

# How Wide Convolutional Neural Networks Learn Hierarchical Tasks

Francesco Cagnetta<sup>1\*</sup> Alessandro Favero<sup>1\*</sup> Matthieu Wyart<sup>1</sup>

<sup>1</sup>EPFL, <sup>\*</sup>Equal contribution.

## Abstract

Despite their success, understanding how convolutional neural networks (CNNs) can efficiently learn high-dimensional functions remains a fundamental challenge. A popular belief is that these models harvest the compositional and hierarchical structure of natural data such as images. Yet, we lack a quantitative understanding of how such structure affects performance, e.g. the rate of decay of the generalisation error with the number of training samples. In this paper, we study deep CNNs in the kernel regime: *i)* we show that the spectrum of the corresponding kernel and its asymptotics inherit the hierarchical structure of the network; *ii)* we use generalisation bounds to prove that deep CNNs adapt to the spatial scale of the target function; *iii)* we illustrate this result by computing the rate of decay of the error in a teacher-student setting, where a deep CNN is trained on the output of another deep CNN with randomly-initialised parameters. We find that, if the teacher function depends on certain low-dimensional subsets of the input variables, then the rate is controlled by the effective dimensionality of these subsets. Conversely, if the teacher function depends on the full set of input variables, then the error rate is inversely proportional to the input dimension. Interestingly, this implies that despite their hierarchical nature, the data generated by deep CNNs are too rich to be efficiently learnable in high dimensions.

## How can CNNs learn in high-dimensions?

- Learning a generic high-dimensional function is plagued by the **curse of dimensionality**: the rate at which the generalisation error  $\epsilon$  decays with the number of training examples  $n$  vanishes when the input dimension  $d$  grows, i.e.  $\epsilon = O(n^{-\beta})$  with  $\beta = 1/d$  [1].
- Convolutional neural networks (CNNs) can learn data whose dimensions can be in the hundreds or more. This points to the existence of some underlying structure that CNNs can leverage. What is this structure?
- A popular hypothesis is that learnable tasks are compositional and hierarchical: features at any scale are made of sub-features at smaller scales. Can we quantify the effect of this structure on the generalisation error's decay  $\beta$ ?

## Deep CNNs and the kernel regime

**$L$ -hidden-layers hierarchical CNN.** Denote with  $\sigma$  the normalised ReLU function,  $\sigma(x) = \sqrt{2} \max(0, x)$ . For each input  $\mathbf{x} \in \mathbb{R}^d$ , the output of a  $L$ -hidden-layers hierarchical neural network can be defined recursively as follows.

$$\begin{aligned} f_{h,i}^{(1)}(\mathbf{x}) &= \sigma(\mathbf{w}_h^{(1)} \cdot \mathbf{x}_i), \forall h \in [1 \dots H_1], \forall i \in [1 \dots p_1]; \\ f_{h,i}^{(l)}(\mathbf{x}) &= \sigma\left(\frac{1}{\sqrt{H_{l-1}}} \sum_{h'} \frac{\mathbf{w}_{h,h'}^{(l)} \cdot (\mathbf{f}_{h'}^{(l-1)})_i}{\sqrt{s_l}}\right), \forall h \in [1 \dots H_l], i \in [1 \dots p_l], l \in [2 \dots L]; \\ f(\mathbf{x}) &= f^{(L+1)}(\mathbf{x}) = \frac{1}{\sqrt{H_L}} \sum_{h=1}^{H_L} \sum_{i=1}^{p_L} \frac{w_{h,i}^{(L+1)} f_{h,i}^{(L)}(\mathbf{x})}{\sqrt{p_L}}. \end{aligned} \quad (1)$$

$H_l$  denotes the width of the  $l$ -th layer,  $s_l$  the filter size (with  $s_1 \equiv s$ , the size of the input patches),  $p_l$  the number of patches (with  $p_1 \equiv p$ , the number of input patches). The  $p_l$ 's and  $s_l$ 's satisfy  $s_1 p_1 = d$  and  $p_{l-1} = s_l p_l$ .  $\mathbf{w}_h^{(1)} \in \mathbb{R}^{s_1}$ ,  $\mathbf{w}_{h,h'}^{(l)} \in \mathbb{R}^{s_l}$ ,  $w_{h,i}^{(L+1)} \in \mathbb{R}$ .

