
What Breaks the Curse of Dimensionality in Deep Learning?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Although learning in high dimensions is commonly believed to suffer from the
2 curse of dimensionality, modern machine learning methods often exhibit an as-
3 tonishing power to tackle a wide range of challenging real-world learning prob-
4 lems without using abundant amounts of data. How exactly these methods break
5 this curse remains a fundamental open question in the theory of deep learning.
6 While previous efforts have investigated this question by studying the data (\mathcal{D}),
7 model (\mathcal{M}), and inference algorithm (\mathcal{I}) as independent modules, in this paper
8 we analyze the triple (\mathcal{D} , \mathcal{M} , \mathcal{I}) as an integrated system. We examine the basic
9 symmetries of such systems, focusing on four of the main architectures in deep
10 learning: fully-connected networks (FCN), locally-connected networks (LCN), and
11 convolutional networks with and without pooling (GAP/VEC). By computing an
12 eigen-decomposition of the infinite-width limits (aka Neural Kernels) of these
13 architectures, we characterize how inductive biases (locality, weight-sharing, pool-
14 ing, etc) and the breaking of spurious symmetries can affect the performance of
15 these learning systems. Our theoretical analysis shows that for many real-world
16 tasks it is locality rather than symmetry that provides the first-order remedy to the
17 curse of dimensionality. Empirical results on state-of-the-art models on ImageNet
18 corroborate our results.

19 1 Introduction

20 Statistical problems with high-dimensional data are frequently plagued by the *curse of dimensionality*,
21 in which the number of samples required to solve the problem grows rapidly with the dimensionality
22 of the input. Classical theory explains this phenomenon as the consequence of basic geometric and
23 algebraic properties of high-dimensional spaces; for example, the number of ϵ -cubes inside a unit
24 cube in \mathbb{R}^d grows exponentially like ϵ^{-d} , and the number of degree r polynomials in \mathbb{R}^d grows like a
25 power-law d^r . Since for real-world problems d is typically in the hundreds or thousands, classical
26 wisdom suggests that learning is likely to be infeasible. However, starting from the groundbreaking
27 work AlexNet [1], practitioners in deep learning have tackled a wide range of difficult real-world
28 learning problems ([2–6]) in high dimensions, once believed by many to be out-of-scope of current
29 techniques. The astonishing success of modern machine learning methods clearly contradicts the
30 curse of dimensionality and therefore poses the fundamental question: mathematically, how do
31 modern machine learning methods break the curse of dimensionality?

32 To answer this question, we must trace back to the most fundamental ingredients of machine learning
33 methods. They are the data (\mathcal{D}), the model (\mathcal{M}), and the inference algorithm (\mathcal{I}).

34 Data (\mathcal{D}) is of course central in machine learning. In the classical learning theory setting, the learning
35 objective usually has a power-law decay $m^{-\beta}$ as the function of the number of training samples
36 m . The theoretical bound on β is usually tiny, owing to the curse of dimensionality, and is of

37 limited practical utility for high-dimensional data. On the other hand, empirical measurements of
 38 β in state-of-the-art deep learning models typically reveal values of β that are not at all small (e.g.
 39 $\beta = 0.43$ for ResNet in Fig.S2) even though d is quite large (e.g. $d \sim 10^5$ for ImageNet). This
 40 example suggests that the learning curve must have important functional dependence on \mathcal{M} and \mathcal{I} .
 41 Indeed, as we will observe later, many of the best performing methods exhibit learning curves for
 42 which $\beta = \beta(m)$ actually *increases* as m becomes larger, i.e. data makes the usage of data more
 43 efficient. We call this phenomenon DIDE, for **d**ata **i**mproves **d**ata **e**fficiency.

44 Designing machine learning models (\mathcal{M}) that maximize data-efficiency is critical to the success
 45 of solving real-world tasks. Indeed, breakthroughs in machine learning are often driven by novel
 46 architectures LeNet [7], AlexNet[1], Transformer [2], etc. While some of the inductive biases of these
 47 methods are clear (e.g. translation symmetries of CNNs), others tend to build off of prior empirical
 48 success and are less well-understood (e.g. the implicit bias of SGD). To build our understanding of
 49 these biases and how they affect learning, we conduct a theoretical analysis of them in the infinite-
 50 width setting [8–12], which preserves most salient aspects of the architecture while enabling tractable
 51 calculations. We classify all phenomena that could be explained by infinite networks alone as the
 52 consequences of inductive biases.

53 The inference procedure (\mathcal{I}) is what enables *learning* in machine learning methods. It is widely
 54 believed that modern inference methods, specifically gradient descent and variants, ‘implicitly’ bias
 55 the solutions of the networks towards those that generalize well and away from those that generalize
 56 poorly [13–15]. The effects of the inference algorithm are intimately tied to the specifics of the model
 57 (e.g. weight-sharing) and the data (e.g. augmentation), and might not be fully understood with a
 58 fixed-data, fixed-model analysis. Indeed, good performance may derive from interactions between
 59 $(\mathcal{M}, \mathcal{I})$, or $(\mathcal{D}, \mathcal{I})$, or even $(\mathcal{D}, \mathcal{M}, \mathcal{I})$. In Sec. 3.1, we demonstrate the DIDE effect for a particular
 60 choice of $(\mathcal{D}, \mathcal{M}, \mathcal{I})$ and show that this effect disappears if any one of \mathcal{D} , \mathcal{M} , or \mathcal{I} is altered.

61 The above discussion highlights the insufficiency of treating \mathcal{D} , \mathcal{M} , and \mathcal{I} as separate non-interacting
 62 modules. They must be considered as an integrated system. Throughout this paper, we will refer to
 63 the triplet $(\mathcal{D}, \mathcal{M}, \mathcal{I})$ as a (machine) learning system and the tuple $(\mathcal{M}, \mathcal{I})$ as the learning algorithm
 64 of the system that operates on \mathcal{D} . We summarize our contributions below.

- 65 1. We surface the basic symmetries of various $(\mathcal{D}, \mathcal{M}, \mathcal{I})$ associated to four of the main arch-
 66 itectures in deep learning FCN_{*n*} (fully-connected networks), LCN_{*n*} (locally-connected
 67 networks), VEC_{*n*}/GAP_{*n*} (convolution networks with a flattening /a global average pooling
 68 readout layer), their infinite width counterparts FCN _{∞} /LCN _{∞} /VEC _{∞} /GAP _{∞} . Treating
 69 FCN_{*n*/ ∞} as the baseline model, we show that the locality from LCN_{*n*} and the weight-sharing
 70 from VEC_{*n*}/GAP_{*n*} break spurious symmetries and lead to better systems. Empirically, we
 71 examine the relation between the symmetries and the performance of the systems in the
 72 infinite width setting and finite width setting with various of interventions. Surprisingly,
 73 we observe that state-of-the-art learning system (EfficientNet[16]) on ImageNet can learn
 74 almost equally well even the coordinate of the data are transformed by the symmetry group
 75 defined by LCN_{*n*}.
- 76 2. We show that although the weight-sharing from VEC_{*n*} provides coordinate information of
 77 the data to the system, as the width gets larger, it becomes harder for the learning algorithm
 78 to explore such information and at infinite width, the system restores the symmetry group
 79 that is identical to LCN_{*n*}, and is completely unaware of the coordinate information. As a
 80 consequence, the performance of the network, as a function of width, monotonically decays
 81 [12]. This is in stark contrast to recent finding that the performance of network is positively
 82 correlated to its width. We show that this phenomenon continues to hold even with various
 83 interventions (larger learning rate and l2 regularization) to the training procedures. However,
 84 with more data (e.g. data augmentation) VEC_{*n*} can be on par with GAP_{*n*}.
- 85 3. The function space defined by LCN_{*n*} is a super set of that defined by VEC_{*n*}. We prove the
 86 opposite is true. Therefore, VEC_{*n*} is able to express functions in the space with a stronger
 87 inductive bias GAP_{*n*} (translation invariance) and functions in a seemingly much larger
 88 class LCN_{*n*}. We hypothesize that as the dataset grows, the learned functions using VEC_{*n*} is
 89 transitioned away from those learned using LCN_{*n*} and become closer to those learned using
 90 GAP_{*n*}. This suggests, even though the prior (provided by human) is not 100% correct, with
 91 the help of more data, gradient descent might be able to correct it, a possible explanation of
 92 DIDE.

