

# Characterizing Well-behaved vs. Pathological Deep Neural Network Architectures

Antoine Labatie

antoine.labatie@centraliens.net

## Abstract

We introduce a principled approach, requiring only mild assumptions, for the characterization of deep neural networks at initialization. Our approach applies both to fully-connected and convolutional networks and incorporates the commonly used techniques of batch normalization and skip-connections. Our key insight is to consider the evolution with depth of statistical moments of signal and sensitivity, thereby characterizing the well-behaved or pathological behaviour of input-output mappings encoded by different choices of architecture. We establish: (i) for feedforward networks with and without batch normalization, depth multiplicativity inevitably leads to ill-behaved moments and distributional pathologies; (ii) for residual networks, on the other hand, the mechanism of identity skip-connection induces power-law rather than exponential behaviour, leading to well-behaved moments and no distributional pathology.<sup>1</sup>

## 1 Introduction

Advances in the design of neural network architectures have more often come from the relentless race of practical applications rather than by principled approaches. Even after wide adoption, many common choices and rules of thumbs still await for theoretical validation. This is unfortunate since the emergence of a theory to fully characterize deep neural networks in terms of wanted and unwanted properties would enable to guide further improvements.

An important branch of research towards this theoretical characterization has focused on random networks at initialization, i.e. before any training has occurred. The characterization of random networks offers valuable insights into the hypothesis space of input-output mappings reachable during training. In a Bayesian view, it also unveils the prior on the input-output mapping encoded by the choice of architecture (Neal, 1996a; Williams, 1997; Lee et al., 2017). Experimentally, well-behaved input-output mappings at initialization were extensively found to be predictive of model trainability and even post-training performance (Schoenholz et al., 2016; Balduzzi et al., 2017; Pennington et al., 2017; Yang & Schoenholz, 2017; Chen et al., 2018; Xiao et al., 2018).

Unfortunately, even this simplifying case of random networks is still challenging in various respects: (i) there is a complex interplay of different sources of randomness from input data and from model parameters; (ii) the finite number of units in each layer leads to additional complexity; (iii) convolutional neural networks – widely used in many applications – also lead to additional complexity. The difficulty (i) is typically circumvented by restricting input data to simplifying cases: a one-dimensional input manifold (Poole et al., 2016; Raghu et al., 2017), two input data points (Schoenholz et al., 2016; Balduzzi et al., 2017; Yang & Schoenholz, 2017; Chen et al., 2018; Xiao et al., 2018), a batch of input data points (Anonymous, 2019). The difficulty (ii) is commonly circumvented by either considering the simplifying case of infinite width with the convergence of neural networks to Gaussian processes

<sup>1</sup>Code to reproduce all results will be made available upon publication.

(Neal, 1996a; Daniely et al., 2016; Lee et al., 2017; Matthews et al., 2018) or the simplifying case of typical activation patterns for ReLU networks (Balduzzi et al., 2017). Finally the difficulty (iii) is rarely tackled as most analyses only apply to fully-connected networks. To the best of our knowledge, all attempts have thus far been limited in their scope or simplifying assumptions.

In this paper, we introduce a principled approach to characterize random deep neural networks at initialization. Our approach does not require any of the usual simplifications of infinite width, gaussianity, typical activation patterns, restricted input data. It further applies both to fully-connected and convolutional networks, and incorporates batch normalization and skip-connections. Our key insight is to consider statistical moments of signal and sensitivity with respect to input data as random variables which depend on model parameters. By studying the evolution of these moments with depth, we characterize the well-behaved or pathological behaviour of input-output mappings encoded by different choices of architectures. Our findings span the topics of one-dimensional degeneracy, pseudo-linearity, exploding sensitivity, exponential and power-law evolution with depth.

## 2 Propagation

We start by formulating the propagation for neural networks with neither batch normalization nor skip-connections, that we refer as *vanilla networks*. The formulation will be slightly adapted in Section 6 with *batch-normalized feedforward nets*, and in Section 7 with *batch-normalized resnets*.

**Clean propagation.** We consider a random tensorial input  $\mathbf{x} \in \mathbb{R}^{n \times \dots \times n \times N_0}$ , spatially  $d$ -dimensional with extent  $n$  in all spatial dimensions and  $N_0$  channels. This input  $\mathbf{x}$  is fed into a  $d$ -dimensional convolutional neural network with periodic boundary conditions and fixed spatial extent equal to  $n$ .<sup>1</sup> For each layer  $l$ , we denote  $N_l$  the number of channels or *width*,  $K_l$  the convolutional spatial extent,  $\omega^l \in \mathbb{R}^{K_l \times \dots \times K_l \times N_{l-1} \times N_l}$  the weight tensors,  $\mathbf{b}^l \in \mathbb{R}^{N_l}$  the biases,  $\mathbf{x}^l, \mathbf{y}^l \in \mathbb{R}^{n \times \dots \times n \times N_l}$  the tensors of post-activations and pre-activations. We further denote  $\alpha$  the spatial position,  $c$  the channel, and  $\phi$  the activation function. Adopting the convention  $\mathbf{x}^0 \equiv \mathbf{x}$ , the propagation at each layer is given by

$$\mathbf{y}^l = \omega^l * \mathbf{x}^{l-1} + \beta^l, \quad \mathbf{x}^l = \phi(\mathbf{y}^l),$$

where  $\beta^l \in \mathbb{R}^{n \times \dots \times n \times N_l}$  is the tensor with repeated version of  $\mathbf{b}^l$  at each spatial position. From now on, we refer to the propagated tensor  $\mathbf{x}^l$  as the *signal*.

**Noisy propagation.** Next we suppose that the input signal  $\mathbf{x}$  is corrupted by a small *white noise* tensor  $\mathrm{d}\mathbf{x} \in \mathbb{R}^{n \times \dots \times n \times N_0}$  with independent and identically distributed components such that  $\mathbb{E}_{\mathrm{d}\mathbf{x}}[\mathrm{d}\mathbf{x}_i \mathrm{d}\mathbf{x}_j] = \sigma_{\mathrm{d}\mathbf{x}}^2 \delta_{ij}$  ( $\sigma_{\mathrm{d}\mathbf{x}} \ll 1$ ), with  $\delta_{ij}$  the Kronecker delta. The noisy signal is propagated into the same neural network and we keep track of the noise corruption with the tensor  $\mathrm{d}\mathbf{x}^l$  defined as  $\mathrm{d}\mathbf{x}^0 \equiv \mathrm{d}\mathbf{x}$  and  $\mathrm{d}\mathbf{x}^l \equiv \Phi_l(\mathbf{x} + \mathrm{d}\mathbf{x}) - \Phi_l(\mathbf{x})$ , with  $\mathbf{x}^l = \Phi_l(\mathbf{x})$  the neural network mapping from layer 0 to  $l$ . The simultaneous propagation of the *signal*  $\mathbf{x}^l$  and the *noise*  $\mathrm{d}\mathbf{x}^l$  is given by

$$\mathbf{y}^l = \omega^l * \mathbf{x}^{l-1} + \beta^l, \quad \mathbf{x}^l = \phi(\mathbf{y}^l), \quad (1)$$

$$\mathrm{d}\mathbf{y}^l = \omega^l * \mathrm{d}\mathbf{x}^{l-1}, \quad \mathrm{d}\mathbf{x}^l = \phi'(\mathbf{y}^l) \odot \mathrm{d}\mathbf{y}^l, \quad (2)$$

where  $\odot$  denotes the element-wise tensor multiplication and Eq. (2) is obtained by taking the derivative in Eq. (1). For given  $\mathbf{x}$ , the mapping from  $\mathrm{d}\mathbf{x}$  to  $\mathrm{d}\mathbf{x}^l$  induced by Eq. (2) is linear. The noise  $\mathrm{d}\mathbf{x}^l$  thus stays centered with respect to  $\mathrm{d}\mathbf{x}$  during propagation such that  $\forall \mathbf{x}, \alpha, c: \mathbb{E}_{\mathrm{d}\mathbf{x}}[\mathrm{d}\mathbf{x}_{\alpha,c}^l] = 0$ .

To get rid of the dependence on  $\sigma_{\mathrm{d}\mathbf{x}}$ , the random *sensitivity tensor* is introduced as the rescaling  $\mathbf{s}^0 \equiv \mathbf{s} \equiv \mathrm{d}\mathbf{x}/\sigma_{\mathrm{d}\mathbf{x}}$  and  $\mathbf{s}^l \equiv \mathrm{d}\mathbf{x}^l/\sigma_{\mathrm{d}\mathbf{x}}$ . By linearity of Eq. (2), the sensitivity tensor  $\mathbf{s}^l$  is the result of the simultaneous propagation of  $\mathbf{x}^l$  and  $\mathbf{s}^l$  in Eq. (1) and (2). We also have  $\mathbb{E}_{\mathbf{s}}[\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$  and  $\forall \mathbf{x}, \alpha, c: \mathbb{E}_{\mathbf{s}}[\mathbf{s}_{\alpha,c}^l] = 0$ . The sensitivity tensor encodes *derivative information* while avoiding the burden of increased dimensionality (see calculations in Appendix D.3 and Appendix D.2). This will prove very useful.

<sup>1</sup>The assumptions of periodic boundary conditions and fixed spatial extent  $n$  are made for simplicity of the analysis. Possible relaxations are discussed in Section C.2.

**Scope.** We require two very mild assumptions: (i)  $\mathbf{x}$  is not trivially 0 such that  $\mathbb{E}_{\mathbf{x}, \alpha, c}[\mathbf{x}_{\alpha, c}^2] > 0$ ;<sup>2</sup> (ii) the width  $N_l$  is bounded.

More importantly, we restrict to the – most commonly used – ReLU activation function  $\phi(\cdot) = \max(\cdot, 0)$ . Even though  $\phi$  is not differentiable at 0, we still define  $d\mathbf{x}^l$  and  $\mathbf{s}^l$  as the result of the simultaneous propagation of Eq. (1) and Eq. (2) with the convention  $\phi'(0) \equiv 1/2$ . Section 3 implicitly relies on the assumption that  $\Phi_l(\mathbf{x})$  remains differentiable, implying  $d\mathbf{x}^l = \Phi_l(\mathbf{x} + d\mathbf{x}) - \Phi_l(\mathbf{x})$ , almost surely (a.s.) with respect to  $\mathbf{x}$  (details and justification of this assumption in Appendix D.1).

Note that fully-connected neural networks are still included in our analysis as the subcase  $n = 1$ .

### 3 Input data randomness

We now turn our attention to the distributions with respect to input data of signal and sensitivity  $P_{\mathbf{x}}(\mathbf{x}^l)$ ,  $P_{\mathbf{x}, s}(\mathbf{s}^l)$ . To outline the importance of these distributions, we may express by layer composition the output of an  $L$ -layer neural network as  $\mathbf{x}^L = \Phi_{L, L}(\mathbf{x}^l)$ , with  $l < L$  and  $\Phi_{l, L}$  the *upper neural network* mapping from layer  $l$  to  $L$ . It means that  $\mathbf{x}^l$  and  $d\mathbf{x}^l$  can be seen as input signal and noise of this upper neural network. It also means that the distributions  $P_{\mathbf{x}^l}(\mathbf{x}^l) = P_{\mathbf{x}}(\mathbf{x}^l)$  and  $P_{d\mathbf{x}^l}(d\mathbf{x}^l) = P_{\mathbf{x}, d\mathbf{x}}(d\mathbf{x}^l)$ , or equivalently  $P_{\mathbf{s}^l}(\mathbf{s}^l) = P_{\mathbf{x}, s}(\mathbf{s}^l)$ , must remain well-behaved for this upper network to have a chance to accomplish its task successfully. We will return to this argument in Section 3.2 by detailing the penalizing effects of ill-behaved  $P_{\mathbf{x}}(\mathbf{x}^l)$ ,  $P_{\mathbf{x}, s}(\mathbf{s}^l)$ .

#### 3.1 Characterizing distributions with respect to input data

First we need to introduce the following statistical quantities to characterize  $P_{\mathbf{x}}(\mathbf{x}^l)$ ,  $P_{\mathbf{x}, s}(\mathbf{s}^l)$ :

- The **feature map vector** and **centered feature map vector** associated with any random tensor  $\mathbf{v}$ ,

$$\varphi(\mathbf{v}, \alpha) \equiv \mathbf{v}_{\alpha, :}, \quad \hat{\varphi}(\mathbf{v}, \alpha) \equiv \mathbf{v}_{\alpha, :} - \mathbb{E}_{\mathbf{v}, \alpha}[\mathbf{v}_{\alpha, :}],$$

where  $\alpha$  is uniformly sampled in  $\{1, \dots, n\}^d$  and the denotations  $\varphi(\mathbf{v}, \alpha)$ ,  $\hat{\varphi}(\mathbf{v}, \alpha)$  remind the randomness of both  $\mathbf{v}$  and  $\alpha$ .

In a similar vein as batch normalization and Xiao et al. (2018), feature maps and centered feature maps at different spatial positions  $\varphi(\mathbf{v}, \alpha)$  and  $\hat{\varphi}(\mathbf{v}, \alpha)$  are implicitly seen as sub-signals propagated together as part of tensorial meta-signals. Statistically, we work at the granularity of sub-signals and we treat equally the randomness from  $\mathbf{v}$  and  $\alpha$ .<sup>3</sup>

- The **non-central moment** and **central moment** of order  $p$  associated with any random tensor  $\mathbf{v}$  per-channel and averaged over channels,

$$\begin{aligned} \nu_{p, c}(\mathbf{v}) &\equiv \mathbb{E}_{\mathbf{v}, \alpha}[\varphi(\mathbf{v}, \alpha)_c^p], & \mu_{p, c}(\mathbf{v}) &\equiv \mathbb{E}_{\mathbf{v}, \alpha}[\hat{\varphi}(\mathbf{v}, \alpha)_c^p], \\ \nu_p(\mathbf{v}) &\equiv \mathbb{E}_{\mathbf{v}, \alpha, c}[\varphi(\mathbf{v}, \alpha)_c^p], & \mu_p(\mathbf{v}) &\equiv \mathbb{E}_{\mathbf{v}, \alpha, c}[\hat{\varphi}(\mathbf{v}, \alpha)_c^p]. \end{aligned}$$

- The **effective rank** associated with any random tensor  $\mathbf{v}$  (Vershynin, 2010),

$$r_{\text{eff}}(\mathbf{v}) \equiv \frac{\text{Tr } C[\varphi(\mathbf{v}, \alpha)]}{\|C[\varphi(\mathbf{v}, \alpha)]\|} = \frac{\sum_i \lambda_i}{\max_i \lambda_i},$$

where  $C[\cdot]$  denotes the covariance matrix,  $\|\cdot\|$  the spectral norm, and  $(\lambda_i)$  the eigenvalues of  $C[\varphi(\mathbf{v}, \alpha)]$ . The effective rank is always greater than 1, and it intuitively measures the number of effective directions which concentrate the variance of  $\varphi(\mathbf{v}, \alpha)$ .

<sup>2</sup>Whenever  $\alpha$  and  $c$  are considered as random variables they are supposed uniformly sampled among all spatial positions  $\{1, \dots, n\}^d$  and all channels  $\{1, \dots, N_l\}$ .

<sup>3</sup>In the context e.g. of  $\mathbf{x}_{\alpha, c}^l$  note that propagation from layer 1 to  $l$  solely depends on the randomness of  $\mathbf{x}$ , the randomness of  $\alpha$  being only involved in the selection of the final position at layer  $l$ . For simplicity,  $\alpha$  will still be referred as input data.

- The **normalized sensitivity** – our key metric – derived from the moments of  $\mathbf{x}^l$  and  $\mathbf{s}^l$ ,

$$\chi^l \equiv \left( \frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2}. \quad (3)$$

To better understand the definition of  $\chi^l$  let us consider the classification task. Since the goal is to set apart different signals, the mean signal over  $\mathbf{x}$ ,  $\alpha$  is uninformative. Centered feature maps  $\hat{\varphi}(\mathbf{x}^l, \alpha)$  thus constitute the informative part of the signal and  $\mu_2(\mathbf{x}^0)$ ,  $\mu_2(\mathbf{x}^l)$  normalize by signal ‘informative content’ in Eq. (3). The normalized sensitivity  $\chi^l$  exactly measures a root mean square sensitivity when neural network input and output signals  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{x}}^l$  are rescaled to have unit signal ‘informative content’:  $\mu_2(\tilde{\mathbf{x}}) = 1$  and  $\mu_2(\tilde{\mathbf{x}}^l) = 1$  (proof in Appendix D.2). If we denote the rescaled input-output mapping as  $\tilde{\mathbf{x}}^l = \Psi_l(\tilde{\mathbf{x}})$ , then this property is simply expressed as  $\chi^l = \mathbb{E}_{\tilde{\mathbf{x}}} [\Psi_l'(\tilde{\mathbf{x}})^2]^{1/2}$  in the fully-connected case with 1-dimensional rescaled input and output signals  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{x}}^l$  ( $N_0 = 1$ ,  $N_l = 1$ ). We provide an illustration in Fig. 1.

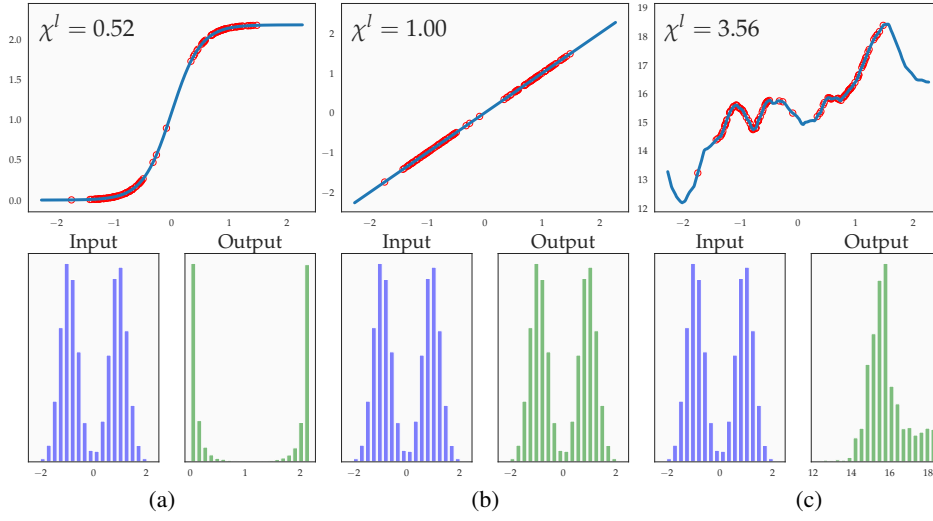


Figure 1: Illustration of the normalized sensitivity  $\chi^l$  for fully-connected networks of  $l$  layers with 1-dimensional input and output rescaled signals  $\tilde{\mathbf{x}}$  and  $\mathbf{x}^l$ . The distribution of the input  $\tilde{\mathbf{x}}$  is a mixture of two Gaussians. We show the result of the propagation in three different cases: (a)  $l = 1$  layer with sigmoid activation; (b)  $l = 1$  layer with linear activation; (c)  $l = 25$  randomly initialized layers with  $N_k = 100$  channels and ReLU activation for  $1 \leq k < l$  and linear activation for  $k = l$ . *Top*: full input-output mapping (blue) and randomly sampled input-output data points (red circles). *Bottom*: histograms of inputs and outputs.

The normalized sensitivity  $\chi^l$  is directly connected to the Jacobian norm in the fully-connected case and to sensitivity to signal perturbation in the general case, known in turn for being connected to generalization (Sokolic et al., 2017; Arora et al., 2018; Morcos et al., 2018; Novak et al., 2018; Philipp & Carbonell, 2018).<sup>4</sup> The notion of *sharpness* which quantifies sensitivity to weight perturbation was similarly shown to be connected to generalization (Hochreiter & Schmidhuber, 1997; Keskar et al., 2016; Smith & Le, 2017; Neyshabur et al., 2017). Sensitivity to signal perturbation and weight perturbation are tightly connected since the introduction of a noise  $d\omega^l$  on the weights amounts to the introduction of a noise  $d\mathbf{y}^l = d\omega^l * \mathbf{x}^{l-1}$  and  $d\mathbf{x}^l = \phi'(\mathbf{y}^l) \odot d\mathbf{y}^l$  on the signal in Eq. (1) and Eq. (2).

<sup>4</sup>The coefficient defined in Philipp & Carbonell (2018) is equivalent to the normalized sensitivity  $\chi^l$  in the fully-connected case. Section D.3 provides details on this equivalence and the reasons for our change of terminology.

### 3.2 Distributional pathologies

Using these statistical tools, we are able to characterize the distributional pathologies – with ill-behaved  $P_{\mathbf{x}}(\mathbf{x}^l)$ ,  $P_{\mathbf{x},\mathbf{s}}(\mathbf{s}^l)$  – that we will encounter:

- **Zero-dimensional signal:**  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 0$ . To understand this pathology, let us consider the following mean vectors and rescaling of the signal:

$$\boldsymbol{\nu}^l \equiv (\nu_{1,c}(\mathbf{x}^l))_{1 \leq c \leq N_l}, \quad \tilde{\mathbf{x}}^l \equiv \frac{1}{\|\boldsymbol{\nu}^l\|_2} \mathbf{x}^l, \quad \tilde{\boldsymbol{\nu}}^l \equiv (\nu_{1,c}(\tilde{\mathbf{x}}^l))_{1 \leq c \leq N_l}.$$

Then  $\|\tilde{\boldsymbol{\nu}}^l\|_2 = 1$  and the pathology  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 0$  implies  $\mu_2(\tilde{\mathbf{x}}^l) \xrightarrow{l \rightarrow \infty} 0$  (proof in Appendix D.4). It follows that  $\varphi(\tilde{\mathbf{x}}^l, \boldsymbol{\alpha})$  becomes *point-like* concentrated at the point  $\tilde{\boldsymbol{\nu}}^l$  of unit  $L^2$  norm. In the limit of strict point-like concentration, the upper neural network from layer  $l$  to  $L$  is limited to *random guessing* since it ‘sees’ all inputs the same and cannot distinguish between them.

- **One-dimensional signal:**  $r_{\text{eff}}(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 1$ . This pathology implies that  $\varphi(\mathbf{x}^l, \boldsymbol{\alpha})$  has its variance which becomes concentrated in only one direction, and thus that  $\varphi(\mathbf{x}^l, \boldsymbol{\alpha})$  becomes *line-like* concentrated. In the limit of strict line-like concentration, the upper neural network from layer  $l$  to  $L$  only ‘sees’ a *single feature* from  $\mathbf{x}$ .
- **Exploding sensitivity:**  $\chi^l \xrightarrow{l \rightarrow \infty} \infty$ . To understand this pathology, let us push further the view of noisy propagation with the *signal-to-noise ratio*  $SNR^l$  and the *noise factor*  $F^l$ :

$$SNR^l \equiv \frac{\mu_2(\mathbf{x}^l)}{\mu_2(\mathbf{d}\mathbf{x}^l)}, \quad F^l \equiv \frac{SNR^0}{SNR^l} = \frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} = (\chi^l)^2, \quad (4)$$

where Eq. (4) follows from  $\mu_2(\mathbf{d}\mathbf{x}^l) = \sigma_{\mathbf{d}\mathbf{x}}^2 \mu_2(\mathbf{s}^l)$  and  $\mu_2(\mathbf{d}\mathbf{x}^0) = \sigma_{\mathbf{d}\mathbf{x}}^2$ . In logarithmic decibel scale  $SNR_{\text{dB}}^l = SNR_{\text{dB}}^0 - 20 \log_{10} \chi^l$ , indicating that  $\chi^l$  measures how the neural network from layer 0 to  $l$  degrades ( $\chi^l > 1$ ) or enhances ( $\chi^l < 1$ ) the input signal-to-noise ratio. The pathology  $\chi^l \xrightarrow{l \rightarrow \infty} \infty$  implies  $F^l \xrightarrow{l \rightarrow \infty} \infty$  and  $SNR^l \xrightarrow{l \rightarrow \infty} 0$ , meaning that the noisy signal  $\mathbf{x}^l + \mathbf{d}\mathbf{x}^l$  becomes completely random for any level of input noise  $\mathbf{d}\mathbf{x}$ .<sup>5</sup> The upper neural network from layer  $l$  to  $L$  is then limited to *random guessing*.

## 4 Model parameters randomness

We now introduce model parameters as the second source of randomness. We consider *random networks* at initialization, which we suppose is *standard*: (i) weights and biases are initialized following He et al. (2015); (ii) when pre-activations are batch-normalized, scale and shift batch normalization parameters are initialized with ones and zeros respectively.

Considering random networks at initialization is justified in two respects:

- From a *Bayesian perspective*, the random distribution on the weights encodes a prior distribution on input-output mappings for a given choice of architecture (Neal, 1996a; Williams, 1997). This prior distribution can be combined with the likelihood of model parameters given training data in order to obtain a posterior distribution on input-output mappings, which in turn enables Bayesian predictions (Lee et al., 2017).
- If distributional pathologies arise at initialization, it is likely that training will be difficult. Indeed if the distributional pathology is severe, it becomes very difficult for the neural network to adjust its parameters and undo the pathology. As an illustration, consider the pathology of zero-dimensional signal:  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 0$ . In this case, the upper neural network from layer  $l$  to  $L$  must adjust its bias parameters very precisely in order to center the signal and distinguish between different inputs  $\mathbf{x}$ . In support of this argument, several works have shown that deep neural networks are precisely trainable when they are not subject to the pathology

<sup>5</sup>In this case, Eq. (2) eventually does not hold since the noise  $\mathbf{d}\mathbf{x}^l$  at layer  $l$  becomes non-infinitesimal even for infinitesimal input noise  $\mathbf{d}\mathbf{x}$ .

of different inputs  $\mathbf{x}$  being strictly correlated at initialization, i.e. nearly undistinguishable (Schoenholz et al., 2016; Xiao et al., 2018).

From now on, our methodology is to consider all moment-related quantities, e.g.  $\mu_p(\mathbf{x}^l)$ ,  $\mu_p(\mathbf{s}^l)$ ,  $\nu_p(\mathbf{x}^l)$ ,  $\nu_p(\mathbf{s}^l)$ ,  $r_{\text{eff}}(\mathbf{x}^l)$ ,  $r_{\text{eff}}(\mathbf{s}^l)$ ,  $\chi^l$ , as random variables depending on  $(\omega^1, \beta^1, \dots, \omega^l, \beta^l)$ . We introduce the notation  $\Theta^l = (\omega^1, \beta^1, \dots, \omega^l, \beta^l)$  for the full set of parameters, and the notation  $\theta^l = \Theta^l | \Theta^{l-1}$  for the conditional set of parameters when  $(\omega^l, \beta^l)$  are considered as random and  $(\omega^1, \beta^1, \dots, \omega^{l-1}, \beta^{l-1})$  as given. We further denote the geometric increment  $\delta\mu_2(\mathbf{x}^l) \equiv \mu_2(\mathbf{x}^l)/\mu_2(\mathbf{x}^{l-1})$ .

**Evolution with Depth.** The evolution with depth of  $\mu_2(\mathbf{x}^l)$  can be written as

$$\begin{aligned} \log \mu_2(\mathbf{x}^l) - \log \mu_2(\mathbf{x}^0) &= \sum_{k \leq l} \underbrace{\log \mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)]}_{\overline{m}[\mu_2(\mathbf{x}^k)]} + \sum_{k \leq l} \underbrace{\mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)] - \log \mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)]}_{\underline{m}[\mu_2(\mathbf{x}^k)]} \\ &\quad + \sum_{k \leq l} \underbrace{\log \delta\mu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)]}_{\underline{s}[\mu_2(\mathbf{x}^k)]}, \end{aligned} \quad (5)$$

where Eq. (5) is obtained using  $\log \mu_2(\mathbf{x}^l) - \log \mu_2(\mathbf{x}^0) = \sum_{k \leq l} \log \delta\mu_2(\mathbf{x}^k)$  and by expressing  $\log \delta\mu_2(\mathbf{x}^k)$  with telescoping terms. Denoting  $\underline{\delta}\mu_2(\mathbf{x}^k) \equiv \delta\mu_2(\mathbf{x}^k)/\mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)]$  the multiplicatively centered increments, and using  $\mathbb{E}_{\theta^k}[\log \mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)]] = \log \mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)]$  and  $\mathbb{E}_{\theta^k}[\mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)]] = \mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)]$ , we obtain

$$\overline{m}[\mu_2(\mathbf{x}^k)] = \log \mathbb{E}_{\theta^k}[\delta\mu_2(\mathbf{x}^k)], \quad (6)$$

$$\underline{m}[\mu_2(\mathbf{x}^k)] = \mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)], \quad (7)$$

$$\underline{s}[\mu_2(\mathbf{x}^k)] = \log \underline{\delta}\mu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k}[\log \delta\mu_2(\mathbf{x}^k)]. \quad (8)$$

**Discussion.** First we note that  $\overline{m}[\mu_2(\mathbf{x}^k)]$  and  $\underline{m}[\mu_2(\mathbf{x}^k)]$  are random variables which depend on  $\Theta^{k-1}$ , while  $\underline{s}[\mu_2(\mathbf{x}^k)]$  is a random variable which depends on  $\Theta^k$ . We also note that  $\underline{m}[\mu_2(\mathbf{x}^k)] < 0$  by log-concavity, and that  $\underline{s}[\mu_2(\mathbf{x}^k)]$  is centered:  $\mathbb{E}_{\Theta^k}[\underline{s}[\mu_2(\mathbf{x}^k)]] = 0$ .

