 5: **List of all linear primitives and currently implemented Structured derivatives rules.** In the future, more primitives and more rules can be supported, yet at the time of writing even the small set currently covered enables dramatic speed-up and memory savings in contemporary ImageNet models as in Fig. 2 and Fig. 4.

time performance and memory requirements. We believe this could correspond to CSE not always working perfectly, and are looking into a more explicitly efficient implementation.

F COMPLEXITY ANALYSIS FOR CONVOLUTIONAL NETWORKS

Here we go through the same analysis as in §3 for the case of convolution, where before the top layer \mathbf{L} global average pooling is applied. In this case the weights of the network θ are expanded by the total filter size \mathbf{F} , and inputs x , pre-activations y^l and post-activations x^l become matrices of shape $\mathbf{D} \times \mathbf{W}$, where \mathbf{D} is the total number of pixels. See Fig. 5 for visual depiction. We will again assume that $\mathbf{O} = \mathcal{O}(\mathbf{LW})$.

F.1 JVP AND VJP

Forward pass, JVP, and VJP costs [cost of all intermediate layers] + [cost of the top layer] = $[\mathbf{LDFW}^2] + [\mathbf{OW}] \sim \mathbf{LDFW}^2$ time. Forward pass and JVP require [size of all weights] + [size of activations at a single layer] = $[\mathbf{LFW}^2 + \mathbf{OW}] + [\mathbf{DW} + \mathbf{O}] \sim \mathbf{LFW}^2 + \mathbf{DW}$ memory. VJP requires [size of all weights] + [size of activations in all layers] + [size of a single weight matrix] = $[\mathbf{LFW}^2 + \mathbf{OW}] + [\mathbf{LDW} + \mathbf{O}] + [\mathbf{FW}^2 + \mathbf{OW}] \sim \mathbf{LFW}^2 + \mathbf{LDW}$ memory.

Batched inputs. Time cost of JVP and VJP increase linearly in \mathbf{N} up to \mathbf{NLDFW}^2 . JVP memory cost becomes [size of all weights] + \mathbf{N} [size of activations at a single layer] = $[\mathbf{LFW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{DW} + \mathbf{O}] \sim \mathbf{LFW}^2 + \mathbf{NDW} + \mathbf{NO}$. VJP memory cost becomes [size of all weights] + \mathbf{N} [size of activations in all layers] + \mathbf{N} [size of a single weight matrix] = $[\mathbf{LFW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LDW} + \mathbf{O}] + \mathbf{N}[\mathbf{FW}^2 + \mathbf{OW}] \sim \mathbf{LFW}^2 + \mathbf{NLDW} + \mathbf{NFW}^2 + \mathbf{NOW}$.

- JVP costs \mathbf{NLDFW}^2 time and $\mathbf{LFW}^2 + \mathbf{NDW} + \mathbf{NO}$ memory.
- VJP costs \mathbf{NLDFW}^2 time and $\mathbf{LFW}^2 + \mathbf{NLDW} + \mathbf{NFW}^2 + \mathbf{NOW}$ memory.

F.2 JACOBIAN

Computing the Jacobian costs \mathbf{O} times the cost of VJP, hence time is $\mathbf{ON}([\text{cost of all intermediate layers}] + [\text{cost of the top layer}]) = \mathbf{ON}([\mathbf{LDFW}^2] + [\mathbf{OW}]) \sim \mathbf{NLODFW}^2 + \mathbf{NO}^2\mathbf{W}$. Memory is [size of all weights] + \mathbf{N} [size of activations in all layers] + \mathbf{ON} [size of a single weight matrix] + \mathbf{ON} [activations in a single layer] = $[\mathbf{LFW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LDW} + \mathbf{O}] + \mathbf{ON}[\mathbf{FW}^2 + \mathbf{OW}] + \mathbf{ON}[\mathbf{DW}] \sim \mathbf{LW}^2 + \mathbf{NLDW} + \mathbf{NODW} + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W}$

Jacobian costs $\mathbf{NLODFW}^2 + \mathbf{NO}^2\mathbf{W}$ time and $\mathbf{LW}^2 + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLDW} + \mathbf{NODW}$ memory.

F.3 JACOBIAN CONTRACTION

Since weight matrices are increased by \mathbf{F} , the contraction cost goes up to $\mathbf{N}^2\mathbf{LO}^2\mathbf{FW}^2$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOFW}^2$ memory. The cost of computing the Jacobian is also modified (§F.2), which results in $\mathbf{N}^2\mathbf{LO}^2\mathbf{FW}^2 + \mathbf{NLODFW}^2 + \mathbf{NO}^2\mathbf{W} \sim \mathbf{N}^2\mathbf{LO}^2\mathbf{FW}^2 + \mathbf{NLODFW}^2$ time and $(\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOFW}^2) + (\mathbf{LFW}^2 + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLDW} + \mathbf{NODW}) \sim \mathbf{N}^2\mathbf{O}^2 + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLDW} + \mathbf{NODW} + \mathbf{LFW}^2$ memory.

Jacobian contraction costs $\mathbf{N}^2\mathbf{LO}^2\mathbf{FW}^2 + \mathbf{NLODFW}^2$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLDW} + \mathbf{NODW} + \mathbf{LFW}^2$ memory.

Structure of $\partial y / \partial \theta \downarrow$	Outside-in	Left-to-right	Inside-out
Constant block-diagonal	$\mathbf{N} \mathbf{O} \mathbf{F} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 \mathbf{F} \mathbf{W}^2$	$\mathbf{N}^2 \mathbf{O} \mathbf{F} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 \mathbf{D} \mathbf{W}$	$\mathbf{N}^2 \mathbf{D}^2 \mathbf{F} \mathbf{W} + \mathbf{N}^2 \mathbf{O} \mathbf{D}^2 \mathbf{W} + \mathbf{N}^2 \mathbf{O}^2 \mathbf{D} \mathbf{W}$

Table 6: **Asymptotic time complexity of contracting** $\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2)$ **corresponding to a CNN primitive** obtained by substituting $\mathbf{Y} = \mathbf{D}\mathbf{W}$, $\mathbf{C} = \mathbf{W}$, and $\mathbf{P} = \mathbf{F}\mathbf{W}^2$ into Table 3. The time cost of **Structured derivatives** is the minimum of the three entries due to using optimal contraction path by `np.einsum`.