93 4. When the input space is the product of hyperspheres, we eigendecompose the kernels
 94 associated to one-hidden layer infinite width network, FCN_∞ , $\text{VEC}_\infty = \text{LCN}_\infty$ and GAP_∞ .
 95 We treat FCN_∞ as the baseline, whose order r eigenspace has dimension of order d^r
 96 and eigenvalues of order d^{-r} for $r \geq 0$ [17]. We show that locality alone (i.e. VEC_∞)
 97 dramatically reduces the dimension of the r -eigenspace for $r \geq 2$ and the spectral gap
 98 between all r -eigenspaces but $r = 0$ and $r = 1$, making learning of higher order eigenspaces
 99 feasible with dramatically fewer samples and gradient steps. In addition, pooling (i.e.
 100 GAP_∞) reduces the dimension of r -eigenspace for $r \geq 1$ by a factor equal to the size of the
 101 pooling window, but it does not change the spectra in an essential way.

102 Our empirical and theoretical results surface the importance of locality which, we believe, provides
 103 the first-order remedy to the curse of dimensionality for many real-world tasks and which has been
 104 largely overlooked.

105 2 Preliminary and Notation

106 2.1 Neural Networks

107 We focus our presentation on the supervised learning setting and more concretely, on image
 108 recognition. Let $\mathcal{D} \subseteq (\mathbb{R}^d)^3 \times \mathbb{R}^k \equiv \mathbb{R}^{3d} \times \mathbb{R}^k$ denote the data set (training and test) and
 109 $\mathcal{X} = \{x : (x, y) \in \mathcal{D}\}$ and $\mathcal{Y} = \{y : (x, y) \in \mathcal{D}\}$ denote the input space (images) and label space,
 110 respectively. Here d is the spatial dimension (e.g. $d = 32 \times 32$ for CIFAR-10) of the images and 3 is
 111 the total number of channels (i.e. RGB). We use FCN_n to denote a L -hidden layer fully-connected
 112 network with identical hidden widths $n_l = n \in \mathbb{N}$ for $l = 1, \dots, L$ and with readout width $n_{L+1} = k$
 113 (the number of logits). For each $x \in \mathbb{R}^{3d} = (\mathbb{R}^d)^3$, we use $h^l(x), x^l(x) \in \mathbb{R}^{n_l}$ to represent the pre-
 114 and post-activation functions at layer l with input x . The recurrence relation FCN is given by

$$\begin{cases} h^{l+1} &= x^l W^{l+1} \\ x^{l+1} &= \phi(h^{l+1}) \end{cases} \text{ and } W_{i,j}^l = \frac{1}{\sqrt{n_l}} \omega_{ij}^l, \quad \omega_{ij}^l \sim \mathcal{N}(0, 1) \quad (1)$$

115 where ϕ is a point-wise activation function, $W^{l+1} \in \mathbb{R}^{n_l \times n_{l+1}}$ are the weights and ω_{ij}^l are the
 116 trainable parameters, drawn i.i.d. from a standard Gaussian $\sim \mathcal{N}(0, 1)$ at initialization. For simplicity
 117 of the presentation, the bias terms and the hyperparameters (the variances of the weights) are omitted.
 118 Adding them back won't affect the conclusion of the paper.

119 For convolutional networks or locally-connected networks, the inputs are treated as tensors in $(\mathbb{R}^d)^3$.
 120 The recurrent relation of convolutional networks can be written as

$$x_{\alpha,j}^{l+1} = \phi(h_{\alpha,j}^{l+1}) \quad \text{and} \quad h_{\alpha,j}^{l+1} \equiv \frac{1}{\sqrt{(2k+1)n^l}} \sum_{j=1}^{n^l} \sum_{\beta=-k}^k x_{\alpha+\beta,i}^l \omega_{ij,\beta}^l \quad (2)$$

121 Here $\alpha \in [d]$ denote the spatial location, $i/j \in [n]$ denotes the fanin/fanout channel indices. For
 122 notational convenience, we assume circular padding and stride equal to 1 for all layers. The features
 123 of the penultimate layer are 2D tensors and there are two commonly used approaches to map them
 124 to the logit layer: stack a dense layer after either vectorizing the 2D tensor to a 1D vector or
 125 applying a global average pooling layer to each channel. We use $\text{VEC}_n/\text{GAP}_n$ to denote the network
 126 obtain from the former/latter, which are known to be equipped with the inductive biases translation
 127 equivariant/invariant. The readout layer of $\text{VEC}_n/\text{GAP}_n$ could be written as

$$x_j^{L+1} = \frac{1}{\sqrt{dn}} \sum_{\alpha \in [d]} x_{\alpha,i}^L w_{\alpha,ij}^{L+1}, \quad x_j^{L+1} = \frac{1}{\sqrt{n}} \sum_{i \in [n]} \left(\frac{1}{d} \sum_{\alpha \in [d]} x_{\alpha,i}^L \right) w_{ij}^{L+1} \quad (3)$$

128 We briefly remark the the key difference between the two. In VEC_n , each pixel in the penultimate
 129 layer has its own (independent random) variable while pixels within the same channel shared the
 130 same (random) variable in GAP_n . It is clear that the function space of VEC_n contains that of GAP_n .
 131 Locally Connected Networks LCN_n [18, 19] are convolutional network without weight sharing
 132 between spatial locations. LCN_n preserve the connectivity pattern, and thus topology, of a convnet.
 133 Mathematically, the current formula is defined as in Equation 2 with all the *shared* parameters $\omega_{ij,\beta}^l$
 134 replaced by *unshared* $\omega_{ij,\alpha,\beta}^l \sim \mathcal{N}(0, 1)$

135 In this note, we assume that the LCN_n are always associated with a vectorization readout layer and it
 136 is clear, as a function space, LCN_n is a super set of VEC_n . Interestingly, the opposite is also true.

137 **Theorem 2.1** (Sec. B). *Let $\text{VEC}_n/\text{LCN}_n/\text{GAP}_n$ denote the set of functions that can be represented
 138 by L -hidden layer $\text{VEC}_n/\text{LCN}_n/\text{GAP}_n$ networks with hidden width n . Then*

$$\text{GAP}_n \subseteq \text{VEC}_n \subseteq \text{LCN}_n \subseteq \text{VEC}_{dn} \quad (4)$$

139 The significance of this theorem is that if we consider the function space VEC_n as a soft *prior*,
 140 gradient descent could move it closer to a *better* prior GAP_n (translation invariance) if the average
 141 pooling is (approximately) learned in the readout layer or it might remain close to LCN_n .

