

Fast Finite Width Neural Tangent Kernel

Roman Novak, Jascha Sohl-Dickstein, Samuel S. Schoenholz

Google Brain, {romann, jaschasd, schsam}@google.com,

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. However, it is notoriously expensive to compute, severely limiting its practical utility. We perform the first in-depth analysis of the compute and memory requirements for NTK computation in finite NNs. Leveraging their structure, we propose two novel algorithms that change the *exponent* of the compute and memory requirements of the finite width NTK, dramatically improving efficiency in a wide range of NN architectures on all hardware platforms. We open-source [\[github.com/iclr2022anon/fast_finite_width_ntk\]](https://github.com/iclr2022anon/fast_finite_width_ntk) our two algorithms as general-purpose JAX function transformations that apply to any differentiable computation and introduce no hyperparameters.

Notation. Consider a NN $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$ with \mathbf{O} outputs (logits) per input x and a total number \mathbf{P} of trainable parameters $\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}]$, with each θ^l of size \mathbf{P}^l , $\mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$. Also assume the network has \mathbf{K} intermediate pre-activations y^k of size \mathbf{Y}^k each, $\mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$. The NTK is

$$\underbrace{\Theta_\theta}_{\mathbf{O} \times \mathbf{O}} := \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{O} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{\mathbf{P} \times \mathbf{O}} = \sum_{l=0}^{\mathbf{L}} \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta^l}}_{\mathbf{O} \times \mathbf{P}^l} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta^l}^T}_{\mathbf{P}^l \times \mathbf{O}} \quad (1)$$

We denote \mathbf{FP} to be the (time or memory, depending on the context) cost of a single forward pass $f(\theta, x)$. For memory, we exclude the cost of storing all \mathbf{P} weights in memory, but rather define it to be the cost of evaluating f one JAX [1] primitive y^k at a time, amounting to no more than $\max_l \mathbf{P}^l + \max_k \mathbf{Y}^k$, which we denote as simply $\mathbf{P}^l + \mathbf{Y}^k$ for brevity. Finally, we will consider x_1 and x_2 to be batches of \mathbf{N} inputs each, in which case the NTK will be a $\mathbf{NO} \times \mathbf{NO}$ matrix.

Jacobian-vector products (JVP) and vector-Jacobian products (VJP). We define

$$\text{JVP}_{(f, \theta, x)} : \theta_t \in \mathbb{R}^{\mathbf{P}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta} \theta_t \in \mathbb{R}^{\mathbf{O}}; \quad \text{VJP}_{(f, \theta, x)} : f_c \in \mathbb{R}^{\mathbf{O}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta}^T f_c \in \mathbb{R}^{\mathbf{P}}. \quad (2)$$

In JAX the time cost of both is comparable to \mathbf{FP} . The memory cost of a JVP is \mathbf{FP} as well, while the memory cost of a VJP is generally $\mathbf{Y} + \mathbf{P}$, since it requires storing all \mathbf{K} intermediate pre-activations for efficient backprop and all \mathbf{L} output cotangents. However, for the purpose of computing the NTK, we never need to store the whole Jacobian $\partial f/\partial \theta$, but only individual cotangents like $\partial f/\partial \theta^l$ to compute the sum in Eq. (1). Hence we consider VJP to cost $\mathbf{Y} + \mathbf{P}^l$ memory. Finally, for a batch of \mathbf{N} inputs x , JVP and VJP cost $\mathbf{N}[\mathbf{FP}]$ time; $\mathbf{N}[\mathbf{FP}] + \mathbf{P}$ and $\mathbf{N}[\mathbf{Y} + \mathbf{P}^l] + \mathbf{P}$ memory respectively.

Jacobian. For NNs, the Jacobian $\partial f/\partial \theta$ is most often computed via \mathbf{O} VJP calls on rows of the identity matrix $\mathbf{I}_\mathbf{O}$, i.e. costs $\mathbf{O}[\text{VJP}]$ time and memory less network weights and pre-activations that can be reused across VJP calls, resulting in $\mathbf{NO}[\mathbf{FP}]$ time and $\mathbf{NO}[\mathbf{Y}^k + \mathbf{P}^l] + \mathbf{NY} + \mathbf{P}$ memory.

1 Jacobian contraction

This baseline method of evaluating the NTK consists in computing the Jacobians $\partial f/\partial \theta$ and contracting them as in Eq. (1). The contraction costs $\mathbf{N}^2 \mathbf{O}^2 \mathbf{P}$ time and $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{NO} \mathbf{P}^l$ memory. Adding up the cost of computing the Jacobian $\partial f/\partial \theta$ we arrive at

Jacobian contraction: $\mathbf{NO}[\mathbf{FP}] + \mathbf{N}^2 \mathbf{O}^2 \mathbf{P}$ time; $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{NO}[\mathbf{Y}^k + \mathbf{P}^l] + \mathbf{NY} + \mathbf{P}$ memory.

2 NTK-vector products – our first contribution

Consider the NTK-vector product function: $\Theta \text{VP} : v \in \mathbb{R}^{\mathbf{O}} \mapsto \Theta_\theta v \in \mathbb{R}^{\mathbf{O}}$. Applying it to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$ allows to compute the NTK, i.e. $\Theta_\theta I_{\mathbf{O}} = \Theta_\theta$. Expand $\Theta \text{VP}(v) = \Theta_\theta v$ as

$$\frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)], \quad (3)$$

where we have observed that the NTK-vector product can be expressed as a composition of a JVP and a VJP. The cost of computing Θ_θ is then equivalent to the cost of **Jacobian**, since it consists of \mathbf{O} VJPs followed by \mathbf{O} (cheaper) JVPs, therefore $\mathbf{O}[\mathbf{FP}]$ time and $\mathbf{O}[\mathbf{Y}^k + \mathbf{P}^k] + \mathbf{Y} + \mathbf{P}$ memory. In the batched setting Eq. (3) is repeated for each pair of inputs, and therefore time increases by a factor of \mathbf{N}^2 to become $\mathbf{N}^2 \mathbf{O}[\mathbf{FP}]$. However, the memory cost grows linearly in \mathbf{N} (except for the cost of storing the NTK of size $\mathbf{N}^2 \mathbf{O}^2$), since intermediate pre-activations and tangents/cotangents necessary to compute the JVP and VJP can be computed for each batch x_1 and x_2 separately, and then reused for every pairwise combination. Therefore memory cost is equivalent to **Jacobian**, and we arrive at

NTK-vector products: $\mathbf{N}^2 \mathbf{O}[\mathbf{FP}]$ time; $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{NO}[\mathbf{Y}^k + \mathbf{P}^l] + \mathbf{NY} + \mathbf{P}$ memory.