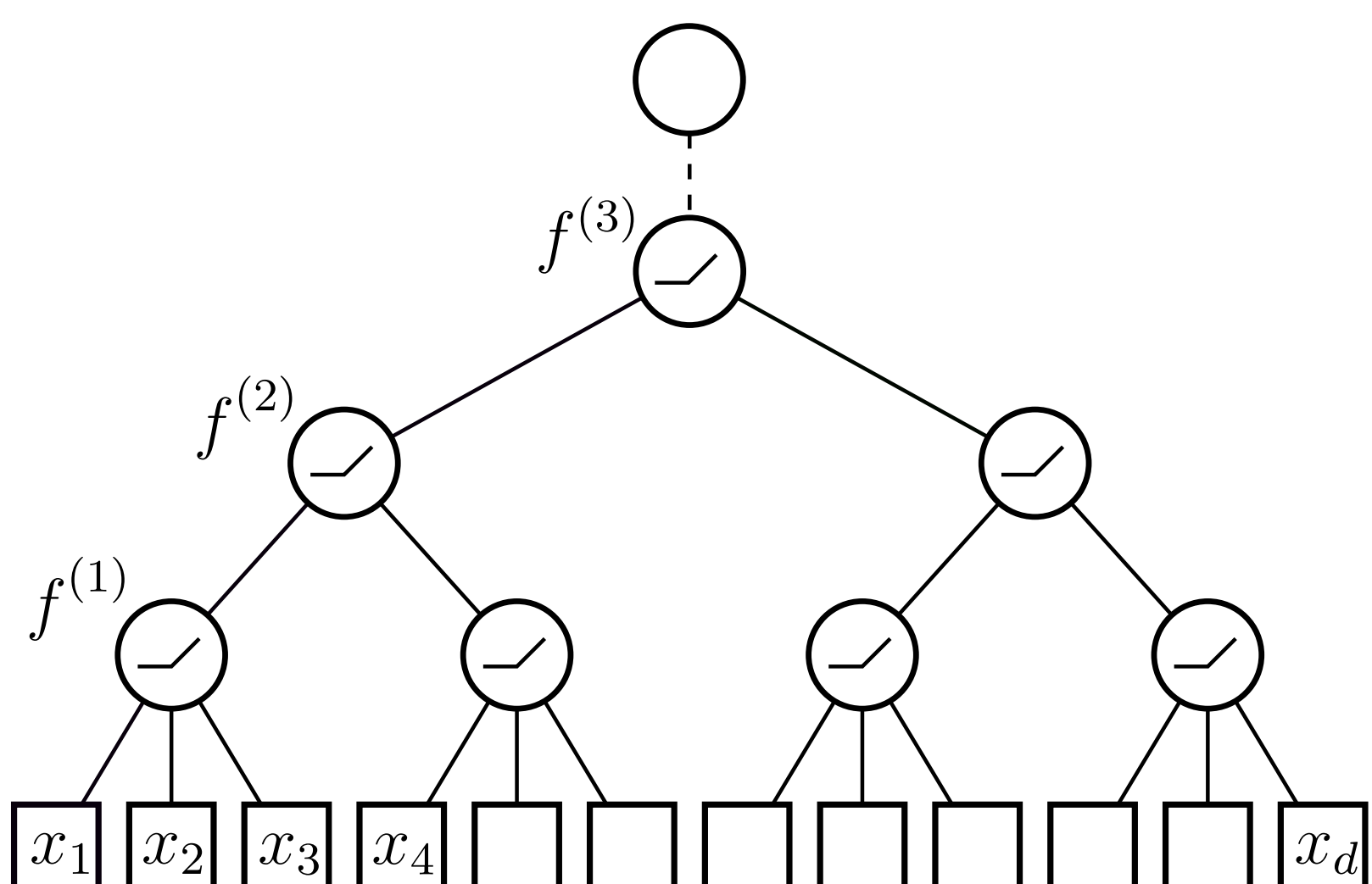


Figure 1. DAG of a hierarchical CNN with  $L = 4$ . The leaves represent the input coordinates, the root the output and all other nodes (channels of) hidden neurons.

