
BEYOND IID WEIGHTS: SPARSE AND LOW-RANK DEEP NEURAL NETWORKS ARE ALSO GAUSSIAN PROCESSES

Thiziri Nait Saada, Alireza Naderi & Jared Tanner

Mathematical Institute

University of Oxford

{naitsaadat, naderi, tanner}@maths.ox.ac.uk

ABSTRACT

The infinitely wide neural network has been proven a useful and manageable mathematical model that enables the understanding of many phenomena appearing in deep learning. One example is the convergence of random deep networks to Gaussian processes that allows a rigorous analysis of the way the choice of activation function and network weights impacts the training dynamics. In this paper, we extend the seminal proof of Matthews et al. (2018) to a larger class of initial weight distributions (which we call PSEUDO-IID), including the established cases of IID and orthogonal weights, as well as the emerging low-rank and structured sparse settings celebrated for their computational speed-up benefits. We show that fully-connected and convolutional networks initialized with PSEUDO-IID distributions are all effectively equivalent up to their variance. Using our results, one can identify the Edge-of-Chaos for a broader class of neural networks and tune them at criticality in order to enhance their training.

1 INTRODUCTION

Deep neural networks are often studied at random initialization, in the limit of infinite width, where they have been shown to generate intermediate entries which approach Gaussian processes. Seemingly this was first studied for one-layer networks in Neal (2012) when the weight matrices have identically and independently distributed (IID) entries and became a popular model for deep networks following the seminal results for fully-connected networks in Lee et al. (2017) where orthogonal matrices were also studied. Specifically, Lee et al. (2017) gave a framework to compute the Gaussian process behaviour as a function of the nonlinear activation as well as the variance of the network weights and biases. This model was then used to explain the exploding and vanishing gradient phenomenon Schoenholz et al. (2017) and Pennington et al. (2018) amongst other network properties. The Gaussian process limit has since been extended to a broad class of network architectures and scaling limits, see Section 1.1.

In this paper, we extend the simultaneous scaling proof of the Gaussian process in Matthews et al. (2018) to a larger class of initial weight distributions (which we call PSEUDO-IID). The PSEUDO-IID distribution includes structured low-dimensional matrices such as low-rank and structured sparse settings celebrated for their computational efficiency and regularizing properties, as well as the already established cases of IID and orthogonal weights. The PSEUDO-IID distribution is defined in Definition 1 from the exchangeable distribution Definition 2 combined with a specified variance and a bounded high order moment.

Definition 1 (PSEUDO-IID). *Let m, n be two integers. We will say that the random matrix $W = (W_{ij}) \in \mathbb{R}^{m \times n}$ is in the PSEUDO-IID distribution with parameter σ^2 if*

- (i) *the matrix is row-exchangeable and column-exchangeable,*
- (ii) *its entries are centered, uncorrelated, with variance $\mathbb{E}(W_{ij}^2) = \frac{\sigma^2}{n}$,*
- (iii) *$\mathbb{E} \left| \sum_{j=1}^n a_j W_{ij} \right|^8 = K \|\mathbf{a}\|_2^8 n^{-4}$ for some constant K ,*
- (iv) *and $\lim_{n \rightarrow \infty} \frac{n^2}{\sigma^4} \mathbb{E}(W_{i_a, j} W_{i_b, j} W_{i_c, j'} W_{i_d, j'}) = \delta_{i_a, i_b} \delta_{i_c, i_d}$, for all $j \neq j'$.*

When $W^{(1)}$ has IID Gaussian entries and the other weight matrices $W^{(l)}$; $2 \leq l \leq L + 1$, of a neural network (see 1) are drawn from PSEUDO-IID distribution, we will say that the network is under the PSEUDO-IID regime.

Definition 2 (Exchangeability) Let $X_1; \dots; X_n$ be scalar or vector-valued random variables. We say $(X_i)_{i=1}^n$ are exchangeable if their joint distribution is invariant under permutations, i.e. $(X_1; \dots; X_n) \stackrel{d}{=} (X_{(1)}; \dots; X_{(n)})$ for all permutations $\sigma: [n] \rightarrow [n]$. A random matrix is called row- (column-) exchangeable if its rows (columns) are exchangeable random vectors, respectively.

A row-exchangeable and column-exchangeable matrix $M \in \mathbb{R}^{m \times n}$ is not in general entrywise exchangeable, which means its distribution is not typically invariant under arbitrary permutations of its entries; particularly, out of $(mn)!$ possible permutations of its entries, only needs to be invariant under $n!$ of them — an exponentially smaller number. Matrices drawn uniformly from the Grassmanian of orthogonal matrices are a special case of random matrices satisfying row- and column-exchangeability and also entrywise exchangeability but are not independently distributed.

The conditions of PSEUDO-IID are sufficient to prove that fully-connected and convolutional networks are Gaussian processes, Theorems 1 and 2 respectively. Moreover, we are able to verify PSEUDO-IID conditions for common initializations as well as parsimonious (e.g. sparse and/or low-rank) ones (see Section 3.1). We will illustrate these conditions further in Section 3.1 via some examples. However, the sharpest variant of Definition 1 condition (iii) remains an open question; its importance is expanded upon in Appendix D.

1.1 RELATED WORK

To the best of our knowledge, the Gaussian Process behaviour in the infinite width regime was first established by Neal (2012) in the case of one-layer fully-connected networks when the weights are sampled from standard distributions. The result has then been extended in Matthews et al. (2018) to deep fully-connected networks, where the depth is fixed, the weights distributed Gaussians and the hidden layers widths growing jointly to infinity. Jointly scaling the network width substantially distinguishes their method of proof from the approach taken in Lee et al. (2017), where the authors considered a sequential limit analysis through layers. That is, analyzing the limiting distribution at one layer when the previous ones have already converged to their limiting distributions as in Lee et al. (2017), is significantly different from analyzing the limiting distribution at a current layer when the previous layers are jointly converging to their limits at the same time, as is done in Matthews et al. (2018).

Since this Gaussian Process behaviour has been established, two main themes of research have further been developed. The first one consists in the extension of such results for more general and complex architectures such as convolutional networks with many channels [Novak et al. (2020), Garriga-Alonso et al. (2019)] or any modern architectures composed of fully-connected, convolutional or residual connections, as summarized in Yang (2021), using the Tensor Program terminology. The second research theme concerns the generalization of this Gaussian Process behaviour to other possible weight distributions, such as orthogonal weights in Huang et al. (2021) or, alternatively, any IID weights with finite moments as derived in Hanin (2021). Note that the orthogonal case does not fit into the latter as entries are exchangeable but not independent (the first column of an orthogonal matrix cannot be independent from the second one in order to satisfy the orthogonality constraints). The same kind of results for general architectures in the setting have been derived by Yang in his Tensor Program framework Golikov & Yang (2022) and Yang (2021). Our contributions into this line of research, where we relax the independence requirement of the weight matrix entries and consider instead PSEUDO-IID distribution of uncorrelated and exchangeable random variables. This broader distribution, Definition 1, enables us to present a unified proof that generalizes the approaches taken so far and encompasses all of them, for two types of architectures, namely, fully-connected and convolutional networks.

We conclude this section by mentioning that this Gaussian Process behaviour is a special case of a more general result about the convergence of wide neural networks towards a symmetric stochastic process Peluchetti et al. (2020) where independent, but not identically distributed, entries were considered in contrast to the setting of the present paper with identically distributed but not independent entries.

1.2 ORGANIZATION OF THE PAPER

In Section 2, we focus on fully-connected neural networks, formally stating the associated Gaussian Process in Theorem 1, and outlining its proof in Section 2.2 with further technical details of its proof relegated to Appendix A. Our PSEUDO-IID regime unifies the previously studied settings of and orthogonal weights, while also allowing for novel settings such as low-rank Gaussian weights. We extend our results to the convolutional neural networks (CNNs) in Section 2.3. In Section 3, we provide examples of PSEUDO-IID distributions in practice, supporting our theoretical results with numerical simulations. Moreover, we explore the problem of stable initialization of deep networks, on the so-called Edge-of-Chaos (EOC), in our more expansive PSEUDO-IID regime. Lastly, in Section 4, we review our main contributions and put forward some further research directions.

2 GAUSSIAN PROCESS BEHAVIOUR IN THE PSEUDO-IID REGIME

We consider an untrained fully-connected neural network with width n at layer $\ell \in \{1, \dots, L+1\}$. Its weights $W^{(\ell)} \in \mathbb{R}^{n \times (n-1)}$ and biases $b^{(\ell)} \in \mathbb{R}^n$ at layer ℓ are sampled from a centered probability distribution, respectively $\mathcal{W}^{(\ell)}$ and $\mathcal{B}^{(\ell)}$. Most commonly, the weights and biases are sampled i.i.d. Gaussian. Starting with such a network, with nonlinear activation $\sigma: \mathbb{R} \rightarrow \mathbb{R}$, the propagation of any input data vector $x^{(0)} := x \in \mathcal{X} \subset \mathbb{R}^{n_0}$ through the network is given by the following equations,

$$h_i^{(\ell)}(x) = \sum_{j=1}^{n_{\ell-1}} W_{ij}^{(\ell)} z_j^{(\ell-1)}(x) + b_i^{(\ell)}; \quad z_j^{(\ell)}(x) = \sigma(h_j^{(\ell)}(x)); \quad (1)$$

where $h^{(\ell)}(x) \in \mathbb{R}^n$ is referred to as the pre-activation vector at layer ℓ or the feature maps.

Throughout this paper, we will consider a specific set of activation functions that satisfy the so-called linear envelope property, Definition 3, which is satisfied by most activation functions used in practice (ReLU, Softmax, Tanh, HTanh, etc.).

Definition 3. (Linear envelope property) A function: $\mathbb{R} \rightarrow \mathbb{R}$ is said to satisfy the linear envelope property if there exist $c, M \geq 0$ such that, for any $x \in \mathbb{R}$,

$$|\sigma(x)| \leq c + M|x|; \quad (2)$$

2.1 THE PSEUDO-IID REGIME FOR FULLY-CONNECTED NETWORKS

Our proofs of PSEUDO-IID networks converging to Gaussian processes are done in the more sophisticated simultaneous width growth limit as pioneered by Matthews et al. (2018). For a review of the literature on deep networks with sequential vs. simultaneous scaling see Section 1.1. One way of characterizing such a simultaneous convergence over all layers is to consider that all widths n_ℓ are increasing functions of one parameter, let us say n , such that, as n grows, all layers' widths grow: $n_\ell \in \{2, \dots, L+1\}; n_\ell := n_\ell[n]$. We emphasize this dependence only by adding a suffix $x[n]$ to the random variables x when is finite and denote by $X[\cdot]$ its limiting distribution, corresponding to $n \rightarrow \infty$. The width of the first layer n_0 is fixed by the input dimension and the final output layer dimension n_{L+1} do not scale with n . Moreover, the input data are assumed to come from a countably infinite input space \mathcal{X} . Equation 1 can thus be rewritten, for any $x \in \mathcal{X}$, as

$$h_i^{(\ell)}(x)[n] = \sum_{j=1}^{n_{\ell-1}[n]} W_{ij}^{(\ell)} z_j^{(\ell-1)}(x)[n] + b_i^{(\ell)}; \quad z_j^{(\ell)}(x)[n] = \sigma(h_j^{(\ell)}(x)[n]); \quad (3)$$

and the associated Gaussian process limit is given in Theorem 1.