Method	Time	Memory	
Jacobian contraction	$\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \mathbf{F} \mathbf{W}^2 + \mathbf{D} \mathbf{F} \mathbf{N} \mathbf{L} \mathbf{O} \mathbf{W}^2$	$\mathbf{N} \mathbf{O} \mathbf{F} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{D} \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{D} \mathbf{N} \mathbf{O} \mathbf{W} + \mathbf{L} \mathbf{F} \mathbf{W}^2$	$\mathbf{D} > \mathbf{O} \mathbf{W}$
NTK-vector products	$\mathbf{N}^2 \mathbf{O}^2 \mathbf{W} + \mathbf{D} \mathbf{F} \mathbf{N}^2 \mathbf{L} \mathbf{O} \mathbf{W}^2$	$\mathbf{N} \mathbf{O} \mathbf{F} \mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{D} \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{D} \mathbf{N} \mathbf{O} \mathbf{W} + \mathbf{L} \mathbf{F} \mathbf{W}^2$	$\mathbf{N} = 1$
Structured derivatives	$\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \min \left(\mathbf{F} \mathbf{W}^2, \mathbf{D} \mathbf{W} + \frac{\mathbf{D} \mathbf{F} \mathbf{W}^2}{\mathbf{O}}, \mathbf{D} \mathbf{W} + \frac{\mathbf{D}^2 \mathbf{W}}{\mathbf{O}} + \frac{\mathbf{D}^2 \mathbf{F} \mathbf{W}}{\mathbf{O}^2} \right) + \mathbf{D} \mathbf{F} \mathbf{N} \mathbf{L} \mathbf{O} \mathbf{W}^2$	$\mathbf{N} \mathbf{D} \mathbf{F} \mathbf{W} + \mathbf{N}^2 \mathbf{O}^2 + \mathbf{D} \mathbf{N} \mathbf{L} \mathbf{W} + \mathbf{D} \mathbf{N} \mathbf{O} \mathbf{W} + \mathbf{L} \mathbf{F} \mathbf{W}^2$	$\mathbf{D} < \mathbf{O} \mathbf{W}$

Table 7: **Asymptotic time and memory cost of computing the NTK for a CNN with global average pooling.** Costs are for a pair of batches of inputs of size \mathbf{N} each, and for \mathbf{L} -deep, \mathbf{W} -wide CNN with \mathbf{O} outputs and \mathbf{D} pixels in each layer, with filter size \mathbf{F} . Resulting NTK has shape $\mathbf{N} \mathbf{O} \times \mathbf{N} \mathbf{O}$. **Structured derivatives** reduce both time and memory complexity (memory under a mild condition of $\mathbf{D} < \mathbf{O} \mathbf{W}$), and are asymptotically beneficial over **Jacobian contraction** under a set of conditions on \mathbf{F} , \mathbf{W} , \mathbf{D} , and \mathbf{O} . **Note:** presented are asymptotic cost estimates; in practice, all methods incur large constant multipliers (e.g. at least 3x for time; see §3.1). However, this generally does not impact the relative performance of different methods. See §F for discussion, Table 1 for FCN, and Table 2 for generic cost analysis.

F.4 STRUCTURED DERIVATIVES

Convolution is **Constant block-diagonal** along the output channel axis with $\mathbf{C} = \mathbf{W}$, $\mathbf{P} = \mathbf{F}\mathbf{W}^2$, $\mathbf{Y} = \mathbf{D}\mathbf{W}$. Substituting this in Table 3, the cost of contraction is the minimum of the costs from Table 6. If we exclude the **Inside-out** contraction path from `np.einsum` (in practice it will always select the best out of three) for simplicity, we can and conclude that for \mathbf{L} layers, the time cost of the contraction is at most $\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \min(\mathbf{F}\mathbf{W}^2, \mathbf{D}\mathbf{W}) + \mathbf{D}\mathbf{F}\mathbf{N}\mathbf{L}\mathbf{O}\mathbf{W}^2$, as the minimum cost between the **Outside-in** and **Left-to-right**. Note that this dominates the time cost of the Jacobian from §F.2, so we don't need to modify it further. Memory due to Jacobian computation is [size of all weights] + \mathbf{N} [size of activations in all layers] + $\mathbf{N}\mathbf{O}$ [activations in a single layer] + [size of primitive derivatives] = $[\mathbf{L}\mathbf{F}\mathbf{W}^2 + \mathbf{O}\mathbf{W}] + \mathbf{N}[\mathbf{L}\mathbf{D}\mathbf{W} + \mathbf{O}] + \mathbf{N}\mathbf{O}[\mathbf{D}\mathbf{W}] + \mathbf{N}[\mathbf{D}\mathbf{W}] \sim \mathbf{L}\mathbf{F}\mathbf{W}^2 + \mathbf{N}\mathbf{L}\mathbf{D}\mathbf{W} + \mathbf{N}\mathbf{O}\mathbf{D}\mathbf{W}$. For **Constant block-diagonal** structure, $\mathbf{J} = \mathbf{Y}\mathbf{P}/\mathbf{C}^2 = \mathbf{D}\mathbf{F}\mathbf{W}$, negligible under a mild condition of $\mathbf{D} < \mathbf{O}\mathbf{W}$ compared to the memory savings (due to not needing to compute or store $\partial f / \partial \theta^l$ derivatives) of $\mathbf{N}\mathbf{O}\mathbf{F}\mathbf{W}^2 + \mathbf{N}\mathbf{O}^2 \mathbf{W}$.

F.5 NTK-VECTOR PRODUCTS

The cost of this approach is asymptotically equivalent to the cost of Jacobian (§F.2), since it consists of \mathbf{O} VJPs followed by \mathbf{O} (cheaper) JVPs. Therefore it costs $\mathbf{L}\mathbf{O}\mathbf{D}\mathbf{F}\mathbf{W}^2 + \mathbf{O}^2 \mathbf{W}$ time and $\mathbf{L}\mathbf{F}\mathbf{W}^2 + \mathbf{O}\mathbf{F}\mathbf{W}^2 + \mathbf{O}^2 \mathbf{W} + \mathbf{L}\mathbf{D}\mathbf{W} + \mathbf{O}\mathbf{D}\mathbf{W}$ memory.