142 2.2 Gradient Descent Training

143 We use f to denote any functions defined by the architectures above and θ to denote the collection
 144 of all parameters. Denote by θ_t the time-dependence of the parameters and by θ_0 their initial
 145 values. We use $f_t(x) \equiv f(x, \theta_t) \in \mathbb{R}^k$ to denote the output (or logits) of the neural network at
 146 time t . Let $\ell(\hat{y}, y) : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}$ denote the loss function where the first/second argument is
 147 the prediction/true label. By applying continuous time gradient descent to minimize the objective
 148 $\mathcal{L} = \sum_{(x,y) \in \mathcal{D}} \ell(f_t(x, \theta), y)$, the evolution of the parameters θ and the logits f can be written as

$$\dot{\theta}_t = -\nabla_{\theta} f_t(\mathcal{X}_T)^T \nabla_{f_t(\mathcal{X}_T)} \mathcal{L}, \quad \dot{f}_t(\mathcal{X}_T) = \nabla_{\theta} f_t(\mathcal{X}_T) \dot{\theta}_t = -\hat{\Theta}_t(\mathcal{X}_T, \mathcal{X}_T) \nabla_{f_t(\mathcal{X}_T)} \mathcal{L} \quad (5)$$

149 where $f_t(\mathcal{X}_T) = \text{vec}([f_t(x)]_{x \in \mathcal{X}_T})$, the $k|\mathcal{D}| \times 1$ vector of concatenated logits for all examples, and
 150 $\nabla_{f_t(\mathcal{X}_T)} \mathcal{L}$ is the gradient of the loss with respect to the model's output, $f_t(\mathcal{X}_T)$. $\hat{\Theta}_t \equiv \hat{\Theta}_t(\mathcal{X}_T, \mathcal{X}_T)$
 151 is the tangent kernel at time t , which is a $k|\mathcal{D}| \times k|\mathcal{D}|$ kernel matrix

$$\hat{\Theta}_t = \nabla_{\theta} f_t(\mathcal{X}_T) \nabla_{\theta} f_t(\mathcal{X}_T)^T \quad (6)$$

152 One can define the tangent kernel for general arguments, e.g. $\hat{\Theta}_t(x, \mathcal{X}_T)$ where x is test input. At
 153 finite-width, $\hat{\Theta}$ will depend on the specific random draw of the parameters and evolve with time. As
 154 such, for a test point x the prediction $f_t(x)$ depends on the random initialization and is also stochastic.
 155 Note that the parameters are initialized randomly and the randomness will be carried out through the
 156 training procedure. As a consequence, the prediction functions are stochastic.

157 2.3 Infinite Network: Gaussian Processes and the Neural Tangent Kernels

158 **Neural Networks as Gaussian Processes (NNGP).** As the width $n \rightarrow \infty$, at initialization the
 159 output $f_0(\mathcal{X})$ forms a Gaussian Process $f_0(\mathcal{X}) \sim \mathcal{GP}(0, \mathcal{K}(\mathcal{X}, \mathcal{X}))$, known as the NNGP [8, 20, 21].
 160 Here \mathcal{K} is the GP kernel and can be computed in closed form for a variety of architectures. By treating
 161 this infinite width network as a Bayesian model (aka Bayesian Neural Networks) and applying
 162 Bayesian inference, the posterior is also a GP

$$\mathcal{N}(\mathcal{K}(\mathcal{X}_*, \mathcal{X}_T) \mathcal{K}^{-1}(\mathcal{X}_T, \mathcal{X}_T) \mathcal{Y}_T, \mathcal{K}(\mathcal{X}_*, \mathcal{X}_*) - \mathcal{K}(\mathcal{X}_*, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \mathcal{K}(\mathcal{X}, \mathcal{X}_*)^T) \quad (7)$$

163 **Neural Tangent Kernel (NTK).** Recent advance in global convergence theory of over-
 164 parameterized networks [22–25, 12] has shown that under certain assumptions, the tangent kernels is
 165 almost stationary over the course of training and is concentrated on its infinite width limit Θ in the
 166 sense there is a constant C independent of t and the network's width n such that

$$\sup_{t \geq 0} \|\hat{\Theta}_t(\mathcal{X}_T, \mathcal{X}_T) - \Theta(\mathcal{X}_T, \mathcal{X}_T)\|_F + \|\hat{\Theta}_t(\mathcal{X}_T, \mathcal{X}_*) - \Theta(\mathcal{X}_T, \mathcal{X}_*)\|_F \leq \frac{C}{\sqrt{n}}. \quad (8)$$

167 where is the infinite width limit of Θ at initialization, whose existence has been proved in [22, 26].
 168 As such, when the loss is the mean squared error (MSE), the mean prediction (marginalized over
 169 random initialization) has the following closed form

$$f(\mathcal{X}_*) = \Theta(\mathcal{X}_*, \mathcal{X}_T) \Theta^{-1}(\mathcal{X}_T, \mathcal{X}_T) \left(I - e^{-\eta \Theta(\mathcal{X}_T, \mathcal{X}_T) t} \right) \mathcal{Y}, \quad (9)$$

170 Letting $t \rightarrow \infty$, the above solution is the same as that of the kernel ridgeless regression using the
 171 infinite width tangent kernel Θ . We use $\text{FCN}_{\infty}(x)$, $\text{LCN}_{\infty}(x)$, $\text{VEC}_{\infty}(x)$ and $\text{GAP}_{\infty}(x)$ to denote
 172 the infinite width solutions (either the GP inference or the NTK regression) for the corresponding
 173 architectures, where we have suppressed the dependence on the training data $(\mathcal{X}_T, \mathcal{Y}_T)$.

174 **3 Symmetries of Machine Learning Systems**

175 Symmetry is fundamental in physical systems. So is it in machine learning systems. We explore
 176 symmetries of various machine learning systems in this section. Given $\mathcal{D} = (\mathcal{X}, \mathcal{Y})$ and a transforma-
 177 tion on the input space $\tau : \mathbb{R}^{3d} \rightarrow \mathbb{R}^{3d}$, we set $\tau(\mathcal{D}) = (\tau(\mathcal{X}), \mathcal{Y})$. Let $O(3d)$ denote the orthogonal
 178 group on the flattened input space \mathbb{R}^{3d} . The subgroup $O(3)^d \leq O(3d)$ operates on the un-flattened
 179 input $(\mathbb{R}^d)^3$, whose element rotates each pixel $x_\alpha \in \mathbb{R}^3$ by an independent element $\tau_\alpha \in O(3)$. The
 180 smaller subgroup $O(3) \otimes \mathbf{I}_d \leq O(3)^d$ applies the *shared* rotation (i.e. $\tau_\alpha = \tau$ to all x_α for $\alpha \in [d]$).
 181 We use $P(3d)$ to denote the permutation group on \mathbb{R}^{3d} and $P(3)^d$ and $P(3) \otimes \mathbf{I}_d$ are defined similarly.
 182 Note that rotating \mathcal{X} by τ is equivalent to transfer the original coordinate system by the adjoint
 183 tranformation $\tau^* = \tau^{-1}$.