Theorem 1 (GP limit for fully-connected PSEUDO-IID networks) Suppose a fully-connected neural network as in equation 3 is under the PSEUDO-IID regime with parameter σ and the activation satisfies the linear envelope property Def. 3. Let be a countably-infinite set of inputs. Then, for every layer $\ell \in \{1, \dots, L+1\}$, the sequence of random elements $\{h_i^{(\ell)}(x)[n]\}_{n \in \mathbb{N}}$ converges in distribution to a centered Gaussian process $\{h_i^{(\ell)}(x)[\cdot]\}_{n \in \mathbb{N}}$, whose covariance function is given by

$$\mathbb{E} h_i^{(\ell)}(x)[\cdot] h_j^{(\ell)}(x^0)[\cdot] = \delta_{ij} K^{(\ell)}(x; x^0); \quad (4)$$

where

$$K^{(\ell)}(x; x^0) = \begin{pmatrix} \frac{2}{b} + \frac{2}{W} E_{(u,v) \sim N(0; K^{(\ell-1)}(x; x^0))} [(u) (v)]; & \dots & 1 \\ \frac{2}{b} + \frac{2}{n_0} h; & x^0; & \dots & 0 \end{pmatrix}; \quad (5)$$

2.2 SKETCH OF THE PROOF OF THEOREM 1: GP LIMIT OF FULLY-CONNECTED PSEUDO-IID NETWORKS

This section includes the outline of our proof, which closely follows the steps taken in the original proof in Matthews et al. (2018) in the Gaussian setting for fully-connected networks. We refer the reader to Appendix A for its complete proof. Some of the results we use are shown in Matthews et al. (2018), so we intentionally choose our notation similarly to aid readers familiar with the prior paper. These steps are as follows:

1. The first step consists in reducing the problem of showing the convergence of the stochastic process $\{h_i^{(\ell)}(x)[n]\}_{i \in [N], x \in \mathcal{X}}$ to a Gaussian Process, defined on a countably-infinite input space $\mathcal{N} \times \mathcal{X}$ (e.g. $\mathcal{X} = \{x_j\}_{j \in [2N]}$), to the convergence of a finite-dimensional vector $(h_i^{(\ell)}(x)[j])_{(i;x) \in \mathcal{L}}$, where $|\mathcal{L}| < \infty$, to a multidimensional Gaussian. This is possible as the convergence towards the Gaussian Process is ensured with respect to the topology generated by a specific metric introduced in Matthews et al. (2018). We will therefore restrict our attention to showing the convergence in distribution of a finite-dimensional vector $(h_i^{(\ell)}(x)[n])_{(i;x) \in \mathcal{L}}$, where i refers to the neuron index and x to the input data.
2. Given this finite-dimensional random vector, the problem is once again reduced to proving the convergence in distribution of any of its linear projections, which are scalars, to the corresponding linearly projected limiting Gaussians (see Cramér & Wold (1936)). We will thus consider the totality of these one-dimensional linear projections of the unbiased feature map $sh_i^{(\ell)}(x) - b_i^{(\ell)}$ onto $(i;x)$,

$$T^{(\ell)}(; L)[n] := \sum_{(i;x) \in \mathcal{L}} X_{(i;x)} h_i^{(\ell)}(x)[n] - b_i^{(\ell)}; \quad (6)$$

where the suffix $x[n]$ emphasizes on the joint limit to be taken simultaneously, as detailed in Section 2.1. Given the recursion formulae (equation 3) satisfied by the feature maps, the latter can be rewritten as

$$T^{(\ell)}(; L)[n] = \frac{1}{N^{(\ell-1)[n]}} \sum_{j=1}^{N^{(\ell-1)[n]}} T_j^{(\ell)}(; L)[n]; \quad (7)$$

where, considering the renormalization $z_j^{(\ell)} := \frac{1}{W^{(\ell-1)}} \frac{1}{N^{(\ell-1)[n]}} W_{ij}^{(\ell)}$, the summands are defined as,

$$T_j^{(\ell)}(; L)[n] := \sum_{(i;x) \in \mathcal{L}} \frac{1}{W^{(\ell-1)}} z_j^{(\ell-1)}(x)[n]; \quad (8)$$

3. The third step is to apply a version of the Central Limit Theorem (CLT) in the case where the random variables are exchangeable rather than independent. This CLT is often attributed to Blum & Rosenblatt (1956). The reason we need such an extension theorem is that, as opposed to taking the width limit sequentially, the distribution at the previous layer has not reached its limit, so when taking the joint limit we cannot claim the independence of the previous activities anymore and the standard CLT is no longer valid.
4. The last step is to proceed by induction through layers, and verifying at each layer that the moment assumptions of the exchangeable CLT hold in PSEUDO-IID regime.¹

Note that these steps can also be found in Garriga-Alonso et al. (2019), where the proof technique of Matthews et al. (2018) is adapted to CNNs, with weights drawn Gaussian. In the following section, we define the PSEUDO-IID regime for the CNNs and present an analogue of Theorem 1.

¹The base case is carried out by the Gaussian condition on the weights of the first layer, guaranteed by our PSEUDO-IID regime (see Definition 1).

2.3 THE PSEUDO-IID REGIME FOR CNNs

We consider a CNN with C number of channels at layer l , $l = 1, \dots, L+1$ and two-dimensional convolutional filters $U_{ij}^{(l)} \in \mathbb{R}^{k \times k}$ mapping the input channel $l-1$ to the output channel l . The input signal X (also two-dimensional) has C_0 channels and its propagation through the network is given by

$$h_i^{(l)}(X)[n] = \sum_{j=1}^C b_j^{(l)} + \sum_{j=1}^C U_{ij}^{(l)} * z_j^{(l-1)}(X)[n]; \quad z_i^{(l)}(X)[n] = \sigma(h_i^{(l)}(X)[n]); \quad (9)$$

In equation 9, $\sigma(\cdot)$ should be understood as having the same size as the convolution output, non-linearity $\sigma(\cdot)$ is applied entrywise, and we emphasize the simultaneous scaling with the addition of $[n]$ to the feature maps $h_i^{(l)}(X)[n]$. We denote spatial (multi-) indices by boldface Greek letters, etc., that are ordered pairs of integers taking values in the range of the size of the array. For example, if X is an RGB ($C_0 = 3$) image of $H \times D$ pixels, $i = 2$, and $\mathbf{i} = (i_1; i_2)$, then $X_{i_1; i_2}$ returns the Green intensity of the $(i_1; i_2)$ pixel. Moreover, we denote \mathbf{K} to be the patch centered at the pixel covered by the filter, e.g. if $\mathbf{i} = (i_1; i_2)$ and the filter covers $k \times k = (2k_0 + 1) \times (2k_0 + 1)$ pixels, then $\mathbf{J} = \mathbf{K} + \mathbf{i}$, with the usual convention of zero-padding for the out-of-range indices. Sufficient conditions for PSEUDO-IID CNNs to converge to a Gaussian process in the simultaneous scaling limit are given in Definition 4.

Definition 4 (PSEUDO-IID for CNNs). Consider a CNN with random filters and biases $b_j^{(l)}$ and $U_{ij}^{(l)}$ as in equation 9. It is said to be in the PSEUDO-IID regime with parameter $\frac{2}{C_0}$ if $U^{(1)}$ has IID $N(0; \frac{2}{C_0})$ entries and, for $l = 2, \dots, L+1$,

- (i) the convolutional kernel $U^{(l)} \in \mathbb{R}^{C \times C_0 \times k \times k}$ is row-exchangeable and column-exchangeable, that is its distribution is invariant under permutations of first and second indices,
- (ii) filters' entries are centered, uncorrelated, with variance $\mathbb{E}[U_{ij}^{(l)}]^2 = \frac{2}{C_0}$,
- (iii) $\mathbb{E} \sum_{j=1}^C a_j U_{ij}^{(l)} = 0$ for some constant \mathbf{a} ,
- (iv) and $\lim_{n \rightarrow \infty} \frac{C_0^{-1} [n]^2}{4} \mathbb{E} U_{i_a; j_a}^{(l)} U_{i_b; j_b}^{(l)} U_{i_c; j_c}^{(l)} U_{i_d; j_d}^{(l)} = \delta_{i_a; i_b} \delta_{i_c; i_d} \delta_{j_a; j_b} \delta_{j_c; j_d}$ for all $j \in \mathbf{j}^0$.

Theorem 2 (GP limit for CNN PSEUDO-IID networks) Suppose a CNN as in equation 9 is under the PSEUDO-IID regime with parameter $\frac{2}{C_0}$ and the activation satisfies the linear envelope property Def. 3. Let X be a countably-infinite set of inputs and \mathbf{l} denote a spatial (multi-) index. Then, for every layer $l = 2, \dots, L+1$, the sequence of random fields $(h_i^{(l)}(X)[n])_{n \in \mathbb{N}}$ converges in distribution to a centered Gaussian process $(h_i^{(l)}(X)[\cdot])_{\cdot \in \mathbb{N}}$, whose covariance function is given by

$$\mathbb{E} h_i^{(l)}(X)[\cdot] h_j^{(l)}(X^0)[\cdot] = \delta_{ij} \left(\frac{2}{C_0} + \frac{2}{C_0} \sum_{\mathbf{J} \in \mathbf{K}} K^{(l)}(X; X^0) \right); \quad (10)$$

where

$$K^{(l)}(X; X^0) = \begin{cases} \mathbb{E} \left(\sum_{i=1}^{C_0} X_i^{(l-1)}(X; X^0) \right) [u] \left(\sum_{i=1}^{C_0} X_i^{(l-1)}(X; X^0) \right) [v]; & \mathbf{l} = \mathbf{1} \\ 0 & \mathbf{l} = \mathbf{0} \end{cases} \quad (11)$$

Our proof of Theorem 2 for PSEUDO-IID CNNs is derived similarly to that of fully-connected networks, see Appendix B; Garriga-Alonso et al. (2019) developed a more restrictive proof to the setting of IID Gaussian CNNs. Note the extra index on the patch and how the fully-connected case is recovered when the filter kernel size is reduced to 1.

3 PSEUDO-IID IN PRACTICE

Untrained networks are typically initialized with IID weights, for example, the Gaussian and the uniform distributions used as default by TORCH. An important non-IID case that also leads to the Gaussian process limit is the random orthogonal initialization Huang et al. (2021). Our proposed PSEUDO-IID regime encompasses both and orthogonal weights as special cases and also allows for a broader class of weight distributions such as random low-rank orthogonal matrices Saada & Tanner (2023) and structured sparse matrices Dao et al. (2022a) and Dao et al. (2022b).

3.1 EXAMPLES OF PSEUDO-IID DISTRIBUTIONS

Here we give examples of PSEUDO-IID distributions that do not fit in the setting of the prior proofs of Gaussian process reviewed in Section 1.1.

IID weights. If $A = (A_{ij}) \in \mathbb{R}^{m \times n}$ has IID entries $A_{ij} \stackrel{\text{iid}}{\sim} D$, then it is automatically row- and column-exchangeable, and the entries are uncorrelated. Therefore, as long as the distribution of the weights D satisfies the moment conditions of Definition 1, then the network is in the PSEUDO-IID regime. A sufficient condition is that A_{ij} be sub-gaussian with parameter $O(n^{-1/2})$. Then, given the independence of entries, the random variable $\sum_{j=1}^n a_j A_{ij}$ would be sub-gaussian with parameter $\text{kek}_2 n^{-1/2}$, and its p th moment is known to be $O(\text{kek}_2^p n^{-p/2})$; see Vershynin (2018). Appropriately scaled Gaussian and uniform weights, for example, meet the sub-gaussianity criterion and therefore fall in the PSEUDO-IID class. Condition (iv) is trivial in this case.