Batched inputs. In a batched setting Eq. (10) is repeated for each pair of inputs, and therefore time increases by a factor of \mathbf{N}^2 to become $\mathbf{N}^2 \mathbf{L}\mathbf{O}\mathbf{D}\mathbf{F}\mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 \mathbf{W}$. Memory only grows linearly in \mathbf{N} (except for storing the result of size $\mathbf{N}^2 \mathbf{O}^2$), by similar argument to §3.5, i.e. becomes $\mathbf{N}^2 \mathbf{O}^2 + (\mathbf{L}\mathbf{F}\mathbf{W}^2 + \mathbf{N}\mathbf{O}\mathbf{F}\mathbf{W}^2 + \mathbf{N}\mathbf{O}^2 \mathbf{W} + \mathbf{N}\mathbf{L}\mathbf{D}\mathbf{W} + \mathbf{N}\mathbf{O}\mathbf{D}\mathbf{W})$ total memory.

NTK computation as a sequence of **NTK-vector products** costs $\mathbf{N}^2 \mathbf{L}\mathbf{O}\mathbf{D}\mathbf{F}\mathbf{W}^2 + \mathbf{N}^2 \mathbf{O}^2 \mathbf{W}$ time and $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N}\mathbf{O}\mathbf{F}\mathbf{W}^2 + \mathbf{L}\mathbf{F}\mathbf{W}^2 + \mathbf{N}\mathbf{L}\mathbf{D}\mathbf{W} + \mathbf{N}\mathbf{O}\mathbf{D}\mathbf{W}$ memory.

G ADDITIONAL DERIVATIONS

Below we derive Eq. (9) in §3.4:

$$\Theta_{\theta}^l(x_1, x_2) = \frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \left(I_W \otimes x_1^{lT} \right) \left(I_W \otimes x_2^{lT} \right)^T \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}^T = \quad (32)$$

$$= \frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \left(I_W \otimes \begin{bmatrix} x_1^{lT} & x_2^{lT} \end{bmatrix} \right) \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}^T = \quad (33)$$

$$= \begin{pmatrix} x_1^{lT} & x_2^{lT} \end{pmatrix} \begin{bmatrix} \frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} & \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}^T \end{bmatrix}, \quad (34)$$

where we were able to pull out $\begin{pmatrix} x_1^{lT} & x_2^{lT} \end{pmatrix}$ since it is a scalar.

Note that similarly to K-FAC (Martens & Grosse, 2015), in this example we leverage the structure in the FCN pre-activation derivative w.r.t. parameters, and we use the mixed-product property, i.e. $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$. However, in the general case this is not enough, and **Structured derivatives** rely on three components: (1) a direct sum linear algebra Eq. (19), (2) symbolic simplification of expressions with an identity matrix (§C.3), and (3) optimal order of contractions in Eq. (12) (e.g. “Inside-out” (Table 4), which is not possible to achieve with standard AD tools). To our knowledge, all three components are necessary to achieve our asymptotic complexities, and cannot be achieved by leveraging the mixed-product property alone.

H EXPERIMENTAL DETAILS

All experiments were performed in JAX (Bradbury et al., 2018) using 32-bit precision.

Throughout this work we assume the cost of multiplying two matrices of shapes (M, K) and (K, P) to be MKP . While there are **faster algorithms** for very large matrices, the XLA compiler (used by JAX, among other libraries) does not implement them, so our assumption is accurate in practice.

Hardware. CPU experiments were run on Dual 28-core Intel Skylake CPUs with at least 240 GiB of RAM. NVIDIA V100 and NVIDIA P100 used a respective GPU with 16 GiB GPU RAM. TPUv3 and TPUv4 have 8 and 32 GiB of RAM respectively, and use the default 16/32-bit mixed precision.

Fig. 5 is adapted from Novak et al. (2019) with authors’ permission.

Fig. 1 and Fig. 3: a 10-layer, ReLU FCN was constructed with the Neural Tangents (Novak et al., 2020) `nt.stax` API. Default settings (weight variance 1, no bias) were used. Individual inputs x had size 3. **Jacobian contraction** was evaluated using `nt.empirical_ntk_fn` with `trace_axes=()`, `diagonal_axes=()`, `vmap_axes=0`. **Jacobian** was evaluated using `jax.jacobian` with a `vmap` over inputs x . For time measurements, all functions were `jax.jit`ted, and timing was measured as the average of 100 random samples (compilation time was not included). For FLOPs, the function was not JITted, and FLOPs were measured on CPU using the `utils.get_flops` function that is released together with our code.⁴

Fig. 2 and Fig. 4: for ResNets, implementations from Flax (Heek et al., 2020) were used, specifically `flax.examples.imagenet.models`. For WideResNets, the **code sample** from Novak et al. (2020) was used.⁵ For all other models, we used implementations from <https://github.com/google-research/vision-transformer>. Inputs were random arrays of shapes $224 \times 224 \times 3$. All models were JITted. All reported values are averages over 10 random samples. For each setting, we ran a grid search over the batch size N in $\{2^k\}_{k=0}^9$, and reported the best time divided by N^2 , i.e. best possible throughput in each setting.

⁴The XLA team has let us know that if JITted, the FLOPs are currently correctly computed only on TPU, but are incorrect on other platforms. Therefore we compute FLOPs of non-JITted functions.

⁵We replaced `stax.AvgPool((8, 8))`, `stax.Flatten()` with `stax.GlobalAvgPool()`.

I APPLICATIONS WITH A LIMITED COMPUTE BUDGET

While our methods allow to dramatically speed-up the computation of NTK, all of them still scale as $\mathbf{N}^2\mathbf{O}^2$ for both time and memory, which can be intractable for large datasets and/or large outputs.

Here we present several settings in which our proposed methods still provide substantial time and memory savings, even when instantiating the entire $\mathbf{NO} \times \mathbf{NO}$ NTK is not feasible or not necessary.

- **NTK-vector products.** In many applications one only requires computing the NTK-vector product linear map

$$\Theta_\theta : v \in \mathbb{R}^{\mathbf{NO}} \mapsto \Theta_\theta v \in \mathbf{NO}, \quad (35)$$

without computing the entire NTK matrix Θ_θ . One common setting is using the power iteration method (Müntz, 1913) to compute NTK condition number and hence trainability of the respective NN (Lee et al., 2019; Chen et al., 2021a;b). Another setting is using conjugate gradients to compute $\Theta_\theta^{-1}\mathcal{Y}$ when doing kernel ridge regression with the NTK (Jacot et al., 2018; Lee et al., 2019; Zhou et al., 2021).