184 For a deterministic (stochastic) learning algorithm $\mathcal{A} = (\mathcal{M}, \mathcal{I})$, we use $\mathcal{A}(\mathcal{D}_T)$ to denote the learned
 185 function (distribution of the learned functions) using training set \mathcal{D}_T . We use $\mathcal{A}^\tau(\mathcal{D}_T)$ to denote
 186 the learned function(s) using $\tau(\mathcal{D}_T)$ and makes prediction on the transformed test point $\tau(\mathcal{X}_*)$. In
 187 another word, the learning algorithm is conducted in the input space whose coordinate system is
 188 transformed by τ^{-1} .

189 **Definition 1.** Let \mathcal{G} be a group of transformations $\mathbb{R}^{3d} \rightarrow \mathbb{R}^{3d}$. We say a deterministic (stochastic)
 190 learning algorithm $\mathcal{A} = (\mathcal{M}, \mathcal{I})$ is g -invariant if $\mathcal{A} = \mathcal{A}^g$ ($\mathcal{A} \stackrel{d}{=} \mathcal{A}^g$). In this case, we say the
 191 system $(\mathcal{D}, \mathcal{M}, \mathcal{I})$ is g -invariant and use the notation $(\mathcal{D}, \mathcal{M}, \mathcal{I}) = (g\mathcal{D}, \mathcal{M}, \mathcal{I})$. If this holds for all
 192 $g \in \mathcal{G}$, then we say the algorithm and the system are \mathcal{G} -invariant.

193 If $(\mathcal{M}, \mathcal{I})$ is the algorithm of minimum norm linear regressor, then $(\mathcal{D}, \mathcal{M}, \mathcal{I})$ is $O(3)^d$ -invariant;
 194 see Sec.G for more details. Note that the symmetry (invariance) in our definition is a property of a
 195 system and is different from the notion of symmetry that are commonly used in the machine learning
 196 community, which is a property of a function (e.g. translation invariance).

197 **Theorem 3.1** (Sec.C). *If the parameters of the networks are initialized with iid $\mathcal{N}(0, 1)$, then*

- | | | | |
|-----|--|-----|---|
| 198 | • FCN _{n/∞} are O(3d)-invariant. | 200 | • VEC _n is O(3) ⊗ I _d -invariant and VEC _∞
is O(3) ^d -invariant. |
| | | 201 | |
| 199 | • LCN _{n/∞} are O(3) ^d -invariant. | 202 | • GAP _{n/∞} are O(3) ⊗ I _d -invariant. |

203 The O(3d)-invariant of FCN_∞ is because the NTK/NNGP kernel is an inner product kernel, namely,
 204 there is a function k such that the kernels have the form $k(\langle x, x' \rangle)$. The O(3d)-invariant of finite
 205 width FCN_n is due to the Gaussian initialization of the first layer which was first observed and
 206 proved in [27]. Rotating the input by $\tau \in O(3d)$ is equivalent to rotating the weight matrix ω of
 207 the first layer by τ^* . Since for $\omega \in \mathcal{N}(0, 1)^{3d}$ $\tau^* \omega \stackrel{d}{=} \omega$, at random initialization, the distribution
 208 of the output functions (the prior) are unchanged if all inputs are rotated by the same element in
 209 $O(3d)$. This property continues to hold throughout the course of (continue/discrete) gradient descent
 210 training with/without L^2 -regularization and Bayesian posterior inference. For the same reason, LCN_n
 211 is $O(3)^d$ -invariant because each patch of the image uses independent Gaussian random variables.
 212 However, weight-sharing in VEC_n and GAP_n breaks the $O(3)^d$ symmetry, reducing it to $O(3) \otimes \mathbf{I}_d$.

213 For infinite networks, LCN_∞ = VEC_∞ [28–31]. The kernels of VEC_∞ and GAP_∞ are of the forms

$$\Theta_{\text{VEC}}(x, x') = k(\{\langle x_\alpha, x'_\alpha \rangle\}_{\alpha \in [d]}) \quad \text{and} \quad \Theta_{\text{GAP}}(x, x') = k(\{\langle x_\alpha, x'_{\alpha'} \rangle\}_{\alpha, \alpha' \in [d]}), \quad (10)$$

214 resp. The former depends only on the inner product between pixels in the same spatial location,
 215 breaking the O(3d) symmetry and reducing it to O(3)^d. In addition, the latter depends also on the
 216 inner products of pixels across different spatial locations due to pooling, which breaks the O(3)^d
 217 symmetry and reduces it to O(3) ⊗ I_d.

218 Note that $\dim(O(3d)) = 3d(3d - 1)/2$, $\dim(O(3)^d) = 3d$ and $\dim(O(3) \otimes \mathbf{I}_d) = 3$. LCN_n/VEC_∞
 219 dramatically reduces the dimension of the symmetry group. It is worth mentioning that while
 220 $\dim(O(3d))$ many pairs of rotated and unrotated images are needed to recover the exact rotation in
 221 $O(3d)$, only 3 pairs are sufficient for $O(3)^d$, same as that of $O(3) \otimes \mathbf{I}_d$. The results of the paper
 222 are presented in the most *vanilla* setting. Our methods can easily extend to more complicated
 223 architectures like ResNet[32], MLP-Mixer[33] and etc. The symmetry groups of such systems
 224 need to be computed in a case-by-case manner by identifying the invariant group of the random
 225 initialization and training procedures.

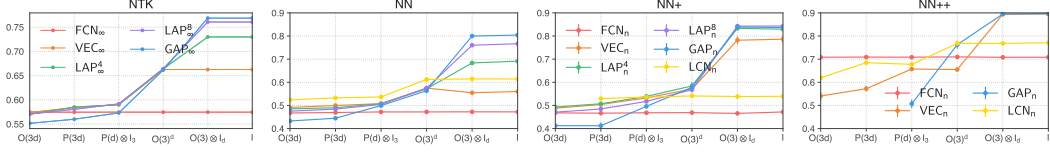


Figure 1: **Performance vs Symmetry.** Machine learning systems are equipped with various kinds of symmetries. Transforming the system by the associated symmetry does not affect the performance of the system. However, injecting spurious symmetries beyond the associated symmetries could dramatically degrade their performance for both finite and infinite networks.

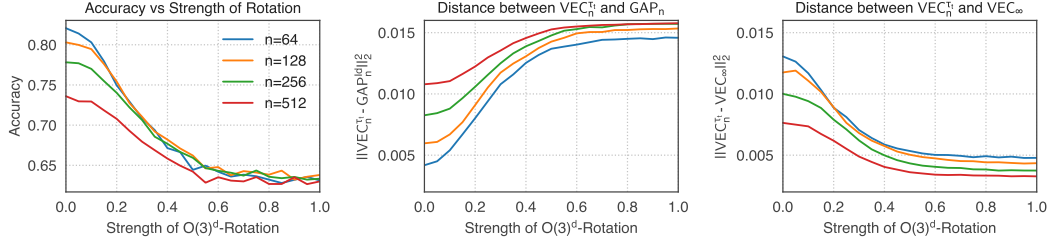


Figure 2: Even in the NN+ setting, VEC_n is closer to GAP_n for small n and moves towards VEC_∞ with more symmetries and/or larger n and accuracy drops.