3 Structured derivatives – our second contribution

Rewrite Θ_θ from Eq. (1) using the chain rule and pre-activation y notation:

$$\Theta_\theta = \sum_{l, k_1, k_2} \left(\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial f_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T = \sum_{l, k_1, k_2} \frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l} \frac{\partial f_2}{\partial y_2^{k_2}}^T, \quad (4)$$

where $f_i = f(\theta, x_i)$, and we only consider $\partial y_i^{k_i} / \partial \theta^l$ to be non-zero if θ^l is a direct input to $y_i^{k_i}$.

Both **Jacobian contraction** and **NTK-vector products** perform this sum of contractions, albeit implicitly via VJPs and JVPs, without explicit instantiation of primitive Jacobians $\partial y / \partial \theta$. However, while VJPs and JVPs themselves are guaranteed to be computationally optimal, higher order computations like their composition (**NTK-vector products**) or contraction (**Jacobian contraction**) are not. The idea of **Structured derivatives** is to design rules for efficient computation of such contractions, similarly to how JAX and other AD packages have rules for JVPs and VJPs.

Specifically, our rules identify a few simple types of structure (e.g. block diagonal, constant-block diagonal, tiling) in $\partial y_i^{k_i} / \partial \theta^l$, that allow us to simplify the contraction in Eq. (4). 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 (much) smaller subarrays and modifying the contraction. In §E we provide specific descriptions of our rules and their impact on the computational complexity of Eq. (4). Notably, the contraction is never slower than **Jacobian contraction**, and is at most $\mathbf{N}^2 \mathbf{O}^2 \min[\mathbf{Y}, \mathbf{P}]$. For a simple concrete example, see §I.4.

The remaining cost to compute the factors $\partial f_i / \partial y_i^{k_i}$, and $\partial y_i^{k_i} / \partial \theta^l$ depends on the specific pair of primitives $y_1^{k_1}$ and $y_2^{k_2}$, but is generally similar to the cost of **Jacobian** except for (1) we don't need to compute and store \mathbf{NO} final weight space cotangents $\partial f_i / \partial \theta^l$, but (2) we do have to instead process \mathbf{N} small subarrays of primitive Jacobians $\partial y_i^{k_i} / \partial \theta^l$, which we consider to cost $\mathbf{J}_i^{k_i}$. We summarize generic cost estimates below and in Table 1, and show next that they end up beneficial (asymptotically and practically) in most common settings.

$\mathbf{NO}[\mathbf{FP}] + \mathbf{N}^2 \mathbf{O}^2 \min[\mathbf{Y}, \mathbf{P}] + \mathbf{N}[\mathbf{J} - \mathbf{OP}]$ time; $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{NOY}^k + \mathbf{NJ}_i^k + \mathbf{NY} + \mathbf{P}$ memory.

Application to FCNs and CNNs. We consider $\mathbf{K} = \mathbf{L}$ -layer CNNs with channel count \mathbf{W} , pixel count \mathbf{D} , filter size \mathbf{F} , and global average pooling before the top FC layer. Plugging $\mathbf{P}^l = \mathbf{FW}^2$ (\mathbf{OW} for $k = \mathbf{K}$), $\mathbf{Y}^k = \mathbf{DW}$ (\mathbf{O} for $k = \mathbf{K}$), $\mathbf{FP} = \mathbf{LDFW}^2 + \mathbf{OW}$, $\mathbf{J}_i^k = \mathbf{DFW}$ (\mathbf{W} for $k = \mathbf{K}$; convolutions and matrix multiplications have the **Constant block-diagonal** structure – see §E.3) we arrive at Table 2. For FCNs we simply put $\mathbf{D} = \mathbf{F} = 1$, and obtain Table 3 (and Fig. 1, Fig. 3). Notably, in both cases **Structured derivatives** are asymptotically better than **Jacobian contraction** in time and memory, under a mild condition of $\mathbf{D} \leq \mathbf{OW}$. Finally, we also confirm that our methods are practically beneficial in a much wider set of operations used by contemporary ImageNet models in Fig. 2 and Fig. 4.

Method	Time	Memory	Use when
Jacobian contraction	$N \cdot O \cdot FP + N^2 O^2 P$	$N^2 O^2 + NO \cdot Y^k + P^l + NY + P$	$P < Y$, small O
NTK-vector products	$N^2 O \cdot FP$	$N^2 O^2 + NO \cdot Y^k + P^l + NY + P$	$FP < OP$, large O , small N
Structured derivatives	$N \cdot O \cdot FP + N^2 O^2 \min[Y, P] + N \cdot [J - OP]$	$N^2 O^2 + NOY^k + NJ^l + NY + P$	$FP > OP$, large O , large N

Table 1: **Generic NTK computation cost.** NTK-vector products trade-off contractions for more FP. Structured derivatives make the contraction cheaper, and usually also reduce memory.

Method	Time	Memory	Use when
Jacobian contraction	$N \cdot O \cdot LDFW^2 + OW + N^2 O^2 [LFW^2 + OW]$	$N^2 O^2 + NO \cdot DW + FW^2 + OW + N \cdot LDW + LFW^2 + OW^2$	$D > OW$
NTK-vector products	$N^2 O \cdot LDFW^2 + OW$	$N^2 O^2 + NO \cdot DW + FW^2 + OW + N \cdot LDW + LFW^2 + OW^2$	$N = 1$
Structured derivatives	$N \cdot O \cdot LDFW^2 + OW + N^2 O^2 L \min(FW^2, DW) + O$	$N^2 O^2 + NO \cdot DW + NDFW + N \cdot LDW + LFW^2 + OW^2$	$D < OW$

Table 2: **CNN NTK computation cost.** Structured derivatives reduce time complexity, and have lower memory cost if $D < OW$, which is a common setting.

Method	Time	Memory	Use when
Jacobian contraction	$N^2 O^2 LW^2$	$N^2 O^2 + NOW^2 + NLW + LW^2$	Don't
NTK-vector products	$N^2 OLW^2 + N^2 O^2 W$	$N^2 O^2 + NOW^2 + NLW + LW^2$	$O > W$ or $N = 1$
Structured derivatives	$N OLW^2 + N^2 O^2 LW$	$N^2 O^2 + NOW + NLW + LW^2$	$O < W$ or $L = 1$

Table 3: **FCN NTK computation cost.** NTK-vector products allow a reduction of the time complexity, while Structured derivatives reduce both time and memory complexity. For brevity $O = \mathcal{O}(LW)$ is assumed in this table.