Eq. (35) is the same map as the one we considered in §3.5, and naturally, **NTK-vector products** can provide a substantial speed-up over **Jacobian contraction** in this setting. Precisely, a straightforward application of **Jacobian contraction** yields

$$\underbrace{\Theta_\theta v}_{\mathbf{NO} \times \mathbf{1}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{NO}} \underbrace{v}_{\mathbf{NO} \times \mathbf{1}}. \quad (36)$$

Combined with the cost of computing the weight space cotangents $\partial f/\partial \theta$, such evaluation costs $\mathbf{NO}[\mathbf{FP}]$ time, i.e. the cost of instantiating the entire **Jacobian**. Alternatively, one could store the entire Jacobians of sizes \mathbf{NOP} in memory, and compute a single NTK-vector product in \mathbf{NOP} time.

In contrast, **NTK-vector products** allow to compute an NTK-vector product at a cost asymptotically equivalent to a single VJP call (§3.5), i.e. $\mathbf{N}[\mathbf{FP}]$, \mathbf{O} times faster than **Jacobian contraction** without caching. With caching, fastest method will vary based on the cost of $[\mathbf{FP}]$ relative to \mathbf{OP} , as discussed in §4, but **NTK-vector products** will remain substantially more memory-efficient due to not caching the entire \mathbf{NOP} Jacobians.

- **Batching.** In many applications it suffices to compute the NTK over small batches of the data. For example Dauphin & Schoenholz (2019); Chen et al. (2021a;b) estimate the conditioning by computing an approximation to the NTK on \mathbf{N} equal to 128, 32, and 48 examples respectively. Similarly, Zhou et al. (2021) use a small batch size of $\mathbf{N} = 25$ to meta-learn the network parameters by replacing the inner SGD training loop with NTK regression.
- **Pseudo-NTK.** Many applications (§2) compute a pseudo-NTK of size $\mathbf{N} \times \mathbf{N}$, which is commonly equal to one of its \mathbf{O} diagonal blocks, or to the mean of all \mathbf{O} blocks. The reason for considering such approximation is that in the infinite width limit, off-diagonal entries often converge to zero, and for wide-enough networks this approximation can be justified. Compute-wise, these approximations are equivalent to having $\mathbf{O} = 1$. While an important contribution of our work is to enable computing the full $\mathbf{NO} \times \mathbf{NO}$ NTK quickly, if necessary, **Structured derivatives** can be combined with the $\mathbf{O} = 1$ approximations, and still provide an asymptotic speed-up and memory savings relative to prior works.

J FINITE AND INFINITE WIDTH NTK

In this work we focus on computing the finite width NTK $\Theta_\theta(x_1, x_2)$, defined in Eq. (1), that we repeat below with an addition of a batch size \mathbf{N} :

$$\mathbf{F}\text{-NTK (finite width): } \underbrace{\Theta_\theta(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} := \underbrace{\left[\frac{\partial f(\theta, x_1)}{\partial \theta} \right]}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\left[\frac{\partial f(\theta, x_2)}{\partial \theta} \right]^T}_{\mathbf{P} \times \mathbf{NO}}. \quad (37)$$

Another highly important object in deep learning theory is the *infinite width* NTK $\Theta(x_1, x_2)$, introduced in the seminal work of [Jacot et al. \(2018\)](#):

$$\text{I-NTK (infinite width): } \underbrace{\Theta(x_1, x_2)}_{\text{NO} \times \text{NO}} := \lim_{\mathbf{W} \rightarrow \infty} \mathbb{E}_{\theta \sim \mathcal{N}(\mathbf{0}, I_{\mathbf{P}})} \left[\underbrace{\Theta_{\theta}(x_1, x_2)}_{\text{NO} \times \text{NO}} \right]. \quad (38)$$

A natural question to ask is what are the similarities and differences of F- and I-NTK, when is one more applicable than the other, and what are their implementation and computational costs.

Applications. At a high level, F-NTK describes the local/linearized behavior of the finite width NN $f(\theta, x)$ ([Lee et al., 2019](#)). In contrast, I-NTK is an approximation that is exact only in the infinite width \mathbf{W} limit, and only at initialization ($\theta \sim \mathcal{N}(\mathbf{0}, I_{\mathbf{P}})$). As such, the resulting I-NTK has no notion of width \mathbf{W} , parameters θ , and cannot be computed during draining, or in a transfer or meta-learning setting, where the parameters θ are updated. As a consequence, any application to finite width networks (§2) is better served by the F-NTK, and often impossible with the I-NTK.

In contrast, I-NTK describes the behavior of an infinite ensemble of infinitely wide NNs. In certain settings this can be desirable, such as when studying the inductive bias of certain NN architectures ([Xiao et al., 2020](#)) or uncertainty ([Adlam et al., 2020](#)), marginalizing away the dependence on specific parameters θ . However, care should be taken when applying I-NTK findings to the finite width realm, since many works have demonstrated substantial finite width effects that cannot be captured by the I-NTK ([Novak et al., 2019](#); [Arora et al., 2019b](#); [Lee et al., 2019](#); [Yaida, 2020](#); [Hanin & Nica, 2020](#); [Lee et al., 2020](#)).

Mathematical scope. Another significant difference between F- and I-NTK is the scope of their definitions in [Eq. \(37\)](#) and [Eq. \(38\)](#) and mathematical tractability.

The F-NTK is well-defined for any differentiable (w.r.t. θ) function f , and our methods are respectively applicable to any differentiable functions. In fact, our work supports any Tangent Kernels (not necessarily “Neural”), and is not specific to NNs at all.

In contrast, the I-NTK requires the function f to have the concept of width \mathbf{W} (that can be meaningfully taken to infinity) to begin with, and further requires f and θ to satisfy many conditions in order for the I-NTK to be well-defined ([Yang, 2019](#)). In order for I-NTK to be well-defined *and computable in closed-form*, f needs to be built out of a relatively small, hand-selected number of primitives that admit certain Gaussian integrals to have closed-form solutions. Examples of ubiquitous primitives that *don’t* allow a closed-form solution include attention with standard parameterization ([Hron et al., 2020](#)); max-pooling; sigmoid, (log-)softmax, tanh, and many other nonlinearities; various kinds of normalization ([Yang et al., 2019](#)); non-trivial weight sharing ([Yang, 2020](#)); and many other settings. Going forward, it is unlikely that the I-NTK will scale to the enormous variety of architectures introduced by the research community each year.