Under standard initialization, each channel provides an independent contribution to  $\mu_2(\mathbf{x}^k) = \frac{1}{N_k} \sum_c \mu_{2,c}(\mathbf{x}^k)$ . As a consequence, for large  $N_k$  the relative increment  $\underline{\delta}\mu_2(\mathbf{x}^k)$  has low expected deviation to 1, meaning with high probability that  $|\log \underline{\delta}\mu_2(\mathbf{x}^k)| \ll 1$ ,  $|\underline{m}[\mu_2(\mathbf{x}^k)]| \ll 1$ ,  $|\underline{s}[\mu_2(\mathbf{x}^k)]| \ll 1$ . In addition,  $\underline{s}[\mu_2(\mathbf{x}^k)]$  is centered and *non-correlated* at different  $k$  so its sum scales as  $\sqrt{l}$ , whereas the sums of  $\overline{m}[\mu_2(\mathbf{x}^k)]$  and  $\underline{m}[\mu_2(\mathbf{x}^k)]$  scale as  $l$  (see Lemma 10 in Appendix E.1). The term  $\underline{s}[\mu_2(\mathbf{x}^k)]$  is thus doubly negligible. In summary, the evolution with depth is dominated by  $\overline{m}[\mu_2(\mathbf{x}^k)]$  when this term is non-vanishing and by  $\underline{m}[\mu_2(\mathbf{x}^k)]$  otherwise. The same analysis can be applied to other positive moments such as  $\nu_2(\mathbf{x}^l)$ ,  $\nu_2(\mathbf{s}^l)$ ,  $\mu_2(\mathbf{s}^l)$ . It can also be applied to the ratio  $\mu_2(\mathbf{s}^l)/\mu_2(\mathbf{x}^l)$  and thus to  $\chi^l$ .

**Further notation.** From now on, the geometric increment of any quantity is denoted with  $\delta$ . The definitions of  $\overline{m}$ ,  $\underline{m}$  and  $\underline{s}$  in Eq. (6), (7) and (8) are extended to other central and non-central moments of signal and sensitivity, as well as  $\chi^l$  with  $\overline{m}[\chi^l] = \frac{1}{2}(\overline{m}[\mu_2(\mathbf{s}^l)] - \overline{m}[\mu_2(\mathbf{x}^l)])$ ,  $\underline{m}[\chi^l] = \frac{1}{2}(\underline{m}[\mu_2(\mathbf{s}^l)] - \underline{m}[\mu_2(\mathbf{x}^l)])$ ,  $\underline{s}[\chi^l] = \frac{1}{2}(\underline{s}[\mu_2(\mathbf{s}^l)] - \underline{s}[\mu_2(\mathbf{x}^l)])$ .

We introduce the notation  $a \simeq b$  when  $a(1 + \epsilon_a) = b(1 + \epsilon_b)$ , with  $|\epsilon_a| \ll 1$ ,  $|\epsilon_b| \ll 1$  with high probability. And the notation  $a \lesssim b$  when  $a(1 + \epsilon_a) \leq b(1 + \epsilon_b)$ , with  $|\epsilon_a| \ll 1$ ,  $|\epsilon_b| \ll 1$  with high probability. From now on, we assume that the *width is large*, implying  $\exp(\underline{m}[\chi^l] + \underline{s}[\chi^l]) \simeq 1$  and  $\delta\chi^l \simeq \exp(\overline{m}[\chi^l])$ . We stress that this assumption is milder than the commonly used *mean-field* assumption of infinite width:  $N_l \rightarrow \infty$ . Indeed we still consider  $\overline{m}[\chi^l]$  as random and  $\mathbf{x}^l$  as possibly non-Gaussian, whereas mean-field would consider  $\overline{m}[\chi^l]$  as non-random and  $\mathbf{x}^l$  as Gaussian.

## 5 Vanilla Networks

We are fully equipped to characterize neural network architectures at initialization. We start by analyzing vanilla networks corresponding to the equations of propagation introduced in Section 2.



**Theorem 1. Moments of vanilla networks.** (proof in Appendix E.2) *There exist positive constants  $m_{\min}, m_{\max}, v_{\min}, v_{\max} > 0$ , random variables  $(m_l), (m'_l), (s_l), (s'_l)$  and events  $(A_l)$  with probability  $\mathbb{P}_{\Theta^l}[A_l] \geq \prod_{k=1}^l (1 - 2^{-N_k+1})$  such that under  $A_l$ :  $s_l$  and  $s'_l$  are centered and*

$$\begin{aligned} \log \nu_2(\mathbf{x}^l) &= -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), & m_{\min} \leq m_l \leq m_{\max}, & v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= -lm'_l + \sqrt{l}s'_l, & m_{\min} \leq m'_l \leq m_{\max}, & v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max}. \end{aligned}$$

**Discussion.** First we discuss the conditionality on  $A_l$  which is necessary to exclude the collapse  $\nu_2(\mathbf{x}^l) = 0$  and  $\mu_2(\mathbf{s}^l) = 0$  (with undefined  $\log \nu_2(\mathbf{x}^l)$  and  $\log \mu_2(\mathbf{s}^l)$ ), occuring e.g. when all elements of  $\omega^l$  are strictly negative (Lu et al., 2018). The complementary event  $A^c$  has probability exponentially small in the width due to

$$-\mathbb{P}_{\Theta^l}[A_l^c] \simeq \log(1 - \mathbb{P}_{\Theta^l}[A_l]) = \log \mathbb{P}_{\Theta^l}[A_l] \geq \sum_{k=1}^l \log(1 - 2^{-N_k+1}) \simeq -\sum_{k=1}^l 2^{-N_k+1},$$

implying  $\mathbb{P}_{\Theta^l}[A_l^c] \lesssim \sum_{k=1}^l 2^{-N_k+1}$ . The conditionality on  $A_l$  thus has highly negligible effect.

Next we discuss the evolution of  $\log \nu_2(\mathbf{x}^l)$  and  $\log \mu_2(\mathbf{s}^l)$  under  $A_l$ . The particularity of the initialization He et al. (2015) is to keep stable  $\mathbb{E}_{\Theta^l} \mathbb{E}_{\mathbf{x}, \alpha, c}[(\mathbf{x}_{\alpha, c}^l)^2] = \mathbb{E}_{\Theta^l}[\nu_2(\mathbf{x}^l)]$  and  $\mathbb{E}_{\Theta^l} \mathbb{E}_{\mathbf{x}, \alpha, c}[(\mathbf{s}_{\alpha, c}^l)^2] = \mathbb{E}_{\Theta^l}[\mu_2(\mathbf{s}^l)]$  during propagation. To do so, it enforces  $\mathbb{E}_{\theta^l}[\nu_2(\mathbf{x}^l)] = \nu_2(\mathbf{x}^{l-1})$  and  $\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)] = \mu_2(\mathbf{s}^{l-1})$  such that  $\overline{m}[\nu_2(\mathbf{x}^l)], \overline{m}[\mu_2(\mathbf{s}^l)]$  vanish in Eq. (5) and  $\log \nu_2(\mathbf{x}^l), \log \mu_2(\mathbf{s}^l)$  are subject to a *slow diffusion* with small negative drift terms:  $\underline{m}[\nu_2(\mathbf{x}^l)] < 0, \underline{m}[\mu_2(\mathbf{s}^l)] < 0$ , and small diffusion terms:  $\underline{s}[\nu_2(\mathbf{x}^l)], \underline{s}[\mu_2(\mathbf{s}^l)]$  (see Section E.3 for details).

As evidenced by Fig. 2a and Fig. 2b, the small negative drift and diffusion terms cause slowly decreasing negative expectation and increasing variance of  $\log \nu_2(\mathbf{x}^l), \log \mu_2(\mathbf{s}^l)$ . The decreasing negative expectation and increasing variance of  $\log \nu_2(\mathbf{x}^l), \log \mu_2(\mathbf{s}^l)$  can be seen as opposing forces which compensate in order to keep stable  $\mathbb{E}_{\Theta^l}[\nu_2(\mathbf{x}^l)], \mathbb{E}_{\Theta^l}[\mu_2(\mathbf{s}^l)]$  during propagation. Indeed Fig. 2a and Fig. 2b show that  $\log \nu_2(\mathbf{x}^l), \log \mu_2(\mathbf{s}^l)$  are nearly Gaussian, implying that  $\nu_2(\mathbf{x}^l)$  and  $\mu_2(\mathbf{s}^l)$  are nearly lognormal. And the expectation of a lognormal variable  $\exp(X)$  with  $X \sim \mathcal{N}(\mu, \sigma)$  and  $\mu < 0$  is equal to  $\exp(\mu + \sigma^2/2)$ . Note that the diffusion happens in log-space since layer composition amounts to a multiplicative random effect in real space. Also note that this is a finite-width effect since the terms  $\underline{m}[\nu_2(\mathbf{x}^l)], \underline{m}[\mu_2(\mathbf{s}^l)], \underline{s}[\nu_2(\mathbf{x}^l)], \underline{s}[\mu_2(\mathbf{s}^l)]$  also vanish in the limit of infinite width.

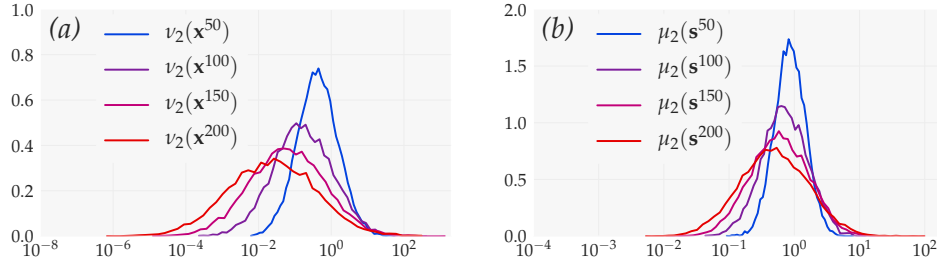


Figure 2: **Slowly diffusing moments of vanilla networks** with  $N_l = 128$  and  $L = 200$  layers. (a) Distributions of  $\nu_2(\mathbf{x}^l)$  for  $l = 50, 100, 150, 200$ . (b) Distributions of  $\mu_2(\mathbf{s}^l)$  for  $l = 50, 100, 150, 200$ . Both  $\nu_2(\mathbf{x}^l)$  and  $\mu_2(\mathbf{s}^l)$  show clear lognormality and are subject to a slow diffusion with small negative drift in log-space.

**Theorem 2. Normalized Sensitivity increments of vanilla networks.** (proof in Appendix E.4) *Denote  $\mathbf{y}^{l,+} = \max(\mathbf{y}^l, 0)$  and  $\mathbf{y}^{l,-} = \max(-\mathbf{y}^l, 0)$ . The dominating term under  $A_l$  in the evolution of  $\chi^l$  is*

$$\delta \chi^l \simeq \exp\left(\overline{m}_{\text{vanilla}}[\chi^l]\right) = \left(1 - \mathbb{E}_{c, \theta^l|A_{l-1}}\left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})}\right]\right)^{-1/2}. \quad (9)$$

**Discussion.** An immediate consequence is that  $\delta\chi^l \gtrsim 1$ , i.e. that normalized sensitivity always increases with depth for random ReLU vanilla networks.

Now let us show that Theorem 1 and Theorem 2 only allow two possibilities of evolution, which both converge to distributional pathology:

- (i) If sensitivity is exploding:  $\chi^l \rightarrow \infty$ , then  $\mu_2(\mathbf{x}^l)/\mu_2(\mathbf{s}^l) \rightarrow 0$  by definition of  $\chi^l$ . Given that  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) = (\mu_2(\mathbf{s}^l)/\nu_2(\mathbf{x}^l))(\mu_2(\mathbf{x}^l)/\mu_2(\mathbf{s}^l))$  and that both  $\mu_2(\mathbf{s}^l)$  and  $\nu_2(\mathbf{x}^l)$  are limited to the slow diffusion of Theorem 1, this qualitatively means that  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l)$  becomes  $\ll 1$ . Rigorously, if  $\chi^l$  has an exponential drift stronger than diffusion and if we assume that  $\nu_2(\mathbf{x}^l)$  and  $\mu_2(\mathbf{s}^l)$  are lognormally-distributed as supported by Fig. 2, then  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \rightarrow 0$  a.s. (proof in Appendix E.5). Consequently, this case leads to double pathologies: *exploding sensitivity* and *zero-dimensional signal*.
- (ii) Otherwise, geometric increments  $\delta\chi^l$  are strongly limited and we may assume  $\delta\chi^l \simeq \exp(\overline{m}[\chi^l]) \rightarrow 1$ . If we further assume that moments of the unit-variance rescaled signal  $\tilde{\mathbf{x}}^l = \mathbf{x}^l/\sqrt{\mu_2(\mathbf{x}^l)}$  are bounded, then Theorem 2 implies the convergence to the pathology of *one-dimensional signal*:  $r_{\text{eff}}(\mathbf{x}^l) \rightarrow 1$  (proof in Appendix E.6), as well as the convergence to neural network *pseudo-linearity* where each layer  $l$  becomes arbitrary well approximated by a linear function (proof in Appendix E.7)

**Experimental verification.** The evolution with depth of vanilla networks is shown in Fig. 3. From the possibilities (i) and (ii), it is case (ii) occurring with subexponential normalized sensitivity:  $\delta\chi^l \simeq \exp(\overline{m}[\chi^l]) \rightarrow 1$ , the convergence to the pathology of one-dimensional signal:  $r_{\text{eff}}(\mathbf{x}^l) \rightarrow 1$ , and to neural network pseudo-linearity. A typical symptom of neural network pseudo-linearity is that signals ‘cross’ the ReLU on the same side. Our analysis offers a novel insight into this *coactivation* phenomenon, previously observed experimentally (Balduzzi et al., 2017).

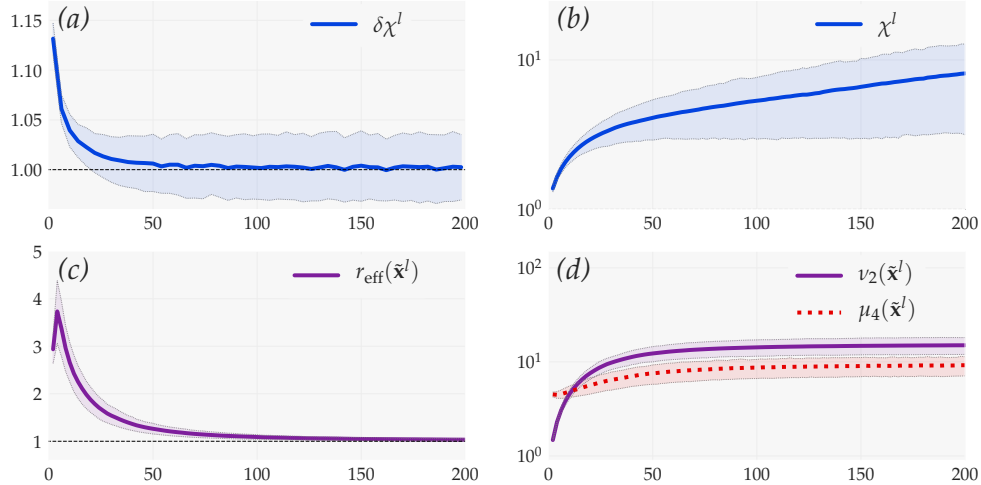


Figure 3: **Pathology of one-dimensional signal for vanilla networks** with  $N_l = 384$  and  $L = 200$  layers. (a) Geometric increments  $\delta\chi^l$  with  $\delta\chi^l \rightarrow 1$ . (b) The growth of the normalized sensitivity is subexponential. (c) The effective rank  $r_{\text{eff}}(\mathbf{x}^l)$  indicates one-dimensional pathology:  $r_{\text{eff}}(\mathbf{x}^l) \rightarrow 1$ . (d) Moments  $\nu_2(\tilde{\mathbf{x}}^l)$  and  $\mu_4(\tilde{\mathbf{x}}^l)$  of the unit-variance rescaled signal  $\tilde{\mathbf{x}}^l = \mathbf{x}^l/\sqrt{\mu_2(\mathbf{x}^l)}$  remain well-behaved due to neural network pseudo-linearity.

## 6 Batch-normalized feedforward nets

Next we incorporate batch normalization (Ioffe & Szegedy, 2015). For simplicity, we only consider the *test mode* which consists in subtracting  $\nu_{1,c}(\mathbf{y}^l)$  and dividing by  $\sqrt{\mu_{2,c}(\mathbf{y}^l)}$  for each channel  $c$  in  $\mathbf{y}^l$ . The propagation is given by



$$\mathbf{y}^l = \boldsymbol{\omega}^l * \mathbf{x}^{l-1} + \boldsymbol{\beta}^l, \quad \mathbf{z}^l = BN(\mathbf{y}^l), \quad \mathbf{x}^l = \phi(\mathbf{z}^l), \quad (10)$$

$$\mathbf{t}^l = \boldsymbol{\omega}^l * \mathbf{s}^{l-1}, \quad \mathbf{u}^l = BN'(\mathbf{y}^l) \odot \mathbf{t}^l, \quad \mathbf{s}^l = \phi'(\mathbf{z}^l) \odot \mathbf{u}^l, \quad (11)$$

where  $BN$  denotes batch normalization. Note that Eq. (10) and (11) explicitly formulate a finer-grained subdivision of three different steps between layers  $l-1$  and  $l$  in the simultaneous propagation of  $(\mathbf{x}^l, \mathbf{s}^l)$ .

**Theorem 3. Normalized Sensitivity increments of batch-normalized feedforward nets.** (proof in Appendix F.1) *The dominating term in the evolution of  $\chi^l$  can be decomposed as the sum of a term  $\overline{m}_{BN}[\chi^l]$  due to batch normalization and a term  $\overline{m}_\phi[\chi^l]$  due to the nonlinearity  $\phi$ :*

$$\exp(\overline{m}_{BN}[\chi^l]) = \left( \frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \mathbb{E}_{\mathbf{c}, \theta^l} \left[ \frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)} \right]^{1/2}, \quad (12)$$

$$\exp(\overline{m}_\phi[\chi^l]) = \left( 1 - 2\mathbb{E}_{\mathbf{c}, \theta^l} [\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-})] \right)^{-1/2}, \quad (13)$$

$$\delta\chi^l \simeq \exp(\overline{m}_{BN-FF}[\chi^l]) = \exp(\overline{m}_{BN}[\chi^l] + \overline{m}_\phi[\chi^l]). \quad (14)$$

**Effect of batch normalization.** The term of Eq. (12) corresponds to the evolution of  $\chi^l$  from  $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$  at layer  $l-1$  to  $(\mathbf{z}^l, \mathbf{u}^l)$  just after  $BN$ . To understand this term qualitatively, the pre-activation tensor  $\mathbf{y}^l$  can be seen as  $N_l$  random projections of  $\mathbf{x}^{l-1}$ , while batch normalization can be seen as an alteration of the magnitude for each projection. Given that batch normalization uses  $\sqrt{\mu_{2,c}(\mathbf{y}^l)}$  as normalization factor, directions of high signal variance are dampened while directions of low signal variance are amplified. This *preferential exploration* of low signal directions naturally deteriorates the signal-to-noise ratio and amplifies  $\chi^l$  due to the *noise factor equivalence* of Eq. (4).

Now let us look directly at the quantity inside the expectation in Eq. (12). By spherical symmetry under standard initialization, geometric increments from  $\mathbf{x}^{l-1}$  to  $\mathbf{y}^l$  for the signal and  $\mathbf{s}^{l-1}$  to  $\mathbf{t}^l$  for the sensitivity have the same expectations  $\mathbb{E}_{\mathbf{c}, \theta^l} [\mu_{2,c}(\mathbf{y}^l)] / \mu_{2,c}(\mathbf{x}^{l-1}) = \mathbb{E}_{\mathbf{c}, \theta^l} [\mu_{2,c}(\mathbf{t}^l)] / \mu_{2,c}(\mathbf{s}^{l-1})$ . On the other hand, the fluctuation of these increments depends on the fluctuation of signal and sensitivity in the  $N_l$  random projections, i.e. on whether directions of signal and sensitivity variances are rare in the ambient space. This effect of *conditioning* precisely depends on whether the effective ranks  $r_{\text{eff}}(\mathbf{x}^{l-1})$  and  $r_{\text{eff}}(\mathbf{s}^{l-1})$  are small compared to the ambient space dimension  $N_{l-1}$ .

If we assume that  $\mathbf{s}^{l-1}$  is well-conditioned, then  $\mu_{2,c}(\mathbf{t}^l)$  has small relative deviation to its expectation  $\mathbb{E}_{\mathbf{c}, \theta^l} [\mu_{2,c}(\mathbf{t}^l)]$  and this term can be treated as a constant. In turn, this implies by convexity of  $x \mapsto 1/x$  that  $\exp(\overline{m}_{BN}[\chi^l]) \gtrsim 1$ . The worse the conditioning of  $\mathbf{x}^{l-1}$ , i.e. the smaller  $r_{\text{eff}}(\mathbf{x}^{l-1})$ , the larger the variance of  $\mu_{2,c}(\mathbf{y}^l) / \mathbb{E}_{\mathbf{c}, \theta^l} [\mu_{2,c}(\mathbf{y}^l)]$  at the denominator and the impact of the convexity. Thus the smaller  $r_{\text{eff}}(\mathbf{x}^{l-1})$  and the larger  $\exp(\overline{m}_{BN}[\chi^l])$ . This argument is strictly valid for the first step of the propagation where sensitivity has perfect conditioning, resulting in  $\exp(\overline{m}_{BN}[\chi^1]) \geq 1$  (proof in Section F.2).

**Effect of the nonlinearity  $\phi$ .** The term of Eq. (13) corresponds to the evolution of  $\chi^l$  from  $(\mathbf{z}^l, \mathbf{u}^l)$  after  $BN$  to  $(\mathbf{x}^l, \mathbf{s}^l)$  after  $\phi$ . The quantity inside the expectation can also be expressed as  $\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) / \mu_{2,c}(\mathbf{z}^l)$  due to  $\mu_{2,c}(\mathbf{z}^l) = 1$  after batch normalization. We then find a very similar expression as Eq. (9) for vanilla networks.

Note that the separate attribution of  $\exp(\overline{m}_{BN}[\chi^l])$  and  $\exp(\overline{m}_\phi[\chi^l])$  to respectively batch normalization and the nonlinearity  $\phi$  is still simplistic. In reality, both are working together since batch normalization maintains the exponential effect of  $\phi$  by centering signal around 0. This is indeed necessary to ensure that, unlike  $\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})$  in Eq. (9) for vanilla networks,  $\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-})$  in Eq. (13) remains non-vanishing for batch-normalized feedforward nets.

**Experimental verification.** In Fig. 4, we confirm experimentally the distributional pathology of *exploding sensitivity*:  $\chi^l \rightarrow 0$ . We also confirm that  $\mathbf{s}^l$  is well-conditioned, while  $\mathbf{x}^l$  becomes ill-conditioned. There is a clear inverse correlation between  $r_{\text{eff}}(\mathbf{x}^l)$  and the effect of batch normalization on  $\delta\chi^l$ .

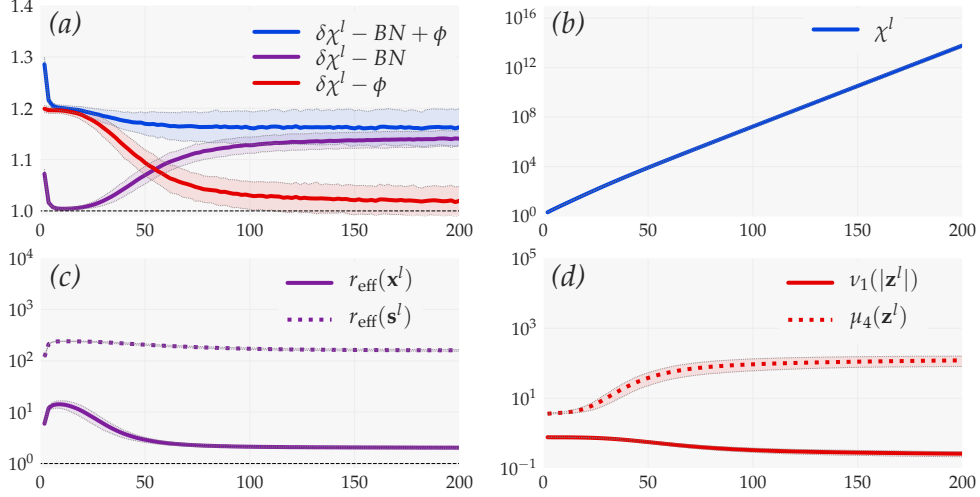


Figure 4: **Pathology of exploding sensitivity for batch-normalized feedforward nets** with  $N_l = 384$  and  $L = 200$  layers. (a) Geometric increments  $\delta\chi^l$  decomposed as the evolution from  $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$  to  $(\mathbf{z}^l, \mathbf{u}^l)$  corresponding to batch normalization, and the evolution from  $(\mathbf{z}^l, \mathbf{u}^l)$  to  $(\mathbf{x}^l, \mathbf{s}^l)$  corresponding to the nonlinearity  $\phi$ . (b) The growth of  $\chi^l$  indicates exploding sensitivity pathology:  $\chi^l \rightarrow \infty$ . (c) Effective ranks  $r_{\text{eff}}(\mathbf{x}^l)$  and  $r_{\text{eff}}(\mathbf{s}^l)$  of signal and sensitivity confirm that sensitivity is well-conditioned. There is a clear inverse correlation between  $r_{\text{eff}}(\mathbf{x}^l)$  and the effect of batch normalization on  $\delta\chi^l$ . (d)  $\mathbf{z}^l$  becomes ill-behaved with small  $\nu_1(|\mathbf{z}^l|)$  and large  $\mu_4(\mathbf{z}^l)$ . This explains the decay of  $\bar{m}_\phi[\chi^l]$  in Eq. (13), and thus of the effect of  $\phi$  on  $\delta\chi^l$ .

Interestingly, the effect of the nonlinearity  $\phi$  becomes subdominant as  $l$  grows. This is explained by the fact that  $\mathbf{z}^l$  becomes heavy-tailed, with large  $\mu_4(\mathbf{z}^l)$  and small  $\nu_1(|\mathbf{z}^l|)$ . Combined with  $\nu_1(\mathbf{z}^{l,+}) \leq \nu_1(|\mathbf{z}^l|)$  and  $\nu_1(\mathbf{z}^{l,-}) \leq \nu_1(|\mathbf{z}^l|)$ , this explains the decay of  $\bar{m}_\phi[\chi^l]$  in Eq. (13).