**Convolutional neural tangent kernel (CNTK).** Assume that the weights of the network are initialised from a Gaussian distribution  $\mathcal{N}(0, 1)$  and the number of channels diverges, i.e.  $H_l \rightarrow \infty$ . Then, denoting with  $\theta$  the collection of the network's parameter, the network's function  $f(\mathbf{x}; \theta)$  learned by gradient descent on a square loss to zero training error is equivalent to the function obtained from ridgeless kernel regression with the neural tangent kernel (NTK) [2, 3]:

$$\mathcal{K}_{\text{CNTK}}(\mathbf{x}, \mathbf{x}') = \nabla_{\theta} f(\mathbf{x}; \theta_0) \cdot \nabla_{\theta} f(\mathbf{x}'; \theta_0), \quad (2)$$

where  $\theta_0$  denotes the values of the parameters at initialisation. Interestingly, the NTK inherits the hierarchical structure of the network's function. In this limit, studying the generalisation properties of CNNs becomes equivalent to studying those of their NTK.

## Statistical properties and spatial adaptation

**Supervised learning setup.** Consider a data distribution  $\rho$  on data-label pairs  $(\mathbf{x}, y)$  such that  $\mathbb{E}[y|\mathbf{x}] = f^*(\mathbf{x})$  for a target function  $f^*$ . Given  $n$  examples i.i.d. from  $\rho$ , the KRR estimator is given by

$$f_{\lambda}^n(\mathbf{x}) = \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{\mu=1}^n (f(\mathbf{x}_{\mu}) - y_{\mu})^2 + \lambda \|f\|_{\mathcal{H}}^2, \quad (3)$$

where  $\mathcal{H}$  is the RKHS of the kernel  $\mathcal{K}$ . The generalisation error of this estimator is defined as

$$\epsilon(f_{\lambda}^n) = \mathbb{E}_{(\mathbf{x}, y) \sim \rho} [(f_{\lambda}^n(\mathbf{x}) - y)^2]. \quad (4)$$

**Bound on the excess risk.** The statistical properties of a kernel  $\mathcal{K}$  are deeply connected with the eigendecomposition of its integral operator

$$\mathcal{T}_{\mathcal{K}} f(\mathbf{x}) = \int \mathcal{K}(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') d\tau(\mathbf{x}'), \quad (5)$$

with  $d\tau$  the probability measure that generates the inputs. Let  $\Lambda_{\rho}$  denote its eigenvalues. If one can find  $\alpha \geq 1$  and  $r \geq 1 - 1/\alpha$  satisfying  $\text{Tr}(\mathcal{T}_{\mathcal{K}}^{1/\alpha}) < \infty$

(capacity condition) and  $\|T_{\mathcal{K}}^{1-r} f^*\|_{\mathcal{H}}^2 < \infty$  (source condition), then, by choosing a  $n$ -dependent optimal  $\lambda_n$ , one gets the following bound on the excess risk [4, 5]:

$$\mathbb{E}_n[\epsilon(f_{\lambda}^n)] - \epsilon(f^*) \leq C n^{-\frac{\alpha r}{\alpha r + 1}}. \quad (6)$$

**Spectral analysis of the CNTK.** Let  $d\tau$  be the uniform measure on the product of  $p$   $s$ -dimensional unit spheres  $\prod_i \mathbb{S}^{s-1}$ . Then, the eigenfunctions of  $\mathcal{T}_{\mathcal{K}_{\text{CNTK}}}$  are  $\tilde{Y}_{\mathbf{k}, \ell} = \prod_i Y_{k_i, \ell_i}(\mathbf{x}_i)$ , where  $Y_{k, \ell}$  denotes a spherical harmonic. In the paper, we show that the eigenfunctions of  $\mathcal{T}_{\mathcal{K}_{\text{CNTK}}}$  can be organised into groups associated with the hidden layers of the network. The eigenfunctions of each group have an *effective dimensionality*  $d_{\text{eff}}$ , which coincides with the receptive field of the neurons of the associated hidden layer. The corresponding eigenvalues scale as a power of the polynomial degree  $k$  of the eigenfunctions, with an exponent depending on  $d_{\text{eff}}$ . For  $s = 2$ , spherical harmonics reduce to Fourier atoms and the eigenfunctions reduce to the multidimensional Fourier basis. Then, for  $\|\mathbf{k}\| \rightarrow \infty$ ,

$$\Lambda_{\mathbf{k}}^{(l)} = C_{2,l} \|\mathbf{k}\|^{-1-d_{\text{eff}}(l)} + o(\|\mathbf{k}\|^{-1-d_{\text{eff}}(l)}), \quad (7)$$

with  $d_{\text{eff}} = s_2 s_3 \dots s_L$ .