Orthogonal weights. Let $A = (A_{ij}) \in \mathbb{R}^{n \times n}$ be drawn from the uniform (Haar) measure on the group of orthogonal matrices $\mathcal{O}(n)$. While the entries are not independent, they are uncorrelated, and the rows and columns are exchangeable. To bound the moments, we may employ the concentration of the Lipschitz functions on the sphere, since one individual row $A_{i; \cdot}$ or say the i -th one $(A_{i;1}; \dots; A_{i;n})$, is drawn uniformly from S^{n-1} . Let $f(A_{i;1}; \dots; A_{i;n}) := \sum_{j=1}^n a_j A_{ij}$, then f is Lipschitz with constant kek_2 . By Theorem 5.1.4 of Vershynin (2018), the random variable $f(A_{i;1}; \dots; A_{i;n}) = \sum_{j=1}^n a_j A_{ij}$ is sub-gaussian with parameter $\text{kek}_{\text{Lip}} n^{-1/2} = \text{kek}_2 n^{-1/2}$, which also implies $\mathbb{E} \sum_{j=1}^n a_j A_{ij}^p = O(\text{kek}_2^p n^{-p/2})$. Moreover, the exact expressions of the moments are known for orthogonal matrices (see Collins et al. (2021) for an introduction to the calculus of Weingarten functions) and satisfy condition (iv).

Low-rank weights. Low-rank structures are widely recognized for speeding up matrix multiplications and can be used to reduce memory requirements of feature maps Price & Tanner (2023). Whilst such structures inevitably impose dependencies between the weight matrix entries $A_{ij} \in \mathbb{R}^{m \times n}$, thus breaking the IID assumption, Saada & Tanner (2023) introduced a low-rank framework that falls within our PSEUDO-IID regime. Let $C := [C_1; \dots; C_r] \in \mathbb{R}^{m \times r}$ be a uniformly drawn orthonormal basis for a random-dimensional subspace. Let $P = (P_{ij}) \in \mathbb{R}^{r \times n}$ has IID entries $P_{ij} \stackrel{\text{iid}}{\sim} D$. If we set $A := CP$, it is easy to see that each column $A_{i; \cdot}$ is a linear combination of the columns given by C , with coefficients given by the matrix P . The row and column exchangeability of A follows immediately from that of C and P and the moment conditions are controlled by the choice of distribution D . Direct computation of the four-cross product that appears in condition (iv) gives us

$$\mathbb{E}(A_{i_a;1} A_{i_b;1} A_{i_c;2} A_{i_d;2}) = s^2 \sum_{1 \leq k, k^0 \leq r} \mathbb{E} C_{i_a;k} C_{i_b;k} C_{i_c;k^0} C_{i_d;k^0};$$

wheres $s := \mathbb{E}(P_{1;1}^2)$. Using the expression in (Huang et al., 2021, Lemma 3) we can calculate the above expectation and deduce condition (iv) which is linearly proportional to kek_2^4 .

Permuted block-sparse weights Block-wise pruned networks have recently been under extensive study for their efficient hardware implementation (Dao et al. (2022b), Dao et al. (2022a)). Once the sparsity pattern is fixed, we may apply random row and column permutations on the weight matrices without compromising the accuracy or the computational benefit. Let $A = (A_{ij}) \in \mathbb{R}^{m \times n}$ has IID entries $A_{ij} \stackrel{\text{iid}}{\sim} D$ and $B \in \mathbb{R}^{m \times n}$ be the binary block-sparse mask. Let $A \leftarrow P_m (A \odot B) P_n$, where P_m and P_n are random permutation matrices of size m and n respectively, and \odot represents entrywise multiplication. Then, by construction, A is row- and column-exchangeable and,

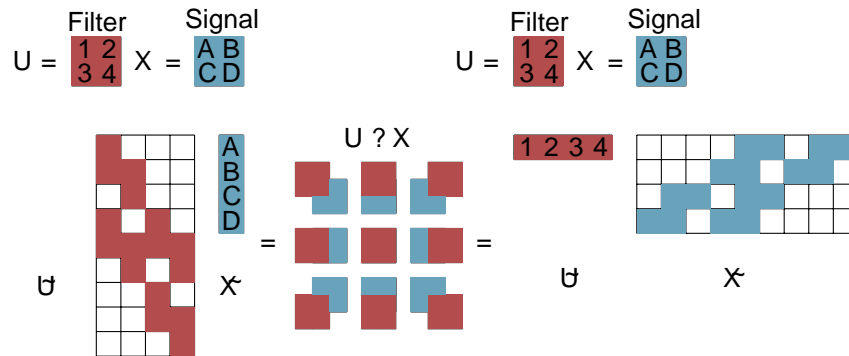


Figure 1: There exist multiple ways to compute convolutions between a tensor \mathcal{U} and a 2-dimensional signal X (shown in the middle) based on matrix multiplications. We illustrate the approach taken in Garriga-Alonso et al. (2019) on the left, where the reshaping procedure is applied to the filter before computing a matrix multiplication, whilst the method we followed, shown in the right hand side of the figure consists in reshaping the signal in order to define special structures on the CNN filters such as orthogonality, sparsity and low-rank.

for suitable choices of underlying distribution \mathcal{D} , it satisfies the moment conditions of Definition 1. Some examples of realisations from such distribution are provided in Appendix E.

Orthogonal CNN filters. Unlike the fully-connected case, it is not obvious how to define the orthogonality of a convolutional layer and, once defined, how to randomly generate such layers for initialization. Xiao et al. (2018) defines an orthogonal convolutional kernel $\mathcal{U} \in \mathbb{R}^{c_{out} \times c_{in} \times k \times k}$ made of c_{out} filters of size k by k via the energy preserving property $\|\mathcal{U} * X\|_2 = \|X\|_2$, for any signal X with c_{in} input channels. Wang et al. (2020) requires the matricised version of the kernel to be orthogonal, while Qi et al. (2020) gives a more stringent definition imposing isometry, i.e.

$$\sum_{i=1}^{c_{out}} U_{ij} * U_{ij}^0 = \begin{cases} 1; & j = j^0 \\ 0; & \text{otherwise} \end{cases}$$

Another definition in Huang et al. (2021) calls for orthogonality of “spatial” slices $\mathcal{U} \in \mathbb{R}^{c_{out} \times c_{in} \times k \times k}$, for all positions \mathbf{s} .

We take a different approach than Wang et al. (2020) for matricising the tensor convolution operator, setting the stride to 1 and padding to 0: reshape the kernel to a matrix $\mathcal{U} \in \mathbb{R}^{c_{out} \times k^2 c_{in}}$ and unfold the signal X into $X' \in \mathbb{R}^{k^2 c_{in} \times d}$, where d is the number of patches depending on the sizes of the signal and the filter. This allows \mathcal{U} to be an arbitrary unstructured matrix rather than the doubly block-Toeplitz matrix in Wang et al. (2020), as shown in Figure 1. Matricising the tensor convolution operator, imposes the structure on the signal rather than the filter \mathcal{U} . Orthogonal (i.e. energy-preserving) kernels can then be drawn uniformly random with orthogonal columns, such that

$$\mathcal{U}^T \mathcal{U} = \frac{1}{k^2} I \tag{12}$$

and then reshaped into the original tensor kernel. Note that this construction is only possible when \mathcal{U} is a tall matrix with trivial null space, that is when $c_{out} \leq k^2 c_{in}$, otherwise the transpose might be considered instead. We emphasize that equation 12 is a sufficient (and not necessary) condition for \mathcal{U} to be energy-preserving, since, by construction, \mathcal{U} has a very specific structure $\mathcal{U} \in \mathbb{R}^{c_{out} \times k^2 c_{in} \times d}$, and, therefore, \mathcal{U} only needs to preserve norm $\|\cdot\|_2$ (and not everywhere). Therefore, we do not claim the generated orthogonal convolutional kernels are “uniformly distributed” over the set of all such kernels.

Now let us verify the conditions of Definition 4. Each $U_{i,j}$ is attained as $U_{i,j} \in \mathbb{R}^{1 \times k^2}$ and forms part of a row of U as shown below:

$$U = \begin{matrix} & \begin{matrix} 1 & 2 & \dots & k^2 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ c_{out} \end{matrix} & \begin{bmatrix} U_{1;1} & U_{1;2} & \dots & U_{1;c_{in}} \\ U_{2;1} & U_{2;2} & \dots & U_{2;c_{in}} \\ \vdots & \vdots & \ddots & \vdots \\ U_{c_{out};1} & U_{c_{out};2} & \dots & U_{c_{out};c_{in}} \end{bmatrix} \end{matrix} \quad (13)$$

Applying permutations on the indices i and j translates to permuting rows and ‘‘column blocks’’ of the orthogonal matrix U , which does not affect the joint distribution of its entries. Hence, the kernel’s distribution is unaffected, that is U is row- and column-exchangeable. The moment conditions are both straightforward to check as $U_{i,j} = U_{i,(j-1)k^2+1} \dots U_{i,(j-1)k^2+k^2}$; where i is the counting number of the pixel. To check condition (iv), note that $U_{i_a;1}^{(1)} U_{i_b;1}^{(1)} U_{i_c;2}^{(1)} U_{i_d;2}^{(1)} = E[U_{i_a;1} U_{i_b;1} U_{i_c;2} U_{i_d;2}]$, that is a four-cross product of the entries of an orthogonal matrix, whose expectation is explicitly known to be $\frac{C+1}{(C-1)C(C+2)} \delta_{i_a i_b} \delta_{i_c i_d} \delta_{a;b} \delta_{c;d}$ (Huang et al., 2021, Lemma 3).

3.2 SIMULATION OF THE GAUSSIAN PROCESSES IN THEOREM 1 FOR FULLY-CONNECTED NETWORKS WITH PSEUDO-IID

Theorem 1 establishes that in the infinite width simultaneous scaling the fully-connected PSEUDO-IID networks converge to Gaussian processes. Here we conduct numerical simulations which validate this for modest dimensions of width $n = n$ for $n = 3, 30, \text{ and } 300$. Fig. 2 shows histograms of entries equation 3 weight matrices with uniform entries, Gaussian with dropout, PSEUDO-IID low-rank matrices, and PSEUDO-IID structured sparse matrices. Even at $n = 30$ there is excellent agreement of the histogram and the variance σ^2 in the infinite width limit. These histograms in Fig. 2 are quantified with Q-Q plots in Appendix F.

Fig. 3 explores the rate with which two independent inputs x_a and x_b generate uncorrelated Gaussian processes for the same neuron. This is done by plotting the joint distribution $h_1^{(1)}(x_a, x_b)[n]$ and $h_1^{(1)}(x_a)[n] h_1^{(1)}(x_b)[n]$ for the same value of n and with the same network. Convergence to the limiting correlation between $h_1^{(1)}(x_a)[n]$ and $h_1^{(1)}(x_b)[n]$ given by Theorem 1 is also shown with the overlaid level curves. These experiments are conducted for weight matrices with uniform entries with dropout as well as PSEUDO-IID orthogonal and PSEUDO-IID Gaussian low-rank and structured sparse matrices. Interestingly the PSEUDO-IID orthogonal converge to the large width distribution most quickly with good agreement at even $n = 3$. The other distributions considered show good agreement at $n = 30$ which improves at $n = 300$. The horizontal and vertical axis in each subplot of Fig. 7 are $h_1^{(5)}(x_a)$ and $h_1^{(5)}(x_b)$ respectively with x_a and x_b drawn independently.