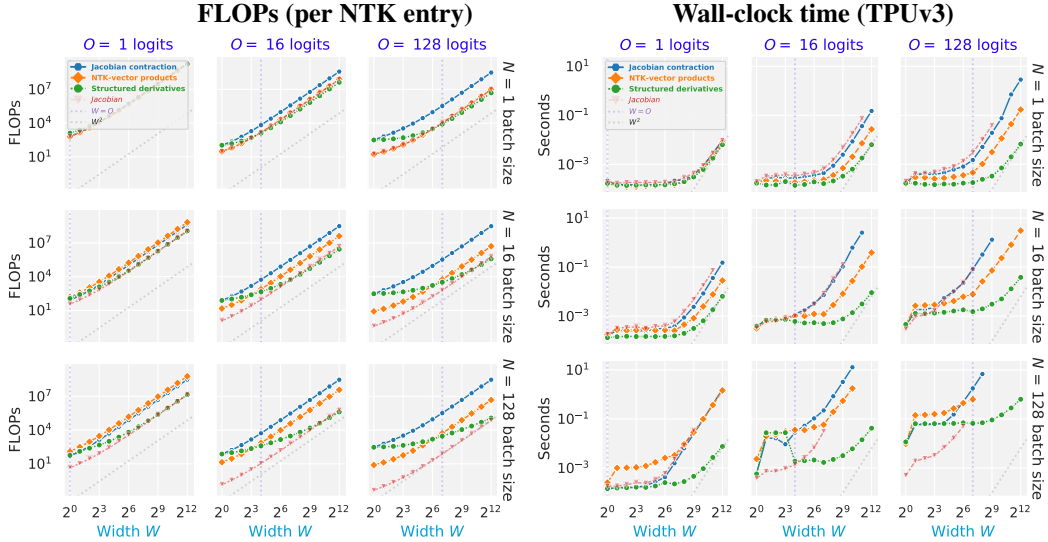


Figure 1: **FLOPs (left) and wall-clock time (right) of computing the NTK for a 10-layer ReLU FCN.** As predicted by Table 3, our methods almost always outperform Jacobian contraction, allowing orders of magnitude speed-ups and memory improvements (missing points are out-of-memory).

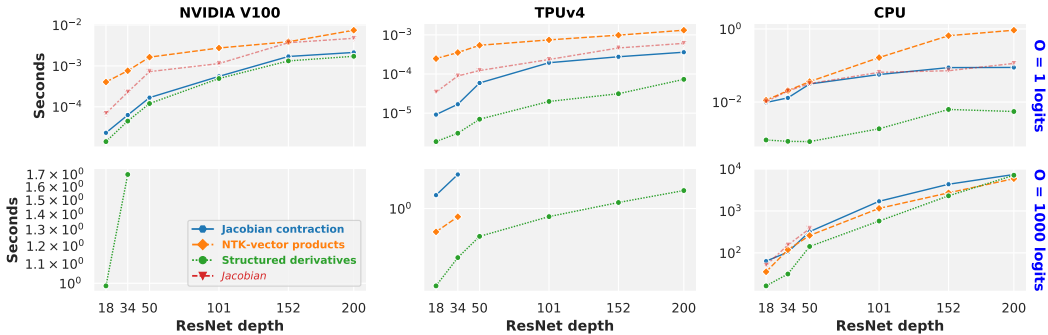


Figure 2: **Wall-clock time cost of computing an NTK for several ResNet sizes on a pair of ImageNet inputs.** Structured derivatives allow the NTK to be computed faster and for larger models (see bottom row – missing points indicate out-of-memory error). NTK-vector products, as predicted by Table 1, are advantageous for large O (bottom row), but also scale worse with FP than other methods, which is especially noticeable in CNNs.

References

- [1] 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>.
- [2] 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>.
- [3] 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.
- [4] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.
- [5] 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.
- [6] Radford M. Neal. Priors for infinite networks (tech. rep. no. crg-tr-94-1). *University of Toronto*, 1994.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] Adrià Garriga-Alonso, Laurence Aitchison, and Carl Edward Rasmussen. Deep convolutional networks as shallow gaussian processes. In *International Conference on Learning Representations*, 2019.
- [11] 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.
- [12] 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.
- [13] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- [14] 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.
- [15] 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.

- [16] 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.
- [17] Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. *International Conference on Learning Representations*, 2017.
- [18] 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.
- [19] Lechao Xiao, Jeffrey Pennington, and Samuel S Schoenholz. Disentangling trainability and generalization in deep learning. In *International Conference on Machine Learning*, 2020.
- [20] Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv preprint arXiv:1901.09321*, 2019.
- [21] Yann Dauphin and Samuel S Schoenholz. Metainit: Initializing learning by learning to initialize. 2019.
- [22] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. *arXiv preprint arXiv:2101.08692*, 2021.
- [23] Andrew Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. *arXiv preprint arXiv:2102.06171*, 2021.
- [24] 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.
- [25] Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets without pretraining or strong data augmentations, 2021.
- [26] 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.
- [27] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *arXiv preprint arXiv:2102.06701*, 2021.
- [28] 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.
- [29] Sanjeev Arora, Simon S. Du, Zhiyuan Li, Ruslan Salakhutdinov, Ruosong Wang, and Dingli Yu. Harnessing the power of infinitely wide deep nets on small-data tasks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkl8sJBYvH>.
- [30] Timothy Nguyen, Zhouong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-regression. *arXiv preprint arXiv:2011.00050*, 2020.
- [31] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. *arXiv preprint arXiv:2107.13034*, 2021.
- [32] 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, editors, *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>.
- [33] 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.

- [34] 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*, pages 8141–8150. Curran Associates, Inc., 2019.
- [35] 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>.
- [36] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/finn17a.html>.
- [37] 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>.
- [38] Boris N. Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *NeurIPS*, 2018.
- [39] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. 2017. URL <https://arxiv.org/abs/1611.01578>.
- [40] 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*, 2021.
- [41] 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*, pages 248–255. Ieee, 2009.
- [42] 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*, pages 770–778, 2016.
- [43] Roy Frostig, Matthew J Johnson, Dougal Maclaurin, Adam Paszke, and Alexey Radul. Decomposing reverse-mode automatic differentiation. *arXiv preprint arXiv:2105.09469*, 2021.
- [44] Dougal Maclaurin, David Duvenaud, and Ryan P Adams. Autograd: Effortless gradients in numpy.
- [45] 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>.