## Adaptivity to spatial structure.

(Informal) If the support of a target function  $f^*$  is entirely contained in a set of contiguous patches with dimension  $d_{\text{eff}}(l)$  ( $l \leq L$ ), then only the first  $l$  groups of the spectrum of  $\mathcal{T}_{\mathcal{K}_{\text{CNTK}}}$  contribute to the source condition. E.g., if  $f^*$  is Lipschitz, one obtains a fast rate of

$$n^{-2/(3+d_{\text{eff}}(l))} \quad (8)$$

as  $d_{\text{eff}}(L) \rightarrow \infty$ , which does not vanish with the input dimensionality!

**Statistical mechanics of generalisation for the ridgeless limit.** In the ridgeless case, where the correspondence between kernel methods and neural networks actually holds, there are no rigorous bounds. However, one can derive similar results using the replica methods of statistical physics [6].

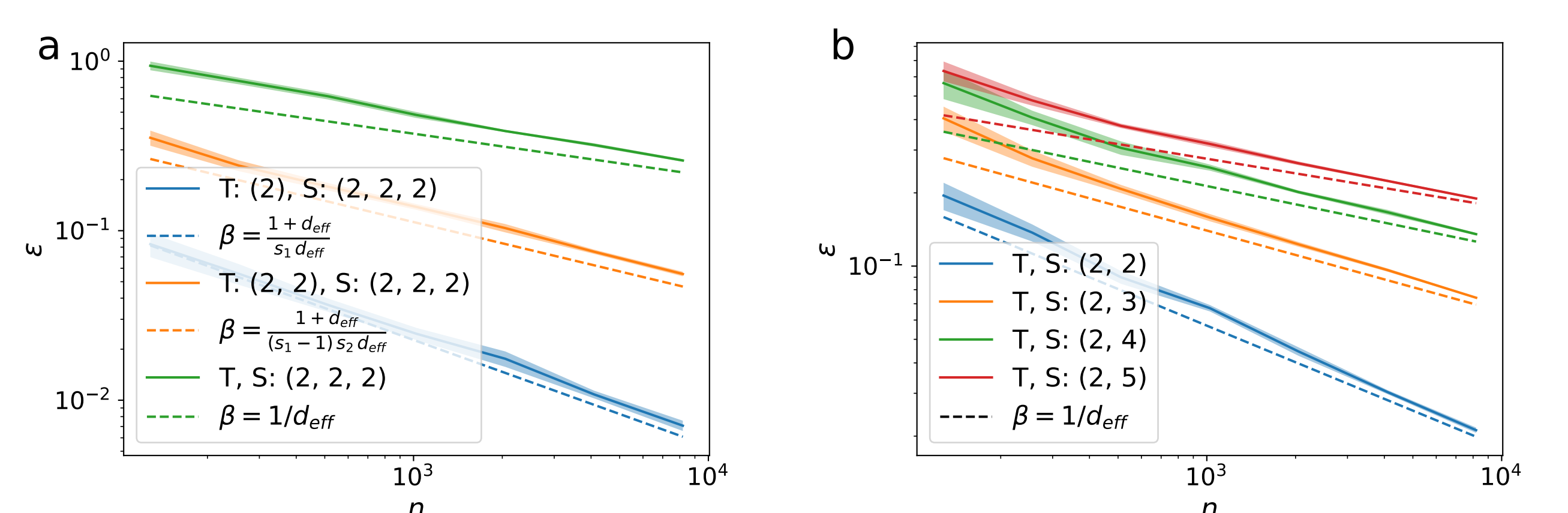


Figure 2. Agreement of theoretical predictions (dashed) with numerical experiments in a teacher (T) student (S) setting where data are generated and learned by CNNs in the kernel limit.

## References

- [1] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. 2019.
- [2] Arthur Jacot, Franck Gabriel, and Clément Hongler. *Neural tangent kernel: Convergence and generalization in neural networks*. 2018.
- [3] Sanjeev Arora et al. *On exact computation with an infinitely wide neural net*. 2019.
- [4] Andrea Caponnetto and Ernesto De Vito. *Optimal rates for the regularized least-squares algorithm*. 2007.
- [5] Francis Bach. *Learning Theory from First Principles*. 2021.
- [6] Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. *Spectrum dependent learning curves in kernel regression and wide neural networks*. 2020.