3.3 PROPAGATION OF GAUSSIAN PROCESSES THROUGH DEEP NETWORKS

Poole et al. (2016) and Xiao et al. (2018) developed formulae for the dynamics of the co-variance matrix of a Gaussian process through layers of fully-connected and convolutional networks respectively. These formulae determine fixed points of the covariance matrices along with sensitivity of the networks to small perturbations of inputs. Specifically, they derive the Edge-of-Chaos (EoC) condition for which the network is stable to perturbations. This same EoC condition was subsequently shown in Schoenholz et al. (2017) and Pennington et al. (2018) to avoid the exploding and vanishing gradient phenomenon.

All of the aforementioned results require the network to generate a Gaussian process in the large width limit. Theorems 1 and 2 extend this theory to networks in the PSEUDO-IID regime, allowing for the first initialization conditions for the first time a rigorous EoC analysis of random networks with

²Experiments conducted for Fig. 2-7 used a fully-connected network with activation $\sigma = \tanh(x)$, weight variance $\sigma_w = 2$ and without bias. Dropout used probability $p = 0.2$ of setting an entry to zero, low-rank used rank $r = 2e$, and block-sparsity used randomly permuted block-diagonal matrices with block size $s = 2e$. The code to reproduce all these figures can be found at <https://shorturl.at/gNOQ0>.

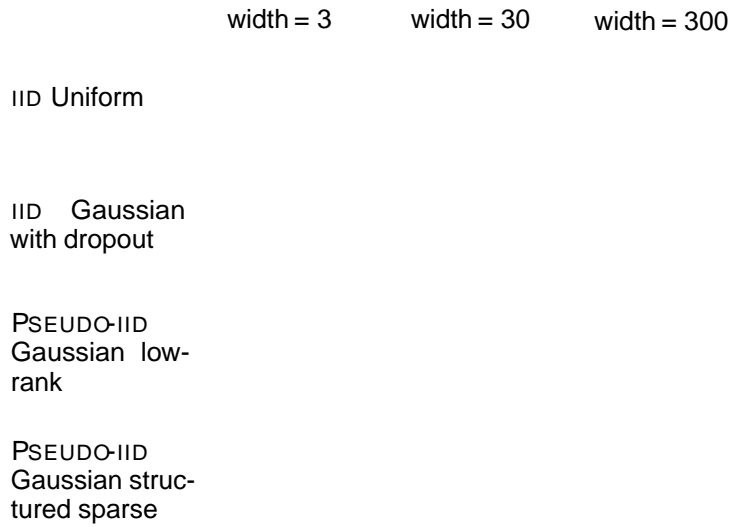


Figure 2: For different instances of the PSEUDO-IID regime, as the width of a fully-connected network grows, the pre-activation given in the i th neuron at the l th layer tend to a Gaussian whose moments are given by Theorem 1. The experiments were conducted 1000 times on a network whose depth is set to be 7 and the input data is sampled from \mathcal{S}^8 .

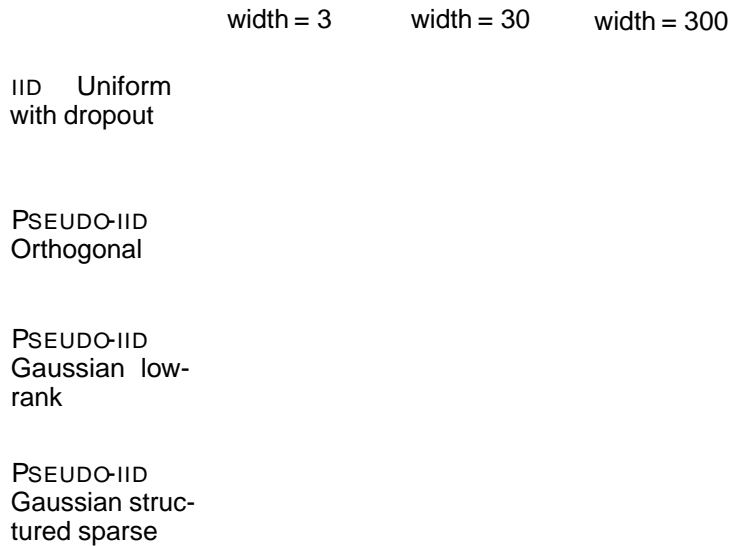


Figure 3: The empirical joint distribution of two pre-activation values at the same neuron resulting from two distinct inputs. Different PSEUDO-IID regimes for fully-connected networks are simulated and their large width limiting distribution from Theorem 1 is included as level curves. The input data are taken from \mathcal{S}^9 and 1000 experiments were conducted on a 7-layer deep network.

parsimonious initialization as developed in Saada & Tanner (2023). More precisely, Theorems 1 and 2 can be expected to form the foundation of developing initialization theory for networks designed for greater efficiency (and accuracy) by using structured sparse and low-rank weight matrices.

4 CONCLUSION

Here we have presented the first Gaussian process theory for deep networks with weight matrices having low-dimensional dependent entries. Theorems 1 and 2 for fully-connected and convolutional networks allow calculation of conditions necessary for initialization of these networks which train efficiently Saada & Tanner (2023). We anticipate these theorems will be extended to additional network architectures and may serve as a roadmap for yet other network regularizers. Moreover, the theory presented here can be further refined to determine finite dimensional corrections, following the approach of Roberts et al. (2022), such as the rate of convergence of these quantities and the variance of these quantities for finite dimensions.

ACKNOWLEDGMENTS

Thiziri Nait Saada is supported by the UK Engineering and Physical Sciences Research Council (EPSRC) through the grant EP/W523781/1. Jared Tanner is supported by the Hong Kong Innovation and Technology Commission (InnoHK Project CIMDA).

REFERENCES

- Patrick Billingsley. Convergence of probability measures; 2nd Wiley series in probability and statistics. Wiley, Hoboken, NJ, 1999. URL <https://cds.cern.ch/record/1254129>.
- Julius R. Blum and Murray Rosenblatt. A class of stationary processes and a central limit theorem. Proceedings of the National Academy of Sciences of the United States of America 41:12–3, 1956.
- Benoit Collins, Sho Matsumoto, and Jonathan Novak. The weingarten calculus, 2021.
- H. Cramér and H. Wold. Some theorems on distribution functions. Journal of the London Mathematical Society 11(4):290–294, 1936. doi: <https://doi.org/10.1112/jlms/s1-11.4.290>. URL <https://londmathsoc.onlinelibrary.wiley.com/doi/abs/10.1112/jlms/s1-11.4.290>.
- Tri Dao, Beidi Chen, Kaizhao Liang, Jiaming Yang, Zhao Song, Atri Rudra, and Christopher Ré. Pixelated butter y: Simple and efficient sparse training for neural network models, 2022a.
- Tri Dao, Beidi Chen, Nimit Sohoni, Arjun Desai, Michael Poli, Jessica Grogan, Alexander Liu, Aniruddh Rao, Atri Rudra, and Christopher Ré. Monarch: Expressive structured matrices for efficient and accurate training, 2022b.
- Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional networks as shallow gaussian processes, 2019.
- Eugene Golikov and Greg Yang. Non-gaussian tensor programs. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds), Advances in Neural Information Processing Systems 2022. URL <https://openreview.net/forum?id=AchUIG2wA8->.
- Boris Hanin. Random neural networks in the infinite width limit as gaussian processes, 2021.
- Wei Huang, Weitao Du, and Richard Yi Da Xu. On the neural tangent kernel of deep networks with orthogonal initialization, 2021.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes, 2017. URL <https://arxiv.org/abs/1711.00165>.
- Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks, 2018. URL <https://arxiv.org/abs/1804.11271>.
- Radford M Neal. Bayesian learning for neural networks, volume 118. Springer Science & Business Media, 2012.
- Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A. Abolaj, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are gaussian processes, 2020.
- Stefano Peluchetti, Stefano Favaro, and Sandra Fortini. Stable behaviour of infinitely wide deep neural networks. In Silvia Chiappa and Roberto Calandra (eds), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics, volume 108 of Proceedings of Machine Learning Research, pp. 1137–1146. PMLR, 26–28 Aug 2020. URL <https://proceedings.mlr.press/v108/peluchetti20b.html>.
- Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. The emergence of spectral universality in deep networks. International Conference on Artificial Intelligence and Statistics, pp. 1924–1932. PMLR, 2018.
- Ben Poole, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos, 2016. URL <https://arxiv.org/abs/1606.05340>.

-
- Ilan Price and Jared Tanner. Improved projection learning for lower dimensional feature maps. In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023.
- Haozhi Qi, Chong You, Xiaolong Wang, Yi Ma, and Jitendra Malik. Deep isometric learning for visual recognition. In International conference on machine learning, pp. 7824–7835. PMLR, 2020.
- Daniel A. Roberts, Sho Yaida, and Boris Hanin. The Principles of Deep Learning Theory. Cambridge University Press, 2022. <https://deeplearningtheory.com>
- Thiziri Nait Saada and Jared Tanner. On the initialisation of wide low-rank feedforward neural networks, 2023. URL <https://arxiv.org/abs/2301.13710>
- Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. In International Conference on Learning Representations, 2017. URL <https://openreview.net/forum?id=H1W1UN9gg>
- Roman Vershynin. High-dimensional probability: An introduction with applications in data science volume 47. Cambridge university press, 2018.
- Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X Yu. Orthogonal convolutional neural networks. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition pp. 11505–11515, 2020.
- 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 Jennifer Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pp. 5393–5402. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/xiao18a.html>
- Greg Yang. Tensor programs i: Wide feedforward or recurrent neural networks of any architecture are gaussian processes, 2021.

A PROOF OF THEOREM 1: GAUSSIAN PROCESS BEHAVIOUR IN FULLY-CONNECTED NETWORKS IN THE PSEUDO-IID REGIME

A.1 STEP 1: REDUCTION OF THE PROBLEM FROM COUNTABLY INFINITE TO FINITE DIMENSIONAL

Firstly, we must clarify in what sense a sequence of stochastic processes $\{X_n\}_{n \in \mathbb{N}}$ converges in distribution to its limit X . For a sequence of real-valued random variables, we can define convergence in distribution $X_n \xrightarrow{d} X$ by the following condition $E f(X_n) \rightarrow E f(X)$ for all continuous functions $f: \mathbb{R} \rightarrow \mathbb{R}$. Similarly, we can define weak convergence for random objects taking values in \mathbb{R}^N (countably-indexed stochastic processes), provided that we equip it with a “good” topology³ that maintains a sufficiently rich class of continuous functions. The metric on the space of real sequences $\mathbb{R}^{\mathbb{N}}$ defined by

$$(h; h^0) := \sum_{i=1}^{\infty} 2^{-i} \min(1; |h_i - h_i^0|) \quad (14)$$

is one example. Thus, we can speak of the weak convergence of $\{h_i[n]\}_{n \in \mathbb{N}}$ in the sense that $E f(h_i[n]) \rightarrow E f(h_i)$ for all $f: \mathbb{R}^N \rightarrow \mathbb{R}$ continuous with respect to the metric

Fortunately, to prove the weak convergence of finite-dimensional distributions, it is sufficient to show the convergence of their finite-dimensional marginals Billingsley (1999).