## 7 Batch-normalized resnets

We finish our exploration of DNN architectures with the incorporation of skip-connections. We now suppose that the width is constant  $N_l = N$ ,<sup>6</sup> and following He et al. (2016b) we adopt the perspective of *pre-activation units*. The propagation is given by

$$(\mathbf{y}^l, \mathbf{t}^l) = (\mathbf{y}^{l-1}, \mathbf{t}^{l-1}) + (\mathbf{y}^{l,H}, \mathbf{t}^{l,H}), \quad (15)$$

$$\mathbf{z}^{l,h} = BN(\mathbf{y}^{l,h-1}), \quad \mathbf{x}^{l,h} = \phi(\mathbf{z}^{l,h}), \quad \mathbf{y}^{l,h} = \omega^{l,h} * \mathbf{x}^{l,h} + \beta^{l,h}, \quad (16)$$

$$\mathbf{u}^{l,h} = BN'(\mathbf{y}^{l,h-1}) \odot \mathbf{t}^{l,h-1}, \quad \mathbf{s}^{l,h} = \phi'(\mathbf{z}^{l,h}) \odot \mathbf{u}^{l,h}, \quad \mathbf{t}^{l,h} = \omega^{l,h} * \mathbf{s}^{l,h}. \quad (17)$$

for  $1 \leq h \leq H$  with  $H$  the number of layers inside residual units, and with  $(\mathbf{y}^{l,0}, \mathbf{t}^{l,0}) \equiv (\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$

If we adopt the additional convention  $(\mathbf{y}^{0,H}, \mathbf{t}^{0,H}) \equiv (\mathbf{y}^0, \mathbf{t}^0)$ , then Eq. (15) can be expanded as

$$(\mathbf{y}^l, \mathbf{t}^l) = \sum_{k=0}^l (\mathbf{y}^{k,H}, \mathbf{t}^{k,H}). \quad (18)$$

For consistency reasons, we redefine the inputs of the propagation as  $(\mathbf{y}^0, \mathbf{t}^0) \equiv (\mathbf{y}, \mathbf{t})$  and the normalized sensitivity and its increments as

$$\chi^l \equiv \left( \frac{\mu_2(\mathbf{t}^l) \mu_2(\mathbf{y}^0)}{\mu_2(\mathbf{y}^l)} \right)^{1/2}, \quad \delta\chi^l \equiv \frac{\chi^l}{\chi^{l-1}}, \quad \chi^{l,h} \equiv \left( \frac{\mu_2(\mathbf{t}^{l,h}) \mu_2(\mathbf{y}^0)}{\mu_2(\mathbf{y}^{l,h})} \right)^{1/2}, \quad \delta\chi^{l,h} \equiv \frac{\chi^{l,h}}{\chi^{l,h-1}}.$$

<sup>6</sup>Again this assumption is only made for simplicity of the analysis. In practice, it holds at least approximately since  $N_l$  is only modified by very few units.

**Theorem 4. Normalized Sensitivity increments of batch-normalized resnets.** (proof in Appendix G.3) Suppose that for all depth  $l$  we can bound the second-order central moments  $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$  and the feedforward increments inside residual units  $\delta_{\min} \lesssim \delta\chi^{l,h} \lesssim \delta_{\max}$ .<sup>7</sup> Denote  $\eta_{\min} = ((\delta_{\min})^{2H} \mu_{2,\min} - \mu_{2,\max})/\mu_{2,\max}$  and  $\eta_{\max} = ((\delta_{\max})^{2H} \mu_{2,\max} - \mu_{2,\min})/\mu_{2,\min}$ , as well as  $\tau_{\min} = \eta_{\min}/2$  and  $\tau_{\max} = \eta_{\max}/2$ . Then there exist positive constants  $C_{\min}, C_{\max} > 0$  such that

$$\left(1 + \frac{\eta_{\min}}{l+1}\right)^{1/2} \lesssim \delta\chi^l \lesssim \left(1 + \frac{\eta_{\max}}{l+1}\right)^{1/2}, \quad (19)$$

$$C_{\min} l^{\tau_{\min}} \lesssim \chi^l \lesssim C_{\max} l^{\tau_{\max}}. \quad (20)$$

**Discussion.** The evolution in Eq. (19) remains exponential inside residual units since  $\eta_{\min}$  and  $\eta_{\max}$  have an exponential dependence in  $H$ . However, it is slowed down by the factor  $1/(l+1)$  between successive residual units. This comes from the *dilution* of the residual path  $(\mathbf{y}^{l,H}, \mathbf{t}^{l,H})$  into the skip-connection path  $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$  with ratio of signal variances  $\mu_2(\mathbf{y}^{l,H})/\mu_2(\mathbf{y}^{l-1})$  scaling as  $1/l$ . If we set  $\mu_{2,\min} = \mu_{2,\max}$  and we remove the dilution effect by replacing  $l+1$  by 1, then Eq. (19) recovers the feedforward evolution with  $(\delta_{\min})^H \lesssim \delta\chi^l \lesssim (\delta_{\max})^H$ . The dilution is clearly visible as a side effect of layer aggregation in Eq. (18): each residual unit  $l$  adds a new term of increased sensitivity but its relative contribution to the aggregation becomes smaller and smaller with  $l$ , so it gets harder and harder for the model to grow  $\chi^l$ .

Given that  $\log(1 + \frac{1}{x}) \simeq \frac{1}{x}$  for  $x \gg 1$  and that  $\int \frac{1}{x} dx = \log x$ , the bounds in Eq. (20) are obtained by integration of the bounds in Eq. (19). A direct consequence of the dilution is thus the power-law evolution of  $\chi^l$  in Eq. (20) instead of the exponential evolution for feedforward nets. Equivalently, when Eq. (20) is written as  $C_{\min} \exp(\tau_{\min} \log l) \lesssim \chi^l \lesssim C_{\max} \exp(\tau_{\max} \log l)$ , the evolution of  $\chi^l$  for resnets is the same as the evolution of  $\chi^{\log l}$  for feedforward nets. In words, the evolution with depth of resnets is the *logarithmic* version of the evolution with depth of feedforward nets. Up to some factor, an evolution from 100, to 1000, and to 10000 layers for resnets is equivalent to an evolution from 20, to 30, and to 40 layers for feedforward nets.

It is noteworthy that the benefits of adding skip connections (though in less sophisticated forms) were already known before the seminal works of He et al. (2016a) and He et al. (2016b). Both Duvenaud et al. (2014) and Neal (1996b) noted that the introduction of direct connections from the input to each layer alleviated distributional pathologies caused by layer composition. More recently, Philipp et al. (2017) also outlined the role of the dilution to avoid exploding gradients.

**Experimental verification.** The slow power-law evolution of  $\chi^l$  is shown in Fig. 5. Notably, the exponent in the power-law fit of Fig. 5b is set to  $\tau = \eta/2 = (\langle \delta\chi^{l,1} \rangle^{2H} - 1)/2$ , with the feedforward increment averaged over the whole evolution  $\langle \delta\chi^{l,1} \rangle$ . This shows that in practice the evolution of  $\chi^l$  is very well described by Eq. (20).

Contrary to batch-normalized feedforward nets, the signal remains well-behaved with: (i) many directions of signal variance preserved in  $r_{\text{eff}}(\mathbf{x}^{l,1})$ ; (ii) close to Gaussian signal distribution as indicated by the moments  $\nu_1(|\mathbf{z}^{l,1}|)$  and  $\mu_4(\mathbf{z}^{l,1})$ . No distributional pathology occurs.

## 8 Discussion and summary

This paper introduced a novel approach to characterize deep convolutional and fully-connected neural networks at initialization. The usual assumptions of infinite width, gaussianity, typical activation patterns, restricted input data were not required, and the commonly used techniques of batch normalization and skip connections were incorporated. The main scope restriction comes from our focus on the ReLU activation function. We leave for future work the application of our methodology to other activation functions – and potentially to recurrent neural networks.

It is interesting to look a posteriori at the commonly used assumption – not required in our analysis – that  $\mathbf{x}^l$  is Gaussian with respect to input data  $\mathbf{x}$ ,  $\alpha$  and model parameters  $\Theta^l$ . Usually, the maximal

<sup>7</sup>The assumption  $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$  is very reasonable since batch normalization controls signal variance at the beginning of layer  $H$ :  $\mu_2(\mathbf{z}^{l,H}) = 1$ .

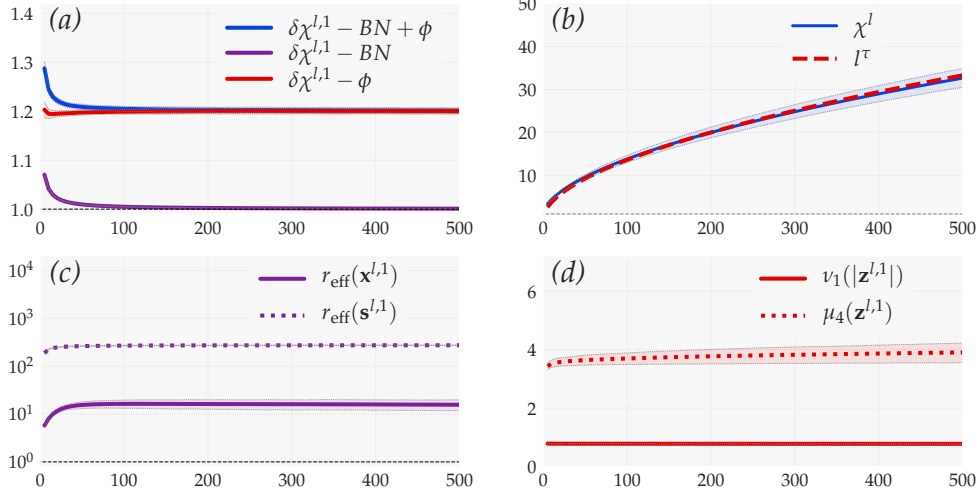


Figure 5: **Well-behaved evolution of batch-normalized resnets** with  $N_l = 384$  and  $L = 500$  residual units of  $H = 2$  layers. (a) Geometric feedforward increments  $\delta\chi^{l,1}$  decomposed as the product of a term corresponding to batch normalization and a term corresponding to  $\phi$ . (b)  $\chi^l$  has power-law evolution. (c) Effective ranks  $r_{\text{eff}}(\mathbf{x}^{l,1})$  and  $r_{\text{eff}}(\mathbf{s}^{l,1})$  of signal and sensitivity indicate that many directions of signal variance are preserved. (d) Moments  $\nu_1(|\mathbf{z}^{l,1}|)$  and  $\mu_4(\mathbf{z}^{l,1})$  indicate that  $\mathbf{z}^{l,1}$  has close to Gaussian distribution.

depth  $L$  is considered constant and the width in the limit  $N_l \rightarrow \infty$  for  $l \leq L$ . Our analysis has reversed this perspective by considering the depth in the limit  $L \rightarrow \infty$  and the width  $N_l$  as bounded. Then the Gaussian assumption becomes invalid in two different contexts:

- (i) In the case of vanilla networks, with e.g. an input  $\varphi(\mathbf{x}, \alpha)$  reduced to a single point of  $\mathbb{R}^{N_0}$  such that  $\varphi(\mathbf{x}^l, \alpha)$  remains a single point of  $\mathbb{R}^{N_l}$ . Given the evolution of Fig. 3, the  $L^2$  norm  $\|\varphi(\mathbf{x}^l, \alpha)\|_2^2 = N_l \nu_2(\mathbf{x}^l)$  becomes lognormally-distributed with diverging variance with respect to  $\Theta^l$ . This means that  $\mathbf{x}^l$  becomes highly non-Gaussian.
- (ii) In the case of batch-normalized feedforward nets, since batch normalization imposes  $\mathbb{E}_{\mathbf{x}, \alpha}[\mathbf{z}_{\alpha, c}^l] = \nu_{1, c}(\mathbf{z}^l) = 0$  and  $\mathbb{E}_{\mathbf{x}, \alpha}[(\mathbf{z}_{\alpha, c}^l)^2] = \mu_{2, c}(\mathbf{z}^l) = 1$ , it follows that  $\mathbb{E}_{\mathbf{x}, \alpha, \Theta^l}[\mathbf{z}_{\alpha, c}^l] = 0$  and  $\mathbb{E}_{\mathbf{x}, \alpha, \Theta^l}[(\mathbf{z}_{\alpha, c}^l)^2] = 1$ . Given the evolution of Fig. 4, the fourth-order standardized moment with respect to  $\mathbf{x}, \alpha, \Theta^l$ , i.e. the kurtosis, of the pre-activation signal  $\mathbf{z}^l$  then becomes  $\mathbb{E}_{\mathbf{x}, \alpha, \Theta^l}[(\mathbf{z}_{\alpha, c}^l)^4] = \mathbb{E}_{\mathbf{x}, \alpha, c, \Theta^l}[(\mathbf{z}_{\alpha, c}^l)^4] = \mathbb{E}_{\Theta^l}[\mu_4(\mathbf{z}^l)] \gg 1$ . Again this means that  $\mathbf{z}^l$  becomes highly non-Gaussian.

Similar observations were made in previous works. Duvenaud et al. (2014) showed that the composition of Gaussian processes, referred as deep Gaussian processes, leads to lognormal and eventually ill-behaved distributions. Matthews et al. (2018) found experimentally that convergence to gaussianity as  $N_l \rightarrow \infty$  becomes slower with respect to  $N_l$  as  $l$  grows. This is explained by the fact that the *affine transform at each layer is an additive process* with respect to the width dimension, but *layer composition is a multiplicative process* with respect to the depth dimension. Qualitatively, the Central Limit Theorem implies that  $\mathbf{x}^l$  becomes normally-distributed as  $N_l \rightarrow \infty$  but log-normally distributed as  $l \rightarrow \infty$ .

Beside from this insight, our findings can be summarized as follows:

- In the case of vanilla networks, the initialization He et al. (2015) limits the evolution of second-order moments of signal and sensitivity. Combined with the limited growth of the normalized sensitivity  $\chi^l$ , this results in the convergence to a distributional pathology where the neural network becomes linear and its signal shrunk to a single dimension:  $r_{\text{eff}}(\mathbf{x}^l) \rightarrow 1$ .
- In the case of batch-normalized feedforward nets, the pathology of exploding sensitivity with  $\chi^l \rightarrow \infty$  has two origins: on the one hand batch normalization which upweights low-signal

pre-activation directions; on the other hand the nonlinearity  $\phi$ . Beside from its direct effect of noise amplification, batch normalization also supports the exponential effect of  $\phi$  on  $\chi^l$  by keeping the signal centered around 0, where  $\phi$  is the most nonlinear.

- Finally in the case of resnets, the normalized sensitivity  $\chi^l$  only grows as a power-law. Equivalently, the evolution with depth of resnets is the logarithmic version of the evolution with depth of feedforward nets. The underlying phenomenon is the dilution of the residual path into the skip-connection path with ratio of signal variances decaying as  $1/l$ . This ingenious mechanism is responsible for breaking the circle of depth multiplicativity which causes distributional pathology for vanilla networks and batch-normalized feedforward nets.

We hope that our findings will open new perspectives in the statistical understanding of deep neural network architectures.

## References

- Anonymous. A mean field theory of batch normalization. In *Submitted to International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SyMDXnCcF7>. under review.
- Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pp. 254–263, 2018. URL <http://proceedings.mlr.press/v80/arora18b.html>.
- David Balduzzi, Marcus Frean, Lennox Leary, J. P. Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pp. 342–350, 2017. URL <http://proceedings.mlr.press/v70/balduzzii17b.html>.
- Minmin Chen, Jeffrey Pennington, and Samuel S. Schoenholz. Dynamical isometry and a mean field theory of rnns: Gating enables signal propagation in recurrent neural networks. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pp. 872–881, 2018. URL <http://proceedings.mlr.press/v80/chen18i.html>.
- Marco Chiani, Davide Dardari, and Marvin K. Simon. New exponential bounds and approximations for the computation of error probability in fading channels. *IEEE Trans. Wireless Communications*, 2(4):840–845, 2003. doi: 10.1109/TWC.2003.814350. URL <https://doi.org/10.1109/TWC.2003.814350>.
- Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 29*, pp. 2253–2261. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6427-toward-deeper-understanding-of-neural-networks-the-power-of-initialization-and-a-dual-view-on-expressivity.pdf>.
- Richard Durrett. *Probability: theory and examples*. Duxbury Press, second edition, 1996. ISBN 0-534-24318-5. URL [https://services.math.duke.edu/~rtd/PTE/PTE4\\_1.pdf](https://services.math.duke.edu/~rtd/PTE/PTE4_1.pdf).
- David Duvenaud, Oren Rippel, Ryan Adams, and Zoubin Ghahramani. Avoiding pathologies in very deep networks. In Samuel Kaski and Jukka Corander (eds.), *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33 of *Proceedings of Machine Learning Research*, pp. 202–210, Reykjavik, Iceland, 22–25 Apr 2014. PMLR. URL <http://proceedings.mlr.press/v33/duvenaud14.html>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV), ICCV ’15*, pp. 1026–1034, Washington, DC, USA, 2015. IEEE Computer Society. ISBN 978-1-4673-8391-2. doi: 10.1109/ICCV.2015.123. URL <http://dx.doi.org/10.1109/ICCV.2015.123>.



- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778, 2016a. doi: 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pp. 630–645, 2016b. doi: 10.1007/978-3-319-46493-0\_38. URL [https://doi.org/10.1007/978-3-319-46493-0\\_38](https://doi.org/10.1007/978-3-319-46493-0_38).
- Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Comput.*, 9(1):1–42, January 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.1.1. URL <http://dx.doi.org/10.1162/neco.1997.9.1.1>.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pp. 448–456, 2015. URL <http://jmlr.org/proceedings/papers/v37/lofffe15.html>.
- Nitish S. Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping T. P. Tang. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. *ArXiv e-prints*, September 2016. URL <http://adsabs.harvard.edu/abs/2016arXiv160904836S>.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep Neural Networks as Gaussian Processes. *ArXiv e-prints*, October 2017. URL <http://adsabs.harvard.edu/abs/2017arXiv171100165L>.
- Lu Lu, Yanhui Su, and George E. Karniadakis. Collapse of deep and narrow neural nets. *CoRR*, abs/1808.04947, 2018. URL <http://arxiv.org/abs/1808.04947>.
- Alexander G. de G. Matthews, Mark Rowland, Jiri Hron, Richard E. Turner, and Zoubin Ghahramani. Gaussian Process Behaviour in Wide Deep Neural Networks. *ArXiv e-prints*, April 2018. URL <http://adsabs.harvard.edu/abs/2018arXiv180411271M>.
- Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance of single directions for generalization. *ArXiv e-prints*, March 2018. URL <http://adsabs.harvard.edu/abs/2018arXiv180306959M>.
- Radford M. Neal. *Priors for Infinite Networks*, pp. 29–53. Springer New York, New York, NY, 1996a. ISBN 978-1-4612-0745-0. doi: 10.1007/978-1-4612-0745-0\_2. URL [https://doi.org/10.1007/978-1-4612-0745-0\\_2](https://doi.org/10.1007/978-1-4612-0745-0_2).
- Radford M. Neal. *Bayesian Learning for Neural Networks*. Springer-Verlag, Berlin, Heidelberg, 1996b. ISBN 0387947248.
- Behnam Neyshabur, Srinadh Bhojanapalli, David Mcallester, and Nati Srebro. Exploring generalization in deep learning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 5947–5956. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7176-exploring-generalization-in-deep-learning.pdf>.
- Roman Novak, Yasaman Bahri, Daniel A. Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. volume abs/1802.08760, 2018. URL <http://arxiv.org/abs/1802.08760>.
- Jeffrey Pennington, Samuel S. Schoenholz, and Surya Ganguli. Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pp. 4788–4798, 2017. URL <http://papers.nips.cc/paper/7064-resurrecting-the-sigmoid-in-deep-learning-through-dynamical-isometry-theory-and-practice>.
- George Philipp and Jaime G. Carbonell. The nonlinearity coefficient - predicting overfitting in deep neural networks. *CoRR*, abs/1806.00179, 2018. URL <http://arxiv.org/abs/1806.00179>.



- George Philipp, Dawn Song, and Jaime G. Carbonell. Gradients explode - deep networks are shallow - resnet explained. *CoRR*, abs/1712.05577, 2017. URL <http://arxiv.org/abs/1712.05577>.
- Ben Poole, Subhaneil Lahiri, Maithreyi Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 3360–3368, 2016. URL <http://papers.nips.cc/paper/6322-exponential-expressivity-in-deep-neural-networks-through-transient-chaos>.
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2847–2854, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL <http://proceedings.mlr.press/v70/raghu17a.html>.
- Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. *CoRR*, abs/1611.01232, 2016. URL <http://arxiv.org/abs/1611.01232>.
- Samuel L. Smith and Quoc V. Le. A bayesian perspective on generalization and stochastic gradient descent. *CoRR*, abs/1710.06451, 2017. URL <http://arxiv.org/abs/1710.06451>.
- Jure Sokolic, Raja Giryes, Guillermo Sapiro, and Miguel R. D. Rodrigues. Robust large margin deep neural networks. *IEEE Trans. Signal Processing*, 65(16):4265–4280, 2017. doi: 10.1109/TSP.2017.2708039. URL <https://doi.org/10.1109/TSP.2017.2708039>.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. volume abs/1011.3027, 2010. URL <http://arxiv.org/abs/1011.3027>.
- Christopher K. I. Williams. Computing with infinite networks. In M. C. Mozer, M. I. Jordan, and T. Petsche (eds.), *Advances in Neural Information Processing Systems 9*, pp. 295–301. MIT Press, 1997. URL <http://papers.nips.cc/paper/1197-computing-with-infinite-networks.pdf>.
- 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, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/xiao18a.html>.
- Greg Yang and Samuel S. Schoenholz. Mean field residual networks: On the edge of chaos. *CoRR*, abs/1712.08969, 2017. URL <http://arxiv.org/abs/1712.08969>.

## A Details of the experiments

All four experiments of Fig. 2, 3, 4, 5 were made on `cifar10` with a random initial convolution of stride 2 reducing the spatial dimension to  $n = 16$ . In each case, we considered the convolutional extent  $K_l = 3$  and periodic boundary conditions. A few experiments with approximately statistics-preserving boundary conditions such as *symmetric mirroring* indicated qualitatively equivalent behaviour.

In Fig. 2, we considered the width  $N_l = 128$  and the total depth  $L = 200$ . For 10000 random realizations of He et al. (2015), we randomly sampled 1024 images and computed the evolution with depth of  $\nu_2(\mathbf{x}^l)$  and  $\mu_2(\mathbf{s}^l)$ . The distributions of  $\nu_2(\mathbf{x}^l)$  and  $\mu_2(\mathbf{s}^l)$  were estimated with the empirical distributions on the 10000 realizations. The limited width – slightly less than standard values – had two purposes: (i) outlining the diffusion process, stronger for smaller width; (ii) limiting computation time in order to gather more realizations.

In Fig. 3, 4, 5, we increased the width to a more realistic value  $N_l = 384$  and considered 1000 realization with single batches of 64 randomly sampled images. For each realization we computed the empirical expectation over all realizations as well as  $1\sigma$  intervals shown as shaded areas. Decreasing the number of inputs for each realization enabled to reduce computation time, mostly penalized by the computation of  $r_{\text{eff}}(\mathbf{x}^l)$  and  $r_{\text{eff}}(\mathbf{s}^l)$  involving an eigenvalue decomposition. We found that this reduction had very little impact for vanilla networks. As for batch-normalized feedforward nets and batch-normalized resnets, it enabled to match the usual setup of batch normalization in training mode.

We plan to release `jupyter` notebooks to enable replication of our results upon publication.

## B Complementary definitions and notations

**Receptive field mapping.** Here we temporarily need to handle the mechanics of convolution. Let us consider the convolution at layer  $l$  of an input  $\mathbf{v} \in \mathbb{R}^{n \times \dots \times n \times N_{l-1}}$  from layer  $l - 1$ . The output feature map of the convolution  $(\omega^l * \mathbf{v})_{\alpha,:}$  at position  $\alpha \in \{1, \dots, n\}^d$  is obtained by the application of the convolution kernel  $\omega^l$  over a local input region of size  $(K_l^d N_{l-1})$ , with  $K_l^d$  the spatial extent and  $N_{l-1}$  the extent in the channel dimension. The local input region is called the *receptive field* of  $(\omega^l * \mathbf{v})$  at spatial position  $\alpha$ .

The *receptive field mapping*  $RF(\cdot)$  associates an input  $\mathbf{v}$  from layer  $l - 1$  to the tensor  $RF(\mathbf{v}) \in \mathbb{R}^{n \times \dots \times n \times K_l^d N_{l-1}}$ , defined such that  $RF(\mathbf{v})_{\alpha,:}$  is the reshaped vectorial form of the receptive field of  $(\omega^l * \mathbf{v})$  at spatial position  $\alpha$ . We denote  $R_l = K_l^d N_{l-1}$  the dimensionality of  $RF(\mathbf{v})_{\alpha,:}$  and  $\mathcal{I}_c^l$  the set of indices in  $RF(\mathbf{v})_{\alpha,:}$  corresponding to elements in channel  $c$  in the input  $\mathbf{v}$ . Strictly speaking,  $RF$  depend on  $l$ , but this is implied by the argument so we write  $RF$  for simplicity.

**Receptive field vectors.** The *receptive field vector*  $\rho$  and *centered receptive field vector*  $\hat{\rho}$  associated with an input  $\mathbf{v}$  from layer  $l$  are random vectors which depend on  $\mathbf{v}$ ,  $\alpha$  such that

$$\rho(\mathbf{v}, \alpha) \equiv RF(\mathbf{v})_{\alpha,:} \quad \text{and} \quad \hat{\rho}(\mathbf{v}, \alpha) \equiv \rho(\mathbf{v}, \alpha) - \mathbb{E}_{\mathbf{v}, \alpha}[\rho(\mathbf{v}, \alpha)],$$

where we kept the same denotation  $\mathbf{v}$  for the variable in the expectation, as a slight abuse of notation. Again  $\rho$  and  $\hat{\rho}$  are strictly speaking dependent on  $l$ , but this is implied by the argument.

**Statistics-preserving property.**  $RF$  is *statistics-preserving* with respect to  $\mathbf{v}$  if for any channel  $c$  and any index  $i_c \in \mathcal{I}_c^l$ , the random variables  $RF(\mathbf{v})_{\alpha, i_c}$  and  $\mathbf{v}_{\alpha, c}$  which depend on  $\mathbf{v}$ ,  $\alpha$  have the same distribution:  $RF(\mathbf{v})_{\alpha, i_c} \sim_{\mathbf{v}, \alpha} \mathbf{v}_{\alpha, c}$ .