Appendix

A Additional figures

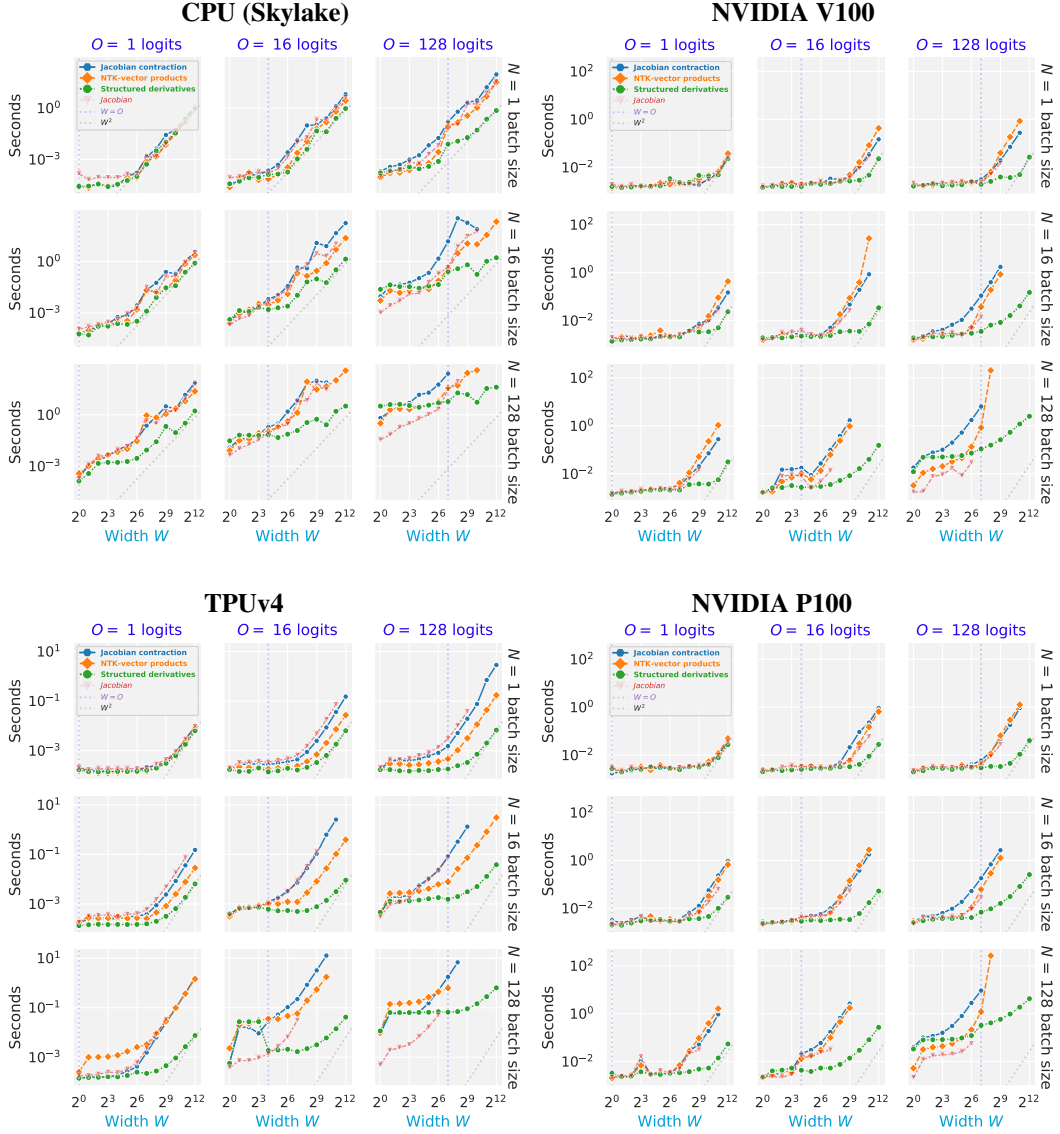


Figure 3: **Wall-clock time of computing NTK of a 10-layer ReLU FCN on different platforms.** In all settings, **Structured derivatives** allow orders of magnitude improvement in wall-clock time and memory (missing points indicate out-of-memory error). However, we remark that on GPU platforms (right), **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 and **TPUv3** platform. See §K for details.

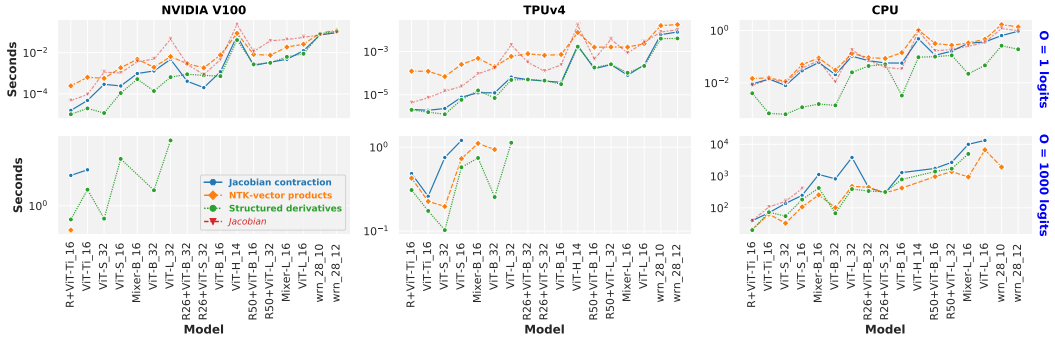


Figure 4: **Wall-clock time per input pair of computing NTK on various ImageNet models** like Vision Transformers and hybrids [2, 3], WideResNets [4] and MLP-Mixers [5].

Structured derivatives generally allow fastest computation, but are also able to process more models due to lower memory requirements (lower left; missing points indicate out-of-memory error). For the case of single output logit $O = 1$ (top row), **NTK-vector products** are generally detrimental due to costly forward pass **FP** relative to the size of parameters **P** (i.e. a lot of weight sharing; see Table 1). However, since **NTK-vector products** scale well with output size, for $O = 1000$ (bottom row), they perform comparably or better than other methods.

Finally, we remark that **Jacobian** not only runs out of memory faster, but can also take more time to compute. We conjecture that due to a larger memory footprint, **XLA** can sometimes perform optimizations that trade off speed for memory, and therefore compute the **Jacobian** in a less optimal way than if it had more memory available. Alternatively, **XLA** could also be performing simplifications of the NTK expression in these cases, such that those would not be possible in **Jacobian** computation alone.

See Fig. 2 for ResNets, and §K for details.

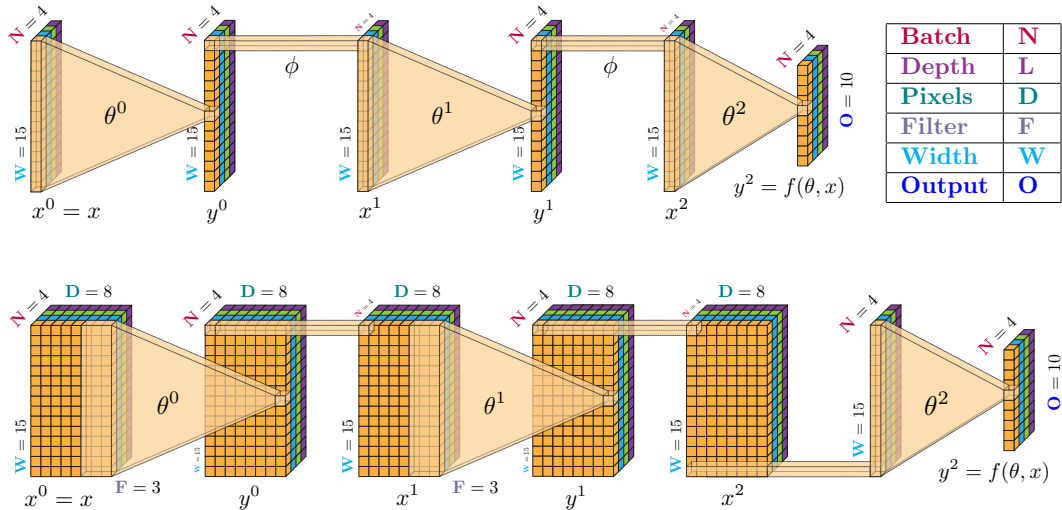


Figure 5: **Notation used in main text, §I (FCN, top) and §J (CNN, bottom).** For FCN, $D = F = 1$. For CNN, $D = 8$, $F = 3$, and the penultimate layer is global average pooling.