A.2 STEP 2: REDUCTION OF THE PROBLEM FROM MULTIDIMENSIONAL TO ONE-DIMENSIONAL

Let $L = \{(i_1; x_1); \dots; (i_P; x_P)\}$ be a finite subset of the index set $\mathbb{N} \times \mathbb{X}$. We need to show that the vector $h_{i_p}^{(\cdot)}(x_p)[n] \in \mathbb{R}^P$ converges in distribution to $h_{i_p}^{(\cdot)}(x_p) \in \mathbb{R}^P$. By the Cramér-Wold theorem (Cramér & Wold (1936)), we may equivalently show the weak convergence of an arbitrary linear projection

$$T^{(\cdot)}(\cdot; L)[n] = \sum_{(i;x) \in L} h_i^{(\cdot)}(x)[n] b^{(\cdot)} \quad (15)$$

$$= \sum_{(i;x) \in L} \sum_{j=1}^N X_j^{[n]} W_{ij}^{(\cdot)} Z_j^{(\cdot-1)}(x)[n] \quad (16)$$

$$= \frac{1}{N} \sum_{j=1}^N X_j^{[n]} t_j^{(\cdot)}(\cdot; L)[n]; \quad (17)$$

where

$$t_j^{(\cdot)}(\cdot; L)[n] := \sum_{(i;x) \in L} W_{ij}^{(\cdot)} Z_j^{(\cdot-1)}(x)[n]; \quad (18)$$

and W_{ij} is centered and normalized, i.e. $E W_{ij} = 0$; $E W_{ij}^2 = 1$. We will be using a suitable version of Central Limit Theorem (CLT) to prove the weak convergence of the series in equation 17 to a Gaussian random variable.

A.3 STEP 3: USE OF AN EXCHANGEABLE CENTRAL LIMIT THEOREM

The classical Central Limit Theorem (CLT) establishes that for a sequence of random variables, the properly scaled sample mean converges to a Gaussian random variable in distribution. Here we recall an extension of the Central Limit Theorem introduced in Blum & Rosenblatt (1956), where the independence assumption on the summands is relaxed and replaced by an exchangeability condition. The following statement is an adapted version derived in Matthews et al. (2018), which is more suited to our case.

³In fact, it needs to be a Polish space, i.e. a complete separable metric space.

Theorem 3 (Matthews et al. (2018), Lemma 10) For each positive integer n , let $(X_{n,j}; j \in \mathbb{N})$ be an infinitely exchangeable process with mean zero, finite variance, and finite absolute third moment. Suppose also that the variance has a limit $\lim_{n \rightarrow \infty} \frac{1}{N[n]} = \sigma^2$. Define

$$S_n := \frac{1}{\sqrt{N[n]}} \sum_{j=1}^{N[n]} X_{n,j}^{[n]}; \quad (19)$$

where $N : \mathbb{N} \rightarrow \mathbb{N}$ is a strictly increasing function. If

- (a) $E[X_{n,1} X_{n,2}] = 0$,
- (b) $\lim_{n \rightarrow \infty} E[X_{n,1}^2 X_{n,2}^2] = \sigma^4$,
- (c) $E[j X_{n,1}^3] = o(n) \left(\frac{1}{N[n]} \right)$,

then S_n converges in distribution to $\mathcal{N}(0; \sigma^2)$.

Comparing equation 17 with equation 19, we need to check if the summands $z_j^{(\cdot)}(\cdot; L)[n]$ satisfy the conditions of Theorem 3. We will carefully verify each condition in the following sections.

A.3.1 EXCHANGEABILITY OF THE SUMMANDS

To apply Theorem 3, we must first prove that the random variables $z_j^{(\cdot)}(\cdot; L)[n]$ are exchangeable. Let us expand the expression for $z_j^{(\cdot)}(\cdot; L)[n]$ further and write it in terms of the pre-activations of layer $\ell = 2$, i.e.

$$\begin{aligned} z_j^{(\cdot)}(\cdot; L)[n] &= \sum_{(i;x) \in \mathcal{L}} w_{(i;x)}^{(\cdot)} \sum_{ij}^{(\cdot)} (h_j^{(\ell-1)}(x)[n]) \\ &= \sum_{(i;x) \in \mathcal{L}} w_{(i;x)}^{(\cdot)} \sum_{ij}^{(\cdot)} \sum_{k=1}^{N_{\ell-1}} X_{jk}^{2[n]} W_{jk}^{(\ell-1)} z_k^{(\ell-2)}(x)[n] + b_j^{(\ell-1)}; \end{aligned}$$

If we apply a random permutation on the indices let us say $\sigma : \{1; \dots; N_{\ell-1}\} \rightarrow \{1; \dots; N_{\ell-1}\}$, then the row-exchangeability and column-exchangeability of PSEUDO-IID weights ensure that $W_{i(j)}^{(\cdot)}$ has same distribution as $W_{ij}^{(\cdot)}$ and that $W_{(j)k}^{(\ell-1)}$ has same distribution as $W_{jk}^{(\ell-1)}$, for any $i \in \{1; \dots; N_{\ell-1}\}; k \in \{1; \dots; N_{\ell-2}\}$. Note that this extends easily to the normalized versions of the weights $\frac{(\cdot)}{ij}$. Additionally, as the biases are set to be Gaussians, they are a fortiori exchangeable and their distributions remain unchanged when considering random permutations of indices. Therefore, $z_j^{(\cdot)}(\cdot; L)[n]$ is equal to $z_{\sigma(j)}^{(\cdot)}(\cdot; L)[n]$ in distribution, hence exchangeability.

A.3.2 MOMENT CONDITIONS

We mean by moment conditions the existence of a limiting variance as well as the conditions (a)-(c) in Theorem 3. We will prove the moment conditions by induction on the layer number

Existence of the limiting variance of the summands. To show the existence of the limiting variance of the summands, let us first write down such a variance at a finite width. Since they are exchangeable (see A.3.1), their distribution is identical and we may simply calculate the variance of

(\cdot) as follows.

$$\begin{aligned}
(\cdot)^2(\cdot; L)[n] &:= E \left(\cdot(\cdot; L)[n] \right)^2 = E \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;x_a)} \sum_{(i_b;x_b)} h^{(i;x)} z_1^{(\cdot)}(x)[n] z_1^{(\cdot)}(x)[n] \\
&= \sum_{\substack{(i_a;x_a)2L \\ (i_b;x_b)2L}} \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;i_b)} h^{(i;x)} E z_1^{(\cdot)}(x_a)[n] z_1^{(\cdot)}(x_b)[n] \\
&= \sum_{\substack{(i_a;x_a)2L \\ (i_b;x_b)2L}} \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;i_b)} h^{(i;x)} E z_1^{(\cdot)}(x_a)[n] z_1^{(\cdot)}(x_b)[n]; \quad (20)
\end{aligned}$$

where we first considered the independence between the normalized weights at layer i and the activations at layer $i-1$; then used the fact that the normalized weights are uncorrelated in the PSEUDO-IID regime.

The convergence of the second moment of the summands is thus dictated by the convergence of the covariance of the activations of the last layer. By the induction hypothesis, the feature maps $h_j^{(\cdot)}(x)[n]$ converge in distribution, so the continuous mapping theorem guarantees the existence of a limiting distribution for the activations $z_1^{(\cdot)}(x_a)[n]$ and $z_1^{(\cdot)}(x_b)[n]$ as n tends to infinity. Note that this result holds even if the activation function has a set of discontinuity points of Lebesgue measure zero, e.g. the step function. Thus, the product inside the expectation in equation 20 converges in distribution to a limiting random variable $z_1^{(\cdot)}(x_a)[\cdot] z_1^{(\cdot)}(x_b)[\cdot]$.

From Billingsley (1999), one knows that if a sequence weakly converges to a limiting distribution and is uniformly integrable, then we can swap the order of taking the limit and the expectation. Thus, by Proposition 1, we have

$$\begin{aligned}
(\cdot)^2(\cdot; L)[\cdot] &:= \lim_{n \rightarrow \infty} (\cdot)^2(\cdot; L)[n] \\
&= \sum_{\substack{(i_a;x_a)2L \\ (i_b;x_b)2L}} \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;i_b)} h^{(i;x)} E z_1^{(\cdot)}(x_a)[\cdot] z_1^{(\cdot)}(x_b)[\cdot]; \quad (21)
\end{aligned}$$

Condition (a). At a given layer i , we need to show that $X_{n;1} := \sum_{(i;x)} w^{(i;x)} z_1^{(\cdot)}(x)[n]$ and $X_{n;2} := \sum_{(i;x)} w^{(i;x)} z_2^{(\cdot)}(x)[n]$ are uncorrelated. We have

$$\begin{aligned}
E X_{n;1} X_{n;2} &= E \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;x_a)} \sum_{(i_b;x_b)} h^{(i;x)} z_1^{(\cdot)}(x)[n] \sum_{(i;x)} w^{(i;x)} \sum_{(i_c;x_c)} \sum_{(i_d;x_d)} h^{(i;x)} z_2^{(\cdot)}(x)[n] \\
&= \sum_{\substack{(i_a;x_a)2L \\ (i_b;x_b)2L}} \sum_{(i;x)} w^{(i;x)} \sum_{(i_a;i_b)} h^{(i;x)} E z_1^{(\cdot)}(x_a)[n] z_2^{(\cdot)}(x_b)[n] \\
&= 0;
\end{aligned}$$

since $\sum_{(i_a;i_b)} h^{(i;x)} E z_1^{(\cdot)}(x_a)[n] z_2^{(\cdot)}(x_b)[n]$ are uncorrelated by the PSEUDO-IID assumption.

Condition (b).

$$\begin{aligned}
 E X_{n;1}^2 X_{n;2}^2 &= E \int_{(i;x)}^h \int_{(i;x)}^i X_{i;1}^{(c)} z_1^{(c-1)}(x)[n]^2 \int_{(i;x)}^h \int_{(i;x)}^i X_{i;2}^{(c)} z_2^{(c-1)}(x)[n]^2 \\
 &= \frac{4}{W} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f a;b;c;d g}^h \int_{t2f a;b;c;d g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[n]^2 \int_{t2f a;b;c;d g}^h \int_{t2f a;b;c;d g}^i Y_{(i_t;x_t)}^{(c)} z_2^{(c-1)}(x_t)[n]^2 \\
 &= \frac{4}{W} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f a;b;c;d g}^h \int_{t2f a;b;c;d g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[n]^2 F_n(i_a; i_b; i_c; i_d) \int_{t2f c;dg}^h \int_{t2f c;dg}^i Y_{(i_t;x_t)}^{(c)} z_2^{(c-1)}(x_t)[n]^2 ; \quad (22)
 \end{aligned}$$

where

$$F_n(i_a; i_b; i_c; i_d) := E_{i_a;1 i_b;1 i_c;2 i_d;2} \quad (23)$$

We justify the convergence in distribution of the random variable inside the expectation in the exact same way as in the previous section, referring to the continuous mapping theorem and the induction hypothesis. By Proposition 1, as $n \rightarrow \infty$, the above expectation converges to the expectation of the limiting Gaussian process, i.e.

$$\begin{aligned}
 \lim_{n \rightarrow \infty} E \int_{t2f a;b g}^h \int_{t2f a;b g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[n]^2 \int_{t2f c;dg}^h \int_{t2f c;dg}^i Y_{(i_t;x_t)}^{(c)} z_2^{(c-1)}(x_t)[n]^2 \\
 = E \int_{t2f a;b g}^h \int_{t2f a;b g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 \int_{t2f c;dg}^h \int_{t2f c;dg}^i Y_{(i_t;x_t)}^{(c)} z_2^{(c-1)}(x_t)[\]^2 ;
 \end{aligned}$$