Implementation tractability. Above we have demonstrated that the I-NTK is defined for a very small subset of functions admitting the F-NTK. A closed-form solution exists for an even smaller subset. However, even when the I-NTK admits a closed-form solution, it is important to consider the complexity of implementing it.

Our implementation for computing the F-NTK is applicable to any differentiable function f , and requires no extra effort when switching to a different function g . It is similar to JAX’s highly-generic function transformations such as `jax.jit` or `jax.vmap`.

In contrast, there is no known way to compute the I-NTK for an arbitrary function f , even if the I-NTK exists in closed form. The best existing solution to date is provided by [Novak et al. \(2020\)](#), which allows to *construct* f out of the limited set of building blocks provided by the authors. However, one cannot compute the I-NTK for a function implemented in a different library such as Flax ([Heek et al., 2020](#)), or Haiku ([Hennigan et al., 2020](#)), or bare-bone JAX. One would have to re-implement it using the primitives provided by [Novak et al. \(2020\)](#). Further, for a generic architecture, the primitive set is unlikely to be sufficient, and the function will need to be adapted to admit a closed-form I-NTK.

Computational tractability. F-NTK and I-NTK have different time and memory complexities, and a fully general comparison is an interesting direction for future work. Here we provide discussion for deep FCNs and CNNs.

Networks having a fully-connected top (L) readout layer have a constant block-diagonal I-NTK, hence its cost *does not* scale with O . The cost of computing the I-NTK for a deep FCN scales as N^2L for time and N^2 for memory. A deep CNN without pooling costs N^2DL time and N^2D memory (where D is the total number of pixels in a single input/activation; $D = 1$ for FCNs). Finally, a deep CNN with pooling, or any other generic architecture that leverages the spatial structure of inputs/activations, costs N^2D^2L time and N^2D^2 memory. This applies to all models in Fig. 2 and Fig. 4, Graph Neural Networks (Du et al., 2019), and the vast majority of other architectures used in practice.

The quadratic scaling of the I-NTK cost with D is especially burdensome, since, for example, for ImageNet $D^2 = 224^2 = 2, 517, 630, 976$. As a result, it would be impossible to evaluate the I-NTK on even a single ($N = 1$) pair of inputs with a V100 GPU for any model for which we’ve successfully evaluated the F-NTK in Fig. 2 and Fig. 4.

The F-NTK time and memory only scale linearly with D (Table 7). However, the F-NTK cost scales with other parameters such as width W or number of outputs O , and in general the relative F- and I-NTK performance will depend on these parameters. As a rough point of comparison, we consider the cost of evaluating the I-NTK of a 20-layer binary classification ReLU CNN with pooling on a V100 GPU used by Arora et al. (2019b) against the respective F-NTK with $W = 128$ also used by Arora et al. (2019b, Section B). Arora et al. (2019b) and Novak et al. (2020) report from 0.002 to 0.003 seconds per I-NTK entry on a pair of CIFAR-10 inputs. Using Structured derivatives, we can compute the respective F-NTK entry on same hardware in at most 0.000014 seconds, i.e. at least 100 times faster than the I-NTK. In 0.002 – 0.003 seconds per NTK entry, we can compute the F-NTK on a pair of ImageNet inputs (about 50x larger than CIFAR-10) for a 200-layer ResNet (about 10x deeper than the model above) in Fig. 2 (top left).

Finally, we remark that efficient NTK-vector products without instantiating the entire $NO \times NO$ NTK are only possible using the F-NTK (§I).

K LEVERAGING JAX DESIGN FOR EFFICIENT NTK COMPUTATION

At the time of writing, Tensorflow (Abadi et al., 2016) and PyTorch (Paszke et al., 2019) are more widely used than JAX (Bradbury et al., 2018). However, certain JAX features and design choices made it much more suitable, if not indispensable, for our project:

1. Structured derivatives require manual implementation of structure rules for different primitives in the computational graph of a function $f(\theta, x)$. JAX has a small primitive set of about 131 primitives, while PyTorch and Tensorflow have more than 400 (Frostig et al., 2021, Section 2). Further, by leveraging `jax.linearize`, we reduce our task to implementing structure rules for only *linear* primitives, of which JAX has only 54.⁶ To our knowledge neither PyTorch nor Tensorflow have an equivalent transformation, which makes JAX a natural choice due to the very concise set of primitives that we need to handle (Table 5).
2. NTK-vector products critically rely on forward mode AD (JVP), and Structured derivatives also use it (albeit it’s not crucial; see §D). At the time of writing, PyTorch does not implement an efficient forward mode AD.
3. Structured derivatives rely crucially on the ability to traverse the computation graph to rewrite contractions using our substitution rules. JAX provides a highly-convenient graph representation in the form of a `Jaxpr`, as well as tooling and documentation for writing custom Jaxpr interpreters.

⁶This follows from the fact that the NTK of a function is equal to the NTK of the linearized function at the same primal parameters θ . See also (Frostig et al., 2021, Section 1) for how JAX uses the same insight to not implement all 131 VJP rules, but only implement 54 transpose rules for reverse mode AD.

4. All implementations (even [Jacobian contraction](#)) rely heavily on `jax.vmap` (and to our knowledge, in many cases, it is indispensable). While PyTorch has [released](#) a prototype of `vmap` in May 2021, it was not available when we started this project.

For researchers interested in interfacing with our library, we recommend looking into tools facilitating data exchange between different ML frameworks, such as [DLPack](#) and [Jax2TF](#). See §5 for more implementation details.

L COMPLEXITY ANALYSIS WITHOUT THE $\mathbf{O} = \mathcal{O}(\mathbf{LW})$ ASSUMPTION

Here we repeat the same analysis as in §3 without the assumption of $\mathbf{O} = \mathcal{O}(\mathbf{LW})$. This results in [Table 8](#), where [Jacobian contraction](#) and [Structured derivatives](#) gain an extra $\mathbf{N}^2\mathbf{O}^3\mathbf{W}$ and $\mathbf{N}^2\mathbf{O}^3$ time terms respectively. This does not affect our main text conclusions.

L.1 JVP AND VJP

As in §3.1, the time cost of both operations is comparable to the forward pass (**FP**), i.e. $[\mathbf{FP}] = [\text{cost of all intermediate layers}] + [\text{cost of the top layer}] = \mathbf{LW}^2 + \mathbf{OW}$.