**Equation of Propagation.** Using the definition of  $RF$ , the affine transformation from the receptive field  $RF(\mathbf{x}^{l-1})_{\alpha,:}$  to the feature map in the next layer  $\mathbf{y}_{\alpha,:}^l$  can be written as

$$\mathbf{y}_{\alpha,:}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1})_{\alpha,:} + \mathbf{b}^l, \quad (21)$$

where  $\mathbf{W}^l \in \mathbb{R}^{N_l \times R_l}$  is the suitably reshaped matricial form of  $\omega^l$ . To lighten notation, we write  $\mathbf{y}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l$  as a short for the affine transformation of Eq. (21) occurring at all spatial

positions  $\alpha$ . We have the following equivalence between the notations with receptive field and with convolutions:

$$\mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l = \omega^l * \mathbf{x}^{l-1} + \beta^l.$$

For vanilla networks, the simultaneous propagation of  $\mathbf{x}^l$  and  $\mathbf{s}^l$  is then written as

$$\mathbf{y}^l = \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l, \quad \mathbf{x}^l = \phi(\mathbf{y}^l), \quad (22)$$

$$\mathbf{t}^l = \mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \mathbf{s}^l = \phi'(\mathbf{y}^l) \odot \mathbf{t}^l. \quad (23)$$

For batch-normalized feedforward nets, the simultaneous propagation of  $\mathbf{x}^l$  and  $\mathbf{s}^l$  is written as

$$\begin{aligned} \mathbf{y}^l &= \mathbf{W}^l RF(\mathbf{x}^{l-1}) + \mathbf{b}^l, \quad \mathbf{z}^l = BN(\mathbf{y}^l), \quad \mathbf{x}^l = \phi(\mathbf{z}^l), \\ \mathbf{t}^l &= \mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \mathbf{u}^l = BN'(\mathbf{y}^l) \odot \mathbf{t}^l, \quad \mathbf{s}^l = \phi'(\mathbf{z}^l) \odot \mathbf{u}^l. \end{aligned}$$

**Gramian and Covariance.** The *Gramian matrix*  $\mathbf{G}$  and *covariance matrix*  $\mathbf{C}$  associated with random vector  $\mathbf{v} \in \mathbb{R}^N$  are defined as

$$\mathbf{G}[\mathbf{v}] \equiv \mathbb{E}_{\mathbf{v}}[\mathbf{v}\mathbf{v}^T] \quad \text{and} \quad \mathbf{C}[\mathbf{v}] \equiv \mathbb{E}_{\mathbf{v}}[\mathbf{v}\mathbf{v}^T] - \mathbb{E}_{\mathbf{v}}[\mathbf{v}]\mathbb{E}_{\mathbf{v}}[\mathbf{v}]^T.$$

Note that gramian matrix and covariance matrix of feature map vectors and receptive field vectors are related by

$$\mathbf{G}[\hat{\varphi}(\mathbf{v}, \alpha)] = \mathbf{C}[\hat{\varphi}(\mathbf{v}, \alpha)] = \mathbf{C}[\varphi(\mathbf{v}, \alpha)], \quad \mathbf{G}[\hat{\rho}(\mathbf{v}, \alpha)] = \mathbf{C}[\hat{\rho}(\mathbf{v}, \alpha)] = \mathbf{C}[\rho(\mathbf{v}, \alpha)].$$

**Symmetric propagation for vanilla networks.** We define additional tensors obtained by *symmetric propagation* at each layer  $l$ . In the case of vanilla networks they are given by:

$$\begin{aligned} \bar{\mathbf{y}}^l &= -\mathbf{W}^l RF(\mathbf{x}^{l-1}) - \mathbf{b}^l, \quad \bar{\mathbf{x}}^l = \phi(\bar{\mathbf{y}}^l), \\ \bar{\mathbf{t}}^l &= -\mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \bar{\mathbf{s}}^l = \phi'(\bar{\mathbf{y}}^l) \odot \bar{\mathbf{t}}^l. \end{aligned}$$

Under standard initialization, tensor moments have the *same distribution with respect to  $\theta^l$  for both propagations*. Furthermore  $\forall \alpha, c, \mathbf{x}_{\alpha,c}^l + \bar{\mathbf{x}}_{\alpha,c}^l = |\mathbf{y}_{\alpha,c}^l|$  and  $(\mathbf{x}_{\alpha,c}^l)^2 + (\bar{\mathbf{x}}_{\alpha,c}^l)^2 = (\mathbf{y}_{\alpha,c}^l)^2$  since  $\mathbf{x}_{\alpha,c}^l \bar{\mathbf{x}}_{\alpha,c}^l = 0$ . We deduce that

$$\forall c: \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l) = \nu_{2,c}(\mathbf{y}^l). \quad (24)$$

Now we consider the second-order moments of the sensitivity tensor. We have the following identity:

$$\begin{aligned} (\mathbf{s}_{\alpha,c}^l)^2 + (\bar{\mathbf{s}}_{\alpha,c}^l)^2 &= (\mathbf{t}_{\alpha,c}^l)^2 \phi'(\mathbf{y}_{\alpha,c}^l)^2 + (\bar{\mathbf{t}}_{\alpha,c}^l)^2 \phi'(\bar{\mathbf{y}}_{\alpha,c}^l)^2 = (\mathbf{t}_{\alpha,c}^l)^2 [\phi'(\mathbf{y}_{\alpha,c}^l)^2 + \phi'(\bar{\mathbf{y}}_{\alpha,c}^l)^2] \\ &= (\mathbf{t}_{\alpha,c}^l)^2, \end{aligned} \quad (25)$$

where Eq. (25) is obtained using  $\mathbf{y}_{\alpha,c}^l = -\bar{\mathbf{y}}_{\alpha,c}^l$  and the convention  $\phi'(0) \equiv 1/2$ . Since  $\mathbf{s}^l, \bar{\mathbf{s}}^l, \mathbf{t}^l$  are centered, it follows that

$$\forall c: \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\bar{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\bar{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{t}^l) = \mu_{2,c}(\mathbf{t}^l). \quad (26)$$

**Symmetric propagation for batch-normalized feedforward nets.** For batch-normalized feedforward nets, the symmetric propagation at each layer  $l$  is given by

$$\bar{\mathbf{y}}^l = -\mathbf{W}^l RF(\mathbf{x}^{l-1}) - \mathbf{b}^l, \quad \bar{\mathbf{z}}^l = BN(\bar{\mathbf{y}}^l), \quad \bar{\mathbf{x}}^l = \phi(\bar{\mathbf{z}}^l), \quad (27)$$

$$\bar{\mathbf{t}}^l = -\mathbf{W}^l RF(\mathbf{s}^{l-1}), \quad \bar{\mathbf{u}}^l = BN'(\bar{\mathbf{y}}^l) \odot \bar{\mathbf{t}}^l, \quad \bar{\mathbf{s}}^l = \phi'(\bar{\mathbf{z}}^l) \odot \bar{\mathbf{u}}^l. \quad (28)$$

$BN$  in Eq. (27) and (28) uses the statistics of  $\bar{\mathbf{y}}^l$  such that, under standard initialization, tensor moments have the *same distribution with respect to  $\theta^l$  for both propagations*. We then simply have

$$\bar{\mathbf{z}}^l = -\mathbf{z}^l, \quad \bar{\mathbf{u}}^l = -\mathbf{u}^l. \quad (29)$$

The exact same analysis as before gives

$$\forall c: \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l) = \nu_{2,c}(\mathbf{z}^l), \quad (30)$$

$$\forall c: \mu_{2,c}(\mathbf{s}^l) + \mu_{2,c}(\bar{\mathbf{s}}^l) = \mu_{2,c}(\mathbf{u}^l). \quad (31)$$

## C Statistics-preserving property

### C.1 Case of periodic boundary conditions and constant spatial extent $n$

**Lemma 1.** *If convolutions have periodic boundary conditions and the global spatial extent  $n$  is constant, then  $RF$  is statistics-preserving with respect to any input  $\mathbf{v}$ .*

**Proof.** Fix a channel  $c$  in  $\mathbf{v}$ , an index  $i_c \in \mathcal{I}_c^l$ , and consider the tensors  $\mathbf{v}_{:,c}$  and  $RF(\mathbf{v})_{:,i_c} \in \mathbb{R}^{n \times \dots \times n}$ . The index  $i_c$  corresponds to a given convolution kernel position  $\kappa \in \{1, \dots, K_l\}^d$ . Under periodic boundary conditions, this fixed kernel position implies that each position  $\alpha$  in  $RF(\mathbf{v})_{\alpha,i_c}$  originates from a different position  $\alpha'$  in the tensor  $\mathbf{v}_{\alpha',c}$ . Therefore the index mapping  $f : \alpha \rightarrow \alpha'$  from  $\{1, \dots, n\}^d$  to  $\{1, \dots, n\}^d$  is bijective. We then have  $RF(\mathbf{v})_{\alpha,i_c} = \mathbf{v}_{f(\alpha),c} \sim_{\alpha} \mathbf{v}_{\alpha,c}$  when  $\alpha$  is considered as random and  $\mathbf{v}$  as given. In turn, this implies that  $RF(\mathbf{v})_{\alpha,i_c} \sim_{\mathbf{v},\alpha} \mathbf{v}_{\alpha,c}$ , when both  $\mathbf{v}$  and  $\alpha$  are considered as random.  $\square$

**Proposition 2.** *If convolutions have periodic boundary conditions and the global spatial extent  $n$  is constant, then  $RF$  is statistics-preserving with respect to  $\mathbf{x}^{l-1}$  and  $\mathbf{s}^{l-1}$ .*

**Proof.** This follows immediately from Lemma 1.  $\square$

**Corollary 3.** *For any  $c$  and  $i_c \in \mathcal{I}_c^l$ , we have  $\rho(\mathbf{x}^{l-1}, \alpha)_{i_c} \sim_{\mathbf{x},\alpha} \varphi(\mathbf{x}^{l-1}, \alpha)_c$  and  $\rho(\mathbf{s}^{l-1}, \alpha)_{i_c} \sim_{\mathbf{s},\alpha} \varphi(\mathbf{s}^{l-1}, \alpha)_c$ . Since the cardinality  $|\mathcal{I}_c^l| = K_l^d$  is the same for all channels  $c$ , it follows that*

$$\begin{aligned} \nu_2(\mathbf{x}^{l-1}) &= \frac{1}{N_{l-1}} \text{Tr } \mathbf{G}[\varphi(\mathbf{x}^{l-1}, \alpha)] = \frac{1}{R_l} \text{Tr } \mathbf{G}[\rho(\mathbf{x}^{l-1}, \alpha)], \\ \mu_2(\mathbf{x}^{l-1}) &= \frac{1}{N_{l-1}} \text{Tr } \mathbf{C}[\varphi(\mathbf{x}^{l-1}, \alpha)] = \frac{1}{R_l} \text{Tr } \mathbf{C}[\rho(\mathbf{x}^{l-1}, \alpha)], \\ \nu_2(\mathbf{s}^{l-1}) &= \mu_2(\mathbf{s}^{l-1}) = \frac{1}{N_{l-1}} \text{Tr } \mathbf{C}[\varphi(\mathbf{s}^{l-1}, \alpha)] = \frac{1}{R_l} \text{Tr } \mathbf{C}[\rho(\mathbf{s}^{l-1}, \alpha)]. \end{aligned}$$

### C.2 Relaxing the assumptions on boundary conditions and constant spatial extent

In this section, we detail possible relaxations of the assumptions on boundary conditions and constant spatial extent  $n$ . The global spatial extent is denoted  $n_l$  when it is not constant.

#### C.2.1 Case of stationary inputs

**Periodic extension.** The *periodic extension*  $\tilde{\mathbf{v}}$  of a random tensor  $\mathbf{v} \in \mathbb{R}^{n \times \dots \times n \times N}$  is defined as

$$\tilde{\mathbf{v}}_{\alpha_1+k_1n, \dots, \alpha_d+k_dn, c} = \mathbf{v}_{\alpha_1, \dots, \alpha_d, c},$$

with  $(k_1, \dots, k_d) \in \mathbb{Z}^d$ , and where  $\alpha_1, \dots, \alpha_d$  are the  $d$  components of the spatial position  $\alpha = (\alpha_1, \dots, \alpha_d) \in \{1, \dots, n\}^d$ .

**Stationarity.** The distribution of a random vector  $\tilde{\mathbf{v}}$  is defined as *stationary* if, for any  $k$ , any configuration of spatial positions  $(\alpha_1, \dots, \alpha_k)$  and channels  $(c_1, \dots, c_k)$ , the joint distribution of  $\tilde{\mathbf{v}}_{\alpha_1+c_1, \dots, \tilde{\mathbf{v}}_{\alpha_k+c_k}$  does not depend on  $\alpha$ .

**Lemma 4.** *If convolutions have periodic boundary conditions and the inputs  $\mathbf{x}$  and  $\mathbf{s}$  have stationary periodic extensions  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{s}}$ , then their periodic extension  $\tilde{\mathbf{x}}^l$  and  $\tilde{\mathbf{s}}^l$  remain stationary during propagation.*

**Proof.** First note that convolutions with periodic boundary conditions followed by periodic extensions on  $\mathbf{v}$  are equivalent to convolutions on  $\tilde{\mathbf{v}}$ . This is also the case for componentwise operations such as batch normalization, nonlinear activation and their derivatives. So we can restrict our attention to periodic extensions.

Consider a given spatial shift  $\alpha$  and define the translation operator  $T_\alpha$ , such that  $T_\alpha(\tilde{\mathbf{u}}) = \tilde{\mathbf{v}}$  with  $\forall \alpha', c : \tilde{\mathbf{v}}_{\alpha', c} = \tilde{\mathbf{u}}_{\alpha + \alpha', c}$ . It is easy to see that  $T_\alpha$  commutes with convolutions as well as componentwise operations such as batch normalization, nonlinear activation and their derivatives. It follows that  $T_\alpha$  commutes with the mapping  $\Phi_l$  defined as  $(\tilde{\mathbf{x}}^l, \tilde{\mathbf{s}}^l) = \Phi_l(\tilde{\mathbf{x}}, \tilde{\mathbf{s}})$ , and thus

$$T_\alpha(\tilde{\mathbf{x}}^l, \tilde{\mathbf{s}}^l) = \Phi_l(T_\alpha(\tilde{\mathbf{x}}, \tilde{\mathbf{s}})), \quad (32)$$

where we adopted the notation  $T_\alpha(\tilde{\mathbf{u}}, \tilde{\mathbf{v}}) = (T_\alpha(\tilde{\mathbf{u}}), T_\alpha(\tilde{\mathbf{v}}))$ . Now consider a given  $k$  and a given configuration of spatial positions  $(\alpha_1, \dots, \alpha_k)$  and channels  $(c_1, \dots, c_k)$  in  $\tilde{\mathbf{x}}^l$  and  $\tilde{\mathbf{s}}^l$ . Due to limited convolutional spatial extent and due to Eq. (32), there exist a function  $\Phi$  and a configuration of spatial positions  $(\alpha'_1, \dots, \alpha'_{k'})$  and channels  $(c'_1, \dots, c'_{k'})$  in  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{s}}$  such that we can write

$$\tilde{\mathbf{x}}_{\alpha_1, c_1}^l, \dots, \tilde{\mathbf{x}}_{\alpha_k, c_k}^l = \Phi(\tilde{\mathbf{x}}_{\alpha'_1, c'_1}, \dots, \tilde{\mathbf{x}}_{\alpha'_{k'}, c'_{k'}}), \quad (33)$$

$$\tilde{\mathbf{x}}_{\alpha + \alpha_1, c_1}^l, \dots, \tilde{\mathbf{x}}_{\alpha + \alpha_k, c_k}^l = \Phi(\tilde{\mathbf{x}}_{\alpha + \alpha'_1, c'_1}, \dots, \tilde{\mathbf{x}}_{\alpha + \alpha'_{k'}, c'_{k'}}), \quad (34)$$

$$\tilde{\mathbf{s}}_{\alpha_1, c_1}^l, \dots, \tilde{\mathbf{s}}_{\alpha_k, c_k}^l = \Phi(\tilde{\mathbf{s}}_{\alpha'_1, c'_1}, \dots, \tilde{\mathbf{s}}_{\alpha'_{k'}, c'_{k'}}), \quad (35)$$

$$\tilde{\mathbf{s}}_{\alpha + \alpha_1, c_1}^l, \dots, \tilde{\mathbf{s}}_{\alpha + \alpha_k, c_k}^l = \Phi(\tilde{\mathbf{s}}_{\alpha + \alpha'_1, c'_1}, \dots, \tilde{\mathbf{s}}_{\alpha + \alpha'_{k'}, c'_{k'}}). \quad (36)$$

By stationarity of  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{s}}$ , the terms on the right-hand sides of Eq. (33) and (34) have the same distribution, and the terms on the right-hand sides of Eq. (35) and (36) have the same distribution. It follows that the terms on the left-hand sides of both pairs of equations have the same distribution, meaning that  $\tilde{\mathbf{x}}^l$  and  $\tilde{\mathbf{s}}^l$  are stationary.  $\square$

**Lemma 5.** *If convolutions have periodic boundary conditions, then RF is statistics-preserving with respect to any input  $\mathbf{v}$  which has stationary periodic extension  $\tilde{\mathbf{v}}$  for any global spatial extent  $n_l$ .*

**Proof.** If  $\tilde{\mathbf{v}}$  has stationary distribution, it means in particular that for any channel  $c$ , the distribution of  $\mathbf{v}_{\alpha, c}$  is the same for all  $\alpha \in \{1, \dots, n_{l-1}\}^d$ . Fix a channel  $c$  and an index  $i_c \in \mathcal{I}_c^l$  corresponding to a given convolution kernel position  $\kappa \in \{1, \dots, K_l\}^d$ . A given position  $\alpha$  in the receptive field tensor  $RF(\mathbf{v})_{\alpha, i_c}$  then corresponds to a given position  $\alpha'$  in the original tensor, such that  $RF(\mathbf{v})_{\alpha, i_c} = \mathbf{v}_{\alpha', c}$ . Since the distribution of  $\mathbf{v}_{\alpha', c}$  does not depend on  $\alpha'$ , it follows that  $RF(\mathbf{v})_{\alpha, i_c} \sim_{\mathbf{v}, \alpha'} \mathbf{v}_{\alpha', c}$  for given  $\alpha$  and random  $\alpha'$ , and thus that  $RF(\mathbf{v})_{\alpha, i_c} \sim_{\mathbf{v}, \alpha} \mathbf{v}_{\alpha, c}$  for random  $\alpha$ .  $\square$

**Proposition 6.** *If convolutions have periodic boundary conditions and the input  $\mathbf{x}$  has stationary periodic extension  $\tilde{\mathbf{x}}$ , then RF is statistics-preserving with respect to  $\mathbf{x}^l$  and  $\mathbf{s}^l$ , for any global spatial extent  $n_l$ .*

**Proof.** This follows from Lemmas 4 and 5, and from the fact that the input sensitivity tensor  $\mathbf{s}$  has stationary periodic extension  $\tilde{\mathbf{s}}$  due to its definition as a white noise tensor with independent and identically distributed components.  $\square$

### C.2.2 Case $n_l \gg K_l$

**Proposition 7.** *If the convolution stride is one in most layers (i.e.  $n_{l-1} = n_l$  in most layers) and the global spatial extent is much larger than the convolutional spatial extent  $n_l \gg K_l$  in most layers, then RF is approximately statistics-preserving with respect to  $\mathbf{x}^{l-1}$  and  $\mathbf{s}^{l-1}$ , for any boundary conditions.*

**Proof.** Fix a layer  $l-1$  such that  $n_{l-1} = n_l$  and  $n_l \gg K_l$ . Denote  $RF^{(p)}$  the receptive field mapping at layer  $l$  associated with periodic boundary conditions. Since  $n_{l-1} = n_l \gg K_l$  the receptive fields  $RF(\mathbf{x}^{l-1})_{\alpha, :}$ ,  $RF(\mathbf{s}^{l-1})_{\alpha, :}$  and  $RF^{(p)}(\mathbf{x}^{l-1})_{\alpha, :}$ ,  $RF^{(p)}(\mathbf{s}^{l-1})_{\alpha, :}$  do not intersect boundary regions for most  $\alpha$ , implying  $RF(\mathbf{x}^{l-1})_{\alpha, :} = RF^{(p)}(\mathbf{x}^{l-1})_{\alpha, :}$  and  $RF(\mathbf{s}^{l-1})_{\alpha, :} = RF^{(p)}(\mathbf{s}^{l-1})_{\alpha, :}$  for most  $\alpha$ . This implies for any index  $i_c$  that  $P_{\mathbf{x}, \alpha}[RF(\mathbf{x}^{l-1})_{\alpha, i_c}] \simeq P_{\mathbf{x}, \alpha}[RF^{(p)}(\mathbf{x}^{l-1})_{\alpha, i_c}]$  and  $P_{\mathbf{s}, \alpha}[RF(\mathbf{s}^{l-1})_{\alpha, i_c}] \simeq P_{\mathbf{s}, \alpha}[RF^{(p)}(\mathbf{s}^{l-1})_{\alpha, i_c}]$ .

Since  $RF^{(p)}$  is statistics-preserving with respect to  $\mathbf{x}^{l-1}$  and  $\mathbf{s}^{l-1}$  by Lemma 1, it follows that for any channel  $c$  and index  $i_c \in \mathcal{I}_c^l$ , we have  $P_{\mathbf{x}, \alpha}[RF^{(p)}(\mathbf{x}^{l-1})_{\alpha, i_c}] = P_{\mathbf{x}, \alpha}[\mathbf{x}_{\alpha, c}^{l-1}]$

and  $P_{\mathbf{x},\mathbf{s},\alpha}[RF^{(p)}(\mathbf{s}^{l-1})_{\alpha,i_c}] = P_{\mathbf{x},\mathbf{s},\alpha}[\mathbf{s}_{\alpha,c}^{l-1}]$ . We then deduce that  $P_{\mathbf{x},\alpha}[RF(\mathbf{x}^{l-1})_{\alpha,i_c}] \simeq P_{\mathbf{x},\alpha}[\mathbf{x}_{\alpha,c}^{l-1}]$  and  $P_{\mathbf{x},\mathbf{s},\alpha}[RF(\mathbf{s}^{l-1})_{\alpha,i_c}] \simeq P_{\mathbf{x},\mathbf{s},\alpha}[\mathbf{s}_{\alpha,c}^{l-1}]$ , meaning that  $RF$  is approximately statistics-preserving for  $\mathbf{x}^{l-1}$  and  $\mathbf{s}^{l-1}$ .  $\square$

## D Details of Section 3

### D.1 Assumption that $\Phi_l$ is differentiable a.s. with respect to $\mathbf{x}$

Let us detail which parts of Section 3 rely on the assumption that  $\Phi_l$  is differentiable a.s. with respect to  $\mathbf{x}$ .

Firstly the *factor of noise equivalence* detailed in Eq. (4) relies on the assumption that moments with respect to  $\mathbf{x}, \alpha$  of the true noise tensor  $\Phi_l(\mathbf{x} + d\mathbf{x}) - \Phi_l(\mathbf{x})$  at layer  $l$  and  $d\mathbf{x}^l$  defined as the result of the simultaneous propagation of Eq. (1) and Eq. (2) coincide. If  $\Phi_l(\mathbf{x})$  is differentiable a.s. with respect to  $\mathbf{x}$ , then  $d\mathbf{x}^l = \Phi_l(\mathbf{x} + d\mathbf{x}) - \Phi_l(\mathbf{x})$  a.s. with respect to  $\mathbf{x}$ . In that case,  $d\mathbf{x}^l$  and  $\Phi_l(\mathbf{x} + d\mathbf{x}) - \Phi_l(\mathbf{x})$  share the same probability density function, and thus the same moments with respect to  $\mathbf{x}, \alpha$ .

Secondly the *Jacobian equivalence* and the *sensitivity equivalence* detailed respectively in Sections D.3 and D.2 rely on the assumption that  $\Phi_l(\mathbf{x})$  is differentiable *surely* with respect to  $\mathbf{x}$ . This can be relaxed using subdifferentials if  $\Phi_l(\mathbf{x})$  is differentiable a.s. with respect to  $\mathbf{x}$ . Indeed the same argument as before shows that moments with respect to  $\mathbf{x}, \alpha$  are left unchanged when ignoring the probability-zero event such that  $\Phi_l(\mathbf{x})$  is not differentiable with respect to  $\mathbf{x}$ .

Now let us justify the assumption that  $\Phi_l(\mathbf{x})$  is differentiable a.s. with respect to  $\mathbf{x}$  in the case of the propagation of Eq. (1). As in Section 4, we adopt the notation  $\Theta^l = (\omega^1, \beta^1, \dots, \omega^l, \beta^l)$  and we further assume standard initialization. For given  $\mathbf{x}$  such that  $\forall \alpha: \mathbf{x}_{\alpha,:} \neq 0$ , it is easy to see that  $\Phi_l(\mathbf{x}^l)$  is *not differentiable* implies that  $\exists k \leq l, \exists \alpha, c$  such that  $\rho(\mathbf{x}^{k-1}, \alpha) \neq 0$  and  $\mathbf{x}_{\alpha,c}^k = 0$ . Under standard initialization, this corresponds to a zero-probability event. Denoting  $D$  the event such that  $\Phi_l(\mathbf{x}^l)$  is not differentiable, we then have  $\mathbb{P}_{\Theta^l|\mathbf{x}}[D] = 0$ . Now considering  $\mathbf{x}$  again as random, using Fubini's Theorem and making the assumption that  $\mathbf{x}_{\alpha,:} \neq 0$  a.s. with respect to  $\mathbf{x}, \alpha$  (which is the case e.g. if this random vector  $\mathbf{x}_{\alpha,:}$  admits a well-defined probability density function):

$$\mathbb{E}_{\Theta^l} \mathbb{E}_{\mathbf{x}|\Theta^l}[\mathbf{1}_D] = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\Theta^l|\mathbf{x}}[\mathbf{1}_D] = \mathbb{E}_{\mathbf{x}}[\mathbb{P}_{\Theta^l|\mathbf{x}}[D]] = 0. \quad (37)$$

By contradiction, if there would be non-zero probability with respect to  $\Theta^l$  that  $\mathbb{P}_{\mathbf{x}|\Theta^l}[D] > 0$ , then Eq. (37) would not hold. Therefore with probability 1 with respect to  $\Theta^l$ :  $\mathbb{P}_{\mathbf{x}|\Theta^l}[D] = 0$ , meaning that  $\Phi_l(\mathbf{x})$  is differentiable a.s. with respect to  $\mathbf{x}$ .

### D.2 Property of normalized sensitivity

**Proposition 8.** *The sensitivity tensor  $\mathbf{s}^l$  and the tensor  $\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l$  containing for given  $\alpha, c$  the derivatives of  $\mathbf{x}_{\alpha,c}^l$  with respect to  $\mathbf{x}$  are related by  $\mathbb{E}_{\mathbf{s}}[(\mathbf{s}_{\alpha,c}^l)^2] = \|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l\|_2^2$ .*

**Proof of Proposition 8.** Due to the definition of  $d\mathbf{x}^l$  as a small corruption to  $\mathbf{x}^l$ ,  $d\mathbf{x}_{\alpha,c}^l$  can be written as a function of  $\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l$  and the input noise  $d\mathbf{x}$ ,

$$d\mathbf{x}_{\alpha,c}^l = \langle \nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l, d\mathbf{x} \rangle,$$

where  $\langle \cdot, \cdot \rangle$  denotes the standard dot product in input space. It follows from the definition of  $\mathbf{s}^l = d\mathbf{x}^l / \sigma_{d\mathbf{x}}$  that  $\mathbf{s}_{\alpha,c}^l = \langle \nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l, \mathbf{s} \rangle$ . Due to the white noise property  $\mathbb{E}_{\mathbf{s}}[\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$ , we then get

$$\mathbb{E}_{\mathbf{s}}[(\mathbf{s}_{\alpha,c}^l)^2] = \|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha,c}^l\|_2^2.$$

$\square$

**Proposition 9.** *Define  $\tilde{\mathbf{x}}$  the rescaling by a constant factor of  $\mathbf{x}$  with unit variance  $\mu_2(\tilde{\mathbf{x}}) = 1$ , and  $\tilde{\mathbf{x}}^l$  the rescaling by a constant factor of  $\mathbf{x}^l$  with unit variance  $\mu_2(\tilde{\mathbf{x}}^l) = 1$ . When  $\tilde{\mathbf{x}}^l$  is considered*



as a function of  $\tilde{\mathbf{x}}$ , the normalized sensitivity  $\chi^l$  measures a root mean square sensitivity as  $\chi^l = \mathbb{E}_{\tilde{\mathbf{x}}, \alpha, c} [\|\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{x}}_{\alpha, c}^l\|_2^2]^{1/2}$ .