Moreover, condition (iv) of Definition 1 implies that

$$\lim_{n \rightarrow \infty} F_n(i_a; i_b; i_c; i_d) = E_{i_a; i_b; i_c; i_d}$$

Substituting the two limits back in the equation 22 and using the independence of the activations at layer $\ell - 1$ given by the induction hypothesis, we get

$$\begin{aligned}
 \lim_{n \rightarrow \infty} E X_{n;1}^2 X_{n;2}^2 &= \frac{4}{W} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f a;b;c;d g}^h \int_{t2f a;b;c;d g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 E_{i_a; i_b; i_c; i_d} \int_{t2f c;dg}^h \int_{t2f c;dg}^i Y_{(i_t;x_t)}^{(c)} z_2^{(c-1)}(x_t)[\]^2 \\
 &= \frac{4}{W} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f a;b g}^h \int_{t2f a;b g}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 E_{i_c; i_d} \int_{t2f a;b g}^h \int_{t2f a;b g}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 \\
 &= \frac{4}{W} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 E_{i_c; i_d} \int_{t2f c;dg}^h \int_{t2f c;dg}^i Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 \\
 &= \frac{4}{W} (\int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2)^2 ;
 \end{aligned}$$

Condition (c). To show that the third absolute moment of the $\int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2$ grows slower than $\sqrt{N} \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2$, it is sufficient to bound it by a constant. Applying Hölder's inequality on $X = \int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2$; $Y = 1$; $p = 4$; $q = 4$, we obtain

$$E \left(\int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 \right)^3 \leq E \left(\int_{(i_t;x_t)}^{(i_t;x_t)2L} \int_{t2f c;dg}^h \int_{t2f c;dg}^i X_{(i_t;x_t)}^{(c)} Y_{(i_t;x_t)}^{(c)} z_1^{(c-1)}(x_t)[\]^2 \right)^4 \frac{1}{4} = 1;$$

Therefore, condition (c) boils down to showing that $E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4$ is finite, for which we will once again make use of the uniform integrability of the feature maps derived in C. De ne

$$G_n(i_a; i_b; i_c; i_d) := E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 \quad (24)$$

and observe that

$$\begin{aligned} E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 &= E \left(\prod_{(i;x) \in \mathcal{X}} \prod_{t \in \mathcal{T}} z_1^{(\cdot-1)}(x_t) \right) [n]^4 \\ &= E \left(\prod_{(i_t; x_t) \in \mathcal{X}^{(i;x) \in \mathcal{X}} \mathcal{T}^{(i_t; x_t) \in \mathcal{T}}} \prod_{t \in \mathcal{T}} z_1^{(\cdot-1)}(x_t) \right) [n]^4 \end{aligned}$$

Using Cauchy-Schwarz inequality, we can bound G_n by the fourth moment of the normalized weights, i.e.

$$\begin{aligned} G_n &\leq \sqrt{\text{Var} \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4} \\ &= E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 \end{aligned}$$

Then, we use condition (iii) of the Definition 1 with $p = 4$ and $\alpha = (1; 0; \dots; 0)^T$ to bound the fourth moment:

$$\begin{aligned} E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 &= \frac{n^2}{4} E W_{i_a;1}^4 \\ &= \frac{n^2}{4} K_4 k^4 n^{-2} = \frac{K_4}{4} = o_n(1) \end{aligned}$$

Furthermore, the induction hypothesis gives the convergence in distribution of the feature maps from the last layer, and combined with the continuous mapping theorem we get the convergence in distribution of the above product inside expectation. Using Lemma 1, the uniform integrability of the activations follows, and Billingsley's theorem (Lemma 2) enables us to swap the limit and the expectation. Thus,

$$\lim_{n \rightarrow \infty} E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 = E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) []^4$$

To bound the product of four different random variables, it is sufficient to bound the fourth order moment of each (see Lemma 3). We can do so using the linear envelope property (Definition 3) satisfied by the activation function to get, for a fixed $(i;x) \in \mathcal{X}$,

$$E \left(z_1^{(\cdot-1)}(x) \right) []^4 \leq 2^{4-1} E c^4 + M^4 E h_1^{(\cdot-1)}(x) []^4$$

The induction hypothesis indicates that $z_1^{(\cdot-1)}(x) []$ follows a Gaussian distribution, whose fourth moment is bounded. Using the fact that we chose the σ to be finite, we can take the supremum over all x . Therefore, we have

$$E \left(z_1^{(\cdot-1)}(x) \right) []^4 \leq 2^{4-1} \sup_{(i;x) \in \mathcal{X}} E c^4 + M^4 E h_1^{(\cdot-1)}(x) []^4 = o_n(1)$$

Combining the above bounds, we then have

$$\lim_{n \rightarrow \infty} E \left(\prod_{i=1}^4 z_1^{(\cdot-1)}(x) \right) [n]^4 < 1$$

A.3.3 CONCLUSION FROM THE EXCHANGEABLE CLT

In the above sections, we showed by induction that at any depth, if the feature maps from previous layers converge in distribution to Gaussian processes, then the assumptions of Theorem 3 hold and the one-dimensional projection of the feature maps at the current layer also converges in distribution to a Gaussian random variable with a specified variance. More precisely, we showed that for any finite set L and projection vector, any linear one-dimensional projection of the feature maps at the current layer, $T^{(c)}(\cdot; L)[n]$, converges in distribution to a Gaussian $N(0; \text{Cov}^{(c)}(\cdot; L))$ as n grows. This gives the convergence of the feature maps at layer c to Gaussian processes.

Considering the unbiased quantity

$$T^{(c)}(\cdot; L)[n] := \frac{1}{n} \sum_{(i; X) \in \mathcal{X}^{2L}} h_i^{(c)}(X)[n] - b_i^{(c)};$$

we can compute its variance:

$$E T^{(c)}(\cdot; L)[n]^2 = \frac{1}{n^2} \sum_{\substack{(i_a; X_a) \in \mathcal{X}^{2L} \\ (i_b; X_b) \in \mathcal{X}^{2L}}} E h_{i_a}^{(c)}(X_a)[n] h_{i_b}^{(c)}(X_b)[n] - \frac{1}{n} \sum_{i_a, i_b} b_{i_a}^{(c)} b_{i_b}^{(c)};$$

As we saw in the previous sections,

$$h_i^{(c)}(\cdot; L)[n] = \frac{1}{W} \sum_{\substack{(i_a; X_a) \in \mathcal{X}^{2L} \\ (i_b; X_b) \in \mathcal{X}^{2L}}} h_{i_a, i_b}^{(c-1)}(X_a, X_b)[n] z_1^{(c-1)}(X_b)[n];$$

Thus, by identification, and using the inductive hypothesis, one recovers the recursion formula for the variance as described in Theorem 1: For any $i_a, i_b \in \mathcal{N}$, $X_a, X_b \in \mathcal{X}$,

$$\begin{aligned} E h_{i_a}^{(c)}(X_a)[n] h_{i_b}^{(c)}(X_b)[n] &= \frac{1}{W^2} \sum_{i_a, i_b} E z_1^{(c-1)}(X_a)[n] z_1^{(c-1)}(X_b)[n] + \frac{1}{W} \\ &= \frac{1}{W} \sum_{(u; v) \in \mathcal{N} \times \mathcal{N}} E_{(u; v) \sim N(0; K^{(c-1)}(X; X^0))} (u)(v) + \frac{1}{W}; \end{aligned}$$

A.4 IDENTICAL DISTRIBUTION AND INDEPENDENCE OVER NEURONS

As we saw, for any n , the feature maps $h_j^{(c)}(\cdot; L)[n]$ are exchangeable, and, in particular, identically distributed. This still holds after taking the limit, that is $h_j^{(c)}(\cdot; L)$ and $h_k^{(c)}(\cdot; L)$ have the same distribution for any $j, k \in \mathcal{N}$.

It still remains to show the independence between $h_i^{(c)}(\cdot; L)$ and $h_j^{(c)}(\cdot; L)$ for $i \neq j$. As we now know the limiting distribution is Gaussian, it suffices to analyze their covariance to conclude about their independence. As derived in the previous section, for any $X \in \mathcal{X}$, $i \neq j$, $E h_i^{(c)}(X)[n] h_j^{(c)}(X^0)[n] = 0$, hence the independence.

B PROOF OF THEOREM 2: GAUSSIAN PROCESS BEHAVIOUR IN CONVOLUTIONAL NEURAL NETWORKS IN THE PSEUDO-IID REGIME

We apply the same machinery to show the Gaussian Process behaviour in CNNs under the pseudo-IID regime, closely following the steps detailed in the fully-connected case. To reduce the problem to a simpler one, one can proceed as previously by considering a finite subset of the feature maps at layer c , $L = \{f(i_1; X_{i_1}; \cdot); \dots; f(i_P; X_{i_P}; \cdot)\} \subset \mathcal{C} \times \mathcal{X}^{2L}$, where \mathcal{C} consists of all the spatial multi-indices. We will follow the same strategy outlined in Sections A.1-A.3. Given a finite set L and the projection vector $\mathbf{w} \in \mathbb{R}^{|\mathcal{C}|}$, we may form the unbiased one-dimensional projection as

$$T^{(c)}(\cdot; L)[n] := \frac{1}{n} \sum_{(i; X) \in \mathcal{X}^{2L}} h_i^{(c)}(X)[n] - b_i^{(c)}; \quad (25)$$

which can be rewritten, using equation 9, as the sum

$$T^{(j)}(\cdot; L)[n] = \frac{1}{C^{j-1}[n]} \sum_{j=1}^{C^{j-1}[n]} T^{(j)}(\cdot; L)[n]; \quad (26)$$

where the summands are

$$T^{(j)}(\cdot; L)[n] := \frac{1}{w} \sum_{(i; X; \cdot) \in \mathcal{L}} \sum_{2^J K} E_{ij}^{(j)}(z_{j; \cdot}^{(j-1)}(X)[n]); \quad (27)$$

As before, we introduced the renormalized version of the filter $U^{(j)}$ such that $E_{ij}^{(j)} := \frac{1}{w} \sum_{(i; X; \cdot) \in \mathcal{L}} U_{ij}^{(j)}$.

We will proceed once again by induction forward through the network's layer verifying the assumptions of the exchangeable CLT (Theorem 3) and using equation 27 to conclude its convergence in distribution to a Gaussian random variable.

The exchangeability of the summands. Similar to the fully-connected case, we employ the row- and column-exchangeability of the PSEUDO-IID convolutional kernel to show the exchangeability of $T^{(j)}$'s. Let us expand equation 27 and write

$$\begin{aligned} T^{(j)}(\cdot; L)[n] &= \frac{1}{w} \sum_{(i; X; \cdot) \in \mathcal{L}} \sum_{2^J K} E_{ij}^{(j)}(h_{j; \cdot}^{(j-1)}(X)[n]) \\ &= \frac{1}{w} \sum_{(i; X; \cdot) \in \mathcal{L}} \sum_{2^J K} E_{ij}^{(j)} \sum_{k=1}^{C^{j-1}[n]} \sum_{2^J K} U_{j;k}^{(j-1)}(z_{k; \cdot}^{(j-2)}(X)[n]); \end{aligned}$$

The joint distributions of $U_{ij}^{(j)}$ $\sum_{j=1}^{C^{j-1}}$ and $U_{j;k}^{(j-1)}$ $\sum_{j=1}^{C^{j-1}}$ are invariant under any permutation π of (j) , therefore $T^{(j)}$'s are exchangeable.

Moment conditions. Condition (a) is straightforward as the filters' entries are uncorrelated by the PSEUDO-IID assumption $E[U_{ij}^{(j)} U_{i'j'}^{(j)}] = \frac{2}{C^{j-1}} \delta_{ij} \delta_{i'j'} = \delta_{ij} \delta_{i'j'}$.