For a single input, the memory cost of computing both the JVP and the VJP are respectively,

$$\begin{aligned} [\text{size of all weights}] + [\text{size of activations at a single layer}] &= [\mathbf{LW}^2 + \mathbf{OW}] + [\mathbf{W} + \mathbf{O}] \sim \mathbf{LW}^2 + \mathbf{OW}, \\ [\text{size of all weights}] + [\text{size of activations in all layers}] &= [\mathbf{LW}^2 + \mathbf{OW}] + [\mathbf{LW} + \mathbf{O}] \sim \mathbf{LW}^2 + \mathbf{OW}. \end{aligned}$$

As in §3.1, despite the fact that the VJP requires more memory to store intermediate activations (which is necessary for efficient backpropagation), we see that both computations are dominated by the cost of storing the weights.

Batched inputs. If x is a batch of inputs of size \mathbf{N} , the time cost of JVP and VJP increases linearly to $\mathbf{NLW}^2 + \mathbf{NOW}$. The memory cost is more nuanced. Since weights can be shared across inputs, the memory cost of the JVP and VJP are respectively,

$$\begin{aligned} &[\text{size of all weights}] + \mathbf{N}[\text{size of activations at a single layer}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{W} + \mathbf{O}] \sim \mathbf{LW}^2 + \mathbf{OW} + \mathbf{NW} + \mathbf{NO}, \\ &[\text{size of all weights}] + \mathbf{N}[\text{size of activations in all layers}] + \mathbf{N}[\text{size of all weight matrices}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LW} + \mathbf{O}] + \mathbf{N}[\mathbf{LW}^2 + \mathbf{OW}] \sim \mathbf{NLW}^2 + \mathbf{OW}^2 + \mathbf{NOW}. \end{aligned}$$

Recall from §3.1 we only need to store one cotangent weight matrix, $\partial f / \partial \theta^l$, at a time. Therefore

- JVP costs $\mathbf{NLW}^2 + \mathbf{NOW}$ time and $\mathbf{LW}^2 + \mathbf{OW} + \mathbf{NW} + \mathbf{NO}$ memory.
- VJP costs $\mathbf{NLW}^2 + \mathbf{NOW}$ time and $\mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NOW}^2 + \mathbf{NOW}$ memory.

L.2 JACOBIAN

The time and memory costs to compute the Jacobian are identical to §3.2,

$$\begin{aligned} &\mathbf{ON}([\text{cost of all intermediate layers}] + [\text{cost of the top layer}]) \\ &= \mathbf{ON}([\mathbf{LW}^2] + [\mathbf{OW}]) \sim \mathbf{NLOW}^2 + \mathbf{NO}^2\mathbf{W}, \\ &[\text{size of all weights}] + \mathbf{N}[\text{size of activations in all layers}] + \mathbf{ON}[\text{size of a single weight matrix}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LW} + \mathbf{O}] + \mathbf{ON}[\mathbf{W}^2 + \mathbf{OW}] \sim \mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W}. \end{aligned}$$

Therefore, asymptotically, the costs are identical to §3.2:

Jacobian costs $\mathbf{NLOW}^2 + \mathbf{NO}^2\mathbf{W}$ time and $\mathbf{LW}^2 + \mathbf{NLW} + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W}$ memory.

L.3 JACOBIAN CONTRACTION

The time cost of the contraction in Eq. (6) is $\mathbf{O}^2\mathbf{W}^2$ for $l < \mathbf{L}$, but is $\mathbf{O}^3\mathbf{W}$ for the top layer $l = \mathbf{L}$. The memory necessary to instantiate each factor and the result is $\mathbf{OW}^2 + \mathbf{O}^2$ for $l < \mathbf{L}$ and $\mathbf{O}^2\mathbf{W}$ for the top layer $l = \mathbf{L}$.

Accounting for all layers together, we arrive at $\mathbf{LO}^2\mathbf{W}^2 + \mathbf{O}^3\mathbf{W}$ time cost and $\mathbf{OW}^2 + \mathbf{O}^2\mathbf{W}$ memory, due to being able to process summands sequentially.

Batched inputs. If we consider x_1 and x_2 to be input batches of size \mathbf{N} , then the resulting NTK is a matrix of shape $\mathbf{NO} \times \mathbf{NO}$, and the time cost becomes $\mathbf{N}^2 (\mathbf{LO}^2\mathbf{W}^2 + \mathbf{O}^3\mathbf{W})$, while memory grows to $[\text{NTK matrix size}] + [\text{factors size}] = \mathbf{N}^2\mathbf{O}^2 + \mathbf{N}(\mathbf{OW}^2 + \mathbf{O}^2\mathbf{W})$.

Adding the cost of the **Jacobian** described in §L.2, we obtain

Jacobian contraction costs $\mathbf{N}^2\mathbf{LO}^2\mathbf{W}^2 + \mathbf{N}^2\mathbf{O}^3\mathbf{W}$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{LW}^2 + \mathbf{NLW}$ memory.

L.4 STRUCTURED DERIVATIVES

The contraction in Eq. (9) takes $\mathbf{O}^2\mathbf{W}$ time and $\mathbf{OW} + \mathbf{O}^2$ memory for $l < \mathbf{L}$, and $\mathbf{O}^3 + \mathbf{W}$ time and $\mathbf{O}^2 + \mathbf{W}$ memory for $l = \mathbf{L}$.

Accounting for all layers, time cost becomes $\mathbf{LO}^2\mathbf{W} + \mathbf{O}^3$, and memory remains $\mathbf{OW} + \mathbf{O}^2$.

Batched inputs. In the batched setting, the time cost grows quadratically with the size of the NTK to $\mathbf{N}^2\mathbf{LO}^2\mathbf{W} + \mathbf{N}^2\mathbf{O}^3$, while the memory cost increases to $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}$.

Extra memory cost for computing the derivatives is

$$\begin{aligned} & [\text{size of all weights}] + \mathbf{N} [\text{size of activations in all layers}] \\ &= [\mathbf{LW}^2 + \mathbf{OW}] + \mathbf{N} [\mathbf{LW} + \mathbf{O}] \sim \mathbf{LW}^2 + \mathbf{OW} + \mathbf{NLW}. \end{aligned}$$

The extra time cost is asymptotically the cost of \mathbf{O} forward passes, $\mathbf{NLOW}^2 + \mathbf{NO}^2\mathbf{W}$ which is the same as the **Jacobian**. Putting everything together we find the following costs,

By leveraging **Structured derivatives** in NN computations, we have reduced the cost of NTK to $\mathbf{N}^2\mathbf{LO}^2\mathbf{W} + \mathbf{N}^2\mathbf{O}^3 + \mathbf{NLOW}^2$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW} + \mathbf{LW}^2 + \mathbf{NLW}$ memory.