B Glossary

- **N** - batch size of inputs x to the NN $f(\theta, x)$. In a more general setting (§E), 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, 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 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 §E, 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 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 1, Table 2, and Table 3.
- \mathbf{J}_l^k is the cost of evaluating a single primitive Jacobian $\partial y^k / \partial \theta^l$, given the structure present in y^k according to §E.

C Motivation

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 [6–12] established that in the limit of infinite width, randomly initialized Neural Networks (NNs) are Gaussian Processes (called the NNGP). Building on this development, [13] showed that in function space the dynamics under gradient descent could be computed analytically using the so-called Neural Tangent Kernel (NTK) and [14] showed that wide neural networks reduce to their linearization in weight space throughout training. A related set of results [15, 16] 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 [17–19].

This notion inspired initialization schemes like Fixup [20], MetaInit [21], and Normalizer Free networks [22, 23] and has enabled efficient neural architecture search [24, 25]. The NTK has additionally given insight into a wide range of phenomena such as: behavior of Generative Adversarial Networks [26], neural scaling laws [27], and neural irradiance fields [28]. Kernel regression using the NTK has further enabled strong performance on small datasets [29], and applications such as dataset distillation [30, 31] and uncertainty prediction [32, 33].

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 can be computed for finite-width networks as the outer-product of Jacobians using forward- or reverse-mode automatic differentiation,

$$\Theta_{\theta}(x_1, x_2) = [\partial f(\theta, x_1)/\partial \theta] [\partial f(\theta, x_2)/\partial \theta]^T. \quad (5)$$

However, as we have shown in §I.3, this is often infeasible due to computational and memory requirements, and our work presents techniques that improve both.

D 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. [34, 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. [14, Figure S6] compute the NTK of up to 3-layer 2^{12} -wide or 1-layer 2^{14} -wide Fully Connected Networks (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 [4] to its infinite-width limit, Novak et al. [35, 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 [36] setting, Zhou et al. [37, Table 7] replace the inner training loop with NTK-based inference. They use up to 5-layer, 200-channel CNNs on MiniImageNet [38] 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. [24, §4.1] use the NTK to predict the generalization performance of architectures in the context of Neural Architecture Search [39, NAS]; however, the authors comment on its high computational burden and ultimately use a different proxy. In another NAS setting, Chen et al. [40, §3.1.1] use the condition number of NTK to predict a model’s trainability. Chen et al. [25, Table 1] also use the NTK to evaluate the trainability of several ImageNet [41] models such as ResNet 50/152 [42], Vision Transformer [2] and MLP-Mixer [5]. However, in all of the above cases the authors only evaluate a pseudo-NTK, i.e. an NTK of a scalar-valued function¹, which impacts the quality of the respective trainability/generalization proxy. In this work we can compute the full 1000×1000 NTK on the same models, 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.

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).

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	$\mathbf{NOYP} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N OYP} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P} + \mathbf{N}^2\mathbf{OY}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Block-diagonal	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C} + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Constant block-diagonal	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C} + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Input block-tiled	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C} + \mathbf{N}^2\mathbf{OY}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y}$
Output block-tiled	$\mathbf{NOYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{P} + \mathbf{NOY}$	$\mathbf{N OYP/C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C} + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C} + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$
Block-tiled	$\mathbf{NOYP/C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{P/C} + \mathbf{NOY}$	$\mathbf{N OYP/C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C}^2 + \mathbf{NOY}$	$\mathbf{N}^2\mathbf{Y}^2\mathbf{P/C}^3 + \mathbf{N}^2\mathbf{OY}^2/\mathbf{C}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{Y/C} + \mathbf{NOY}$

Table 4: **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. (6), for all n_1 and n_2 from 1 to \mathbf{N} (resulting in a $\mathbf{NO} \times \mathbf{NO}$ 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 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 2 and Table 3 for exact comparison in the case of convolution and matrix multiplication. See §B for legend.

E Types of structured derivatives

Here we continue §3 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 4.

E.1 No structure

We first consider the default cost of evaluating a single summand in Eq. (4), denoting individual matrix shapes underneath:

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) := \frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l} \frac{\partial f_2}{\partial y_2^{k_2}} =: \overbrace{\frac{\partial f_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2^T}{\partial \theta} \frac{\partial f_2^T}{\partial y_2}}^{\mathbf{O} \times \mathbf{O}} \quad (6)$$

$\underbrace{\frac{\partial f_1}{\partial y_1}}_{\mathbf{O} \times \mathbf{Y}} \underbrace{\frac{\partial y_1}{\partial \theta}}_{\mathbf{Y} \times \mathbf{P}} \underbrace{\frac{\partial y_2^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{Y}} \underbrace{\frac{\partial f_2^T}{\partial y_2}}_{\mathbf{Y} \times \mathbf{O}}$

We have dropped indices l, k_1 and k_2 on the right-hand side of Eq. (6) 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. There are 3 ways of contracting Eq. (6) 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 $\partial y / \partial \theta$.

E.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. (6) can be re-written as

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \frac{\partial f_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial f_2}{\partial y_2}^T \quad (7)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\oplus_{c=1}^{\mathbf{C}} \frac{\partial y_1^c}{\partial \theta_c} \right) \left(\oplus_{c=1}^{\mathbf{C}} \frac{\partial y_2^c}{\partial \theta_c} \right)^T \frac{\partial f_2}{\partial y_2}^T \quad (8)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\oplus_{c=1}^{\mathbf{C}} \left[\frac{\partial y_1^c}{\partial \theta_c} \frac{\partial y_2^c}{\partial \theta_c}^T \right] \right) \frac{\partial f_2}{\partial y_2}^T \quad (9)$$

$$= \sum_{c=1}^{\mathbf{C}} \frac{\partial f_1}{\partial y_1^c} \left[\frac{\partial y_1^c}{\partial \theta_c} \frac{\partial y_2^c}{\partial \theta_c}^T \right] \frac{\partial f_2}{\partial y_2^c}^T, \quad (10)$$

where we have applied the block matrix identity

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

We now perform a complexity analysis similar to Eq. (6):

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

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:** $\mathbf{Y}^2\mathbf{P}/\mathbf{C}^2 + \mathbf{OY}^2/\mathbf{C} + \mathbf{O}^2\mathbf{Y}$.

E.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 §E.2 with $\frac{\partial y^c}{\partial \theta_c} = \frac{\partial y^1}{\partial \theta_1}$ for any c . Here an identical analysis applies, yielding

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

and the same contraction complexities as in §E.2.