226 3.1 Empirical Supports and Observations

227 **Performance under Rotations.** We examine the performance of: FCN, VEC, LCN, GAP and
 228 $LAP^{4/8}$, when the coordinates of the data are transformed by six different groups (x -axis in Fig.1)
 229 using the standard dataset CIFAR-10. Here $LAP^{4/8}$ is the same as GAP except the readout layer is
 230 replaced by the **Local Average Pooling** with window size $4 \times 4/8 \times 8$. We consider 4 types of training
 231 methods: (1) NTK, i.e. infinite networks (2) NN, our baseline for finite width neural network which
 232 is trained with momentum using a small learning rate and without L^2 regularizer and the network
 233 is centered (+C) to reduce the variance from random initialization (3) NN+= NN+LR+L2-C, i.e.
 234 using a larger learning rate (+LR), adding L^2 regularization (+L2)) and removing the centering (-C)
 235 (4) NN++=NN++DA, adding MixUp[34] data augmentation (+DA) to NN+. Overall, we observe
 236 that, for most of the cases in NTK/NN/NN+, adding spurious symmetry to a system ($\mathcal{D}, \mathcal{M}, \mathcal{I}$)
 237 degrades the performance towards that of the system invariant to that symmetry. Surprisingly, in
 238 the baseline NN, performance of $VEC_n + O(3) \otimes \mathbf{I}_d$ rotation is slightly worse than that of $VEC_n + O(3)^d$
 239 and than that of LCN_n , indicating that the system with $\mathcal{M} = VEC_n$ is likely operating closely
 240 on the $O(3)^d$ symmetry. The interventions -C+L2+LR in NN+ distinguishes the performance of
 241 $VEC_n + O(3) \otimes \mathbf{I}_d$ from $VEC_n + O(3)^d$ and +DA eventually closes the performance gap between
 242 $VEC_n + O(3) \otimes \mathbf{I}_d$ and $GAP_n + O(3) \otimes \mathbf{I}_d$, helping the system to be aware of the smaller symmetry
 243 $O(3) \otimes \mathbf{I}_d$, escaping from the $O(3)^d$ symmetry.

244 **Symmetry Breaking of VEC_n .** Assuming Equation 8, namely, the network is in the NTK regime,

$$\lim_{n \rightarrow \infty} |\mathbb{E}VEC_n(x) - VEC_\infty(x)| + \lim_{n \rightarrow \infty} |\mathbb{E}VEC_n(x) - \mathbb{E}VEC_n^\tau(x)| \leq Cn^{-\frac{1}{2}} \quad (11)$$

245 where the expectation \mathbb{E} is over random initialization and $VEC_n(x)$ is the prediction of the test point
 246 x when $t = \infty$, i.e. training loss is 0. VEC_n^τ is the prediction of the τ -rotated system, $\tau \in O(3)^d$.
 247 The $O(3)^d$ symmetry is restored as $n \rightarrow \infty$. As such, for large n , the system is approximately $O(3)^d$ -
 248 invariant. We randomly sample $\tau \in O(3)^d$ and use the exponential map to construct a continuous
 249 interpolation $\tau_t \in O(3)^d$ with $\tau_0 = \mathbf{Id}$ and $\tau_1 = \tau$. We train the network as in NN++ (+LR+L2-C)
 250 using different n and τ_t and average the predictions over 10 random initialization as an approximation
 251 of $\mathbb{E}VEC_n^{\tau_t}(x)$. Not surprisingly, as n increases and/or t increases, (1) test performance decays
 252 monotonically (left panel in Fig.2), (2) the distance to $\mathbb{E}GAP_n$ increases monotonically (middle
 253 panel) and (3) distance to VEC_∞ decrease monotonically (right panel). Clearly, the coordinate
 254 information from the data is utilized by smaller width VEC_n .

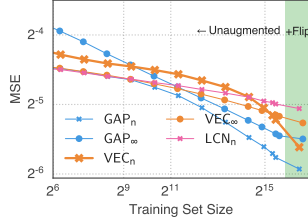


Figure 3: **Data Bends Learning Curve of VEC_n .** We study the effect of training set size to the network’s performance for various models. In the small dataset regime, the slope of the learning curve (in the log-log plot) of VEC_n is similar to that of VEC_∞ and FCN_n . However, as the dataset gets larger, the slope increases significantly. This is hinted by Theorem 2.1.

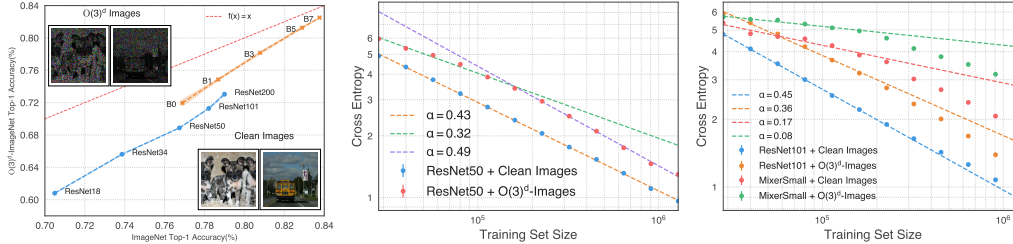


Figure 4: With coordinate of the input data rotated by $O(3)^d$, state of the art models learn as good as without rotation. middle/right: slopes of the learning curves increases due to more data. DIDE

255 **DIDE for VEC_n .** To understand the role of data, we vary the training set size of Cifar10 from about
 256 2^6 to 50k (the whole un-augmented training set) and to 100k (adding left-right flip augmentation) and
 257 plot the learning curves in Fig.3. We observe dramatic speedup of learning for VEC_n in the larger
 258 data set regime, which isn’t the case for VEC_∞ (kernel), LCN_n , GAP_∞ and even for GAP_n
 259 after $m = 2^{12}$. We argue that this is due to the prior (the function space defined by the model) is too large
 260 (and not optimal) for the task and the coupled effect of more data together with inference procedures
 261 corrects the prior, as it is suggested by Theorem 2.1.

262 **DIDE for SOTA models.** In the middle and right panels of Fig.S2, we provide additional evidence
 263 in a larger scale setting. We generate learning curves of ImageNet using ResNet50 and MLP-Mixer,
 264 a very recent architecture that contains no convolution layers except the first layer, which is a
 265 convolution with filter size and stride equal to $(16, 16)$ (patches are disjoint). The symmetry group
 266 associated to ResNet is similar to that of GAP_n which is relatively small. However, the symmetry
 267 group induced by the first layer of the Mixer is $O(3 \times 16^2) \otimes \mathbf{I}_{14^2}$, where 3×16^2 is number of entries
 268 in the $(16, 16, 3)$ patch (RGB channels) and $14^2 = 224^2/16^2$ is the number of patches. Although
 269 the dimension of $O(3 \times 16^2) \otimes \mathbf{I}_{14^2}$ is quite large (about $(3 \times 16^2)^2/2$), it is still dramatically
 270 smaller than that of applying a fully-connected layer to the flatten images, which $O(3 \times 224^2)$ (about
 271 $(3 \times 224^2)^2/2$). In the middle panel of Fig.S2, we observe an almost perfect power-law scaling for
 272 the learning curve for the ResNet50 system with unrotated images. When the images are rotated by
 273 $O(3)^d$ ($d = 224^2$), the learning curve is relatively flat in the smaller data regime (green dashed line).
 274 However, the data set grows, it eventually catches up (purple dashed line) as that of the unrotated
 275 setting; see Sec.E for ResNet34/101. In the third panel, we see the learning curves are much flatter
 276 (red) for the Mixer and even more so for the rotated images (green). Again, these curves are bent
 277 towards that of ResNet50 with unrotated images as data increases, indicating the prior was being
 278 corrected.