L.5 NTK-VECTOR PRODUCTS

The cost analysis of **NTK-vector products** in §3.5 is not impacted by the $\mathbf{O} = \mathcal{O}(\mathbf{LW})$ assumption, hence it remains the same as in §3.5:

NTK computation as a sequence of **NTK-vector products** costs $\mathbf{N}^2\mathbf{LOW}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{W}$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOW}^2 + \mathbf{LW}^2 + \mathbf{NLW}$ memory.

M RELATIONSHIP BETWEEN THE NTK AND THE HESSIAN

Here we briefly touch on the difference between the NTK

$$\text{NTK: } \underbrace{\Theta_\theta(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} := \underbrace{\left[\partial f(\theta, x_1) / \partial \theta \right]}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\left[\partial f(\theta, x_2) / \partial \theta \right]^T}_{\mathbf{P} \times \mathbf{NO}}, \quad (39)$$

and the Hessian:

$$\text{Hessian: } \underbrace{\mathbf{H}_\theta(x)}_{\mathbf{P} \times \mathbf{P}} := \frac{\partial^2 \mathcal{L}(f(\theta, x))}{\partial \theta^2}, \quad (40)$$

Method	Time	Memory	Use when
Jacobian contraction	$\mathbf{N}^2\mathbf{L}\mathbf{O}^2\mathbf{W}^2 + \mathbf{N}^2\mathbf{O}^3\mathbf{W}$	$\mathbf{N}\mathbf{O}\mathbf{W}^2 + \mathbf{N}^2\mathbf{O}^2 + \mathbf{N}\mathbf{L}\mathbf{W} + \mathbf{L}\mathbf{W}^2$	Don't
NTK-vector products	$\mathbf{N}^2\mathbf{O}^2\mathbf{W} + \mathbf{N}^2\mathbf{L}\mathbf{O}\mathbf{W}^2$	$\mathbf{N}\mathbf{O}\mathbf{W}^2 + \mathbf{N}^2\mathbf{O}^2 + \mathbf{N}\mathbf{L}\mathbf{W} + \mathbf{L}\mathbf{W}^2$	$\mathbf{O} > \mathbf{W}$ or $\mathbf{N} = 1$
Structured derivatives	$\mathbf{N}^2\mathbf{L}\mathbf{O}^2\mathbf{W} + \mathbf{N}\mathbf{L}\mathbf{O}\mathbf{W}^2 + \mathbf{N}^2\mathbf{O}^3$	$\mathbf{N}\mathbf{O}\mathbf{W} + \mathbf{N}^2\mathbf{O}^2 + \mathbf{N}\mathbf{L}\mathbf{W} + \mathbf{L}\mathbf{W}^2$	$\mathbf{O} < \mathbf{W}$ or $\mathbf{L} = 1$

Table 8: **Asymptotic time and memory cost of computing the NTK for an FCN without assuming that $\mathbf{O} = \mathcal{O}(\mathbf{L}\mathbf{W})$.** Costs are for a pair of batches of inputs of size \mathbf{N} each, and for \mathbf{L} -deep, \mathbf{W} -wide FCN with \mathbf{O} outputs. Resulting NTK has shape $\mathbf{N}\mathbf{O} \times \mathbf{N}\mathbf{O}$. **NTK-vector products** allow a reduction of the time complexity, while **Structured derivatives** reduce both time and memory complexity. **Note:** presented are asymptotic cost estimates; in practice, all methods incur large constant multipliers (e.g. at least 3x for time; see §3.1). However, this generally does not impact the relative performance of different methods. See §3.6 for discussion, Table 1 for a simplified cost summary under the assumption of $\mathbf{O} = \mathcal{O}(\mathbf{L}\mathbf{W})$ (differing only by lacking the $\mathbf{N}^2\mathbf{O}^3\mathbf{W}$ and $\mathbf{N}^2\mathbf{O}^3$ terms in **Jacobian contraction** and **Structured derivatives** time costs respectively), Table 7 for CNN, and Table 2 for more generic cost analysis.

defined for some differentiable loss function on the output space $\mathcal{L} : \mathbb{R}^{\mathbf{N}\mathbf{O}} \rightarrow \mathbb{R}$.

Both matrices characterize localized training dynamics of a NN, and the NTK can be used as a more tractable quantity in cases where the Hessian is infeasible to instantiate (for example, \mathbf{P} amounts to tens of millions in models considered in Fig. 2 and Fig. 4).

The connection between the NTK and the Hessian can be established for when \mathcal{L} is the squared error (SE), i.e. $\mathcal{L}(y) = \|y - \mathcal{Y}\|_2^2/2$, where $\mathcal{Y} \in \mathbb{R}^{\mathbf{N}\mathbf{O}}$ are the training targets. In this case, as presented in (Pennington & Bahri, 2017, Section 2) and (Grosse, 2021, Equation 13; page 21):

$$\underbrace{\mathbf{H}_\theta(x)}_{\mathbf{P} \times \mathbf{P}} = \underbrace{\left[\frac{\partial f(\theta, x)}{\partial \theta} \frac{\partial f(\theta, x)}{\partial \theta} \right]}_{:= \mathbf{H}_\theta^0(x)} + \underbrace{\sum_{n,o=1}^{\mathbf{N}, \mathbf{O}} (f(\theta, x) - \mathcal{Y})^{n,o} \frac{\partial^2 f(\theta, x)^{n,o}}{\partial^2 \theta}}_{:= \mathbf{H}_\theta^1(x)}, \quad (41)$$

where we have decomposed the Hessian $\mathbf{H}_\theta(x)$ into two summands $\mathbf{H}_\theta^0(x)$ and $\mathbf{H}_\theta^1(x)$ following the notation of Pennington & Bahri (2017).