**Proof of Proposition 9.** First we work with the non-normalized signals  $\mathbf{x}$  and  $\mathbf{x}^l$ . We can express the second-order central moment of  $\mathbf{s}^l$  as

$$\mu_2(\mathbf{s}^l) = \mathbb{E}_{\mathbf{x}, \alpha, c} \mathbb{E}_{\mathbf{s}} [\hat{\varphi}(\mathbf{s}^l, \alpha)_c^2] = \mathbb{E}_{\mathbf{x}, \alpha, c} \mathbb{E}_{\mathbf{s}} [\varphi(\mathbf{s}^l, \alpha)_c^2] = \mathbb{E}_{\mathbf{x}, \alpha, c} \mathbb{E}_{\mathbf{s}} [(\mathbf{s}_{\alpha, c}^l)^2] \quad (38)$$

$$= \mathbb{E}_{\mathbf{x}, \alpha, c} [\|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha, c}^l\|_2^2], \quad (39)$$

where Eq. (38) follows from  $\mathbf{s}^l$  being centered and Eq. (39) follows from Proposition 8. The normalized sensitivity is then given by

$$\chi^l = \left( \frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2} = \left( \frac{\mathbb{E}_{\mathbf{x}, \alpha, c} [\|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha, c}^l\|_2^2] \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2}. \quad (40)$$

Let us define  $\tilde{\mathbf{x}}^l = \mathbf{x}^l / \mu_2(\mathbf{x}^l)^{1/2}$  and  $\tilde{\mathbf{x}} = \mathbf{x} / \mu_2(\mathbf{x})^{1/2} = \mathbf{x} / \mu_2(\mathbf{x}^0)^{1/2}$ . Then we get  $\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{x}}_{\alpha, c}^l = \nabla_{\mathbf{x}} \mathbf{x}_{\alpha, c}^l / \mu_2(\mathbf{x}^l)^{1/2}$  and  $\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{x}}_{\alpha, c}^l = \nabla_{\mathbf{x}} \tilde{\mathbf{x}}_{\alpha, c}^l \mu_2(\mathbf{x}^0)^{1/2}$ . It follows that

$$\begin{aligned} \mu_2(\tilde{\mathbf{x}}^l) &= 1, \\ \mu_2(\tilde{\mathbf{x}}) &= 1, \\ \mathbb{E}_{\tilde{\mathbf{x}}, \alpha, c} [\|\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{x}}_{\alpha, c}^l\|_2^2] &= \frac{\mathbb{E}_{\mathbf{x}, \alpha, c} [\|\nabla_{\mathbf{x}} \mathbf{x}_{\alpha, c}^l\|_2^2] \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)}. \end{aligned} \quad (41)$$

Finally combining Eq. (40) and Eq. (41), we get  $\chi^l = \mathbb{E}_{\tilde{\mathbf{x}}, \alpha, c} [\|\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{x}}_{\alpha, c}^l\|_2^2]^{1/2}$ . □

### D.3 Equivalence between $\chi^l$ and previous definitions

In the fully-connected case  $n = 1$ , Philipp & Carbonell (2018) recently introduced the following coefficient:

$$\left( \frac{\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2]}{N_l} \frac{\mathbb{E}_{\mathbf{c}} [\text{Var}_{\mathbf{x}}[\mathbf{x}_{\mathbf{c}}^0]]}{\mathbb{E}_{\mathbf{c}} [\text{Var}_{\mathbf{x}}[\mathbf{x}_{\mathbf{c}}^l]]} \right)^{1/2}. \quad (42)$$

Let us prove the equivalence between the definitions of Eq. (3) and Eq. (42). In the fully-connected case  $n = 1$ , the spatial position  $\alpha$  can be ignored and tensors and feature map vectors coincide:  $\mathbf{x}^l = \varphi(\mathbf{x}^l)$  and  $\mathbf{s}^l = \varphi(\mathbf{s}^l) = \hat{\varphi}(\mathbf{s}^l)$ . When  $\mathbf{x}$  is fixed, the input-output Jacobian  $\mathbf{J}^l = \frac{\partial \mathbf{x}^l}{\partial \mathbf{x}^0} \in \mathbb{R}^{N_l \times N_0}$  directly summarizes the propagation of the noise  $d\mathbf{x}^l = \mathbf{J}^l d\mathbf{x}$ , and thus the sensitivity  $\mathbf{s}^l = \mathbf{J}^l \mathbf{s}$ . Due to the white noise property  $\mathbb{E}_{\mathbf{s}} [\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$ ,

$$\mathbb{E}_{\mathbf{x}, \mathbf{s}} \left[ \sum_{\mathbf{c}} (\mathbf{s}_{\mathbf{c}}^l)^2 \right] = \mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2], \quad \mathbb{E}_{\mathbf{x}, \mathbf{s}, \mathbf{c}} [(\mathbf{s}_{\mathbf{c}}^l)^2] = \frac{1}{N_l} \mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2].$$

Here we clearly see the advantage of the sensitivity tensor of encoding information on the Jacobian while avoiding increased dimensionality. Going back to our calculation, the definitions

$$\begin{aligned} \mu_2(\mathbf{s}^l) &= \mathbb{E}_{\mathbf{x}, \mathbf{s}, \mathbf{c}} [\hat{\varphi}(\mathbf{s}^l)_{\mathbf{c}}^2] = \mathbb{E}_{\mathbf{x}, \mathbf{s}, \mathbf{c}} [\varphi(\mathbf{s}^l)_{\mathbf{c}}^2] = \mathbb{E}_{\mathbf{x}, \mathbf{s}, \mathbf{c}} [(\mathbf{s}_{\mathbf{c}}^l)^2], \\ \mu_2(\mathbf{x}^l) &= \mathbb{E}_{\mathbf{x}, \mathbf{c}} [\hat{\varphi}(\mathbf{x}^l)_{\mathbf{c}}^2] = \mathbb{E}_{\mathbf{c}} [\text{Var}_{\mathbf{x}}[\mathbf{x}_{\mathbf{c}}^l]], \end{aligned}$$

finally give the equivalence between the two definitions:

$$\chi^l = \left( \frac{\mu_2(\mathbf{s}^l) \mu_2(\mathbf{x}^0)}{\mu_2(\mathbf{x}^l)} \right)^{1/2} = \left( \frac{\mathbb{E}_{\mathbf{x}} [\|\mathbf{J}^l\|_F^2] \mathbb{E}_{\mathbf{c}} [\text{Var}_{\mathbf{x}}[\mathbf{x}_{\mathbf{c}}^0]]}{N_l \mathbb{E}_{\mathbf{c}} [\text{Var}_{\mathbf{x}}[\mathbf{x}_{\mathbf{c}}^l]]} \right)^{1/2}.$$

Philipp & Carbonell (2018) chose the terminology of *nonlinearity coefficient* for this metric. While our analysis unveils a strong connection between  $\chi^l$  and the nonlinearity  $\phi$ , it also reveals a strong connection with batch normalization which is still a linear operation. So we chose instead the terminology of normalized sensitivity.

#### D.4 Distributional pathologies

We consider the following rescaling of the signal:

$$\boldsymbol{\nu}^l \equiv (\nu_{1,c}(\mathbf{x}^l))_{1 \leq c \leq N_l}, \quad \tilde{\mathbf{x}}^l \equiv \frac{1}{\|\boldsymbol{\nu}^l\|_2} \mathbf{x}^l, \quad \tilde{\boldsymbol{\nu}}^l \equiv (\nu_{1,c}(\tilde{\mathbf{x}}^l))_{1 \leq c \leq N_l}.$$

We immediately have  $\|\tilde{\boldsymbol{\nu}}^l\|_2 = 1$ . Furthermore we have

$$\begin{aligned} \nu_2(\mathbf{x}^l) &= \frac{1}{N_l} \sum_c \mathbb{E}_{\mathbf{x}, \alpha} [\varphi(\mathbf{x}^l, \alpha)_c^2] = \frac{1}{N_l} \left( \sum_c \text{Var}_{\mathbf{x}, \alpha} [\varphi(\mathbf{x}^l, \alpha)_c] + \mathbb{E}_{\mathbf{x}, \alpha} [\varphi(\mathbf{x}^l, \alpha)_c]^2 \right) \\ &= \frac{1}{N_l} \left( \sum_c \mu_{2,c}(\mathbf{x}^l) + \nu_{1,c}(\mathbf{x}^l)^2 \right) = \mu_2(\mathbf{x}^l) + \frac{1}{N_l} \|\boldsymbol{\nu}^l\|_2^2. \end{aligned}$$

The pathology  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 0$  then implies  $\nu_2(\mathbf{x}^l)/\|\boldsymbol{\nu}^l\|_2^2 \xrightarrow{l \rightarrow \infty} 1/N_l$  and  $\mu_2(\mathbf{x}^l)/\|\boldsymbol{\nu}^l\|_2^2 \xrightarrow{l \rightarrow \infty} 0$ , i.e.  $\mu_2(\tilde{\mathbf{x}}^l) \xrightarrow{l \rightarrow \infty} 0$ . Thus  $\varphi(\tilde{\mathbf{x}}^l, \alpha)$  becomes *point-like* concentrated around the vector  $\tilde{\boldsymbol{\nu}}^l$  of unit  $L^2$  norm.

### E Details of Section 5

#### E.1 Lemma on the sum of increments

**Lemma 10.** *Consider a sequence  $(X_k)$  of random variables which depend on  $\Theta^k$  and denote  $Y_k = \mathbb{E}_{\Theta^k}[X_k]$  and  $Z_k = X_k - \mathbb{E}_{\Theta^k}[X_k]$ . Then:*

(i) *The random variables  $Z_k$  are centered and non-correlated:*

$$\forall k \leq l : \mathbb{E}_{\Theta^k}[Z_k] = 0, \quad \forall k \neq k' \leq l : \mathbb{E}_{\Theta^{\max(k, k')}}[Z_k Z_{k'}] = 0.$$

(ii) *If there exist constants  $m_{\min}, m_{\max}, v_{\min}, v_{\max}$  with  $\forall k \leq l : m_{\min} \leq Y_k \leq m_{\max}$  and  $v_{\min} \leq \text{Var}_{\Theta^k}[Z_k] \leq v_{\max}$ , then there exist random variables  $m_l$  and  $s_l$  such that  $s_l$  is centered and*

$$\sum_{k=1}^l X_k = l m_l + \sqrt{l} s_l, \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l}[s_l] \leq v_{\max}.$$

**Proof of (i).** First we show that  $Z_k$  is centered:

$$\begin{aligned} \mathbb{E}_{\Theta^k}[Z_k] &= \mathbb{E}_{\Theta^k}[X_k] - \mathbb{E}_{\Theta^k}[X_k] = 0, \\ \mathbb{E}_{\Theta^k}[Z_k] &= \mathbb{E}_{\Theta^{k-1}}[\mathbb{E}_{\Theta^k}[Z_k]] = 0. \end{aligned} \tag{43}$$

Now for  $k < k'$ , we have  $k \leq k' - 1$  and thus  $Z_k$  is a random variable which only depends on  $\Theta^{k'-1}$ . Then we can write

$$\begin{aligned} \mathbb{E}_{\Theta^{k'}}[Z_k Z_{k'}] &= \mathbb{E}_{\Theta^{k'-1}} \mathbb{E}_{\Theta^{k'}}[Z_k Z_{k'}] \\ &= \mathbb{E}_{\Theta^{k'-1}}[Z_k \mathbb{E}_{\Theta^{k'}}[Z_{k'}]] \\ &= 0, \end{aligned} \tag{44}$$

where Eq. (44) follows from Eq. (43). □

**Proof of (ii).** Denote  $M_l = \sum_{k=1}^l Y_k$  and  $S_l = \sum_{k=1}^l Z_k$ . Then we have

$$\begin{aligned} \mathbb{E}_{\Theta^l}[S_l] &= \sum_k \mathbb{E}_{\Theta^l}[Z_k] = 0, \\ \text{Var}_{\Theta^l}[S_l] &= \mathbb{E}_{\Theta^l}[S_l^2] = \sum_{k, k'} \mathbb{E}_{\Theta^l}[Z_k Z_{k'}], \end{aligned}$$

$$= \sum_k \mathbb{E}_{\Theta^k} [Z_k^2] = \sum_k \text{Var}_{\Theta^k} [Z_k], \quad (45)$$

where Eq. (45) follows from (i). The hypothesis further gives  $lm_{\min} \leq M_l \leq lm_{\max}$  and  $lv_{\min} \leq \text{Var}_{\Theta^l}[S_l] \leq lv_{\max}$ . If we now define  $m_l = M_l / l$  and  $s_l = S_l / \sqrt{l}$ , we get that  $s_l$  is centered and the telescoping sum  $\sum_{k=1}^l X_k = \sum_{k=1}^l Y_k + \sum_{k=1}^l Z_k$  can be written as required:

$$\sum_{k=1}^l X_k = M_l + S_l = lm_l + \sqrt{l}s_l, \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l}[s_l] \leq v_{\max}.$$

□

## E.2 Proof of Theorem 1

**Theorem 1. Moments of vanilla networks.** *There exist positive constants  $m_{\min}$ ,  $m_{\max}$ ,  $v_{\min}$ ,  $v_{\max} > 0$ , random variables  $(m_l), (m'_l)$ ,  $(s_l), (s'_l)$  and events  $(A_l)$  with probability  $\mathbb{P}_{\Theta^l}[A_l] \geq \prod_{k=1}^l (1 - 2^{-N_k+1})$  such that under  $A_l$ :  $s_l$  and  $s'_l$  are centered and*

$$\begin{aligned} \log \nu_2(\mathbf{x}^l) &= -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), & m_{\min} \leq m_l \leq m_{\max}, & v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= -lm'_l + \sqrt{l}s'_l, & m_{\min} \leq m'_l \leq m_{\max}, & v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max}. \end{aligned}$$

**Proof.** We use the definitions and notations from Section B and we denote  $(e_1, \dots, e_{R_l})$  and  $(\lambda_1, \dots, \lambda_{R_l})$  respectively the orthogonal eigenvectors and eigenvalues of  $\mathbf{G}[\rho(\mathbf{x}^{l-1}, \alpha)]$ , and  $\hat{\mathbf{W}}^l = \mathbf{W}^l(e_1, \dots, e_{R_l})$ . We then get

$$\begin{aligned} \forall c : \nu_{2,c}(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{x}, \alpha} [(\mathbf{y}_{\alpha,c}^l)^2] = \mathbb{E}_{\mathbf{x}, \alpha} [(\mathbf{W}_{c,:}^l \rho(\mathbf{x}^{l-1}, \alpha))^2] \\ &= \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \lambda_i. \end{aligned} \quad (46)$$

We further define

$$u_c^l = \begin{cases} \frac{\nu_{2,c}(\mathbf{x}^l)}{\nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l)} & \text{if } \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l) > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (47)$$

Combining the definition of  $u_c^l$  with Eq. (24) and Eq. (46), we get that under  $\{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$ :

$$\begin{aligned} \forall c : \nu_{2,c}(\mathbf{x}^l) &= u_c^l \nu_{2,c}(\mathbf{y}^l), \\ \forall c : \nu_{2,c}(\mathbf{x}^l) &= u_c^l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \lambda_i = R_l \nu_2(\mathbf{x}^{l-1}) u_c^l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i, \end{aligned} \quad (48)$$

where we defined  $\hat{\lambda}_i = \lambda_i / \sum_j \lambda_j$  and used  $\sum_j \lambda_j = \text{Tr } \mathbf{G}[\rho(\mathbf{x}^{l-1}, \alpha)] = R_l \nu_2(\mathbf{x}^{l-1})$  by Corollary 3. The symmetric propagation gives that, under  $\{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$ :

$$\begin{aligned} \forall c : \nu_{2,c}(\bar{\mathbf{x}}^l) &= R_l \nu_2(\mathbf{x}^{l-1}) \bar{u}_c^l \sum_i (-\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i, \\ \forall c : \nu_{2,c}(\bar{\mathbf{x}}^l) &= R_l \nu_2(\mathbf{x}^{l-1}) (1 - u_c^l) \sum_i (-\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i, \end{aligned} \quad (49)$$

$$\forall c : \nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l) = R_l \nu_2(\mathbf{x}^{l-1}) \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i, \quad (50)$$

where  $\bar{u}_c^l$  is the symmetric counterpart of  $u_c^l$  obtained by interverting  $\mathbf{x}^l$  and  $\bar{\mathbf{x}}^l$  in Eq. (47), and Eq. (49) follows from  $\bar{u}_c^l = (1 - u_c^l)$  when  $\nu_{2,c}(\mathbf{y}^l) > 0$  and  $\bar{u}_c^l \nu_{2,c}(-\mathbf{y}^l) = (1 - u_c^l) \nu_{2,c}(-\mathbf{y}^l)$  still holding when  $\nu_{2,c}(\mathbf{y}^l) = 0$ . By symmetry of the propagation  $\nu_{2,c}(\mathbf{x}^l) \sim_{\theta^l} \nu_{2,c}(\bar{\mathbf{x}}^l)$ . Combined with Eq. (50) and the assumption of standard initialization  $\mathbf{W}^l \sim_{\theta^l} \hat{\mathbf{W}}^l \sim_{\theta^l} \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I})$ , we deduce that under  $\{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$ :

$$\begin{aligned}
2\mathbb{E}_{\theta^l} [\nu_{2,c}(\mathbf{x}^l)] &= \mathbb{E}_{\theta^l} [\nu_{2,c}(\mathbf{x}^l) + \nu_{2,c}(\bar{\mathbf{x}}^l)] \\
&= \mathbb{E}_{\theta^l} [R_l \nu_2(\mathbf{x}^{l-1}) \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i] \\
&= R_l \nu_2(\mathbf{x}^{l-1}) \frac{2}{R_l} \sum_i \hat{\lambda}_i = 2\nu_2(\mathbf{x}^{l-1}).
\end{aligned}$$

Thus  $\forall c : \mathbb{E}_{\theta^l} [\nu_{2,c}(\mathbf{x}^l)] = \nu_2(\mathbf{x}^{l-1})$  and  $\mathbb{E}_{\theta^l} [\nu_2(\mathbf{x}^l)] = \nu_2(\mathbf{x}^{l-1})$ , meaning that under  $\{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$ :

$$\mathbb{E}_{\theta^l} [\delta\nu_2(\mathbf{x}^l)] = 1. \quad (51)$$

Let us define

$$v_c^l = \begin{cases} 0 & \text{if } (u_c^l < 1/2) \wedge (\nu_{2,c}(\mathbf{y}^l) > 0) \\ 1 & \text{if } (u_c^l > 1/2) \wedge (\nu_{2,c}(\mathbf{y}^l) > 0) \\ b & \text{if } (u_c^l = 1/2) \vee (\nu_{2,c}(\mathbf{y}^l) = 0) \end{cases}$$

with  $\wedge$  the logical *and*,  $\vee$  the logical *or*, and with  $b$  independent of  $\Theta^l$  and following a Bernouilli distribution with probability  $1/2$ :  $b \sim \text{Bernouilli}(1/2)$ . Conditionally on any value of  $\nu_{2,c}(\mathbf{y}^l)$ , we have  $v_c^l \sim \text{Bernouilli}(1/2)$  and thus  $\nu_{2,c}(\mathbf{y}^l)$  and  $v_c^l$  are independent. Let us further denote  $B_l = \{\exists c : v_c^l = 1\}$ , which gives  $\mathbb{P}_{\theta^l}[B_l] = 1 - 2^{-N_l}$ .

Since  $(\nu_{2,c}(\mathbf{y}^l))_{1 \leq c \leq N_l}$  and  $(v_c^l)_{1 \leq c \leq N_l}$  are independent, using Eq. (48) we get that  $\exists(w_{i,0})$  such that under  $B_l \cap \{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$ :

$$(w_{i,0})_{1 \leq i \leq R_l} \sim \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I}) : \quad \frac{R_l}{2N_l} \nu_2(\mathbf{x}^{l-1}) \sum_{i=1}^{R_l} w_{i,0}^2 \hat{\lambda}_i \leq \frac{1}{N_l} \sum_{c=1}^{R_l} \nu_{2,c}(\mathbf{x}^l) = \nu_2(\mathbf{x}^l).$$

Thus  $\exists(w_{i,j})_{i,j}$  such that the increments  $\delta\nu_2(\mathbf{x}^l)$  can be bounded under  $B_l \cap \{\nu_2(\mathbf{x}^{l-1}) \neq 0\}$  as

$$\forall j : (w_{i,j})_{1 \leq i \leq R_l} \sim \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I}), \quad \frac{R_l}{2N_l} \sum_{i=1}^{R_l} w_{i,0}^2 \hat{\lambda}_i \leq \delta\nu_2(\mathbf{x}^l) \leq \frac{R_l}{N_l} \sum_{j=1}^{N_l} \sum_{i=1}^{R_l} w_{i,j}^2 \hat{\lambda}_i. \quad (52)$$

Simply replacing  $\mathbf{x}^l$  by  $\mathbf{s}^l$ ,  $\mathbf{y}^l$  by  $\mathbf{t}^l$ ,  $\mathbf{G}$  by  $\mathbf{C}$ , and using Eq. (26) instead of Eq. (24) and the identity with  $\mu_2(\mathbf{s}^{l-1})$  instead of  $\nu_2(\mathbf{x}^{l-1})$  in Corollary 3, we get that under  $\{\mu_2(\mathbf{s}^{l-1}) \neq 0\}$ :

$$\mathbb{E}_{\theta^l} [\delta\mu_2(\mathbf{s}^l)] = 1. \quad (53)$$

We similarly define  $B'_l : \mathbb{P}_{\theta^l}[B'_l] = 1 - 2^{-N_l}$ ,  $(\hat{\lambda}'_i)_{1 \leq i \leq R_l} : \sum_i \hat{\lambda}'_i = 1$  and  $(w'_{i,j})$  such that, under  $B'_l \cap \{\mu_2(\mathbf{s}^{l-1}) \neq 0\}$ :

$$\forall j : (w'_{i,j})_{1 \leq i \leq R_l} \sim \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I}) : \quad \frac{R_l}{2N_l} \sum_{i=1}^{R_l} w'_{i,0}{}^2 \hat{\lambda}'_i \leq \delta\mu_2(\mathbf{s}^l) \leq \frac{R_l}{N_l} \sum_{j=1}^{N_l} \sum_{i=1}^{R_l} w'_{i,j}{}^2 \hat{\lambda}'_i. \quad (54)$$

Let us denote  $A_l = \bigcap_{k=1}^l (B_k \cap B'_k \cap \{\nu_2(\mathbf{x}^k) \neq 0\} \cap \{\mu_2(\mathbf{s}^k) \neq 0\})$ . It follows from Eq. (52) and Eq. (54) that  $\{\nu_2(\mathbf{x}^k) \neq 0\}$  has probability 1 under  $B_k \cap A_{k-1}$ , and that  $\{\mu_2(\mathbf{s}^k) \neq 0\}$  has probability 1 under  $B'_k \cap A_{k-1}$ . Thus

$$\begin{aligned}
\mathbb{P}_{\theta^l|A_{l-1}}[A_l] &= \mathbb{P}_{\theta^l|A_{l-1}}[B_l \cap B'_l \cap \{\nu_2(\mathbf{x}^l) \neq 0\} \cap \{\mu_2(\mathbf{s}^l) \neq 0\}], \\
&= \mathbb{P}_{\theta^l|A_{l-1}}[B_l \cap B'_l] = 1 - \mathbb{P}_{\theta^l|A_{l-1}}[B_l^c \cup B_l'^c] \geq 1 - 2^{-N_l+1}, \quad (55)
\end{aligned}$$

where  $B_k^c$  and  $B_k'^c$  denote the complementary events of  $B_k$  and  $B_k'$  respectively. We then have for  $\mathbb{P}_{\Theta^l}[A_l]$ :

$$\begin{aligned}\mathbb{P}_{\Theta^l}[A_l] &= \prod_{k=1}^l \mathbb{P}_{\theta^k|A_{k-1}}[A_k] \\ &\geq \prod_{k=1}^l (1 - 2^{-N_k+1}).\end{aligned}$$

From Eq. (55), we deduce that  $\mathbb{P}_{\theta^l|A_{l-1}}[A_l^c] \ll 1$ , where  $A_l^c$  denotes the complementary event of  $A_l$ . Further using  $\mathbb{E}_{\theta^l|A_{l-1}}[\delta\nu_2(\mathbf{x}^l)] = 1$ , and  $\mathbb{E}_{\theta^l|A_{l-1}}[\delta\mu_2(\mathbf{s}^l)] = 1$  by Eq. (51) and Eq. (53), and applying Cauchy-Schwarz inequality:

$$\begin{aligned}& \left| \mathbb{E}_{\theta^l|A_l}[\delta\nu_2(\mathbf{x}^l)] - 1 \right| \\ &= \left| \mathbb{E}_{\theta^l|A_l}[\delta\nu_2(\mathbf{x}^l)] - \mathbb{E}_{\theta^l|A_{l-1}}[\delta\nu_2(\mathbf{x}^l)] \right| \\ &= \left| (\mathbb{P}_{\theta^l|A_{l-1}}[A_l]^{-1} - 1) \mathbb{E}_{\theta^l|A_{l-1}}[\mathbf{1}_{A_l} \delta\nu_2(\mathbf{x}^l)] + \mathbb{E}_{\theta^l|A_{l-1}}[\mathbf{1}_{A_l^c} \delta\nu_2(\mathbf{x}^l)] \right| \\ &\leq \mathbb{P}_{\theta^l|A_{l-1}}[A_l^c] \mathbb{P}_{\theta^l|A_{l-1}}[A_l]^{-1} \left| \mathbb{E}_{\theta^l|A_{l-1}}[\mathbf{1}_{A_l} \delta\nu_2(\mathbf{x}^l)] \right| + \mathbb{P}_{\theta^l|A_{l-1}}[A_l^c]^{1/2} \mathbb{E}_{\theta^l|A_{l-1}}[\delta\nu_2(\mathbf{x}^l)^2]^{1/2} \\ &\leq \left( \frac{\mathbb{P}_{\theta^l|A_{l-1}}[A_l^c]}{\mathbb{P}_{\theta^l|A_{l-1}}[A_l]} + \mathbb{P}_{\theta^l|A_{l-1}}[A_l^c]^{1/2} \right) \mathbb{E}_{\theta^l|A_{l-1}}[\delta\nu_2(\mathbf{x}^l)^2]^{1/2} \ll 1.\end{aligned}\tag{56}$$

Similarly,

$$\left| \mathbb{E}_{\theta^l|A_l}[\delta\mu_2(\mathbf{s}^l)] - 1 \right| \ll 1.\tag{57}$$

By log-concavity, it follows that  $\mathbb{E}_{\theta^l|A_l}[\log \delta\nu_2(\mathbf{x}^l)] \leq 0$ , and that  $\mathbb{E}_{\theta^l|A_l}[\log \delta\mu_2(\mathbf{s}^l)] \leq 0$ . Now since  $\log x$  and  $\log^2 x$  are integrable at 0, it follows from Eq. (52) and Eq. (54) that  $\delta\nu_2(\mathbf{x}^l)$  has well-defined first and second-order moments with respect to  $\theta^l$  under  $A_{l-1} \cap B_l \cap \{\nu_2(\mathbf{x}^l) \neq 0\}$ , and that  $\delta\mu_2(\mathbf{s}^l)$  has well-defined first and second-order moments with respect to  $\theta^l$  under  $A_{l-1} \cap B_l' \cap \{\mu_2(\mathbf{s}^l) \neq 0\}$ . Thus  $\delta\nu_2(\mathbf{x}^l)$  and  $\delta\mu_2(\mathbf{s}^l)$  have well-defined first and second-order moments with respect to  $\theta^l$  under  $A_l$ .

From now on we work under  $A_l$ , and we define under  $A_l$ :

$$\begin{aligned}X_k &= \log \delta\nu_2(\mathbf{x}^k), \quad Y_k = \mathbb{E}_{\theta^k|A_l}[\log \delta\nu_2(\mathbf{x}^k)], \quad Z_k = \log \delta\nu_2(\mathbf{x}^k) - \mathbb{E}_{\theta^k|A_l}[\log \delta\nu_2(\mathbf{x}^k)], \\ X'_k &= \log \delta\mu_2(\mathbf{s}^k), \quad Y'_k = \mathbb{E}_{\theta^k|A_l}[\log \delta\mu_2(\mathbf{s}^k)], \quad Z'_k = \log \delta\mu_2(\mathbf{s}^k) - \mathbb{E}_{\theta^k|A_l}[\log \delta\mu_2(\mathbf{s}^k)].\end{aligned}$$

Note that for  $k < l$ ,  $A_l$  is independent of  $\Theta^k$  under  $A_k$ . This implies

$$\begin{aligned}Y_k &= \mathbb{E}_{\theta^k|A_l}[\log \delta\nu_2(\mathbf{x}^k)] = \mathbb{E}_{\theta^k|A_k}[\log \delta\nu_2(\mathbf{x}^k)] \leq 0, & \text{Var}_{\Theta^k|A_k}[Z_k] &= \text{Var}_{\Theta^k|A_l}[Z_k], \\ Y'_k &= \mathbb{E}_{\theta^k|A_l}[\log \delta\mu_2(\mathbf{s}^k)] = \mathbb{E}_{\theta^k|A_k}[\log \delta\mu_2(\mathbf{s}^k)] \leq 0, & \text{Var}_{\Theta^k|A_k}[Z'_k] &= \text{Var}_{\Theta^k|A_l}[Z'_k].\end{aligned}$$

Let us now apply Lemma 10 for given  $l$  conditionally on  $A_l$ . Suppose there exist positive constants  $m_{\min}, m_{\max}, v_{\min}, v_{\max} > 0$ , such that  $\forall k \leq l$  under  $A_l$ :

$$m_{\min} \leq -Y_k, -Y'_k \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\theta^k|A_k}[Z_k], \text{Var}_{\theta^k|A_k}[Z'_k] \leq v_{\max},$$

implying

$$-m_{\max} \leq Y_k, Y'_k \leq -m_{\min}, \quad v_{\min} \leq \text{Var}_{\Theta^k|A_l}[Z_k], \text{Var}_{\Theta^k|A_l}[Z'_k] \leq v_{\max}.$$

Then by Lemma 10 there exist sequences of random variables  $(m_l), (m'_l), (s_l), (s'_l)$  such that  $\forall l$  under  $A_l$ :  $s_l$  and  $s'_l$  are centered and

$$\begin{aligned}\log \nu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^0) &= lm_l + \sqrt{l}s_l, \quad -m_{\max} \leq m_l \leq -m_{\min}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= lm'_l + \sqrt{l}s'_l, \quad -m_{\max} \leq m'_l \leq -m_{\min}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max},\end{aligned}$$

where we used  $\mu_2(\mathbf{s}^0) = 1$  due to the white noise property:  $\mathbb{E}_{\mathbf{s}}[\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$ . By changing the variable  $m_l$  to  $-m_l$  and the variable  $m'_l$  to  $-m'_l$ , we get that  $\forall l$  under  $A_l$ :  $s_l$  and  $s'_l$  are centered and

$$\begin{aligned}\log \nu_2(\mathbf{x}^l) &= -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), \quad m_{\min} \leq m_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s_l] \leq v_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= -lm'_l + \sqrt{l}s'_l, \quad m_{\min} \leq m'_l \leq m_{\max}, \quad v_{\min} \leq \text{Var}_{\Theta^l|A_l}[s'_l] \leq v_{\max}.\end{aligned}$$

To obtain the bounds  $m_{\min}$ ,  $m_{\max}$ ,  $v_{\min}$ ,  $v_{\max}$ , we consider extreme cases for  $u_c^l$  and  $R_l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i$  in Eq. (48). Let us denote Chi-Square( $N$ ) the chi-square distribution with  $N$  degree of freedom and  $N_{\min}$ ,  $N_{\max}$ ,  $R_{\min}$ ,  $R_{\max}$  the bounds on  $N_l$ ,  $R_l$  such that

$$N_{\min} = \min_l N_l, \quad N_{\max} = \max_l N_l, \quad R_{\min} = \min_l R_l, \quad R_{\max} = \max_l R_l.$$

Then we obtain *minimum bounds* by considering  $N_l = N_{\max}$ ,  $R_l = R_{\max}$  and  $u_c^l \sim_{\theta^l} 1/2$  and  $R_l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i \sim_{\theta^l} 2 \text{Chi-Square}(R_{\max}) / R_{\max}$ . This leads to  $\log \delta \nu_2(\mathbf{x}^l), \log \delta \mu_2(\mathbf{s}^l) \sim_{\theta^l} \text{Chi-Square}(N_{\max} R_{\max}) / (N_{\max} R_{\max})$ . We obtain *maximum bounds* by considering  $N_l = N_{\min}$ , and  $u_c^l \sim_{\theta^l} \text{Bernoulli}(1/2)$  and  $R_l \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \hat{\lambda}_i \sim_{\theta^l} 2 \text{Chi-Square}(1)$ .