279 Finally, in the left panel of Fig.S2, we compare the accuracy of state-of-the-art models trained on both
 280 unrotated and $O(3)^d$ rotated images. Surprisingly, the gap between the two are not large and becomes
 281 smaller for better performant models. For EfficientNet B7¹, the top-1 accuracy of the rotated system
 282 is only 1.2% off from the unrotated one.

283 4 Eigencomposition of Neural Kernels

284 To get insights into the inductive biases, we eigendecompose the kernels using spherical harmonics.
 285 We assume the input space $\mathcal{X} = \{\xi = (\xi_0, \dots, \xi_{p-1}) \in (\sqrt{d_0}S^{(d_0-1)})^p\} \subseteq \mathbb{R}^{d_0 p}$, i.e.

¹Still under training

286 the p -product of $(d_0 - 1)$ -sphere with radius $\sqrt{d_0}$. We call $\xi_i \in \sqrt{d_0}\mathbb{S}^{(d_0-1)}$ a mini-patch and
 287 $(\xi_i, \xi_{i+1}, \dots, \xi_{i+s-1}) \in (\sqrt{d_0}\mathbb{S}^{(d_0-1)})^s$ a patch for $i \in [p]$, where circular boundary condition is
 288 assumed. We consider the asymptotic limit when $d_0 = d^\alpha$, $p = d^{1-\alpha}$ and $d = pd_0 \rightarrow \infty$ and treat
 289 $0 < \alpha < 1$ and s as fixed constant. The input space \mathcal{X} is associated with the product measure
 290 $\mu \equiv \sigma_{d_0}^p$, where σ_{d_0} is the normalized uniform measure on $\sqrt{d_0}\mathbb{S}^{(d_0-1)}$. The kernels associated to
 291 the one-hidden layer infinite networks (either NNGP or NTK) has the following general forms

$$k \left(\frac{1}{p} \sum_{i \in [p]} \xi_i^T \eta_i / d_0 \right) \frac{1}{p} \sum_{i \in [p]} k \left(\frac{1}{s} \sum_{b \in [s]} \xi_{i+b}^T \eta_{i+b} / d_0 \right) \frac{1}{p^2} \sum_{i,j \in [p]} k \left(\frac{1}{s} \sum_{b \in [s]} \xi_{i+b}^T \eta_{j+b} / d_0 \right), \quad (12)$$

292 although that exact form of the (positive definite) kernel function $k : \mathbb{R} \rightarrow \mathbb{R}$ depends on the kernel
 293 types (NNGP vs NTK), activations, hyperparameters and etc. We assume the kernel is sufficiently
 294 smooth in $(-1, 1)$ and the Taylor expansion of $k^{(r)}$ converges uniformly in $[-1, 1]$ for sufficiently
 295 many $r \in \mathbb{N}$. We use the notation that $A \sim B$ if there are positive constants c and C such that
 296 $cA \leq B \leq CA$ for d sufficiently large. We use \mathcal{K} to represent any kernels above and consider it as a
 297 Hilbert–Schmidt operator on $L^2(\mathcal{X}, \mu)$

$$\mathcal{K}f(\xi) = \int_{\mathcal{X}} \mathcal{K}(\xi, \eta) f(\eta) d\mu, \quad f \in L^2(\mathcal{X}, \mu), \quad (13)$$

298 which is well-defined since μ is a probability measure and k is bounded. Let $\vec{r} = (r_0, \dots, r_{p-1}) \in \mathbb{N}^p$,
 299 τ the shifting operator $\tau \vec{r} = (r_{p-1}, r_0, \dots, r_{p-2})$. The s -banded subset of \mathbb{N}^p is defined to be

$$B(\mathbb{N}^p, s) = \{\vec{r} \in \mathbb{N}^p : \text{dist}(\text{argmax}_j r_j \neq 0, \text{argmin}_j r_j \neq 0) \leq s - 1\} \quad (14)$$

300 which is a quantifier used to restrict the support of a function on a patch. Here $\text{dist}(i, j) = \min\{|i - j|, p - |i - j|\}$, a distance defined on the cyclic group $[p] = \mathbb{Z}/p\mathbb{Z}$. The quotient space $B(\mathbb{N}^p, s)/\tau$
 301 denotes a subset of $B(\mathbb{N}^p, s)$ by identifying $\vec{v} = \vec{v}'$ as the same element if $\vec{v} = \tau^a \vec{v}'$ for some $a \in [p]$.
 302 Finally, $Y_{r_j, l_j}(\xi_j)$ is used to denote the l_j -th spherical harmonic of degree r_j in the unit sphere
 303 $\mathbb{S}^{(d_0-1)}$ and has unit norm under the normalized measure on $\mathbb{S}^{(d_0-1)}$. As such $Y_{r_j, l_j}(\xi_j / \sqrt{d_0}) \in$
 304 $L^2(\sqrt{d_0}\mathbb{S}^{(d_0-1)}, \sigma_{d_0})$ has unit norm. Recall that the total number of spherical harmonic of degree
 305 r_j in $\mathbb{S}^{(d_0-1)}$ is $N(d_0, r_j) = (2r_j + d_0 - 2) \binom{d_0 + r_j - 3}{r_j - 1} \sim d_0^{r_j} / r_j!$ as $d_0 \rightarrow \infty$. We use $N(d_0, \vec{r}) =$
 306 $\prod_{j \in [p]} N(d_0, r_j)$ and $[N(d_0, \vec{r})] = \prod_{j \in [p]} [N(d_0, r_j)]$, resp. Let

$$\vec{Y}_{\vec{r}, \vec{l}}(\xi) = \prod_{j \in [p]} Y_{r_j, l_j}(\xi_j) \quad (15)$$

308 The following theorem shows that locality (VEC_∞) dramatically reduces both the dimensions of
 309 Eigendecomposition $r \geq 1$ eigenspaces and the spectral gap between them. In addition, pooling
 310 (i.e. translation symmetry of GAP_n) reduces their dimensions by a factor of p . See Sec.E for the
 311 implication of this theorem to learning.

312 **Theorem 4.1.** [Sec.D] We have the following eigendecomposition for the integral operator \mathcal{K}

$$\mathbb{H} = \bigcup_{r \in \mathbb{N}} \mathbb{H}^{(r)} = \bigcup_{r \in \mathbb{N}} \bigcup_{\vec{r} \in Q(\mathcal{K}, r)} \mathbb{H}^{(\vec{r})}, \quad (16)$$

313 where $Q(\mathcal{K}, r)$ is a quantifier defined below. If $r = 0$, $\mathbb{H}^{(0)}$ is the space of constant functions and the
 314 eigenvalue is $\sim k(0)$. For $r \geq 1$,

315 1. if $\mathcal{K} = \mathcal{K}_{\text{FCN}}$, then $Q(\mathcal{K}, r) = \{\vec{r} \in \mathbb{N}^p : |\vec{r}| = r\}$ and the unit eigenfunctions are

$$\begin{cases} \mathbb{H}^{(\vec{r})} = \text{span} \left\{ Y_{\vec{r}, \vec{l}} \right\}_{\vec{l} \in [B(d_0, \vec{r})]} \\ \dim(\mathbb{H}^{(r)}) \sim d^r \quad \text{and} \quad \lambda(\mathbb{H}^{(\vec{r})}) \sim d^{-r} \delta(k^{(r)}(0)) \end{cases} \quad (17)$$