Notice that if $f(\theta, x) = \mathcal{Y}$, i.e. the SE loss is 0, $\mathbf{H}_\theta^1(x) = 0$, yielding

$$\underbrace{\mathbf{H}_\theta(x)}_{\mathbf{P} \times \mathbf{P}} = \mathbf{H}_\theta^0(x) = \underbrace{\frac{\partial f(\theta, x)}{\partial \theta}}_{\mathbf{P} \times \mathbf{N}\mathbf{O}} \underbrace{\frac{\partial f(\theta, x)}{\partial \theta}}_{\mathbf{N}\mathbf{O} \times \mathbf{P}}, \quad \underbrace{\Theta_\theta(x, x)}_{\mathbf{N}\mathbf{O} \times \mathbf{N}\mathbf{O}} = \underbrace{\frac{\partial f(\theta, x)}{\partial \theta}}_{\mathbf{N}\mathbf{O} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}}_{\mathbf{P} \times \mathbf{N}\mathbf{O}}, \quad (42)$$

and, as a consequence, the Hessian and the NTK have the same eigenvalues (see also Grosse (2021, Page 21)) in this particular case. Moreover, the Hessian (and Hessian-vector products) can be computed very similarly to **NTK-vector products**, by switching the order of VJP and JVP operations in Eq. (10).

However, except for zero SE loss case above, the NTK and the Hessian have different spectra, and their computations share less similarity. Precisely, Hessian-vector products (and consequently the Hessian) are computed in JAX through a composition of JVPs and VJP similar to **NTK-vector products**:

$$\mathbf{H}_\theta(x)v = \left[\frac{\partial^2 \mathcal{L}(f(\theta, x))}{\partial \theta^2} \right] v = \frac{\partial}{\partial \theta} \left[\frac{\partial \mathcal{L}(f(\theta, x))}{\partial \theta} \right] v = \quad (43)$$

$$= \frac{\partial [\text{VJP}_{\mathcal{L} \circ f, \theta, x}(1)]}{\partial \theta} v = \text{JVP}_{[\text{VJP}_{\mathcal{L} \circ f, \cdot, \cdot}(1)], \theta, x}(v). \quad (44)$$

While Eq. (43) is similar to Eq. (10) in that both are compositions of JVPs and VJP, in Eq. (10) the result of a VJP is the input tangent to the JVP of f , while in Eq. (43) it is the function to be differentiated by the JVP (instead of f).

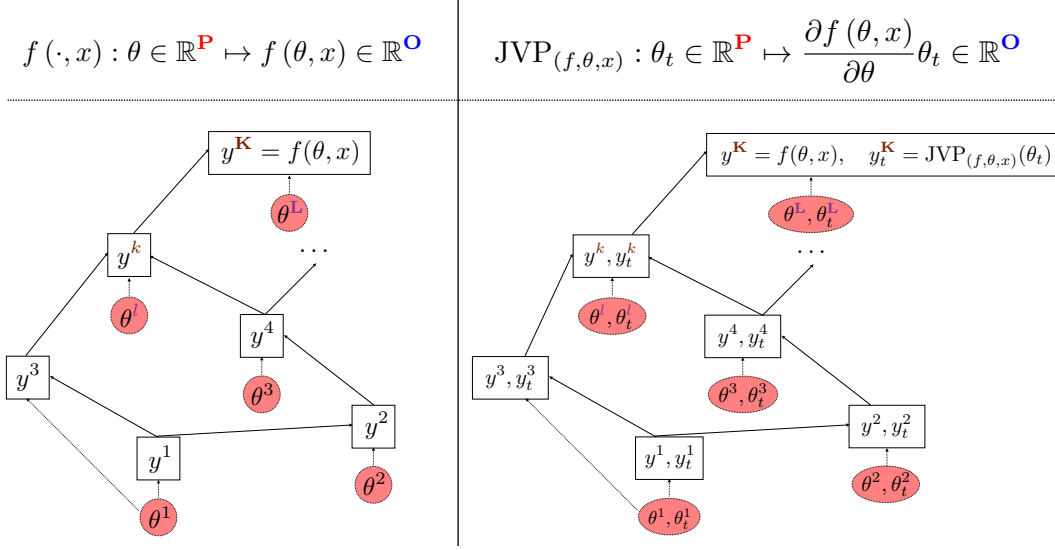


Figure 6: **Visual demonstration for why JVP time and memory costs are asymptotically comparable to the forward pass (FP).** **Left:** computational graph of the forward pass $f(\theta, x)$. **Right:** computational graph of joint evaluation of the forward pass $f(\theta, x)$ along with $\text{JVP}_{(f, \theta, x)}(\theta_t)$. Each node of the JVP graph accepts both primal and tangent inputs, and returns primal and tangent outputs, but the topology of the graph and order of execution remains identical to **FP**. As long as individual nodes of the JVP graph do not differ significantly in time and memory from the **FP** nodes, time and memory of a JVP ends up asymptotically equivalent to **FP** due to identical graph structure. However, in order to create JVP nodes and evaluate them, the time cost does grow by a factor of about 3 compared to **FP**. See §3.1 and §N for discussion.

N GENERIC COMPUTATIONAL COSTS OF JVP AND VJP

In this section we present a brief intuition for why JVP and VJP are asymptotically equivalent in compute time to the forward pass **FP**, as we mentioned in §3.1. We refer the reader to the [JAX Autodiff Cookbook](#) for more details, and (Griewank & Walther, 2008, Section 3) for a rigorous treatment of the subject.

JVP can be computed by traversing a computational graph of the same topology as **FP**, except for primitive nodes in the graph need to be augmented to compute not only the forward pass of the node, but also the JVP of the node (see Fig. 6). Due to identical topology and order of evaluation, asymptotically time and memory costs remain unchanged. However, constructing the augmented nodes in the JVP graph, and their consequent evaluation results in extra time cost proportional to the size of the graph. Therefore in practice JVP costs about $3 \times \text{FP}$ time and $2 \times \text{FP}$ memory.

VJP, as a linear function of cotangents f_c , is precisely the transpose of the linear function JVP. As such, it can be computed by traversing the transpose of the JVP graph (Fig. 6, right), with each JVP node replaced by its transposition as well. This results in identical time and memory costs, as long as node transpositions are implemented efficiently. However, their evaluation requires primal outputs y^k (now inputs to the transpose nodes), which is why VJP necessitates an extra **FP** time cost to compute them (hence costlier than JVP, but still inconsequential asymptotically) and extra memory to store them, which can generally increase asymptotic memory requirements.