As an illustration, in the fully-connected case with constant width  $N_l = 100$ , we numerically find  $m_{\min} \simeq 9.7 \times 10^{-5}$  and  $v_{\min} \simeq 2.0 \times 10^{-4}$  as minimum bounds and  $m_{\max} \simeq 2.5 \times 10^{-2}$  and  $v_{\max} \simeq 5.2 \times 10^{-2}$  as maximum bounds.  $\square$

### E.3 Relation to the terms $\overline{m}$ , $\underline{m}$ , $\underline{s}$ defined in Section 4

Here we relate Theorem 1 to the terms  $\overline{m}$ ,  $\underline{m}$ ,  $\underline{s}$  defined in Section 4, under the conditionality  $A_k$ . By Eq. (56) and Eq. (57), we have  $|\mathbb{E}_{\theta^k|A_k}[\delta \nu_2(\mathbf{x}^k)] - 1| \ll 1$  and  $|\mathbb{E}_{\theta^k|A_k}[\delta \mu_2(\mathbf{s}^k)] - 1| \ll 1$ , which implies

$$\begin{aligned}|\overline{m}[\nu_2(\mathbf{x}^k)]| &= |\log \mathbb{E}_{\theta^k|A_k}[\delta \nu_2(\mathbf{x}^k)]| \simeq |\mathbb{E}_{\theta^k|A_k}[\delta \nu_2(\mathbf{x}^k)] - 1| \ll 1, \\ |\overline{m}[\mu_2(\mathbf{s}^k)]| &\ll 1.\end{aligned}$$

The terms  $\overline{m}[\nu_2(\mathbf{x}^k)]$  and  $\overline{m}[\mu_2(\mathbf{s}^k)]$  are thus vanishing and the evolution is dominated by the terms  $\underline{m}[\nu_2(\mathbf{x}^k)] < 0$ ,  $\underline{m}[\mu_2(\mathbf{s}^k)] < 0$ . These terms correspond to  $Y_k, Y'_k$  in the proof of Theorem 1.

### E.4 Proof of Theorem 2

**Theorem 2. Normalized Sensitivity increments of vanilla networks.** Denote  $\mathbf{y}^{l,+} = \max(\mathbf{y}^l, 0)$  and  $\mathbf{y}^{l,-} = \max(-\mathbf{y}^l, 0)$ . The dominating term under  $A_l$  in the evolution of  $\chi^l$  is

$$\delta \chi^l \simeq \exp\left(\overline{m}_{\text{vanilla}}[\chi^l]\right) = \left(1 - \mathbb{E}_{c,\theta^l|A_{l-1}}\left[\frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})}\right]\right)^{-1/2}. \quad (58)$$

**Proof.** The dominating term in the evolution of  $\chi^l$  is  $\frac{1}{2}(\overline{m}[\mu_2(\mathbf{s}^l)] - \overline{m}[\mu_2(\mathbf{x}^l)])$ . The terms  $\overline{m}[\mu_2(\mathbf{s}^l)]$  and  $\overline{m}[\mu_2(\mathbf{x}^l)]$  are simply obtained by considering  $\mathbb{E}_{\theta^l}[\delta \mu_2(\mathbf{s}^l)]$  and  $\mathbb{E}_{\theta^l}[\delta \mu_2(\mathbf{x}^l)]$ .

By Eq. (53) in the proof of Theorem 1:  $\mathbb{E}_{\theta^l}[\delta \mu_2(\mathbf{s}^l)] = 1$ , and thus  $\overline{m}[\mu_2(\mathbf{s}^l)] = \log \mathbb{E}_{\theta^l|A_{l-1}}[\delta \mu_2(\mathbf{s}^l)] = 0$ .

Next we turn to the term  $\overline{m}[\mu_2(\mathbf{x}^l)]$ . Again we use the definitions and notations from Section B. We further denote  $(\mathbf{e}_1, \dots, \mathbf{e}_{R_l})$  and  $(\lambda_1, \dots, \lambda_{R_l})$  respectively the orthogonal eigenvectors and



eigenvalues of  $C[\rho(\mathbf{x}^{l-1}, \boldsymbol{\alpha})]$  and  $\hat{\mathbf{W}}^l = \mathbf{W}^l(e_1, \dots, e_{R_l})$ . Using these notations, we get

$$\begin{aligned} \forall c : \mu_{2,c}(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [\hat{\varphi}(\mathbf{y}^l, \boldsymbol{\alpha})_c^2] = \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [(\mathbf{W}_{c,:}^l \hat{\rho}(\mathbf{x}^{l-1}, \boldsymbol{\alpha}))^2] \\ &= \sum_i (\hat{\mathbf{W}}_{c,i}^l)^2 \lambda_i. \end{aligned} \quad (59)$$

Then due to  $\mathbf{W}^l \sim_{\theta^l} \hat{\mathbf{W}}^l \sim_{\theta^l} \mathcal{N}(0, \sqrt{2/R_l} \mathbf{I})$ :

$$\begin{aligned} \mathbb{E}_{\theta^l | A_{l-1}} [\mu_{2,c}(\mathbf{y}^l)] &= \frac{2}{R_l} \sum_i \lambda_i = \frac{2}{R_l} \text{Tr } C[\rho(\mathbf{x}^{l-1}, \boldsymbol{\alpha})] \\ &= 2\mu_2(\mathbf{x}^{l-1}). \end{aligned} \quad (60)$$

where Eq. (60) follows from Corollary 3. Furthermore the symmetric propagation gives:

$$\begin{aligned} \mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\bar{\mathbf{x}}^l) &= \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [(\mathbf{y}_{\boldsymbol{\alpha},c}^{l,+})^2] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [\mathbf{y}_{\boldsymbol{\alpha},c}^{l,+}]^2 + \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [(\mathbf{y}_{\boldsymbol{\alpha},c}^{l,-})^2] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\alpha}} [\mathbf{y}_{\boldsymbol{\alpha},c}^{l,-}]^2 \\ &= \nu_{2,c}(\mathbf{y}^{l,+}) - \nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{2,c}(\mathbf{y}^{l,-}) - \nu_{1,c}(\mathbf{y}^{l,-})^2 \\ &= \nu_{2,c}(\mathbf{y}^l) - (\nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{1,c}(\mathbf{y}^{l,-})^2) \end{aligned} \quad (61)$$

We have  $\nu_{1,c}(\mathbf{y}^l) = \nu_{1,c}(\mathbf{y}^{l,+}) - \nu_{1,c}(\mathbf{y}^{l,-})$  and thus  $\nu_{1,c}(\mathbf{y}^l)^2 = \nu_{1,c}(\mathbf{y}^{l,+})^2 + \nu_{1,c}(\mathbf{y}^{l,-})^2 - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})$ . We can then rewrite Eq. (61) as

$$\begin{aligned} \mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\bar{\mathbf{x}}^l) &= \nu_{2,c}(\mathbf{y}^l) - \nu_{1,c}(\mathbf{y}^l)^2 - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) \\ &= \mu_{2,c}(\mathbf{y}^l) - 2\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-}) \end{aligned} \quad (62)$$

Combining Eq. (60) and Eq. (62):

$$\begin{aligned} \mathbb{E}_{\theta^l | A_{l-1}} [\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\bar{\mathbf{x}}^l)] &= 2\mu_2(\mathbf{x}^{l-1}) - 2\mathbb{E}_{\theta^l | A_{l-1}} [\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})], \\ 2\mathbb{E}_{\theta^l | A_{l-1}} [\mu_{2,c}(\mathbf{x}^l)] &= 2\mu_2(\mathbf{x}^{l-1}) - 2\mathbb{E}_{\theta^l | A_{l-1}} [\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})], \\ \mathbb{E}_{\theta^l | A_{l-1}} [\mu_{2,c}(\mathbf{x}^l)] &= \mu_2(\mathbf{x}^{l-1}) \left( 1 - \mathbb{E}_{\theta^l | A_{l-1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right). \end{aligned} \quad (63)$$

where Eq. (63) is obtained by symmetry of the propagation. We finally get

$$\begin{aligned} \mathbb{E}_{\theta^l | A_{l-1}} [\mu_2(\mathbf{x}^l)] &= \mathbb{E}_{\theta^l | A_{l-1}} [\mathbb{E}_c [\mu_{2,c}(\mathbf{x}^l)]] = \mathbb{E}_c [\mathbb{E}_{\theta^l | A_{l-1}} [\mu_{2,c}(\mathbf{x}^l)]] \\ &= \mu_2(\mathbf{x}^{l-1}) \left( 1 - \mathbb{E}_{c, \theta^l | A_{l-1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right), \\ \mathbb{E}_{\theta^l | A_{l-1}} [\delta\mu_2(\mathbf{x}^{l-1})] &= 1 - \mathbb{E}_{c, \theta^l | A_{l-1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right]. \end{aligned}$$

Combined with  $\mathbb{E}_{\theta^l | A_{l-1}} [\delta\mu_2(\mathbf{s}^l)] = 1$ :

$$\begin{aligned} \delta\chi^l &\simeq \exp \left( \overline{m}_{\text{vanilla}}[\chi^l] \right) = \left( \frac{\mathbb{E}_{\theta^l | A_{l-1}} [\delta\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l | A_{l-1}} [\delta\mu_2(\mathbf{x}^l)]} \right)^{1/2} \\ &= \left( 1 - \mathbb{E}_{c, \theta^l | A_{l-1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{l,+})\nu_{1,c}(\mathbf{y}^{l,-})}{\mu_2(\mathbf{x}^{l-1})} \right] \right)^{-1/2}. \end{aligned}$$

□

**E.5** If  $\chi^l$  has exponential drift larger than diffusion and  $\mu_2(\mathbf{x}^l), \nu_2(\mathbf{x}^l)$  are lognormal, then  $\mu_2(\mathbf{x}^l) / \nu_2(\mathbf{x}^l)$  converges a.s. to zero

**Lemma 11. Borel Cantelli implies a.s. convergence** (see [Durrett \(1996\)](#)). For a sequence of random variables  $(X_l)$  and a random variable  $X$ , if  $\forall \epsilon > 0 : \sum_{l=1}^{\infty} \mathbb{P}[|X_l - X| > \epsilon] < \infty$ , then

$$X_l \xrightarrow{l \rightarrow \infty} X \text{ a.s.}$$

**Proof.** For given  $\epsilon > 0$ , denote  $N_\epsilon$  the number of times that the event  $\{|X_l - X| > \epsilon\}$  occurs such that  $N_\epsilon = \sum_{l=1}^{\infty} \mathbf{1}_{\{|X_l - X| > \epsilon\}}$ . Fubini's Theorem implies  $\mathbb{E}[N_\epsilon] = \sum_{l=1}^{\infty} \mathbb{P}[|X_l - X| > \epsilon] < \infty$ , and thus that  $N_\epsilon$  is finite a.s.

Now let us reason by contradiction and suppose  $\exists E$  with  $\mathbb{P}[E] > 0$  such that, under  $E$ :  $X_l \not\xrightarrow{l \rightarrow \infty} X$ . Under  $E$ ,  $\exists \epsilon$  random variable and  $\exists (k_l)$  increasing sequence with  $\forall l: |X_{k_l} - X| > \epsilon$ . This implies in turn  $\exists \epsilon' > 0$  non-random and  $\exists E'$  with  $\mathbb{P}[E'] > 0$  such that under  $E'$ :  $\exists (k_l)$  increasing sequence with  $\forall l: |X_{k_l} - X| > \epsilon'$ . Thus  $N_{\epsilon'}$  thus has non-zero probability to be infinite:  $\mathbb{P}[N_{\epsilon'} = \infty] \geq \mathbb{P}[E'] > 0$ , which is a contradiction. We deduce that  $X_l \xrightarrow{l \rightarrow \infty} X$  a.s. □

**Proposition 12.** Suppose that

- (i) We can neglect the events  $A_l$  with probability exponentially small in the width
- (ii) The event  $D$  under which  $\chi^l$  has drift larger than diffusion has probability  $\mathbb{P}[D] > 0$
- (iii)  $\mu_2(\mathbf{x}^l), \nu_2(\mathbf{x}^l)$  are lognormal

Then, under  $D$ :

$$\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} \xrightarrow{l \rightarrow \infty} 0 \text{ a.s.} \quad (64)$$

**Proof.** Neglecting the events  $A_l$ , Theorem 1 implies that  $\exists (m_l), (m'_l), (s_l), (s'_l)$  such that  $(s_l), (s'_l)$  are centered and

$$\begin{aligned} \log \nu_2(\mathbf{x}^l) &= -lm_l + \sqrt{l}s_l + \log \nu_2(\mathbf{x}^0), & m_{\min} \leq m_l \leq m_{\max}, & v_{\min} \leq \text{Var}[s_l] \leq v_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= -lm'_l + \sqrt{l}s'_l, & m_{\min} \leq m'_l \leq m_{\max}, & v_{\min} \leq \text{Var}[s'_l] \leq v_{\max}. \end{aligned}$$

Given that  $\log \nu_2(\mathbf{x}^l)$  and  $\log \mu_2(\mathbf{s}^l)$  are Gaussian by the assumption of lognormality, and given  $\mathbb{E}_{\Theta^l}[\nu_2(\mathbf{x}^l)] = \nu_2(\mathbf{x}^0)$  and  $\mathbb{E}_{\Theta^l}[\mu_2(\mathbf{s}^l)] = \mu_2(\mathbf{s}^0) = 1$  under standard initialization, we deduce:  $\exists (X_l), (X'_l)$  random variables and  $\exists (M_l), (M'_l) > 0$  constants such that

$$\begin{aligned} \log \nu_2(\mathbf{x}^l) &= X_l - M_l + \log \nu_2(\mathbf{x}^0), & X_l &\sim \mathcal{N}(0, \sqrt{M_l}), & lm_{\min} \leq M_l \leq lm_{\max}, \\ \log \mu_2(\mathbf{s}^l) &= X'_l - M'_l, & X'_l &\sim \mathcal{N}(0, \sqrt{M'_l}), & lm_{\min} \leq M'_l \leq lm_{\max}. \end{aligned}$$

Now let us make more precise the conditionality on  $D$ . We may assume  $\exists m > (m_{\max} - m_{\min})/2$  such that  $\forall l$  under  $D$ :  $\log \chi^l \geq lm$ .

The ratio  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l)$  can be expressed as

$$\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} = \frac{\mu_2(\mathbf{x}^l)}{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)} \frac{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)}{\nu_2(\mathbf{x}^l)} = \frac{1}{(\chi^l)^2} \frac{\mu_2(\mathbf{s}^l)\mu_2(\mathbf{x}^0)}{\nu_2(\mathbf{x}^l)},$$

which gives with logarithms that, under  $D$ :

$$\begin{aligned} \log \mu_2(\mathbf{x}^l) - \log \nu_2(\mathbf{x}^l) &= -2 \log \chi^l + \log \mu_2(\mathbf{s}^l) - \log \nu_2(\mathbf{x}^l) + \log \mu_2(\mathbf{x}^0) \\ &\leq -2lm + (X'_l - M'_l) - (X_l - M_l + \log \nu_2(\mathbf{x}^0)) + \log \mu_2(\mathbf{x}^0) \end{aligned}$$

$$\begin{aligned}
&\leq -2lm + lm_{\max} - lm_{\min} - \log \nu_2(\mathbf{x}^0) + \log \mu_2(\mathbf{x}^0) + X'_l - X_l \\
&\leq -lM + C + X'_l - X_l,
\end{aligned}$$

with  $M = 2m - m_{\max} + m_{\min} > 0$  and  $C = -\log \nu_2(\mathbf{x}^0) + \log \mu_2(\mathbf{x}^0)$ . Thus for given  $\epsilon$ , under  $D$ :

$$\begin{aligned}
\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon &\implies \log \epsilon < -lM + C + X'_l - X_l \\
&\implies \left( X'_l \geq \frac{\log \epsilon + lM - C}{2} \right) \vee \left( -X_l \geq \frac{\log \epsilon + lM - C}{2} \right) \\
&\implies \left( \tilde{X}'_l \geq \frac{\log \epsilon + lM - C}{2\sqrt{M'_l}} \right) \vee \left( -\tilde{X}_l \geq \frac{\log \epsilon + lM - C}{2\sqrt{M_l}} \right) \\
&\implies \left( \tilde{X}'_l \geq \frac{\log \epsilon + lM - C}{2\sqrt{lm_{\max}}} \right) \vee \left( -\tilde{X}_l \geq \frac{\log \epsilon + lM - C}{2\sqrt{lm_{\max}}} \right),
\end{aligned}$$

with  $\vee$  the logical *or*, and with  $\tilde{X}_l \equiv X_l/\sqrt{M_l}$  and  $\tilde{X}'_l \equiv X'_l/\sqrt{M'_l}$ . In turn, this implies  $\exists C_\epsilon$  such that  $\forall l$ , under  $D$ :

$$\begin{aligned}
\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon &\implies \left( \tilde{X}'_l \geq \sqrt{l}C_\epsilon \right) \vee \left( -\tilde{X}_l \geq \sqrt{l}C_\epsilon \right), \\
\mathbb{P}_D \left[ \frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon \right] &\leq \mathbb{P}_D [\tilde{X}'_l \geq \sqrt{l}C_\epsilon] + \mathbb{P}_D [-\tilde{X}_l \geq \sqrt{l}C_\epsilon] \\
&\leq \frac{1}{\mathbb{P}[D]} \mathbb{P} [1_D \mathbf{1}_{\{\tilde{X}'_l \geq \sqrt{l}C_\epsilon\}}] + \frac{1}{\mathbb{P}[D]} \mathbb{P} [1_D \mathbf{1}_{\{-\tilde{X}_l \geq \sqrt{l}C_\epsilon\}}] \\
&\leq \frac{1}{\mathbb{P}[D]} \operatorname{erfc} \left( \sqrt{\frac{l}{2}} C_\epsilon \right) \tag{65}
\end{aligned}$$

$$\leq \frac{1}{\mathbb{P}[D]} \exp \left( -\frac{l}{2} C_\epsilon^2 \right), \tag{66}$$

where Eq. (65) is obtained using  $\tilde{X}_l, \tilde{X}'_l \sim \mathcal{N}(0, 1)$ , while Eq. (66) is obtained using  $\operatorname{erfc}(x) \leq \exp(-x^2)$  (Chiani et al., 2003). It follows from Eq. (66) that

$$\sum_{l=1}^{\infty} \mathbb{P}_D \left[ \frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} > \epsilon \right] < \infty.$$

By Lemma 11, we finally deduce that, under  $D$ :

$$\frac{\mu_2(\mathbf{x}^l)}{\nu_2(\mathbf{x}^l)} \xrightarrow{l \rightarrow \infty} 0 \text{ a.s.}$$

□

**E.6 If  $\exp(\overline{m}[\chi^l]) \rightarrow 1$  and  $\tilde{\mathbf{x}}^l$  has bounded moments, then  $\mathbf{x}^l$  converges to one-dimensional signal pathology**

**Proposition 13.** We recall the notation for the unit-variance rescaled signal:  $\tilde{\mathbf{x}}^l \equiv \mathbf{x}^l/\mu_2(\mathbf{x}^l)^{1/2}$ , and the usual notation:

$$X_l = \mathcal{O}(Y_l) \iff \exists M > 0, \forall l : |X_l| \leq M|Y_l|.$$

We further suppose that

- (i)  $\tilde{\mathbf{x}}^l$  is well-defined with bounded moments:  $\nu_p(|\tilde{\mathbf{x}}^l|) = \mathcal{O}(1)$ , implying in particular  $\nu_2(\mathbf{x}^l)/\mu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} \infty$  and thus  $\mu_2(\mathbf{x}^l)/\nu_2(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 0$ , i.e. that  $\mathbf{x}^l$  does not converge to zero-dimensional signal pathology,

$$(ii) \delta\chi^l \simeq \exp(\overline{m}[\chi^l]) \xrightarrow{l \rightarrow \infty} 1.$$

Then  $\mathbf{x}^l$  converges to one-dimensional signal pathology.

**Proof.** Again we use the notations from Section B and we denote

$$\boldsymbol{\nu}_\varphi^l \equiv \mathbb{E}_{\mathbf{x}, \alpha} [\varphi(\tilde{\mathbf{x}}^l, \alpha)] = (\nu_{1,c}(\tilde{\mathbf{x}}^l))_{1 \leq c \leq N_l}, \quad \boldsymbol{\nu}_\rho^l \equiv \mathbb{E}_{\mathbf{x}, \alpha} [\rho(\tilde{\mathbf{x}}^l, \alpha)].$$

Due to the statistic-preserving property, we have  $\frac{1}{N_l} \|\boldsymbol{\nu}_\varphi^l\|_2^2 = \frac{1}{R_l} \|\boldsymbol{\nu}_\rho^l\|_2^2$ , which implies

$$\begin{aligned} \nu_2(\tilde{\mathbf{x}}^l) &= \frac{1}{N_l} \left( \sum_c \mu_{2,c}(\tilde{\mathbf{x}}^l) + \nu_{1,c}(\tilde{\mathbf{x}}^l)^2 \right) = \mu_2(\tilde{\mathbf{x}}^l) + \frac{1}{N_l} \|\boldsymbol{\nu}_\varphi^l\|_2^2 \\ &= 1 + \frac{1}{N_l} \|\boldsymbol{\nu}_\varphi^l\|_2^2 = 1 + \frac{1}{R_l} \|\boldsymbol{\nu}_\rho^l\|_2^2 \end{aligned}$$

i.e.  $\|\boldsymbol{\nu}_\rho^l\|_2^2 = R_l(\nu_2(\tilde{\mathbf{x}}^l) - 1)$ . Combined with  $\nu_2(\tilde{\mathbf{x}}^l) = \mathcal{O}(1)$ , we deduce that  $\|\boldsymbol{\nu}_\rho^l\|_2 = \mathcal{O}(1)$ .

Now let us reason by contradiction and suppose that  $r_{\text{eff}}(\mathbf{x}^l) = r_{\text{eff}}(\tilde{\mathbf{x}}^l) \xrightarrow{l \rightarrow \infty} 1$ , which implies  $\exists \eta > 0$  and  $\exists(k_l)$  increasing sequence with  $\forall l: r_{\text{eff}}(\tilde{\mathbf{x}}^{k_l}) \geq 1 + \eta$ . In turn this implies  $\exists \eta' > 0$  such that  $\forall l$ :

$$\exists \mathbf{v}_\varphi^{k_l} \in \mathbb{R}^{N_{k_l}} \perp \boldsymbol{\nu}_\varphi^{k_l}, \quad \|\mathbf{v}_\varphi^{k_l}\|_2 = 1 : \quad \text{Var}_{\mathbf{x}, \alpha} [\langle \varphi(\tilde{\mathbf{x}}^{k_l}, \alpha), \mathbf{v}_\varphi^{k_l} \rangle] = \mathbb{E}_{\mathbf{x}, \alpha} [\langle \varphi(\tilde{\mathbf{x}}^{k_l}, \alpha), \mathbf{v}_\varphi^{k_l} \rangle^2] \geq \eta',$$

i.e. that  $\varphi(\tilde{\mathbf{x}}^{k_l}, \alpha)$  necessarily has a direction of variance  $> \eta'$  which is orthogonal to its mean vector  $\boldsymbol{\nu}_\varphi^{k_l}$ . By padding this direction appropriately with zeros, it follows that  $\exists \eta' > 0$  such that  $\forall l$ :

$$\exists \mathbf{v}_\rho^{k_l} \in \mathbb{R}^{R_{k_l}} \perp \boldsymbol{\nu}_\rho^{k_l}, \quad \|\mathbf{v}_\rho^{k_l}\|_2 = 1 : \quad \text{Var}_{\mathbf{x}, \alpha} [\langle \rho(\tilde{\mathbf{x}}^{k_l}, \alpha), \mathbf{v}_\rho^{k_l} \rangle] = \mathbb{E}_{\mathbf{x}, \alpha} [\langle \rho(\tilde{\mathbf{x}}^{k_l}, \alpha), \mathbf{v}_\rho^{k_l} \rangle^2] \geq \eta'.$$

Let us denote  $\tilde{\mathbf{W}}^{k_l+1}$  such that  $\forall c: \tilde{\mathbf{W}}_{c,:}^{k_l+1} = \mathbf{W}_{c,:}^{k_l+1} / \|\mathbf{W}_{c,:}^{k_l+1}\|_2$  and  $\tilde{\boldsymbol{\nu}}_\rho^{k_l} = \boldsymbol{\nu}_\rho^{k_l} / \|\boldsymbol{\nu}_\rho^{k_l}\|_2$ . We further decompose  $\tilde{\mathbf{W}}_{c,:}^{k_l+1}$  as

$$\tilde{\mathbf{W}}_{c,:}^{k_l+1} = w_{\mathbf{v}} \mathbf{v}_\rho^{k_l} + w_{\boldsymbol{\nu}} \tilde{\boldsymbol{\nu}}_\rho^{k_l} + \mathbf{w}', \quad \mathbf{w}' \perp \mathbf{v}_\rho^{k_l}, \quad \mathbf{w}' \perp \tilde{\boldsymbol{\nu}}_\rho^{k_l}.$$

Then we get

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \alpha} \left[ \left( \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right)^2 \right]^2 &\geq \left( w_{\mathbf{v}}^2 \mathbb{E}_{\mathbf{x}, \alpha} [\langle \rho(\tilde{\mathbf{x}}^{k_l}, \alpha), \mathbf{v}_\rho^{k_l} \rangle^2] \right)^2 \geq w_{\mathbf{v}}^4 \eta'^2, \\ \mathbb{E}_{\mathbf{x}, \alpha} \left[ \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right]^2 &= w_{\boldsymbol{\nu}}^2 \|\boldsymbol{\nu}_\rho^{k_l}\|_2^2. \end{aligned}$$

Given that  $\|\boldsymbol{\nu}_\rho^{k_l}\|_2 = \mathcal{O}(1)$ , this implies by spherical symmetry that  $\forall \epsilon > 0, \exists p_\epsilon > 0$  such that  $\forall l$ :

$$\mathbb{P}_{\theta^{k_l+1}} \left[ \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ \left( \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right)^2 \right] \geq \eta'^2 - \epsilon \right) \wedge \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right]^2 \leq \epsilon \right) \right] \geq p_\epsilon, \quad (67)$$

with  $\wedge$  the logical *and*. By Cauchy-Schwarz inequality:

$$\mathbb{E}_{\mathbf{x}, \alpha} \left[ \left( \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right)^2 \right]^2 \leq \mathbb{E}_{\mathbf{x}, \alpha} \left[ \left| \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right| \right] \mathbb{E}_{\mathbf{x}, \alpha} \left[ \left| \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right|^3 \right]. \quad (68)$$