The moment conditions are shown to be satisfied by induction through the network and the proofs boil down to showing the uniform integrability of the activation vectors to be able to swap limit and expectation. This uniform integrability in the CNN case under PSEUDO-IID weights is rigorously demonstrated in proposition 2. One can compute the variance of one representative of the summands, let us say the first one, as follows:

$$\begin{aligned}
(\cdot)^2(\cdot; L)[n] &:= E \left[(\cdot)^2(\cdot; L)[n] \right]^2 = E \int_W \int_{(i; X; \cdot)} \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} E_{i; 1; a}^{(\cdot)} E_{i; 1; b}^{(\cdot)} (X_a)[n] z_{i; b}^{(\cdot)} (X_b)[n] \, d\mu \\
&= \int_W \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} E_{i; 1; a}^{(\cdot)} E_{i; 1; a}^{(\cdot)} E_{i; 1; b}^{(\cdot)} E_{i; 1; b}^{(\cdot)} (X_a)[n] z_{i; a}^{(\cdot)} (X_a)[n] z_{i; b}^{(\cdot)} (X_b)[n] \, d\mu \\
&= \int_W \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} E_{i; 1; a}^{(\cdot)} E_{i; 1; a}^{(\cdot)} E_{i; 1; b}^{(\cdot)} E_{i; 1; b}^{(\cdot)} (X_a)[n] z_{i; a}^{(\cdot)} (X_a)[n] z_{i; b}^{(\cdot)} (X_b)[n] \, d\mu \\
&= \int_W \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} E_{i; 1; a}^{(\cdot)} E_{i; 1; a}^{(\cdot)} E_{i; 1; b}^{(\cdot)} E_{i; 1; b}^{(\cdot)} (X_a)[n] z_{i; a}^{(\cdot)} (X_a)[n] z_{i; b}^{(\cdot)} (X_b)[n] \, d\mu
\end{aligned}$$

Given the uniform integrability of the product the activations in a CNN, we may swap the order of taking the limit and the expectation to have

$$\begin{aligned}
(\cdot)^2(\cdot; L)[\cdot] &:= \lim_{n \rightarrow \infty} (\cdot)^2(\cdot; L)[n] \\
&= \int_W \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2L} \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} E_{i; 1; a}^{(\cdot)} E_{i; 1; a}^{(\cdot)} E_{i; 1; b}^{(\cdot)} E_{i; 1; b}^{(\cdot)} (X_a)[\cdot] z_{i; a}^{(\cdot)} (X_a)[\cdot] z_{i; b}^{(\cdot)} (X_b)[\cdot] \, d\mu
\end{aligned}$$

As in the fully-connected case, we conclude that the feature maps at the next layer converge to Gaussians processes, whose covariance function is given by

$$E \left[h_{i; 1; a}^{(\cdot)} (X_a)[\cdot] h_{i; 1; b}^{(\cdot)} (X_b)[\cdot] \right] = \delta_{i; a; b} + \int_W \int_{(i; X; \cdot)}^{2J} \int_{(i; X; \cdot)}^{2K} K^{(\cdot)}(X_a; X_b) \, d\mu$$

where

$$K^{(\cdot)}(X_a; X_b) = \begin{cases} E_{(u; v)} \left(\prod_{j=1}^N C_0^{(0; K^{(\cdot)}(X_a; X_b))} (X_a)_j (X_b)_j \right) & \text{if } \delta_{i; a; b} = 1 \\ \frac{1}{C_0} & \text{if } \delta_{i; a; b} = 0 \end{cases}$$

C LEMMAS USED IN THE PROOF OF THEOREMS 1 AND 2

We will present in this section the lemmas used in the derivation of our results. The proofs are omitted in cases where they can easily be found in the literature.

Lemma 1 (Hölder's inequality) For a probability space $(\Omega; \mathcal{F}; P)$, let $X; Y$ be two random variables on Ω and $p; q > 1$ such that $\frac{1}{p} + \frac{1}{q} = 1$. Then,

$$E |XY| \leq E |X|^p \frac{1}{p} + E |Y|^q \frac{1}{q}$$

Lemma 2 (Bellingsley's theorem) Let $X[n]$ be a sequence of random variables, $n \in \mathbb{N}$ and X another random variable. $X[n]$ is uniformly integrable and the sequence $X[n]$ converges in distribution to X , then X is integrable and one can swap the limit and the expectation, i.e.

$$\lim_{n \rightarrow \infty} E X[n] = E X$$

We adapt Lemma 18 from Matthews et al. (2018) to our setting and the proof can trivially be obtained from the proof derived in the cited article.

Lemma 3 (Sufficient condition to uniformly bound the expectation of a four-cross product) Let $X_1, X_2, X_3,$ and X_4 be random variables on \mathbb{R} with the usual Borel σ -algebra. Assume that $E |X_i|^{4+\epsilon} < \infty$ for all $i \in \{1, 2, 3, 4\}$. Then, for any choice of $\alpha_i \in [0, 1]$ (where $\sum \alpha_i = 1$), the expectations $E \prod_{i=1}^4 |X_i|^{\alpha_i}$ are uniformly bounded by a polynomial in the eighth moments $E |X_i|^8 < \infty$.

This Lemma can be easily derived from standard Pearson correlation bounds.

We now have the tools to show that all four-cross products of the activations are uniformly integrable in the PSEUDO-IID setting, which enables us to swap the limit and the expectation in the previous sections of the appendix.

Proposition 1 (Uniform integrability in the PSEUDO-IID regime – fully-connected networks) Consider a fully-connected neural network in the PSEUDO-IID regime. Consider the random activations $z_i^{(\cdot)}(x_a)[n]; z_j^{(\cdot)}(x_b)[n]; z_k^{(\cdot)}(x_c)[n]; z_l^{(\cdot)}(x_d)[n]$ with any $i, j, k, l \in \{1, 2, 3, 4\}$; $x_a, x_b, x_c, x_d \in \mathbb{R}^N$, neither necessarily distinct, as in 3. Then, the family of random variables

$$z_i^{(\cdot)}(x_a)[n] z_j^{(\cdot)}(x_b)[n] z_k^{(\cdot)}(x_c)[n] z_l^{(\cdot)}(x_d)[n];$$

indexed by n is uniformly integrable for any $\gamma = 1; \dots; L + 1$.

Proof. The proof is adapted from Matthews et al. (2018) to the PSEUDO-IID regime in fully-connected neural networks described in equation 3.

If a collection of random variables is uniformly p -bounded for $p > 1$, then it is uniformly integrable. So, we will show that our family of random variables is uniformly p -bounded for some $p > 0$, i.e. there exist $K < \infty$ independent of n such that,

$$E |z_i^{(\cdot)}(x_a)[n] z_j^{(\cdot)}(x_b)[n] z_k^{(\cdot)}(x_c)[n] z_l^{(\cdot)}(x_d)[n]|^{1+p} \leq K;$$

which is equivalent to

$$E |z_i^{(\cdot)}(x_a)[n]|^{1+p} |z_j^{(\cdot)}(x_b)[n]|^{1+p} |z_k^{(\cdot)}(x_c)[n]|^{1+p} |z_l^{(\cdot)}(x_d)[n]|^{1+p} \leq K;$$

To do so, Lemma 3 gives us a sufficient condition: bounding the moment of order $4+\epsilon$ of each term in the product by a constant independent of n . For any $x_t \in \mathbb{R}^N$ and $i \in \{1, 2, 3, 4\}$, this moment can be rewritten in terms of the feature maps using the linear envelope property 3 and the convexity of the map $x \mapsto |x|^{4+\epsilon}$ as

$$E |z_i^{(\cdot)}(x_t)[n]|^{4(1+\epsilon)} \leq 2^{4(1+\epsilon)} E |c^{(1+\epsilon)} + M^{(1+\epsilon)} h_i^{(\cdot)}(x_t)[n]|^{4(1+\epsilon)};$$

Thus it is sufficient to show that the absolute feature maps $|h_i^{(\cdot)}(x_t)[n]|$ have a finite moment of order $4+\epsilon$, independent of $x_t, i,$ and n , for some ϵ . For the sake of simplicity, we will show this by induction on n for $\epsilon = 1$.

Base case. For $n = 1$, the feature maps $h_i^{(1)} = \sum_{j=1}^{N_0} W_{ij}^{(1)} x_j + b_i^{(1)}$ are identically distributed for all $i \in \{1, \dots, N_1[n]\}$ from the row-exchangeability of the weights. Moreover, from the moment condition of Definition 1, there exists $p = 8$ such that

$$\begin{aligned} E |h_i^{(1)}(x)[n]|^p &= E \left| \sum_{j=1}^{N_0} x_j W_{ij}^{(1)} + b_i^{(1)} \right|^p \\ &= K \sum_{j=1}^{N_0} |x_j|^{p/2} + E |b_i^{(1)}|^p. \end{aligned}$$

The RHS of the above equality is independent of n (and n). Therefore, for all $n \in \mathbb{N}$ and $i \in \{1, \dots, N_1[n]\}$, we have found a constant bound for the feature map's moment of order

Inductive step. Let us assume that for any $x_t \in \mathcal{X}_t^4$ and $i \geq 2$, the eighth-order moment of $h_i^{(i-1)}(x_t)[n]$ is bounded by a constant independent from

We will show that this implies

$$\mathbb{E} h_i^{(i)}(x_t)[n]^8 < 1 :$$

Considering the vector of activations $h^{(i-1)}(x_t)[n]$, we have, from the moment condition (iii) of the PSEUDO-IID regime the following conditional expectation,

$$\mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} W_{ij}^{(i-1)} h_j^{(i-1)}(x_t)[n]^8 = \sum_{j=1}^{N_X^{i-1}[n]} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 = \sum_{j=1}^{N_X^{i-1}[n]} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 \cdot \frac{1}{N_X^{i-1}[n]} :$$

Thus using the recursion formulae for the data propagation in such architecture given by 3, the convexity of $x \mapsto x^8$ on \mathbb{R}^+ and taking conditional expectations,

$$\mathbb{E} h_i^{(i)}(x_t)[n]^8 \leq 2^{8-i} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 + \sum_{j=1}^{N_X^{i-1}[n]} W_{ij}^{(i-1)} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 \quad (28)$$

$$\begin{aligned} &= 2^{8-i} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 + \mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} W_{ij}^{(i-1)} h_j^{(i-1)}(x_t)[n]^8 \\ &= 2^{8-i} \mathbb{E} h_j^{(i-1)}(x_t)[n]^8 + \sum_{j=1}^{N_X^{i-1}[n]} \mathbb{E} W_{ij}^{(i-1)} h_j^{(i-1)}(x_t)[n]^8 : \end{aligned} \quad (29)$$

As the biases are Gaussian, the first expectation involving the bias is trivially finite and does not depend on n as they are identically distributed. For the second expectation, we compute,

$$\begin{aligned} \mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} W_{ij}^{(i-1)} h_j^{(i-1)}(x_t)[n]^8 &= \mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} h_j^{(i-1)}(x_t)[n]^{2i} \\ &= \mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} (c + M h_j^{(i-1)}(x_t)[n])^{2i} \\ &= \mathbb{E} \sum_{j=1}^{N_X^{i-1}[n]} (c^2 + M^2 h_j^{(i-1)}(x_t)[n]^2 + 2cM h_j^{(i-1)}(x_t)[n])^{i} : \end{aligned} \quad (30)$$

This last expression can be written as a linear combination of n^{i-4} quantities of the form

$$\mathbb{E} h_{j_1}^{(i-1)}(x_t)[n]^{p_1} h_{j_2}^{(i-1)}(x_t)[n]^{p_2} h_{j_3}^{(i-1)}(x_t)[n]^{p_3} h_{j_4}^{(i-1)}(x_t)[n]^{p_4} ;$$

and we bound each of them making use of Lemma 3. The factor n^{i-4} in 29 thus cancels out with the number of terms in the sum 30. As the pre-activations are exchangeable, they are identically distributed so the dependence on the neuron index can be ignored, and taking the supremum over the finite set of input data x_t does not affect the uniformity of the bound; which concludes the proof. □

Proposition 2 (Uniform integrability in the PSEUDO-IID regime – CNN.) Consider a convolutional neural network in the PSEUDO-IID regime. Consider a collection of random variables

$z_l^{(c)}(x_a)[n]; z_l^{(c)}(x_b)[n]; z_l^{(c)}(x_c)[n]; z_l^{(c)}(x_d)[n]$ with any $j; k; l \in \{1, \dots, N\}; X_a; X_b; X_c; X_d \in \mathbb{R}^D$, neither necessarily distinct, obtained by the recursion 9. Then, the family of random variables

$$z_l^{(c)}(X_a)[n]z_l^{(c)}(X_b)[n]z_l^{(c)}(X_c)[n]z_l^{(c)}(X_d)[n];$$

indexed by n is uniformly integrable for any $\gamma = f + 1; \quad L + 1 \leq \gamma \leq L + 1 + g$.