E.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

$$\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2) = \frac{\partial f_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial f_2}{\partial y_2}^T \quad (11)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbb{1}_{(1,\mathbf{C})} \otimes \frac{\partial y_1}{\partial \theta} \right) \left(\mathbb{1}_{(1,\mathbf{C})} \otimes \frac{\partial y_2}{\partial \theta} \right)^T \frac{\partial f_2}{\partial y_2}^T \quad (12)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbf{C} \mathbb{1}_{(1,1)} \otimes \left[\frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \right] \right) \frac{\partial f_2}{\partial y_2}^T \quad (13)$$

$$= \mathbf{C} \frac{\partial f_1}{\partial y_1} \left[\frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \right] \frac{\partial f_2}{\partial y_2}^T. \quad (14)$$

The matrix shapes are

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

Which leads to the following resulting complexities:

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 and inside-out-right:** $\mathbf{Y}^2\mathbf{P}/\mathbf{C} + \mathbf{OY}^2 + \mathbf{O}^2\mathbf{Y}$.

E.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 $(\mathbf{Y}/\mathbf{C}) \times \mathbf{P}$. This occurs during broadcasting or broadcasted arithmetic operations. In this case

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

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C},1)} \otimes \frac{\partial y_1^1}{\partial \theta} \right) \left(\mathbb{1}_{(\mathbf{C},1)} \otimes \frac{\partial y_2^1}{\partial \theta} \right)^T \frac{\partial f_2}{\partial y_2}^T \quad (16)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C},\mathbf{C})} \otimes \left[\frac{\partial y_1^1}{\partial \theta} \frac{\partial y_2^1}{\partial \theta}^T \right] \right) \frac{\partial f_2}{\partial y_2}^T \quad (17)$$

$$= \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial f_1}{\partial y_1^c} \right) \left[\frac{\partial y_1^1}{\partial \theta} \frac{\partial y_2^1}{\partial \theta}^T \right] \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial f_2}{\partial y_2^c} \right)^T, \quad (18)$$

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,

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

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

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

E.6 Block-tiled

Assume $\frac{\partial y}{\partial \theta} = \mathbb{1}_{(\mathbf{C}, \mathbf{C})} \otimes \frac{\partial y_1^1}{\partial \theta_1}$, where $\frac{\partial y_1^1}{\partial \theta_1}$ has shape $(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})$. This occurs for instance when y is a constant. In this case

$$\Theta_{\theta}^{l, k_1, k_2}(f_1, f_2) = \frac{\partial f_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T \frac{\partial f_2}{\partial y_2}^T \quad (19)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbb{1}_{(\mathbf{C}, \mathbf{C})} \otimes \frac{\partial y_1^1}{\partial \theta_1} \right) \left(\mathbb{1}_{(\mathbf{C}, \mathbf{C})} \otimes \frac{\partial y_2^1}{\partial \theta_1} \right)^T \frac{\partial f_2}{\partial y_2}^T \quad (20)$$

$$= \frac{\partial f_1}{\partial y_1} \left(\mathbf{C} \mathbb{1}_{(\mathbf{C}, \mathbf{C})} \otimes \left[\frac{\partial y_1^1}{\partial \theta_1} \frac{\partial y_2^1}{\partial \theta_1}^T \right] \right) \frac{\partial f_2}{\partial y_2}^T \quad (21)$$

$$= \mathbf{C} \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial f_1}{\partial y_1^c} \right) \left[\frac{\partial y_1^1}{\partial \theta_1} \frac{\partial y_2^1}{\partial \theta_1}^T \right] \left(\sum_{c=1}^{\mathbf{C}} \frac{\partial f_2}{\partial y_2^c} \right)^T, \quad (22)$$

This results in the following contraction:

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

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:** $\mathbf{Y}^2\mathbf{P}/\mathbf{C}^3 + \mathbf{OY}^2/\mathbf{C}^2 + \mathbf{O}^2\mathbf{Y}/\mathbf{C} + \mathbf{OY}$.

E.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 $f_1^{n_i}(\theta^0, \dots, \theta^L) := f(\theta^0, \dots, \theta^L, x_i^{n_i})$. 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 4.

E.8 Complex structure cost analysis

In previous sections 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 to be leveraged for NTK computation, and we intend to investigate it in future work.

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

Table 5: **Asymptotic time complexities of computing a single fully-connected layer NTK contribution.** See §E.9 for discussion, Table 4 for a more general setting, Table 3 for the case of deep networks, and §B for detailed legend.

E.9 Example

In §3 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 3. See §J 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, $\mathbf{Y} = \mathbf{W}$. We recognize that individual primitives $y^{k,n}(\theta_k) = \theta_l x^{k,n}$, 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 $\mathbf{C} = \mathbf{Y} = \mathbf{W} = \mathbf{J}$. Finally, **FP** costs \mathbf{W}^2 . Substituting all these equalities into Table 4 we get a simplified Table 5, that confirms the benefits of **NTK-vector products** and **Structured derivatives** for FCNs.

F Implementation

Both algorithms are implemented in JAX [1] 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 [35].

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 $\mathbf{I}_\mathbf{O}$ matrix and to parallelize the computation over pairwise combinations of \mathbf{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. (4). Each summand is a contraction of 4 factors: $\partial f_1 / \partial y_1, \partial y_1 / \partial \theta, \partial y_2 / \partial \theta, \partial f_2 / \partial y_2$.

First, we linearize f to obtain a computational graph constructed out of a limited set (54,² see Table 6) of linear primitives y^1, \dots, y^K . Then, we can obtain two factors $\partial f_1 / \partial y_1, \partial f_2 / \partial y_2$ as part of a backward pass almost identical to calling `jax.jacobian(f)(\theta, 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 (§E.8); for a given pair of primitives, these rules allow us to analytically simplify the contraction and avoid explicitly instantiating the derivatives.

Finally, 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.

²JAX leverages a similar approach to implement only 54 transpose rules for linear primitives for reverse-mode differentiation instead of 131 VJP rules [43].

Transposable primitive in <code>jax.ad.primitive_transposes</code>	CBD	BD	OBT
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 6: **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.

G Jacobian rules for structured derivatives

Here we discuss computing primitive $\partial y / \partial \theta$ Jacobians as part of our implementation in §F. 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 \mathbf{P} columns of the $I_{\mathbf{P}}$ identity matrix. Best for $\mathbf{P} < \mathbf{Y}$.
2. **Reverse mode**, `fwd = False`, is equivalent to `jax.jacrev`, reverse mode Jacobian computation, performed by applying the VJP to \mathbf{Y} columns of the $I_{\mathbf{Y}}$ identity matrix. Best for $\mathbf{P} > \mathbf{Y}$.
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 §L.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 $I_{\mathbf{P}/\mathbf{C}}$ or $I_{\mathbf{Y}/\mathbf{C}}$ respectively, and all methods will return a smaller Jacobian matrix of size $(\mathbf{Y}/\mathbf{C}) \times (\mathbf{P}/\mathbf{C})$. 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.