We can bound the second term on the right-hand side as

$$\mathbb{E}_{\mathbf{x}, \alpha} \left[ \left| \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right|^3 \right]$$

$$\leq \mathbb{E}_{\mathbf{x}, \alpha} \left[ \left( \sum_{i_c=1}^{R_l} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_c}^2 \right)^{3/2} \right] \quad (69)$$

$$\leq \mathbb{E}_{\mathbf{x}, \alpha} \left[ \sum_{i_1, i_2, i_3} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_1}^2 \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_2}^2 \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_3}^2 \right]^{1/2} \quad (70)$$

$$\leq \sum_{i_1, i_2, i_3} \mathbb{E}_{\mathbf{x}, \alpha} [\rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_1}^4]^{1/4} \mathbb{E}_{\mathbf{x}, \alpha} [\rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_2}^8]^{1/8} \mathbb{E}_{\mathbf{x}, \alpha} [\rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_{i_3}^8]^{1/8} \quad (71)$$

$$\leq R_l^3 N_l^{1/2} \nu_4(\tilde{\mathbf{x}}^{k_l})^{1/4} \nu_8(\tilde{\mathbf{x}}^{k_l})^{1/4}, \quad (72)$$

where Eq. (69), Eq. (70) and Eq. (71) are obtained by again applying Cauchy-Schwarz inequality, while Eq. (72) is obtained with  $\forall i, \forall p: \mathbb{E}_{\mathbf{x}, \alpha} [\rho(\tilde{\mathbf{x}}^{k_l}, \alpha)_i^p] \leq \sum_c \mathbb{E}_{\mathbf{x}, \alpha} [\varphi(\tilde{\mathbf{x}}^{k_l}, \alpha)_c^p] = N_l \nu_p(\tilde{\mathbf{x}}^l)$ . It then follows from Eq. (68) and the hypothesis that all moments are bounded  $\nu_p(|\tilde{\mathbf{x}}^l|) = \mathcal{O}(1)$  that

$$\mathbb{E}_{\mathbf{x}, \alpha} \left[ (\tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^2 \right] = \mathcal{O} \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ |\tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)| \right] \right). \quad (73)$$

Combining Eq. (67) and Eq. (73), we deduce that  $\exists \eta'' > 0$  with  $\forall \epsilon > 0, \exists p_\epsilon > 0$  such that  $\forall l$ :

$$\mathbb{P}_{\theta^{k_l+1}} \left[ \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ |\tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)| \right] \geq \eta'' - \epsilon \right) \wedge \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ \tilde{\mathbf{W}}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right]^2 \leq \epsilon \right) \right] \geq p_\epsilon.$$

Due to the assumption of standard initialization  $\mathbf{W}^{k_l+1} \sim_{\theta^{k_l+1}} \mathcal{N}(0, \sqrt{2/R_{k_l}} \mathbf{I})$ , it follows that  $\tilde{\mathbf{W}}^{k_l+1}$  and  $\|\mathbf{W}_{c,:}^{k_l+1}\|_2$  are independent, and we have e.g.  $\mathbb{P}_{\theta^{k_l+1}} [1 \leq \|\mathbf{W}_{c,:}^{k_l+1}\|_2 \leq 2] \geq 0$  independent of  $l$ . Therefore  $\forall \epsilon > 0, \exists p'_\epsilon > 0$  such that  $\forall l$ :

$$\mathbb{P}_{\theta^{k_l+1}} \left[ \mathbb{E}_{\mathbf{x}, \alpha} \left[ |\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)| \right] \geq \eta'' - \epsilon \right) \wedge \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ \mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right]^2 \leq 4\epsilon \right) \right] \geq p'_\epsilon. \quad (74)$$

Let us note that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, \alpha} \left[ |\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha)| \right] &= \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^+ \right] + \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^- \right], \\ \mathbb{E}_{\mathbf{x}, \alpha} \left[ \mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha) \right]^2 &= \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^+ \right] - \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^- \right] \right)^2, \end{aligned}$$

where  $x^+ \equiv \max(x, 0)$  and  $x^- \equiv \max(-x, 0)$ . We then get  $\exists \eta''' > 0, \exists p > 0$  such that  $\forall l$ :

$$\begin{aligned} \mathbb{P}_{\theta^{k_l+1}} \left[ \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^+ \right] \geq \eta''' \right) \wedge \left( \mathbb{E}_{\mathbf{x}, \alpha} \left[ (\mathbf{W}_{c,:}^{k_l+1} \rho(\tilde{\mathbf{x}}^{k_l}, \alpha))^- \right] \geq \eta''' \right) \right] &\geq p, \\ \mathbb{P}_{\theta^{k_l+1}} \left[ \left( \nu_{1,c}(\mathbf{y}^{k_l+1,+}) \geq \eta''' \mu_2(\mathbf{x}^l)^{1/2} \right) \wedge \left( \nu_{1,c}(\mathbf{y}^{k_l+1,-}) \geq \eta''' \mu_2(\mathbf{x}^l)^{1/2} \right) \right] &\geq p, \\ \mathbb{P}_{\theta^{k_l+1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{k_l+1,+}) \nu_{1,c}(\mathbf{y}^{k_l+1,-})}{\mu_2(\mathbf{x}^l)} \geq (\eta''')^2 \right] &\geq p. \end{aligned}$$

We finally get

$$\mathbb{E}_{c, \theta^{k_l+1}} \left[ \frac{\nu_{1,c}(\mathbf{y}^{k_l+1,+}) \nu_{1,c}(\mathbf{y}^{k_l+1,-})}{\mu_2(\mathbf{x}^l)} \right] \geq p(\eta''')^2.$$

Thus by Theorem 2,  $\exists \eta'''' > 0$  such that  $\forall l$ :

$$\exp(\overline{m}[\chi^{k_l+1}]) \geq 1 + \eta'''' ,$$

which contradicts the hypothesis  $\exp(\overline{m}[\chi^{k_l+1}]) \xrightarrow{l \rightarrow \infty} 1$ . We deduce that  $r_{\text{eff}}(\mathbf{x}^l) \xrightarrow{l \rightarrow \infty} 1$ , i.e. that  $\mathbf{x}^l$  converges to one-dimensional signal pathology.  $\square$

**E.7 If  $\exp(\overline{m}[\chi^l]) \rightarrow 1$ , then each new layer  $l$  becomes arbitrary well approximated by a linear function**

If  $\exp(\overline{m}[\chi^l]) \rightarrow 1$ , then Eq. (9) in Theorem 2 implies for the rescaled signal  $\tilde{\mathbf{y}}^l = \mathbf{y}^l / \sqrt{\mu_2(\mathbf{x}^{l-1})}$ :

$$\begin{aligned} \mathbb{E}_{\mathbf{c}, \theta^l} [\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}) \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})] &\rightarrow 0, \\ \mathbb{E}_{\mathbf{c}, \theta^l} [\min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}), \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-}))^2] &\rightarrow 0, \\ \forall \epsilon : \mathbb{P}_{\mathbf{c}, \theta^l} [\min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}), \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) \geq \epsilon] &\rightarrow 0, \\ \forall \epsilon : \mathbb{P}_{\theta^l} [\exists \mathbf{c} : \min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}), \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) \geq \epsilon] &\rightarrow 0, \\ \forall \epsilon : \mathbb{P}_{\theta^l} [\forall \mathbf{c} : \min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}), \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) \leq \epsilon] &\rightarrow 1. \end{aligned} \quad (75)$$

Now let us fix a channel  $\mathbf{c}$  and suppose that  $\min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}) \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) \leq \epsilon$ . Given that  $\tilde{\mathbf{y}}^{l,-} = \tilde{\mathbf{y}}^l - \tilde{\mathbf{y}}^{l,+}$ , we have

$$\min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}) \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) = \min(\nu_{1,\mathbf{c}}(|\tilde{\mathbf{y}}^{l,+} - 0|), \nu_{1,\mathbf{c}}(|\tilde{\mathbf{y}}^{l,+} - \tilde{\mathbf{y}}^l|)) \leq \epsilon.$$

Both  $\nu_{1,\mathbf{c}}(|\tilde{\mathbf{y}}^{l,+} - 0|)$  and  $\nu_{1,\mathbf{c}}(|\tilde{\mathbf{y}}^{l,+} - \tilde{\mathbf{y}}^l|)$  correspond to the mean absolute error of the approximation of the rescaled output  $\mathbf{x}^l / \mu_2(\mathbf{x}^{l-1}) = \tilde{\mathbf{y}}^{l,+} / \mu_2(\mathbf{x}^{l-1}) = \tilde{\mathbf{y}}^{l,+}$  in channel  $\mathbf{c}$  with a linear function of  $\mathbf{x}^{l-1}$ . So there exists a linear function  $f_{\mathbf{c}} : \mathbb{R}^{n \times \dots \times n \times N_{l-1}} \rightarrow \mathbb{R}^{n \times \dots \times n}$  such that

$$\mathbb{E}_{\mathbf{x}, \alpha} [|\tilde{\mathbf{y}}_{\alpha, \mathbf{c}}^{l,+} - f_{\mathbf{c}}(\mathbf{x}^{l-1})_{\alpha}|] \leq \epsilon.$$

If  $\forall \mathbf{c} : \min(\nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,+}) \nu_{1,\mathbf{c}}(\tilde{\mathbf{y}}^{l,-})) \leq \epsilon$ , and if we define the linear function  $f : \mathbb{R}^{n \times \dots \times n \times N_{l-1}} \rightarrow \mathbb{R}^{n \times \dots \times n \times N_l}$  such that  $\forall \alpha, \mathbf{c} : f(\mathbf{x}^{l-1})_{\alpha, \mathbf{c}} = f_{\mathbf{c}}(\mathbf{x}^{l-1})_{\alpha}$ , then we get

$$\nu_1(|\tilde{\mathbf{y}}^{l,+} - f(\mathbf{x}^{l-1})|) = \mathbb{E}_{\mathbf{x}, \alpha, \mathbf{c}} [|\tilde{\mathbf{y}}_{\alpha, \mathbf{c}}^{l,+} - f(\mathbf{x}^{l-1})_{\alpha, \mathbf{c}}|] = \mathbb{E}_{\mathbf{c}} \mathbb{E}_{\mathbf{x}, \alpha} [|\tilde{\mathbf{y}}_{\alpha, \mathbf{c}}^{l,+} - f_{\mathbf{c}}(\mathbf{x}^{l-1})_{\alpha}|] \leq \epsilon.$$

Combined with Eq. (75), this means that  $\mathbf{x}^l / \mu_2(\mathbf{x}^{l-1})$  can be approximated arbitrary well by a linear function of  $\mathbf{x}^{l-1}$  with probability arbitrary close to 1 in  $\theta^l$ .

## F Details of Section 6

### F.1 Proof of Theorem 3

**Theorem 3. Normalized Sensitivity increments of batch-normalized feedforward nets.** *The dominating term in the evolution of  $\chi^l$  can be decomposed as the sum of a term  $\overline{m}_{BN}[\chi^l]$  due to batch normalization and a term  $\overline{m}_{\phi}[\chi^l]$  due to the nonlinearity  $\phi$ :*

$$\begin{aligned} \exp(\overline{m}_{BN}[\chi^l]) &= \left( \frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \mathbb{E}_{\mathbf{c}, \theta^l} \left[ \frac{\mu_{2,\mathbf{c}}(\mathbf{t}^l)}{\mu_{2,\mathbf{c}}(\mathbf{y}^l)} \right]^{1/2}, \\ \exp(\overline{m}_{\phi}[\chi^l]) &= \left( 1 - 2 \mathbb{E}_{\mathbf{c}, \theta^l} [\nu_{1,\mathbf{c}}(\mathbf{z}^{l,+}) \nu_{1,\mathbf{c}}(\mathbf{z}^{l,-})] \right)^{-1/2}, \\ \delta \chi^l &\simeq \exp(\overline{m}_{BN-FF}[\chi^l]) = \exp(\overline{m}_{BN}[\chi^l] + \overline{m}_{\phi}[\chi^l]). \end{aligned}$$

**Proof.** First let us decompose the dominating term as the product of two terms:

$$\begin{aligned} \exp(\overline{m}_{BN}[\chi^l]) &= \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]}{\mu_2(\mathbf{s}^{l-1})} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2}, \\ \exp(\overline{m}_{\phi}[\chi^l]) &= \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]} \right)^{-1/2}, \end{aligned}$$



$$\begin{aligned}\exp(\overline{m}_{BN-FF}[\chi^l]) &= \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mu_2(\mathbf{s}^{l-1})} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \\ &= \exp(\overline{m}_{BN}[\chi^l]) \exp(\overline{m}_\phi[\chi^l]).\end{aligned}$$

$\overline{m}_{BN}[\chi^l]$  is a dominating term in the evolution of  $\chi^l$  from  $(\mathbf{x}^{l-1}, \mathbf{s}^{l-1})$  to  $(\mathbf{z}^l, \mathbf{u}^l)$ , while  $\overline{m}_\phi[\chi^l]$  is a dominating term in the evolution of  $\chi^l$  from  $(\mathbf{z}^l, \mathbf{u}^l)$  to  $(\mathbf{x}^l, \mathbf{s}^l)$ . These terms can be seen – slightly simplistically – as the direct contribution of respectively batch normalization and  $\phi$  to  $\overline{m}_{BN-FF}[\chi^l]$ . Now let us explicitate both terms.

**Term**  $\exp(\overline{m}_{BN}[\chi^l])$ . First note that batch normalization directly gives:  $\mu_2(\mathbf{z}^l) = 1$ , and thus  $\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)] = 1$ . Now let us explicitate  $\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]$ :

$$\begin{aligned}\forall \mathbf{c} : \mathbf{u}_{:,c}^l &= \frac{\mathbf{t}_{:,c}^l}{\sqrt{\mu_{2,c}(\mathbf{y}^l)}}, \quad \forall \mathbf{c} : \mu_{2,c}(\mathbf{u}^l) = \frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)}, \\ \mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)] &= \mathbb{E}_{\mathbf{c}, \theta^l}[\mu_{2,c}(\mathbf{u}^l)] = \mathbb{E}_{\mathbf{c}, \theta^l} \left[ \frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)} \right]\end{aligned}$$

All together, we get

$$\begin{aligned}\exp(\overline{m}_{BN}[\chi^l]) &= \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]}{\mu_2(\mathbf{s}^{l-1})} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \\ &= \left( \frac{\mu_2(\mathbf{s}^{l-1})}{\mu_2(\mathbf{x}^{l-1})} \right)^{-1/2} \mathbb{E}_{\mathbf{c}, \theta^l} \left[ \frac{\mu_{2,c}(\mathbf{t}^l)}{\mu_{2,c}(\mathbf{y}^l)} \right]^{1/2}\end{aligned}$$

**Term**  $\exp(\overline{m}_\phi[\chi^l])$ . We consider the symmetric propagation for batch-normalized feedforward nets. From Eq. (31) we deduce

$$\begin{aligned}\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)] + \mathbb{E}_{\theta^l}[\mu_2(\overline{\mathbf{s}}^l)] &= \mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)], \\ 2\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)] &= \mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)],\end{aligned}\tag{76}$$

where Eq. (76) is obtained by symmetry of the propagation. Now we turn to the symmetric propagation of the signal:

$$\begin{aligned}\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\overline{\mathbf{x}}^l) &= \mathbb{E}_{\mathbf{x}, \alpha}[(\mathbf{z}_{\alpha,c}^{l,+})^2] - \mathbb{E}_{\mathbf{x}, \alpha}[\mathbf{z}_{\alpha,c}^{l,+}]^2 + \mathbb{E}_{\mathbf{x}, \alpha}[(\mathbf{z}_{\alpha,c}^{l,-})^2] - \mathbb{E}_{\mathbf{x}, \alpha}[\mathbf{z}_{\alpha,c}^{l,-}]^2. \\ &= \nu_{2,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{2,c}(\mathbf{z}^{l,-}) - \nu_{1,c}(\mathbf{z}^{l,-})^2 \\ &= \nu_{2,c}(\mathbf{z}^l) - \left( \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2 \right),\end{aligned}\tag{77}$$

where Eq. (77) follows from Eq. (29). Due to the constraints imposed by batch normalization:  $\nu_{1,c}(\mathbf{z}^l) = 0$  and  $\nu_{2,c}(\mathbf{z}^l) = 1$ , it follows that

$$\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\overline{\mathbf{x}}^l) = 1 - \left( \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2 \right).\tag{78}$$

$$\begin{aligned}\nu_{1,c}(\mathbf{z}^l) &= \nu_{1,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,-}) = 0, \\ \left( \nu_{1,c}(\mathbf{z}^{l,+}) - \nu_{1,c}(\mathbf{z}^{l,-}) \right)^2 &= \nu_{1,c}(\mathbf{z}^{l,+})^2 + \nu_{1,c}(\mathbf{z}^{l,-})^2 - 2\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) = 0.\end{aligned}\tag{79}$$

Using Eq. (78), Eq. (79) and the symmetry of the propagation,

$$\begin{aligned}\mu_{2,c}(\mathbf{x}^l) + \mu_{2,c}(\overline{\mathbf{x}}^l) &= 1 - 2\nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}), \\ 2\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)] &= 1 - 2\mathbb{E}_{\mathbf{c}, \theta^l} \left[ \nu_{1,c}(\mathbf{z}^{l,+})\nu_{1,c}(\mathbf{z}^{l,-}) \right].\end{aligned}\tag{80}$$

We finally combine Eq. (76) and Eq. (80):

$$\begin{aligned}\exp(\overline{m}_\phi[\chi^l]) &= \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{s}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{u}^l)]} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{x}^l)]}{\mathbb{E}_{\theta^l}[\mu_2(\mathbf{z}^l)]} \right)^{-1/2}, \\ &= \left( 1 - 2\mathbb{E}_{c,\theta^l} \left[ \nu_{1,c}(\mathbf{z}^{l,+}) \nu_{1,c}(\mathbf{z}^{l,-}) \right] \right)^{-1/2}.\end{aligned}$$

□

**F.2 In the first step of the propagation:**  $\exp(\overline{m}_{BN}[\chi^1]) \geq 1$

Let us explicitate the second-order moment in channel  $c$  of  $\mathbf{t}^1$ :

$$\mu_{2,c}(\mathbf{t}^1) = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} [\hat{\varphi}(\mathbf{t}^1, \alpha)_c^2] = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} [\varphi(\mathbf{t}^1, \alpha)_c^2] = \mathbb{E}_{\mathbf{x},\mathbf{s},\alpha} [(\mathbf{W}_{c,:}^1 \rho(\mathbf{s}, \alpha))^2] \quad (81)$$

$$= \sum_{i,j} \mathbf{W}_{c,i}^1 \mathbf{W}_{c,j}^1 \mathbb{E}_{\mathbf{s},\alpha} [\rho(\mathbf{s}, \alpha)_i \rho(\mathbf{s}, \alpha)_j] = \sum_i (\mathbf{W}_{c,i}^1)^2 = \|\mathbf{W}_{c,:}^1\|_2^2. \quad (82)$$

where Eq. (81) follows from  $\mathbf{t}^1$  being centered and Eq. (82) follows from the white noise property  $\mathbb{E}_{\mathbf{s}}[\mathbf{s}_i \mathbf{s}_j] = \delta_{ij}$ , which implies for any  $\alpha$  that  $\mathbb{E}_{\mathbf{s}}[\rho(\mathbf{s}, \alpha)_i \rho(\mathbf{s}, \alpha)_j] = \delta_{ij}$  under periodic boundary conditions.

Now we turn to the second-order moment in channel  $c$  of  $\mathbf{y}^1$ . Denoting  $(\mathbf{e}_1, \dots, \mathbf{e}_{R_1})$  and  $(\lambda_1, \dots, \lambda_{R_1})$  respectively the orthogonal eigenvectors and eigenvalues of  $\mathbf{C}[\rho(\mathbf{x}, \alpha)]$ , and  $\tilde{\mathbf{W}}^1 = \mathbf{W}^1(\mathbf{e}_1, \dots, \mathbf{e}_{R_1})$ , we get

$$\begin{aligned}\mu_{2,c}(\mathbf{y}^1) &= \mathbb{E}_{\mathbf{x},\alpha} [\hat{\varphi}(\mathbf{y}^1, \alpha)_c^2] = \mathbb{E}_{\mathbf{x},\alpha} [(\mathbf{W}_{c,:}^1 \hat{\rho}(\mathbf{x}, \alpha))^2] = \sum_i (\hat{\mathbf{W}}_{c,i}^1)^2 \lambda_i \\ &= \|\mathbf{W}_{c,:}^1\|_2^2 \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i = \mu_{2,c}(\mathbf{t}^1) \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i,\end{aligned}$$

where we defined  $\tilde{\mathbf{W}}^1$  such that  $\forall c: \tilde{\mathbf{W}}_{c,:}^1 = \hat{\mathbf{W}}_{c,:}^1 / \|\mathbf{W}_{c,:}^1\|$  and we used Eq. (82). Under standard initialization, the distribution of  $\mathbf{W}^1$  is spherically symmetric, implying that for all  $c$  the distribution of  $\tilde{\mathbf{W}}_{c,:}^1$  is uniform on the sphere of  $\mathbb{R}^{R_1}$ . In turn, this implies

$$\begin{aligned}\forall i: \mathbb{E}_{\theta^1} [(\tilde{\mathbf{W}}_{c,i}^1)^2] &= \frac{1}{R_1}, \\ \forall c: \mathbb{E}_{\theta^1} \left[ \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right] &= \frac{1}{R_1} \sum_i \lambda_i, \quad \mathbb{E}_{c,\theta^1} \left[ \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right] = \frac{1}{R_1} \sum_i \lambda_i.\end{aligned} \quad (83)$$

Finally we can write  $\exp(\overline{m}_{BN}[\chi^1])$  as

$$\begin{aligned}\exp(\overline{m}_{BN}[\chi^1]) &= \left( \frac{\mu_2(\mathbf{s}^0)}{\mu_2(\mathbf{x}^0)} \right)^{-1/2} \mathbb{E}_{c,\theta^1} \left[ \frac{\mu_{2,c}(\mathbf{t}^1)}{\mu_{2,c}(\mathbf{y}^1)} \right]^{1/2} \\ &= \left( \frac{1}{R_1} \sum_i \lambda_i \right)^{1/2} \mathbb{E}_{c,\theta^1} \left[ \frac{1}{\sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i} \right]^{1/2},\end{aligned} \quad (84)$$

$$\geq \left( \frac{1}{R_1} \sum_i \lambda_i \right)^{1/2} \left( \mathbb{E}_{c,\theta^1} \left[ \sum_i (\tilde{\mathbf{W}}_{c,i}^1)^2 \lambda_i \right]^{-1} \right)^{1/2} = 1. \quad (85)$$

where Eq. (84) is obtained using  $\mu_2(\mathbf{s}^0) = \mu_2(\mathbf{s}) = 1$  and  $\mu_2(\mathbf{x}^0) = \frac{1}{R_1} \text{Tr } \mathbf{C}[\rho(\mathbf{x}, \alpha)] = \frac{1}{R_1} \sum_i \lambda_i$  by Corollary 3, while Eq. (85) is obtained using the convexity of  $x \mapsto 1/x$  and Eq. (83).

## G Details of Section 7

### G.1 Adaptation of the previous setup to resnets

Before proceeding to the analysis, slight adaptations and forewords are necessary. Let us denote  $\Theta^{l,h} = (\omega^{1,1}, \beta^{1,1}, \dots, \omega^{1,H}, \beta^{1,H}, \dots, \omega^{l,1}, \beta^{l,1}, \dots, \omega^{l,h}, \beta^{l,h})$  the full set of parameters up to

layer  $h$  in residual unit  $l$  and  $\theta^{l,h} = \Theta^{l,h} | \Theta^{l,h-1}$  the conditional set of parameters of layer  $h$  in residual unit  $l$ . We further denote  $\Theta^l = \Theta^{l,H}$  and  $\theta^l = \Theta^{l,H} | \Theta^{l-1,H}$  respectively the full and conditional sets of parameters at the granularity of the residual unit.

In the pre-activation perspective of Eq. (16) and (17) each residual layer starts with  $(\mathbf{y}^{l,h-1}, \mathbf{t}^{l,h-1})$  after the convolution and ends with  $(\mathbf{y}^{l,h}, \mathbf{t}^{l,h})$  again after the convolution. The concrete effect is that batch normalization and  $\phi$  are completely deterministic conditionally on  $\Theta^{l-1}$  in the first layer  $h = 1$  of each residual unit  $l$ . This occurs again for  $h \geq 2$  since batch normalization and  $\phi$  are random conditionally on  $\Theta^{l-1}$  but completely deterministic conditionally on  $\Theta^{l,h-1}$ . At even larger granularity, due to the aggregation  $(\mathbf{y}^l, \mathbf{t}^l) = \sum_{k=0}^l (\mathbf{y}^{k,H}, \mathbf{t}^{k,H})$ , the input  $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$  of each residual unit becomes more and more correlated between successive  $l$  and less and less dependent on the parameters  $\theta^{l-k}$  of previous individual units.

Since the evolution of  $\chi^l$  is mainly influenced by batch normalization and the nonlinearity  $\phi$ , this shift can be thought as attributing the parameters and thus the stochasticity of layer  $h$  to layer  $h - 1$ . A simple strategy to apply the results of Section 6 is thus to shift back to the post-activation perspective by considering the parameters  $\theta^{l,h-1}$  and the evolution from  $(\mathbf{x}^{l,h-1}, \mathbf{s}^{l,h-1})$  to  $(\mathbf{x}^{l,h}, \mathbf{s}^{l,h})$  for layers  $2 \leq h \leq H$ . Theorem 3 strictly applies in this case.

It remains to understand the evolution from  $(\mathbf{y}^{l,0}, \mathbf{t}^{l,0}) = (\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$  to  $(\mathbf{x}^{l,1}, \mathbf{s}^{l,1})$  in layer  $h = 1$  and the evolution from  $(\mathbf{x}^{l,H}, \mathbf{s}^{l,H})$  to  $(\mathbf{y}^{l,H}, \mathbf{t}^{l,H})$  in layer  $h = H$ . By considering the parameter  $\Theta^{l-1}$  and the evolution from  $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1})$  to  $(\mathbf{x}^{l,1}, \mathbf{s}^{l,1})$ , the batch normalization term of Eq. (12) in Theorem 3 becomes

$$\begin{aligned} & \left( \frac{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{u}^{l,1})]}{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{t}^{l-1})]} \right)^{1/2} \left( \frac{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{z}^{l,1})]}{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{y}^{l-1})]} \right)^{-1/2} \\ &= \left( \frac{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{t}^{l-1})]}{\mathbb{E}_{\Theta^{l-1}}[\mu_2(\mathbf{y}^{l-1})]} \right)^{-1/2} \mathbb{E}_{\mathbf{c}, \Theta^{l-1}} \left[ \frac{\mu_{2,\mathbf{c}}(\mathbf{t}^{l-1})}{\mu_{2,\mathbf{c}}(\mathbf{y}^{l-1})} \right]^{1/2}. \end{aligned}$$

Under the assumption of well-conditioned sensitivity  $\mu_{2,\mathbf{c}}(\mathbf{t}^l) \simeq \mu_{2,\mathbf{c}}(\mathbf{t}^{l-1})$ , this term is again  $\gtrsim 1$  by convexity of  $x \mapsto 1/x$ . For the nonlinearity term, symmetric propagation with respect to  $\Theta^{l-1}$  applies for all terms in the sum  $(\mathbf{y}^{l-1}, \mathbf{t}^{l-1}) = \sum_{k=0}^{l-1} (\mathbf{y}^{k,H}, \mathbf{t}^{k,H})$  except for  $(\mathbf{y}^{0,H}, \mathbf{t}^{0,H}) = (\mathbf{y}, \mathbf{t})$ . The expression of the nonlinearity term of Eq. (13) in Theorem 3 thus remains approximately valid.

Finally by spherical symmetry, the evolution from  $(\mathbf{x}^{l,H}, \mathbf{s}^{l,H})$  to  $(\mathbf{y}^{l,H}, \mathbf{t}^{l,H})$  in layer  $h = H$  has dominating term

$$\left( \frac{\mathbb{E}_{\theta^{l,H}}[\mu_2(\mathbf{t}^{l,H})]}{\mathbb{E}_{\theta^{l,H}}[\mu_2(\mathbf{s}^{l,H})]} \right)^{1/2} \left( \frac{\mathbb{E}_{\theta^{l,H}}[\mu_2(\mathbf{y}^{l,H})]}{\mathbb{E}_{\theta^{l,H}}[\mu_2(\mathbf{x}^{l,H})]} \right)^{-1/2} = 1.$$

In summary, Theorem 3 remains approximately valid in the feedforward evolution inside residual units.