316 2. if $\mathcal{K} = \mathcal{K}_{\text{VEC}}$, $Q(\mathcal{K}, r) = \{\vec{r} \in B(\mathbb{N}^p, s) : |\vec{r}| = r\}$ the unit eigenfunctions are

$$\begin{cases} \mathbb{H}_{\text{VEC}}^{(\vec{r})} = \text{span} \left\{ Y_{\vec{r}, \vec{l}} \right\}_{\vec{l} \in [B(d_0, \vec{r})]} \\ \dim(\mathbb{H}_{\text{VEC}}^{(r)}) \sim p(sd_0)^r = s^r d^{1-\alpha+ra} \quad \text{and} \quad \lambda(\mathbb{H}_{\text{VEC}}^{(\vec{r})}) \sim p^{-1}(sd_0)^{-r} \delta(k^{(r)}(0)) \end{cases} \quad (18)$$

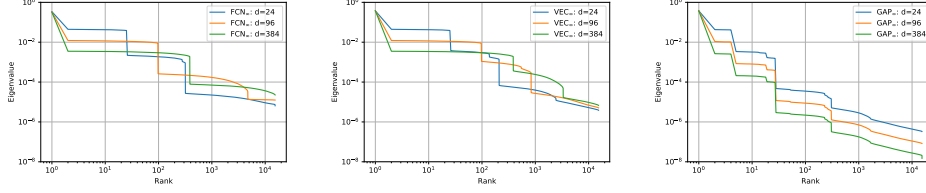


Figure 5: **Eigenvalue Decay of Relu NTK of FCN_∞ , VEC_∞ and GAP_∞ .** $d_0 = s = 3$. The eigenvalues of GAP_∞ decays *faster* because with $m = 15k$ many samples, higher order eigenspace can be covered by GAP_∞ but not FCN_∞/VEC_∞ due to Theorem 4.1.

317 3. and finally, if $\mathcal{K} = \mathcal{K}_{GAP}$, then $Q(\mathcal{K}, r) = \{\vec{r} \in B(\mathbb{N}^p, s)/\tau : |\vec{r}| = r\}$, the unit eigenfunc-
 318 tions are

$$\left\{ \begin{array}{l} H_{GAP}^{(\vec{r})} = \text{span} \left\{ \frac{1}{\sqrt{p}} \sum_{\tau \in [p]} Y_{\vec{r}, i}(\tau \xi) \right\}_{\vec{i} \in [B(d_0, \vec{r})]} \\ \dim(H_{GAP}^{(r)}) \sim (sd_0)^r = s^r d^{r\alpha} \quad \text{and} \quad \lambda(H_{GAP}^{(\vec{r})}) \sim p^{-1} (sd_0)^{-r} \delta(k^{(r)}(0)) \end{array} \right. \quad (19)$$

319 5 Related Work

320 The study of infinite networks dates back to seminal work by Neal [8] who showed the convergence of
 321 single hidden-layer networks to Gaussian Processes (GPs). Recently, there has been renewed interest
 322 in studying random, infinite, networks starting with concurrent work on “conjugate kernels” [10, 35]
 323 and “mean-field theory” [9, 36], taking a statistical learning and statistical physics view of points,
 324 resp. Since then this analysis has been extended to include a wide range for architectures [20, 21, 37,
 325 29, 26, 38]. The inducing kernel is often referred to as the Neural Network Gaussian Process (NNGP)
 326 kernel. The neural tangent kernel (NTK), first introduced in Jacot et al. [22], along with followup
 327 work [12, 39] showed that the distribution of functions induced by gradient descent for infinite-width
 328 networks is a Gaussian Process with NTK as the kernel.

329 The study of implicit bias (regularization) of gradient descent has received considerable interests.
 330 The work [15, 40–43] demonstrate the convergence of SGD to the maximal margin solution for
 331 logistic-type of losses during late time training. [44–50] study the early-time SGD dynamics, spectral
 332 biases of neural networks. These results aim to explain the order of learning of neural networks:
 333 functions of less complexity are usually learned before more complex functions.

334 [27] is the first to show that the prediction functions obtained from training FCN depend, in addition
 335 on the labels, only on the covariance of the input data. This implies our result regarding the $O(3d)$
 336 invariance of FCN. By utilizing this symmetry, recent work [51] constructs a particular task where
 337 the label function is a second order polynomial of the inputs and show that orthogonal invariance
 338 algorithm requires sample size of order d^2 while there is a convnet requires only $O(1)$ samples. Their
 339 convnet essentially corresponds to the $d_0 = s = 1$ and $r = 2$ case of Theorem 4.1, in which the
 340 dimension of this eigenspace (and indeed of all r -eigenspace by treating r as a finite constant as
 341 $d \rightarrow \infty$) of GAP_∞ is $O(1)$ while the dimension of the 2-eigenspace of FCN_∞ is of order d^2 .

342 6 Conclusion

343 In this paper, we consider machine learning methods as an integrated system of data, models and
 344 inference algorithms and study the basic symmetries of various machine learning systems. We surface
 345 the importance of locality in modern machine learning systems through large scale empirical study
 346 and through an eigendecomposition of one-layer infinite networks. However, we haven’t addressed
 347 the two import questions (1) theoretical characterization of the effect of composing locality and (2)
 348 the mathematical understanding of DIDE and how the prior is corrected by the coupled effect of data
 349 and gradient descent. We leave them to future work.

350 Checklist

351 The checklist follows the references. Please read the checklist guidelines carefully for information on
352 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
353 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
354 the appropriate section of your paper or providing a brief inline description. For example:

- 355 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 356 • Did you include the license to the code and datasets? **[No]** The code and the data are
357 proprietary.
- 358 • Did you include the license to the code and datasets? **[N/A]**

359 Please do not modify the questions and only use the provided macros for your answers. Note that the
360 Checklist section does not count towards the page limit. In your paper, please delete this instructions
361 block and only keep the Checklist section heading above along with the questions/answers below.

- 362 1. For all authors...
 - 363 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
364 contributions and scope? **[Yes]**
 - 365 (b) Did you describe the limitations of your work? **[Yes]**
 - 366 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - 367 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
368 them? **[Yes]**
- 369 2. If you are including theoretical results...
 - 370 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 371 (b) Did you include complete proofs of all theoretical results? **[Yes]**
- 372 3. If you ran experiments...
 - 373 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
374 mental results (either in the supplemental material or as a URL)? **[No]**
 - 375 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
376 were chosen)? **[No]**
 - 377 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
378 ments multiple times)? **[Yes]**
 - 379 (d) Did you include the total amount of compute and the type of resources used (e.g., type
380 of GPUs, internal cluster, or cloud provider)? **[No]**
- 381 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 382 (a) If your work uses existing assets, did you cite the creators? **[TODO]**
 - 383 (b) Did you mention the license of the assets? **[TODO]**
 - 384 (c) Did you include any new assets either in the supplemental material or as a URL?
385 **[TODO]**
 - 386 (d) Did you discuss whether and how consent was obtained from people whose data you're
387 using/curating? **[TODO]**
 - 388 (e) Did you discuss whether the data you are using/curating contains personally identifiable
389 information or offensive content? **[TODO]**
- 390 5. If you used crowdsourcing or conducted research with human subjects...
 - 391 (a) Did you include the full text of instructions given to participants and screenshots, if
392 applicable? **[TODO]**
 - 393 (b) Did you describe any potential participant risks, with links to Institutional Review
394 Board (IRB) approvals, if applicable? **[TODO]**
 - 395 (c) Did you include the estimated hourly wage paid to participants and the total amount
396 spent on participant compensation? **[TODO]**