H 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 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.

I Complexity analysis for fully-connected networks

This section presents our contributions in a simplified setting of fully-connected (FCN) networks. For a more general discussion, see main text.

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 $\theta = \text{vec}[\theta^0, \dots, \theta^{\mathbf{L}}] \in \mathbb{R}^{\mathbf{LW}^2 + \mathbf{OW}}$. 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 $\mathbf{O} = \mathcal{O}(\mathbf{LW})$, i.e. the number of logits is dominated by the product of width and depth.

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

$$\Theta_\theta := \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T = \sum_{l=0}^{\mathbf{L}} \frac{\partial f(\theta, x_1)}{\partial \theta^l} \frac{\partial f(\theta, x_2)}{\partial \theta^l}^T =: \sum_{l=0}^{\mathbf{L}} \Theta_\theta^l \in \mathbb{R}^{\mathbf{O} \times \mathbf{O}}, \quad (23)$$

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

In §I.1 and §I.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 [44] 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 §I.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 §I.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 §I.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.

I.1 Jacobian-vector products and vector-Jacobian products

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

$$\text{JVP}_{(f, \theta, x)} : \theta_t \in \mathbb{R}^{\mathbf{LW}^2 + \mathbf{OW}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta} \theta_t \in \mathbb{R}^{\mathbf{O}}, \quad (24)$$

$$\text{VJP}_{(f, \theta, x)} : f_c \in \mathbb{R}^{\mathbf{O}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta}^T f_c \in \mathbb{R}^{\mathbf{LW}^2 + \mathbf{OW}}. \quad (25)$$

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. The time cost³ of both operations is comparable to the forward pass (**FP**), i.e. $[\mathbf{FP}] = [\text{cost of all intermediate layers}] + [\text{cost of the top layer}] = [\mathbf{LW}^2] + [\mathbf{OW}] \sim \mathbf{LW}^2$.

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

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

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 . 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. For the purposes of computing the NTK, we will be contracting Jacobians layerwise and so we will only

³To declutter notation, we omit the \mathcal{O} symbol to indicate asymptotic complexity in this work.

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

I.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.

$$[\partial f(\theta, x) / \partial \theta]^T = [\partial f(\theta, x) / \partial \theta]^T I_{\mathbf{O}} \in \mathbb{R}^{(\text{LW}^2 + \text{OW}) \times \mathbf{O}}. \quad (26)$$

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

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

Therefore, asymptotically,

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

I.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. (1):

$$\underbrace{\Theta_{\theta}^l}_{\mathbf{O} \times \mathbf{O}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta^l}}_{\mathbf{O} \times (\text{W} \times \text{W})} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta^l}}_{(\text{W} \times \text{W}) \times \mathbf{O}}. \quad (27)$$

The time cost of this contraction is $\mathbf{O}^2\text{W}^2$, and the memory necessary to instantiate each factor and the result is $\text{OW}^2 + \mathbf{O}^2$. Repeating the above operation for each θ^l , we arrive at LO^2W^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 N , then the resulting NTK is a matrix of shape $\text{NO} \times \text{NO}$, and the time cost becomes $\text{N}^2\text{LO}^2\text{W}^2$, while memory grows to $[\text{NTK matrix size}] + [\text{factors size}] = \text{N}^2\mathbf{O}^2 + \text{NOW}^2$.

What remains is to account for the cost of computing and storing individual derivatives $\partial f / \partial \theta^l$, which is exactly the cost of computing the **Jacobian** described in §I.2. Adding the costs up we obtain

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

I.4 Leveraging structured derivatives for computing the NTK

We can rewrite Θ_{θ}^l in Eq. (27) using the chain rule and our pre- and post-activation notation as:

$$\Theta_{\theta}^l = \left[\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \frac{\partial y_{x_1}^l}{\partial \theta^l} \right] \left[\frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l} \frac{\partial y_{x_2}^l}{\partial \theta^l} \right]^T = \underbrace{\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l}}_{\mathbf{O} \times \text{W}} \underbrace{\frac{\partial y_{x_1}^l}{\partial \theta^l}}_{\text{W} \times (\text{W} \times \text{W})} \underbrace{\frac{\partial y_{x_2}^l}{\partial \theta^l}}_{(\text{W} \times \text{W}) \times \text{W}} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial y_{x_2}^l}}_{\text{W} \times \mathbf{O}}. \quad (28)$$

At face value, rewriting Eq. (27) in this way is unhelpful as it appears to have introduced additional costly contractions. However, recall that $y^l = \theta^l x^l$, and therefore

$$\frac{\partial y_{x_1}^l}{\partial \theta^l} = I_{\mathbf{W}} \otimes x_1^{lT}, \quad \frac{\partial y_{x_2}^l}{\partial \theta^l} = I_{\mathbf{W}} \otimes x_2^{lT}, \quad (29)$$

where \otimes is the **Kronecker product**. Plugging Eq. (29) into Eq. (28) we get

$$\Theta_{\theta}^l(x_1, x_2) = \frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \left(I_{\mathbf{W}} \otimes x_1^{lT} \right) \left(I_{\mathbf{W}} \otimes x_2^{lT} \right)^T \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l} = \quad (30)$$

$$= \frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \left(I_{\mathbf{W}} \otimes \begin{bmatrix} x_1^{lT} & x_2^{lT} \end{bmatrix} \right) \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l} = \begin{pmatrix} x_1^{lT} & x_2^{lT} \end{pmatrix} \left[\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l} \frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l} \right]^T, \quad (31)$$

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

$$\Theta_{\theta}^l = \begin{pmatrix} \underbrace{x_1^{lT}}_{1 \times \mathbf{W}} & \underbrace{x_2^{lT}}_{\mathbf{W} \times 1} \end{pmatrix} \begin{bmatrix} \underbrace{\frac{\partial f(\theta, x_1)}{\partial y_{x_1}^l}}_{\mathbf{O} \times \mathbf{W}} & \underbrace{\frac{\partial f(\theta, x_2)}{\partial y_{x_2}^l}}_{\mathbf{W} \times \mathbf{O}} \end{bmatrix}, \quad (32)$$

and observe that it takes only $\mathbf{O}^2 \mathbf{W}$ time and $\mathbf{O} \mathbf{W} + \mathbf{O}^2$ memory. Accounting for depth, time cost increases by a factor of depth \mathbf{L} and becomes $\mathbf{L} \mathbf{O}^2 \mathbf{W}$, while memory does not change since the summands can be processed sequentially.

Batched inputs. In the batched setting, the time cost grows quadratically with the size of the NTK to $\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \mathbf{W}$, while the memory cost increases to $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{O} \mathbf{W}$ to store the result, $\Theta_{\theta}^l(x_1, x_2)$, and factors, $\partial f(\theta, x) / \partial y_x^l$, respectively.