## G.2 Lemma on dot-product

**Lemma 14.** *The following hold:*

$$\begin{aligned} & \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\phi}(\mathbf{y}^{l-1})_{\mathbf{c}} \hat{\phi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}]] = 0, \\ & \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\phi}(\mathbf{y}^{l-1})_{\mathbf{c}} \hat{\phi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}]^2] \leq \frac{1}{Nr_{\text{eff}}(\mathbf{y}^{l-1})} \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})], \\ & \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \mathbf{t}, \alpha, \mathbf{c}} [\hat{\phi}(\mathbf{t}^{l-1})_{\mathbf{c}} \hat{\phi}(\mathbf{t}^{l,H}, \alpha)_{\mathbf{c}}]] = 0, \\ & \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \mathbf{t}, \alpha, \mathbf{c}} [\hat{\phi}(\mathbf{t}^{l-1})_{\mathbf{c}} \hat{\phi}(\mathbf{t}^{l,H}, \alpha)_{\mathbf{c}}]^2] \leq \frac{1}{Nr_{\text{eff}}(\mathbf{y}^{l-1})} \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{t}^{l,H})]. \end{aligned}$$

**Proof.** By spherical symmetry, moments of  $\hat{\phi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}$  and  $-\hat{\phi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}} = \hat{\phi}(-\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}$  have the same distribution with respect to  $\theta^l$ . It follows that

$$\begin{aligned}\mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}]] &= \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} (-\hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}})]]], \\ \mathbb{E}_{\theta^l} [\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}]] &= 0.\end{aligned}$$

Next we note that

$$\begin{aligned}\mathbb{E}_{\mathbf{y}, \alpha, \mathbf{c}} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}] &= \frac{1}{N} \sum_{\mathbf{c}} \mathbb{E}_{\mathbf{y}, \alpha} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_{\mathbf{c}} \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_{\mathbf{c}}], \\ &= \frac{1}{N} \mathbb{E}_{\mathbf{y}, \alpha} [\langle \hat{\varphi}(\mathbf{y}^{l-1}, \alpha) \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) \rangle],\end{aligned}\quad (86)$$

where  $\langle \cdot \rangle$  denotes the standard dot product in  $\mathbb{R}^N$ . Let us denote  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$  and  $(\lambda_1, \dots, \lambda_N)$  respectively the orthogonal eigenvectors and eigenvalues of  $\mathbf{C}[\varphi(\mathbf{y}^{l-1}, \alpha)]$ . We further denote  $u_i$  the unit-variance components of  $\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)$  in the basis  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$ , and  $y_i$  the components of  $\hat{\varphi}(\mathbf{y}^{l,H}, \alpha)$  in the basis  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$ . This gives

$$\begin{aligned}\hat{\varphi}(\mathbf{y}^{l-1}, \alpha) &= \sum_i \sqrt{\lambda_i} u_i \mathbf{e}_i, \quad \mathbb{E}_{\mathbf{y}, \alpha} [u_i] = 0, \quad \mathbb{E}_{\mathbf{y}, \alpha} [u_i^2] = 1, \\ \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) &= \sum_i y_i \mathbf{e}_i.\end{aligned}$$

Now we decompose each component  $y_i$  of  $\mathbf{y}^{l,H}$  as:

$$\forall j : \alpha_{i,j} = \mathbb{E}_{\mathbf{y}, \alpha} [y_i u_j], \quad y_i = \sum_j \alpha_{i,j} u_j + z_i,$$

From this definition, we get

$$\begin{aligned}\forall j : \mathbb{E}_{\mathbf{y}, \alpha} [z_i u_j] &= 0, \quad \mathbb{E}_{\mathbf{y}, \alpha} [y_i u_i] = \alpha_{i,i}, \quad \mathbb{E}_{\mathbf{y}, \alpha} [y_i^2] = \sum_j \alpha_{i,j}^2 + \mathbb{E}_{\mathbf{y}, \alpha} [z_i^2], \\ \mu_2(\mathbf{y}^{l,H}) &= \frac{1}{N} \mathbb{E}_{\mathbf{y}, \alpha} [\langle \mathbf{y}^{l,H}, \mathbf{y}^{l,H} \rangle] = \frac{1}{N} \sum_i \mathbb{E}_{\mathbf{y}, \alpha} [y_i^2] = \frac{1}{N} \left( \sum_{i,j} \alpha_{i,j}^2 + \sum_i \mathbb{E}_{\mathbf{y}, \alpha} [z_i^2] \right).\end{aligned}\quad (87)$$

where the dot product in Eq. (87) is computed in the orthogonal basis  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$ . Now computing the dot product of  $\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)$  and  $\hat{\varphi}(\mathbf{y}^{l,H}, \alpha)$  in the orthogonal basis  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$ :

$$\mathbb{E}_{\mathbf{y}, \alpha} [\langle \hat{\varphi}(\mathbf{y}^{l-1}, \alpha) \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) \rangle] = \sum_i \sqrt{\lambda_i} \mathbb{E}_{\mathbf{y}, \alpha} [y_i u_i] = \sum_i \sqrt{\lambda_i} \alpha_{i,i}.$$

Spherical symmetry implies that moments of  $y_1 \mathbf{e}_1 + \dots + y_i \mathbf{e}_i + \dots + y_N \mathbf{e}_N$  and  $y_1 \mathbf{e}_1 + \dots - y_i \mathbf{e}_i + \dots + y_N \mathbf{e}_N$  have the same distribution with respect to  $\theta^l$ . Thus:

$$\begin{aligned}\forall j \neq i : \mathbb{E}_{\mathbf{y}, \alpha} [y_i u_i] \mathbb{E}_{\mathbf{y}, \alpha} [y_j u_j] &\sim_{\theta^l} \mathbb{E}_{\mathbf{y}, \alpha} [-y_i u_i] \mathbb{E}_{\mathbf{y}, \alpha} [y_j u_j], \\ \forall j \neq i : \alpha_{i,i} \alpha_{j,j} &\sim_{\theta^l} (-\alpha_{i,i}) \alpha_{j,j}, \\ \forall j \neq i : \mathbb{E}_{\theta^l} [\alpha_{i,i} \alpha_{j,j}] &= 0.\end{aligned}$$

We deduce that

$$\mathbb{E}_{\theta^l} \left[ \mathbb{E}_{\mathbf{y}, \alpha} [\langle \hat{\varphi}(\mathbf{y}^{l-1}, \alpha) \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) \rangle]^2 \right] = \sum_i \lambda_i \mathbb{E}_{\theta^l} [\alpha_{i,i}^2].$$

Spherical symmetry also implies that the distribution of  $\alpha_{i,j}$  with respect to  $\theta^l$  does not depend on  $i$ . Denoting  $(\beta_j)$  such that  $\forall i, j : \beta_j = \mathbb{E}_{\theta^l} [\alpha_{i,j}^2]$ , we get combined with Eq. (87):

$$\begin{aligned}\mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})] &\geq \frac{1}{N} \sum_{i,j} \mathbb{E}_{\theta^l} [\alpha_{i,j}^2] = \frac{1}{N} \sum_{i,j} \beta_j = \sum_i \beta_i, \\ \mathbb{E}_{\theta^l} \left[ \mathbb{E}_{\mathbf{y},\alpha} \left[ \langle \hat{\varphi}(\mathbf{y}^{l-1}, \alpha) \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) \rangle \right]^2 \right] &= \sum_i \lambda_i \beta_i \leq \lambda_{\max} \left( \sum_i \beta_i \right) \leq \lambda_{\max} \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})].\end{aligned}$$

Finally combining with Eq. (86):

$$\begin{aligned}\mathbb{E}_{\theta^l} \left[ \mathbb{E}_{\mathbf{y},\alpha,c} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_c \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_c]^2 \right] &= \frac{1}{N^2} \mathbb{E}_{\theta^l} \left[ \mathbb{E}_{\mathbf{y},\alpha} \left[ \langle \hat{\varphi}(\mathbf{y}^{l-1}, \alpha) \hat{\varphi}(\mathbf{y}^{l,H}, \alpha) \rangle \right]^2 \right] \\ &\leq \frac{1}{N^2} \lambda_{\max} \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})] \\ &\leq \frac{1}{N r_{\text{eff}}(\mathbf{y}^{l-1})} \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})],\end{aligned}$$

where we used  $\lambda_{\max} r_{\text{eff}}(\mathbf{y}^{l-1}) = \sum_i \lambda_i = N \mu_2(\mathbf{y}^{l-1})$ . The same analysis can be applied to  $\hat{\varphi}(\mathbf{t}^{l-1}, \alpha)$  and  $\hat{\varphi}(\mathbf{t}^{l,H}, \alpha)$ .  $\square$

**Corollary 15.** *Let us denote the dot products:*

$$\begin{aligned}Y_l &= \mathbb{E}_{\mathbf{y},\alpha,c} [\hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_c \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_c], \\ T_l &= \mathbb{E}_{\mathbf{y},\alpha,c} [\hat{\varphi}(\mathbf{t}^{l-1}, \alpha)_c \hat{\varphi}(\mathbf{t}^{l,H}, \alpha)_c], \\ Y_{l,l} &= \mathbb{E}_{\mathbf{y},\alpha,c} [\hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_c \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_c] = \mu_2(\mathbf{y}^{l,H}), \\ T_{l,l} &= \mathbb{E}_{\mathbf{y},\alpha,c} [\hat{\varphi}(\mathbf{t}^{l,H}, \alpha)_c \hat{\varphi}(\mathbf{t}^{l,H}, \alpha)_c] = \mu_2(\mathbf{t}^{l,H}).\end{aligned}$$

Then by spherical symmetry:

$$\begin{aligned}\forall l : \mathbb{E}_{\Theta^l} [Y_l] &= 0, & \forall l \neq l' : \mathbb{E}_{\Theta^{\max(l,l')}} [Y_l Y_{l'}] &= 0, \\ \forall l : \mathbb{E}_{\Theta^l} [T_l] &= 0, & \forall l \neq l' : \mathbb{E}_{\Theta^{\max(l,l')}} [T_l T_{l'}] &= 0.\end{aligned}$$

Furthermore given  $r_{\text{eff}}(\mathbf{y}^{l-1})$ ,  $r_{\text{eff}}(\mathbf{t}^{l-1}) \geq 1$  and given Lemma 14, we deduce the following inequalities in preparation of the proof of Theorem 4:

$$\begin{aligned}\mathbb{E}_{\Theta^l} [Y_l^2] &\leq \frac{1}{N} \mathbb{E}_{\Theta^{l-1}} [\mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} [\mu_2(\mathbf{y}^{l,H})]], \\ \mathbb{E}_{\Theta^l} \left[ \left( \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} T_l \right)^2 \right] &\leq \frac{1}{N} \mathbb{E}_{\Theta^{l-1}} \left[ \frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^{l-1})}{(\chi^{l-1})^2} \mathbb{E}_{\theta^l} \left[ \frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^{l,H})}{(\chi^{l-1})^2} \right] \right], \\ &\leq \frac{1}{N} \mathbb{E}_{\Theta^{l-1}} \left[ \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\theta^l} \left[ \frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^{l,H})}{(\chi^{l-1})^2} \right] \right].\end{aligned}$$

### G.3 Proof of Theorem 4

**Theorem 4. Normalized Sensitivity increments of batch-normalized resnets.** Suppose that for all depth  $l$  we can bound the second-order central moments  $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$  and the feedforward increments inside residual units  $\delta_{\min} \lesssim \delta \chi^{l,h} \lesssim \delta_{\max}$ . Denote  $\eta_{\min} = ((\delta_{\min})^{2H} \mu_{2,\min} - \mu_{2,\max}) / \mu_{2,\max}$  and  $\eta_{\max} = ((\delta_{\max})^{2H} \mu_{2,\max} - \mu_{2,\min}) / \mu_{2,\min}$ , as well as  $\tau_{\min} = \eta_{\min} / 2$  and  $\tau_{\max} = \eta_{\max} / 2$ . Then there exist positive constants  $C_{\min}, C_{\max} > 0$  such that

$$\left( 1 + \frac{\eta_{\min}}{l+1} \right)^{1/2} \lesssim \delta \chi^l \lesssim \left( 1 + \frac{\eta_{\max}}{l+1} \right)^{1/2}, \quad (88)$$

$$C_{\min} l^{\tau_{\min}} \lesssim \chi^l \lesssim C_{\max} l^{\tau_{\max}}. \quad (89)$$

**Proof.** First we introduce the additional constants  $\gamma_{\min} = (\delta_{\min})^{2H}$  and  $\gamma_{\max} = (\delta_{\max})^{2H}$ , so that we can write  $\eta_{\min} = (\gamma_{\min} \mu_{2,\min} - \mu_{2,\max}) / \mu_{2,\max}$  and  $\eta_{\max} = (\gamma_{\max} \mu_{2,\max} - \mu_{2,\min}) / \mu_{2,\min}$ .

We also remind that we write  $a \lesssim b$  when  $a(1 + \epsilon_a) \leq b(1 + \epsilon_b)$  with  $|\epsilon_a| \ll 1, |\epsilon_b| \ll 1$  with high probability. And we write  $a \simeq b$  when  $a(1 + \epsilon_a) = b(1 + \epsilon_b)$  with  $|\epsilon_a| \ll 1, |\epsilon_b| \ll 1$  with high probability. Denoting  $\wedge$  the logical *and*,  $\vee$  the logical *or*, the following rules are easily verified:

$$\begin{aligned} (a \lesssim b) \wedge (a \gtrsim b) &\iff (a \simeq b), \\ (a \lesssim b) &\iff (-a \gtrsim -b), \\ (a \lesssim b) \wedge (c \lesssim d) &\implies (ac \lesssim bd), \\ (a \lesssim b) \wedge (b \lesssim c) &\implies (a \lesssim c), \\ (a \lesssim b) \wedge (a > 0) \wedge (b > 0) &\implies (\sqrt{a} \lesssim \sqrt{b}), \\ (a \lesssim b) \wedge (a > 0) \wedge (b > 0) &\implies (1/a \gtrsim 1/b), \\ (a \lesssim b) \wedge (c \lesssim d) \wedge \left( \mathbb{P} \left[ (|a+c| \ll |a| + |c|) \vee (|b+d| \ll |b| + |d|) \right] \ll 1 \right) \\ &\implies (a+c \lesssim b+d). \end{aligned}$$

Finally consider a constant  $b$  and a random variable  $a$  depending on  $\Theta^l$  with well-defined moments. Let us prove that  $(a \lesssim b) \implies (\mathbb{E}_{\Theta^l}[a] \lesssim b) \wedge (\mathbb{E}_{\Theta^l}[a^2] \lesssim b^2)$ . Denoting  $A$  the event under which  $a(1 + \epsilon_a) \leq b(1 + \epsilon_b)$  with  $|\epsilon_a| \ll 1, |\epsilon_b| \ll 1$ , we can write

$$\frac{1}{\mathbb{E}_{\Theta^l}[a]^2} \left( \mathbb{E}_{\Theta^l}[a] - \mathbb{E}_{\Theta^l}[\mathbf{1}_A a] \right)^2 = \frac{1}{\mathbb{E}_{\Theta^l}[a]^2} \mathbb{E}_{\Theta^l}[\mathbf{1}_{A^c} a]^2 \leq \mathbb{P}_{\Theta^l}[\mathbf{1}_{A^c}] \frac{\mathbb{E}_{\Theta^l}[a^2]}{\mathbb{E}_{\Theta^l}[a]^2}, \quad (90)$$

where  $A^c$  denotes the complementary event of  $A$  and Eq. (90) is obtained using Cauchy-Schwarz inequality. The assumption  $(a \lesssim b)$  further gives  $\mathbb{P}_{\Theta^l}[A] \simeq 1$  and  $\mathbb{P}_{\Theta^l}[A^c] \ll 1$ . By contradiction if we had non negligible probability with respect to  $\Theta^{l-1}$  that  $\mathbb{P}_{\Theta^l}[\mathbf{1}_{A^c}] = \mathbb{P}_{\Theta^{l-1}}[\mathbf{1}_{A^c}]$  is non negligible, then we would not have  $\mathbb{P}_{\Theta^l}[A^c] = \mathbb{E}_{\Theta^{l-1}} \mathbb{E}_{\Theta^l|\Theta^{l-1}}[\mathbf{1}_{A^c}] \ll 1$ . It follows that  $\mathbb{P}_{\Theta^l}[\mathbf{1}_{A^c}] \ll 1$  with high probability with respect to  $\Theta^{l-1}$  and that  $\mathbb{E}_{\Theta^l}[a] \simeq \mathbb{E}_{\Theta^l}[\mathbf{1}_A a] \lesssim b$ .

A similar reasoning gives

$$\frac{1}{\mathbb{E}_{\Theta^l}[a]^2} \left( \mathbb{E}_{\Theta^l}[a] - \mathbb{E}_{\Theta^l}[\mathbf{1}_A a] \right)^2 \leq \mathbb{P}_{\Theta^l}[\mathbf{1}_{A^c}] \frac{\mathbb{E}_{\Theta^l}[a^2]}{\mathbb{E}_{\Theta^l}[a]^2}, \quad \mathbb{E}_{\Theta^l}[a] \simeq \mathbb{E}_{\Theta^l}[\mathbf{1}_A a] \lesssim b.$$

We keep all these rules in mind in the course of this proof.

**Proof of Eq. (88).** Adopting the same notations as Corollary 15 and using  $\mathbf{y}^l = \mathbf{y}^{l-1} + \mathbf{y}^{l,H}$  by Eq. (15), we get

$$\begin{aligned} \mu_2(\mathbf{y}^l) &= \mathbb{E}_{\mathbf{y}, \alpha, c} \left[ \left( \hat{\varphi}(\mathbf{y}^{l-1}, \alpha)_c + \hat{\varphi}(\mathbf{y}^{l,H}, \alpha)_c \right)^2 \right] = \mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2Y_l, \\ \mu_2(\mathbf{t}^l) &= \mathbb{E}_{\mathbf{y}, \mathbf{t}, \alpha, c} \left[ \left( \hat{\varphi}(\mathbf{t}^{l-1}, \alpha)_c + \hat{\varphi}(\mathbf{t}^{l,H}, \alpha)_c \right)^2 \right] = \mu_2(\mathbf{t}^{l-1}) + T_{l,l} + 2T_l. \end{aligned} \quad (91)$$

Due to  $\mu_{2,\min} \lesssim Y_{l,l} = \mu_2(\mathbf{y}^{l,H}) \lesssim \mu_{2,\max}$ , we have  $\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^0) = \mu_2(\mathbf{y}^{0,H}) \lesssim \mu_{2,\max}$ .

Now let us reason by induction and suppose that  $l\mu_{2,\min} \lesssim \mu_2(\mathbf{y}^{l-1}) \lesssim l\mu_{2,\max}$ . Combined with Eq. (91) and Corollary 15 we deduce that

$$\begin{aligned} l\mu_{2,\min} + \mu_{2,\min} + 2Y_l &\lesssim \mu_2(\mathbf{y}^l) \lesssim l\mu_{2,\max} + \mu_{2,\max} + 2Y_l, \\ \mathbb{E}_{\Theta^l}[Y_l^2] &\lesssim \frac{1}{N} l\mu_{2,\max}^2 \leq \frac{1}{N} \frac{1}{l+1} (l+1)^2 \mu_{2,\max}^2. \end{aligned}$$

Further using Chebyshev's inequality, we deduce that

$$\mathbb{P}_{\Theta^l} \left[ |Y_l| > k \frac{1}{\sqrt{N}} \frac{1}{\sqrt{l+1}} (l+1) \mu_{2,\max} \right] \lesssim \frac{1}{k^2}. \quad (92)$$

For large width  $N \gg 1$ , it follows that  $|Y_l| \ll (l+1) \mu_{2,\min}$  and  $|Y_l| \ll (l+1) \mu_{2,\max}$  with high probability, which then gives

$$(l+1) \mu_{2,\min} \lesssim \mu_2(\mathbf{y}^l) \lesssim (l+1) \mu_{2,\max}. \quad (93)$$

Eq. (93) then holds for all  $l$ . Now let us write  $(\chi^l)^2$  as

$$\begin{aligned} (\chi^l)^2 &= \frac{\mu_2(\mathbf{y}^0) \mu_2(\mathbf{t}^l)}{\mu_2(\mathbf{y}^l)} = \mu_2(\mathbf{y}^0) \frac{\mu_2(\mathbf{t}^{l-1}) + T_{l,l} + 2\tilde{T}_l}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2Y_l}, \\ (\chi^l)^2 &= (\chi^{l-1})^2 \frac{\mu_2(\mathbf{y}^{l-1}) + \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} T_{l,l} + 2 \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} \tilde{T}_l}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2Y_l}. \end{aligned}$$

Denoting  $\tilde{T}_{l,l} = \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} T_{l,l}$  and  $\tilde{T}_l = \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} \tilde{T}_l$ , we then get

$$(\delta\chi^l)^2 = \frac{(\chi^l)^2}{(\chi^{l-1})^2} = \frac{\mu_2(\mathbf{y}^{l-1}) + \tilde{T}_{l,l} + 2\tilde{T}_l}{\mu_2(\mathbf{y}^{l-1}) + Y_{l,l} + 2Y_l}. \quad (94)$$

Furthermore we can bound  $Y_{l,l}$  and  $\tilde{T}_{l,l}$  as

$$\mu_{2,\min} \lesssim Y_{l,l} \lesssim \mu_{2,\max}, \quad (95)$$

$$\begin{aligned} \tilde{T}_{l,l} &= \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} \mu_2(\mathbf{t}^{l,H}) = \frac{\mu_2(\mathbf{y}^0)}{(\chi^{l-1})^2} (\chi^{l-1})^2 \prod_h (\delta\chi^{l,h})^2 \frac{\mu_2(\mathbf{y}^{l,H})}{\mu_2(\mathbf{y}^0)}, \\ \gamma_{\min} \mu_{2,\min} &\lesssim \tilde{T}_{l,l} \lesssim \gamma_{\max} \mu_{2,\max}. \end{aligned} \quad (96)$$

By Corollary 15, variances of  $Y_l$  and  $\tilde{T}_l$  are bounded as

$$\begin{aligned} \mathbb{E}_{\Theta^l} [Y_l^2] &\lesssim \frac{1}{N} l \mu_{2,\max}^2, \\ \mathbb{E}_{\Theta^l} [\tilde{T}_l^2] &\lesssim \frac{1}{N} \mathbb{E}_{\Theta^{l-1}} \left[ \mu_2(\mathbf{y}^{l-1}) \mathbb{E}_{\Theta^l} [\tilde{T}_{l,l}] \right] \lesssim \gamma_{\max} \frac{1}{N} l \mu_{2,\max}^2. \end{aligned}$$

The reasoning of Eq. (92) then implies that  $|Y_l| \ll \mu_2(\mathbf{y}^{l-1})$  and  $|\tilde{T}_l| \ll \mu_2(\mathbf{y}^{l-1})$  with high probability. Finally combining Eq. (94), Eq. (95) and Eq. (96):

$$\begin{aligned} \frac{\mu_2(\mathbf{y}^{l-1}) + \gamma_{\min} \mu_{2,\min}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\max}} &\lesssim (\delta\chi^l)^2 \lesssim \frac{\mu_2(\mathbf{y}^{l-1}) + \gamma_{\max} \mu_{2,\max}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\min}}, \\ 1 + \frac{\gamma_{\min} \mu_{2,\min} - \mu_{2,\max}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\max}} &\lesssim (\delta\chi^l)^2 \lesssim 1 + \frac{\gamma_{\max} \mu_{2,\max} - \mu_{2,\min}}{\mu_2(\mathbf{y}^{l-1}) + \mu_{2,\min}}, \\ 1 + \frac{\gamma_{\min} \mu_{2,\min} - \mu_{2,\max}}{(l+1) \mu_{2,\max}} &\lesssim (\delta\chi^l)^2 \lesssim 1 + \frac{\gamma_{\max} \mu_{2,\max} - \mu_{2,\min}}{(l+1) \mu_{2,\min}}, \\ \left(1 + \frac{\eta_{\min}}{l+1}\right)^{1/2} &\lesssim \delta\chi^l \lesssim \left(1 + \frac{\eta_{\max}}{l+1}\right)^{1/2}. \end{aligned}$$

□

**Proof of Eq. (89).** Expanding Eq. (88), we get

$$\prod_{k=1}^l \left(1 + \frac{\eta_{\min}}{k+1}\right)^{1/2} \lesssim \chi^l = \prod_{k=1}^l \delta\chi^k \lesssim \prod_{k=1}^l \left(1 + \frac{\eta_{\max}}{k+1}\right)^{1/2}.$$



We can further explicitate the bounds:

$$\begin{aligned}
& \sum_{k=1}^l \log \left( 1 + \frac{\eta_{\max}}{k+1} \right) \\
& \leq \int_1^{l+1} \log \left( 1 + \frac{\eta_{\max}}{x} \right) dx, \\
& \leq \int_1^{l+1} \log(x + \eta_{\max}) dx - \int_1^{l+1} \log x dx, \\
& \leq \left[ x \log x - x \right]_{1+\eta_{\max}}^{l+1+\eta_{\max}} - \left[ x \log x - x \right]_1^{l+1}, \\
& \leq (l+1+\eta_{\max}) \log(l+1+\eta_{\max}) - (1+\eta_{\max}) \log(1+\eta_{\max}) - (l+1) \log(l+1), \\
& \leq \eta_{\max} \log(l+1+\eta_{\max}) + (l+1) \log \left( 1 + \frac{\eta_{\max}}{l+1} \right) - (1+\eta_{\max}) \log(1+\eta_{\max}), \\
& \leq \eta_{\max} \log(l+1+\eta_{\max}) + \eta_{\max} - (1+\eta_{\max}) \log(1+\eta_{\max}), \tag{97}
\end{aligned}$$

where we used  $\log(1+x) \leq x$  in Eq. (97). Considering the integration between 2 and  $l+2$ , we get with a similar calculation:

$$\begin{aligned}
& \sum_{k=1}^l \log \left( 1 + \frac{\eta_{\min}}{k+1} \right) \\
& \geq \eta_{\min} \log(l+2+\eta_{\min}) + (l+2) \log \left( 1 + \frac{\eta_{\min}}{l+2} \right) - (2+\eta_{\min}) \log(2+\eta_{\min}) + 2 \log 2, \\
& \geq \eta_{\min} \log(l+2+\eta_{\min}) - (2+\eta_{\min}) \log(2+\eta_{\min}) + 2 \log 2.
\end{aligned}$$

Let  $c_{\max} = \exp \left( \eta_{\max} - (1+\eta_{\max}) \log(1+\eta_{\max}) \right)$ ,  $c_{\min} = \exp \left( -(2+\eta_{\min}) \log(2+\eta_{\min}) + 2 \log 2 \right)$ .

We get

$$\begin{aligned}
& \prod_{k=1}^l \left( 1 + \frac{\eta_{\max}}{k+1} \right) \leq c_{\max} (l+1+\eta_{\max})^{\eta_{\max}}, \\
& \prod_{k=1}^l \left( 1 + \frac{\eta_{\min}}{k+1} \right) \geq c_{\min} (l+2+\eta_{\min})^{\eta_{\min}}, \\
& \sqrt{c_{\min}} (l+2+\eta_{\min})^{\eta_{\min}/2} \lesssim \chi^l \lesssim \sqrt{c_{\max}} (l+1+\eta_{\max})^{\eta_{\max}/2}, \\
& \sqrt{c_{\min}} (l+2+\eta_{\min})^{\tau_{\min}} \lesssim \chi^l \lesssim \sqrt{c_{\max}} (l+1+\eta_{\max})^{\tau_{\max}}.
\end{aligned}$$

Since  $x \mapsto \left( \frac{x+2+\eta_{\min}}{x} \right)^{\tau_{\min}}$  and  $x \mapsto \left( \frac{x+1+\eta_{\max}}{x} \right)^{\tau_{\max}}$  are lower-bounded and upper-bounded for  $x \geq 1$ , there exist positive constants  $C_{\min}, C_{\max} > 0$  such that

$$C_{\min} l^{\tau_{\min}} \lesssim \chi^l \lesssim C_{\max} l^{\tau_{\max}}.$$

□