397 **References**

- 398 [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional
399 neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- 400 [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz
401 Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- 402 [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-
403 tional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 404 [4] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian
405 Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go
406 with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 407 [5] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli
408 Qin, Augustin Židek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein structure prediction
409 using potentials from deep learning. *Nature*, 577(7792):706–710, 2020.
- 410 [6] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray,
411 Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint*
412 *arXiv:2001.08361*, 2020.
- 413 [7] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
414 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 415 [8] Radford M. Neal. Priors for infinite networks (tech. rep. no. crg-tr-94-1). *University of Toronto*, 1994.
- 416 [9] Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential
417 expressivity in deep neural networks through transient chaos. In *Advances In Neural Information Processing*
418 *Systems*, 2016.
- 419 [10] Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The
420 power of initialization and a dual view on expressivity. In *Advances In Neural Information Processing*
421 *Systems*, 2016.
- 422 [11] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization
423 in neural networks. *arXiv preprint arXiv:1806.07572*, 2018.
- 424 [12] Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein,
425 and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent.
426 In *Advances in Neural Information Processing Systems*, 2019.
- 427 [13] Behnam Neyshabur. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*, 2017.
- 428 [14] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Implicit bias of gradient descent on
429 linear convolutional networks. *arXiv preprint arXiv:1806.00468*, 2018.
- 430 [15] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias
431 of gradient descent on separable data, 2018.
- 432 [16] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In
433 *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.
- 434 [17] Alex J Smola, Zoltan L Ovari, Robert C Williamson, et al. Regularization with dot-product kernels.
435 *Advances in neural information processing systems*, pages 308–314, 2001.
- 436 [18] Kuniyuki Fukushima. Cognitron: A self-organizing multilayered neural network. *Biological cybernetics*,
437 20(3-4):121–136, 1975.
- 438 [19] Yann Lecun. Generalization and network design strategies. In *Connectionism in perspective*. Elsevier,
439 1989.
- 440 [20] Jaehoon Lee, Yasaman Bahri, Roman Novak, Sam Schoenholz, Jeffrey Pennington, and Jascha Sohl-
441 dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Repre-*
442 *sentations*, 2018.
- 443 [21] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani.
444 Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning*
445 *Representations*, 2018.

- 446 [22] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization
447 in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- 448 [23] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global
449 minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018.
- 450 [24] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-
451 parameterization. In *International Conference on Machine Learning*, 2018.
- 452 [25] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized
453 deep relu networks. *Machine Learning*, 109(3):467–492, 2020.
- 454 [26] Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior,
455 gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- 456 [27] Neha S. Wadia, Daniel Duckworth, Samuel S. Schoenholz, Ethan Dyer, and Jascha Sohl-Dickstein.
457 Whitening and second order optimization both destroy information about the dataset, and can make
458 generalization impossible. *arxiv preprint arXiv:2008.07545*, 2020.
- 459 [28] Lechao Xiao, Yasaman Bahri, Jascha Sohl-Dickstein, Samuel Schoenholz, and Jeffrey Pennington. Dy-
460 namical isometry and a mean field theory of cnns: How to train 10,000-layer vanilla convolutional neural
461 networks. In *International Conference on Machine Learning*, pages 5393–5402, 2018.
- 462 [29] Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A. Abolafia,
463 Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels
464 are gaussian processes. In *International Conference on Learning Representations*, 2019.
- 465 [30] Adrià Garriga-Alonso, Laurence Aitchison, and Carl Edward Rasmussen. Deep convolutional networks as
466 shallow gaussian processes. In *International Conference on Learning Representations*, 2019.
- 467 [31] Roman Novak, Lechao Xiao, Jiri Hron, Jaehoon Lee, Alexander A Alemi, Jascha Sohl-Dickstein, and
468 Samuel S Schoenholz. Neural tangents: Fast and easy infinite neural networks in python. *arXiv preprint*
469 *arXiv:1912.02803*, 2019.
- 470 [32] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
471 In *Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- 472 [33] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner,
473 Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An all-mlp architecture for
474 vision. *arXiv preprint arXiv:2105.01601*, 2021.
- 475 [34] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical
476 risk minimization. In *International Conference on Learning Representations*, 2018. URL [https://
477 openreview.net/forum?id=r1Ddp1-Rb](https://openreview.net/forum?id=r1Ddp1-Rb).
- 478 [35] Amit Daniely. SGD learns the conjugate kernel class of the network. In *Advances in Neural Information*
479 *Processing Systems 30*. 2017.
- 480 [36] Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information
481 propagation. *International Conference on Learning Representations*, 2017.
- 482 [37] Lechao Xiao, Yasaman Bahri, Jascha Sohl-Dickstein, Samuel Schoenholz, and Jeffrey Pennington. Dy-
483 namical isometry and a mean field theory of CNNs: How to train 10,000-layer vanilla convolutional neural
484 networks. In *International Conference on Machine Learning*, 2018.
- 485 [38] Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk for
486 deep attention networks, 2020.
- 487 [39] Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In
488 *Advances in Neural Information Processing Systems*, pages 2937–2947, 2019.
- 489 [40] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In
490 *International Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?
491 id=SJeLIgBKPS](https://openreview.net/forum?id=SJeLIgBKPS).
- 492 [41] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference*
493 *on Learning Theory*, pages 1772–1798, 2019.
- 494 [42] Ziwei Ji and Matus Jan Telgarsky. Gradient descent aligns the layers of deep linear networks. In *7th*
495 *International Conference on Learning Representations, ICLR 2019*, 2019.

- 496 [43] Lénaïc Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
497 trained with the logistic loss. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty*
498 *Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages
499 1305–1338. PMLR, 09–12 Jul 2020. URL <http://proceedings.mlr.press/v125/chizat20a.html>.
- 500 [44] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred Zhang,
501 and Boaz Barak. Sgd on neural networks learns functions of increasing complexity. *arXiv preprint*
502 *arXiv:1905.11604*, 2019.
- 503 [45] Wei Hu, Lechao Xiao, Ben Adlam, and Jeffrey Pennington. The surprising simplicity of the early-time
504 learning dynamics of neural networks. *arXiv preprint arXiv:2006.14599*, 2020.
- 505 [46] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua
506 Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on*
507 *Machine Learning*, pages 5301–5310. PMLR, 2019.
- 508 [47] Zhiqin John Xu. Understanding training and generalization in deep learning by fourier analysis. *arXiv*
509 *preprint arXiv:1808.04295*, 2018.
- 510 [48] Zhi-Qin John Xu, Yaoyu Zhang, Tao Luo, Yanyang Xiao, and Zheng Ma. Frequency principle: Fourier
511 analysis sheds light on deep neural networks. *arXiv preprint arXiv:1901.06523*, 2019.
- 512 [49] Lili Su and Pengkun Yang. On learning over-parameterized neural networks: A functional approximation
513 perspective. *arXiv preprint arXiv:1905.10826*, 2019.
- 514 [50] Greg Yang and Hadi Salman. A fine-grained spectral perspective on neural networks. *arXiv preprint*
515 *arXiv:1907.10599*, 2019.
- 516 [51] Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than
517 fully-connected nets? *arXiv preprint arXiv:2010.08515*, 2020.