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 (in other words, we cut the backpropagation algorithm short) the memory cost will be,

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

The extra time cost is asymptotically the cost of \mathbf{O} forward-passes, $\mathbf{N} \mathbf{L} \mathbf{O} \mathbf{W}^2$ 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. Putting everything together we find the following costs,

By leveraging **Structured derivatives** in NN computations, we have reduced the cost of NTK to $\mathbf{N}^2 \mathbf{L} \mathbf{O}^2 \mathbf{W} + \mathbf{N} \mathbf{L} \mathbf{O} \mathbf{W}^2$ time and $\mathbf{N}^2 \mathbf{O}^2 + \mathbf{N} \mathbf{O} \mathbf{W} + \mathbf{L} \mathbf{W}^2 + \mathbf{N} \mathbf{L} \mathbf{W}$ 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 §I.4 and detail in the §E, similar structure exists for many common operations such as convolutions, pooling, and arithmetic. However, the improvements discussed in this section do not emerge automatically in AD. While JAX and other libraries leverage structures analogous to Eq. (29) to efficiently compute single evaluations of the VJP and JVP, this structure is lost once the (structureless) Jacobian is instantiated (e.g. by composing the VJP with vectorization and contraction). We discuss how we impose this structure to compute the NTK for general neural networks in §E.

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

$$\frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)], \quad (33)$$

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 **Jacobian**, since it consists of \mathbf{O} VJPs followed by \mathbf{O} (cheaper) JVPs. Therefore it costs $\mathbf{LOW}^2 + \mathbf{O}^2\mathbf{W}$ time and $\mathbf{LW}^2 + \mathbf{OW}^2 + \mathbf{O}^2\mathbf{W}$ memory.

Batched inputs. In the batched setting Eq. (33) is repeated for each pair of inputs, and therefore time increases by a factor of \mathbf{N}^2 to become $\mathbf{N}^2\mathbf{LOW}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{W}$. However, the memory cost grows linearly in \mathbf{N} (except for the cost of storing the NTK of size $\mathbf{N}^2\mathbf{O}^2$), 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 cost equal to the cost of computing the **Jacobian** over a batch, i.e. $\mathbf{N}^2\mathbf{O}^2 + (\mathbf{LW}^2 + \mathbf{NOW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLW})$.

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.

J Complexity analysis for convolutional networks

Here we go through the same analysis as in §I 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 $\mathbf{O} = \mathcal{O}(\mathbf{LW})$.

J.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.

J.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] +

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

Table 7: **Time complexity of contracting** $\Theta_{\theta}^{l,k_1,k_2}(f_1, f_2)$ **corresponding to a CNN primitive** obtained by substituting $\mathbf{Y} = \mathbf{DW}$, $\mathbf{C} = \mathbf{W}$, and $\mathbf{P} = \mathbf{FW}^2$ into Table 4. The time cost of **Structured derivatives** are the minimum of the three entries due to using optimal contraction path by `np.einsum`.

$$\mathbf{ON}[\text{size of a single weight matrix}] + \mathbf{ON}[\text{activations in a single layer}] = [\mathbf{LFW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LDW} + \mathbf{O}] + \mathbf{ON}[\mathbf{FW}^2 + \mathbf{OW}] + \mathbf{ON}[\mathbf{DW}] \sim \mathbf{LW}^2 + \mathbf{NLDW} + \mathbf{NODW} + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W}$$

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.

J.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 (§J.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.

J.4 Structured derivatives

Convolution is **Constant block-diagonal** along the output channel axis with $\mathbf{C} = \mathbf{W}$, $\mathbf{P} = \mathbf{FW}^2$, $\mathbf{Y} = \mathbf{DW}$. Substituting this in Table 4, the cost of contraction is the minimum of the costs from Table 7. 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 $\mathbf{N}^2\mathbf{LO}^2 \min(\mathbf{FW}^2, \mathbf{DW}) + \mathbf{DFNLOW}^2$, as the minimum cost between the **Outside-in** and **Left-to-right**. Note that this dominates the time cost of the Jacobian from §J.2, so we don't need to modify it further. Memory due to Jacobian computation is $[\text{size of all weights}] + \mathbf{N}[\text{size of activations in all layers}] + \mathbf{NO}[\text{activations in a single layer}] + [\text{size of primitive derivatives}] = [\mathbf{LFW}^2 + \mathbf{OW}] + \mathbf{N}[\mathbf{LDW} + \mathbf{O}] + \mathbf{NO}[\mathbf{DW}] + \mathbf{N}[\mathbf{DW}] \sim \mathbf{LFW}^2 + \mathbf{NLDW} + \mathbf{NODW}$. Again, as in §I.4, and unlike other methods, we do not need to compute or store $\partial f / \partial \theta^l$ derivatives, allowing to avoid the $\mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W}$ extra memory overhead. However, we need to add the cost of storing the (subarray of) primitive Jacobians $\partial y / \partial \theta$, while have the size of $\mathbf{J} = \mathbf{YP} / \mathbf{C}^2 = \mathbf{DFW}$, hence the extra cost is \mathbf{NDFW} .

Structured derivatives cost $\mathbf{N}^2\mathbf{LO}^2 \min(\mathbf{FW}^2, \mathbf{DW}) + \mathbf{NLODFW}^2$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NDFW} + \mathbf{NLDW} + \mathbf{NODW} + \mathbf{LFW}^2$ memory.

J.5 NTK-vector products

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

Batched inputs. In a batched setting Eq. (33) is repeated for each pair of inputs, and therefore time increases by a factor of \mathbf{N}^2 to become $\mathbf{N}^2\mathbf{LODFW}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{W}$. Memory only grows linearly in \mathbf{N} (except for storing the result of size $\mathbf{N}^2\mathbf{O}^2$), by similar argument to §I.5, i.e. becomes $\mathbf{N}^2\mathbf{O}^2 + (\mathbf{LFW}^2 + \mathbf{NOFW}^2 + \mathbf{NO}^2\mathbf{W} + \mathbf{NLDW} + \mathbf{NODW})$ total memory.

NTK computation as a sequence of **NTK-vector products** costs $\mathbf{N}^2\mathbf{LODFW}^2 + \mathbf{N}^2\mathbf{O}^2\mathbf{W}$ time and $\mathbf{N}^2\mathbf{O}^2 + \mathbf{NOFW}^2 + \mathbf{LFW}^2 + \mathbf{NLDW} + \mathbf{NODW}$ memory.

K Experimental details

All experiments were performed in JAX [1] 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. 1 and **Fig. 3**: a 10-layer, ReLU FCN was constructed with the Neural Tangents [35] `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 [45] were used, specifically `flax.examples.imagenet.models`. For WideResNets, the **code sample** from Novak et al. [35] 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 \mathbf{N} in $\{2^k\}_{k=0}^9$, and reported the best time divided by \mathbf{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()`.