Proof. As previously, this proposition holds some novelty of this paper, extending already known proofs in the standard Gaussian setting to the PSEUDO-IID regime in CNNs. Note that this directly implies the universality of the Gaussian Process behaviour for CNNs in the regime, which has been established so far only by Yang (2021) to the best of our knowledge. Our result goes beyond. We recall that the data propagation is described by Equation 9.

Once again, we observe it is sufficient to show such that the moment of order γ of the feature maps is uniformly bounded. We do this again by induction.

Base case. Since $h_i^{(1)}(X)[n] = b_i^{(1)} + \sum_{j=1}^P c_{0j} \sum_{k=1}^K U_{ij}^{(1)} x_j$, and weights and biases are independent, for some $p \leq 8$ we can write

$$\begin{aligned} E |h_i^{(1)}(X)[n]|^p &= E \left| b_i^{(1)} + \sum_{j=1}^P c_{0j} \sum_{k=1}^K U_{ij}^{(1)} x_j \right|^p + E |b_i^{(1)}|^p \\ &= K_p \sum_{j=1}^P c_{0j}^p \sum_{k=1}^K E |x_j|^{2p} + E |b_i^{(1)}|^p; \end{aligned}$$

where $\sum_{k=1}^K x_j$ is the part of signal around the pixel and we have used the condition (iii) of Definition 4 in the last line. Observe that the RHS is independent of n .

Inductive step. Let us assume that for any $X_t \in \mathbb{R}^D$, $1 \leq t \leq L$ and $i \in \{1, \dots, N\}$, there exists $c_i \in (0; 1)$ such that the eighth moment of the pre-activations from the previous layer $h_j^{(c-1)}(X_t)[n]$ is bounded by a constant independent from $n \in \mathbb{N}$, $1 \leq n \leq n$, and n , for all $X_t \in \mathbb{R}^D$.

Mirroring our proof in the fully-connected case, we will show that this propagates to the next layer,

$$E |h_i^{(c)}(X_t)[n]|^8 < 1 :$$

From the third condition of the PSEUDO-IID regime, we can compute the expectation conditioned on the vector of activations $h^{(c-1)}(X_t)[n]$,

$$\begin{aligned} E \left| \sum_{j=1}^P c_{1j} \sum_{k=1}^K U_{ij}^{(c)} h_j^{(c-1)}(X_t)[n] \right|^8 &= E \left| \sum_{j=1}^P c_{1j} \sum_{k=1}^K U_{ij}^{(c)} h_j^{(c-1)}(X_t)[n] \right|^8 \\ &= K C_{1j}^8 E \left| \sum_{k=1}^K h_j^{(c-1)}(X_t)[n] \right|^8; \end{aligned}$$

where the norm of the activations can be computed using the linear envelope property,

$$\begin{aligned} E \left| \sum_{k=1}^K h_j^{(c-1)}(X_t)[n] \right|^8 &= E \left| \sum_{j=1}^P c_{1j} \sum_{k=1}^K U_{ij}^{(c)} h_j^{(c-1)}(X_t)[n] \right|^8 \\ &= E \left| \sum_{j=1}^P c_{1j} \sum_{k=1}^K U_{ij}^{(c)} h_j^{(c-1)}(X_t)[n] \right|^8 \\ &= c^2 + M^2 E |h_j^{(c-1)}(X_t)[n]|^2 + 2cM E |h_j^{(c-1)}(X_t)[n]|^4; \end{aligned}$$

This last quantity turns out to be the weighted sum of $C_{-1}[n]^4$ terms (recall being the iter size, which is finite) of the form

$$E h_{j_1; \cdot}^{(\cdot-1)}(X_t)^{p_1} [n] h_{j_2; \cdot}^{(\cdot-1)}(X_t)^{p_2} [n] h_{j_3; \cdot}^{(\cdot-1)}(X_t)^{p_3} [n] h_{j_4; \cdot}^{(\cdot-1)}(X_t)^{p_4} [n] ;$$

that can be bounded using Lemma 3 combined with our inductive hypothesis. Observe how the factors $C_{-1}[n]^4$ cancel out and lead to a bound independent from

Using the recursion formulae for the data propagation in the CNN architecture recalled in 9 and the convexity of the map $x \mapsto x^8$ on \mathbb{R}^+ , we have

$$E h_i^{(\cdot)}(X_t)[n]^8 \leq 2^8 E b_i^{(\cdot)}{}^8 + \sum_{j=1}^C \sum_{J=1}^J \sum_{K=1}^K U_{ij}^{(\cdot)} h_j^{(\cdot-1)}(X_t)[n]^8 ;$$

The first expectation is finite and uniformly bounded as the biases are Gaussians in the PSEUDO-IID regime and we have just shown above the boundedness of the second expectation.

Therefore,

$$E h_j^{(\cdot-1)}(X_t)[n]^8 < 1 ; \tag{31}$$

and taking the supremum over the finite input data set does not change the uniformly bounded property, which was needed to be shown. □

D EXAMPLES CONCERNING THE BOUNDED MOMENT CONDITION (III) OF THE PSEUDO-IID DISTRIBUTION

The bounded moment condition (iii) of the PSEUDO-IID distribution in Definition 1 is a key condition in the distinct proof of IID matrices taken in Hanin (2021) (Lemma 2.9). We show some examples of distributions in Figure 4 which verify or violate the conditions we identified as sufficient to rigorously prove the convergence of random neural networks to gaussian processes in the large width limit.

E STRUCTURED SPARSE WEIGHT MATRICES IN THE FULLY-CONNECTED SETTING

Fig. 5 shows examples of permuted block-sparse weight matrices used to initialize a fully-connected network in order to produce the plots given in Fig. 2-7.

F FURTHER NUMERICAL SIMULATIONS VALIDATING THEOREM 1

Fig. 6 shows QQ plots for the histograms in Fig. 2 as compared to their infinite with Gaussian limit. These QQ plots show how the PSEUDO-IID networks approach the Gaussian process at somewhat different rates with uniform approaching the fastest.

Fig. 7 explores the growing independence of entries in $(x)[n]$ for different by showing their joint distributions for two distinct choices of μ and σ . Moreover, convergence to the limiting isotropic Gaussian distribution $h_i^{(\cdot)}(x)[n]$ is overlaid in the same plots. Uniform converges the quickest, while PSEUDO-IID Gaussian low-rank and structured sparse converge towards an isotropic distribution somewhat slower, albeit already showing good agreement at $n=100$. The horizontal and vertical axis in each subplot of Fig. 7 are $h_i^{(\cdot)}(x)$ for $i = 1$ and 2 respectively.

Figure 4: Different cases where conditions (ii) and (iii) of the pseudo IID regime are either satisfied or violated. Condition (iii) is considered with $n = 8$ and $\alpha = (1; \dots; 1)$ such that it becomes $E \sum_{j=1}^n X_j^8 = K$, where $X = (X_1; \dots; X_N)$ is regarded as one row of the weight matrix. The chosen distribution for the vector x impacts whether the network is in the PSEUDO-IID regime. In the identical coordinates case, $x = (X_1; \dots; X_1)$ is the concatenation of the same realisation sampled from a standard normal. Not only the traditional assumption is broken as the coordinates are obviously dependent but also condition (iii) is violated, thus resulting in an unbounded expectation when growing the dimension N . The autoregressive process $(x_1; \dots; x_N)$ shown is obtained by $X_i = \alpha_i + X_{i-1}$, where α_i are IID multivariate gaussian noises of dimension n . The correlation between the coordinates does not decrease fast enough with the dimension to get a bound on the computed expectation and condition (iii) is once again violated. On the contrary, from the plot produced by sampling Cauchy distributions, it is not obvious whether condition (iii) holds. Nonetheless, condition (ii) which ensures the finiteness of the variance is not, thus a random network initialized with IID Cauchy weights falls outside the scope of our identified broad class of distributions to ensure a convergence towards a gaussian process. The last cases of samples taken either from scaled multivariate normals in red or uniformly sampled from the unit sphere in orange (with appropriate scaling such that condition (ii) holds) show that there exists a bound independent from the dimension on the expectation of interest. The empirical expectations are taken considering averages over 100000 samples.

width = 3

width = 30

width = 300

Figure 5: Example of a permuted block-sparse weight matrix at initialization of a fully-connected network with increasing width N . The matrix is initialized with identically and independently sampled diagonal blocks from a scaled Gaussian. Its rows and columns are then randomly permuted in order to satisfy the PSEUDO-IID conditions. The block size is set to $\lfloor \frac{N}{2} \rfloor$. Entries in yellow are zero and entries in black are nonzero..

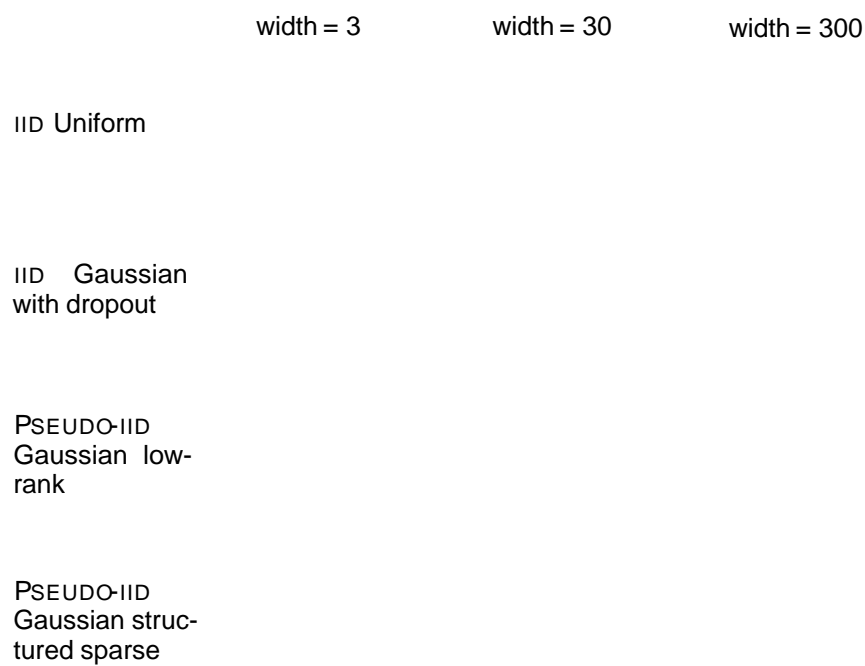


Figure 6: Q-Q plots of the pre-activations values in Fig. 2 as an alternative way of showing the convergence of the pre-activation of a fully-connected network to a Gaussian as fully characterized in Theorem 1. The settings of the experiment are the same as those in Fig. 2.

